

# Implementing the IBM System Storage SAN Volume Controller with IBM Spectrum Virtualize V8.3.1

Jack Armstrong

Tiago Bastos

Pawel Brodacki

Markus Döllinger

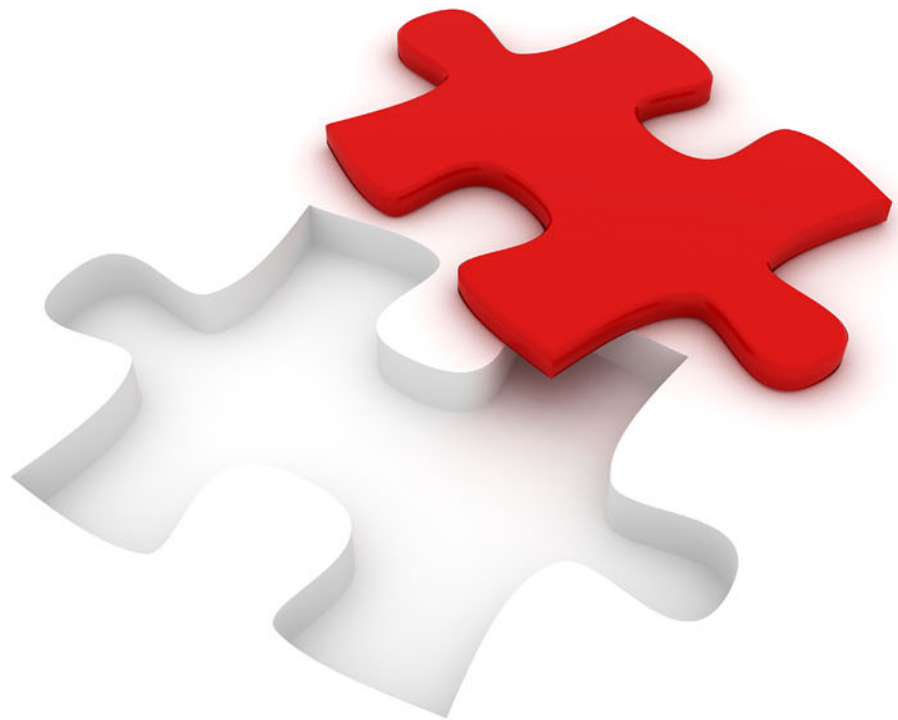
Jon Herd

Sergey Kubin

Carsten Larsen

Hartmut Lonzer

Jon Tate



**Storage**





IBM Redbooks

**Implementing the IBM System Storage SAN Volume  
Controller with IBM Spectrum Virtualize V8.3.1**

January 2021

**Note:** Before using this information and the product it supports, read the information in “Notices” on page xiii.

### **First Edition (January 2021)**

This edition applies to IBM Spectrum Virtualize V8.3.1 and the associated hardware and software that are detailed within. The screen captures might differ from the generally available (GA) version because parts of this book were written with pre-GA code. On 11 February 2020, IBM announced the arrival of IBM SAN Volume Controller SA2 and SV2. This book was written specifically for prior versions of IBM SAN Volume Controller, but most of the general principles apply. If you are in any doubt as to their applicability, then you should work with your local IBM representative.

**© Copyright International Business Machines Corporation 2021. All rights reserved.**

Note to U.S. Government Users Restricted Rights -- Use, duplication or disclosure restricted by GSA ADP Schedule Contract with IBM Corp.



# Contents

<b>Notices</b> .....	xiii
Trademarks .....	xiv
<b>Preface</b> .....	xv
Authors .....	xv
Now you can become a published author, too! .....	xviii
Comments welcome .....	xviii
Stay connected to IBM Redbooks .....	xix
<b>Chapter 1. Introduction and system overview</b> .....	1
1.1 Storage virtualization terminology .....	2
1.2 Latest changes and enhancements .....	5
1.3 IBM SAN Volume Controller architecture .....	7
1.3.1 Brief history of IBM SAN Volume Controller .....	8
1.3.2 IBM SAN Volume Controller architectural overview .....	8
1.3.3 IBM Spectrum Virtualize .....	12
1.3.4 IBM SAN Volume Controller topology .....	12
1.4 IBM SAN Volume Controller models .....	14
1.4.1 IBM SAN Volume Controller DH8 .....	14
1.4.2 IBM SAN Volume Controller SV1 .....	14
1.4.3 IBM SAN Volume Controller SV2 and SA2 .....	15
1.4.4 SAN Volume Controller model comparisons .....	17
1.5 IBM SAN Volume Controller components .....	20
1.5.1 Nodes .....	21
1.5.2 I/O groups .....	21
1.5.3 System .....	22
1.5.4 Expansion enclosures .....	22
1.5.5 Dense expansion drawers .....	23
1.5.6 Flash drives .....	24
1.5.7 MDisks .....	28
1.5.8 Cache .....	30
1.5.9 Quorum disk .....	33
1.5.10 Disk tier .....	35
1.5.11 Storage pool .....	35
1.5.12 Volumes .....	36
1.5.13 Easy Tier .....	38
1.5.14 Hosts .....	39
1.5.15 Host cluster .....	39
1.5.16 RAID .....	40
1.5.17 Encryption .....	40
1.5.18 iSCSI .....	41
1.5.19 IBM Real-time Compression .....	41
1.5.20 Data Reduction Pools .....	42
1.5.21 Deduplication .....	42
1.5.22 IP replication .....	43
1.5.23 IBM Spectrum Virtualize copy services .....	43
1.5.24 Synchronous or asynchronous Remote Copy .....	43
1.5.25 FlashCopy and Transparent Cloud Tiering .....	44
1.6 Business continuity .....	45

1.6.1 Business continuity with Stretched Clusters . . . . .	46
1.6.2 Business continuity with Enhanced Stretched Cluster . . . . .	47
1.6.3 Business continuity with HyperSwap . . . . .	47
1.6.4 Automatic hot spare nodes . . . . .	48
1.7 Management and support tools. . . . .	49
1.7.1 IBM Assist On-site and Remote Support Assistance . . . . .	49
1.7.2 Event notifications. . . . .	50
1.8 Useful IBM SAN Volume Controller web links. . . . .	51
<b>Chapter 2. Planning . . . . .</b>	<b>53</b>
2.1 General planning rules . . . . .	54
2.2 Planning for availability . . . . .	55
2.3 Physical installation planning . . . . .	56
2.4 Planning for system management. . . . .	56
2.5 Connectivity planning . . . . .	57
2.6 Fibre Channel SAN configuration planning. . . . .	58
2.6.1 Physical topology . . . . .	58
2.6.2 Zoning. . . . .	59
2.6.3 N_Port ID Virtualization. . . . .	59
2.6.4 Inter-node zone. . . . .	60
2.6.5 Back-end storage zones . . . . .	60
2.6.6 Host zones . . . . .	60
2.6.7 Zoning considerations for Metro Mirror and Global Mirror . . . . .	61
2.6.8 Port designation recommendations. . . . .	62
2.7 IP SAN configuration planning . . . . .	63
2.7.1 iSCSI and iSER protocols. . . . .	63
2.7.2 Priority flow control . . . . .	64
2.7.3 RDMA clustering . . . . .	65
2.7.4 iSCSI back-end storage attachment . . . . .	65
2.7.5 IP network host attachment . . . . .	66
2.7.6 Native IP replication . . . . .	66
2.7.7 Firewall planning . . . . .	67
2.8 Planning topology . . . . .	67
2.9 Back-end storage configuration . . . . .	68
2.10 Internal storage configuration . . . . .	69
2.11 Storage pool configuration . . . . .	69
2.11.1 The storage pool and cache relationship . . . . .	71
2.12 Volume configuration . . . . .	71
2.12.1 Planning for image mode volumes . . . . .	72
2.12.2 Planning for fully allocated volumes . . . . .	72
2.12.3 Planning for thin-provisioned volumes . . . . .	72
2.12.4 Planning for compressed volumes . . . . .	72
2.12.5 Planning for deduplicated volumes . . . . .	73
2.13 Host attachment planning . . . . .	73
2.13.1 Queue depth . . . . .	74
2.13.2 Microsoft Offloaded Data Transfer . . . . .	74
2.13.3 SAN boot support . . . . .	74
2.13.4 Planning for large deployments . . . . .	74
2.13.5 Planning for SCSI Unmap. . . . .	74
2.14 Planning copy services . . . . .	75
2.14.1 FlashCopy guidelines . . . . .	75
2.14.2 Planning for Metro Mirror and Global Mirror . . . . .	76
2.15 Data migration. . . . .	77

2.16 Performance monitoring with IBM Storage Insights . . . . .	78
2.17 Configuration backup procedure . . . . .	80
<b>Chapter 3. Initial configuration . . . . .</b>	<b>81</b>
3.1 Prerequisites . . . . .	82
3.2 System initialization . . . . .	83
3.2.1 System initialization process . . . . .	84
3.3 System setup . . . . .	87
3.3.1 System setup wizard . . . . .	87
3.4 Base configuration . . . . .	96
3.4.1 Configuring Remote Direct Memory Access clustering . . . . .	96
3.4.2 Adding a node or hot spare node . . . . .	99
3.4.3 Changing the system topology . . . . .	101
3.4.4 Configuring quorum disks or applications . . . . .	104
3.4.5 Configuring the local Fibre Channel port masking . . . . .	106
3.5 Configuring management access . . . . .	108
3.5.1 Configuring secure communications . . . . .	108
3.5.2 Configuring user authentication . . . . .	111
<b>Chapter 4. IBM Spectrum Virtualize GUI . . . . .</b>	<b>119</b>
4.1 Normal operations by using the GUI . . . . .	120
4.1.1 Accessing the GUI . . . . .	120
4.2 Introduction to the GUI . . . . .	124
4.2.1 Task menu . . . . .	124
4.2.2 Suggested tasks . . . . .	125
4.2.3 Notification icons and help . . . . .	126
4.3 System - Overview window . . . . .	129
4.3.1 Content-based organization . . . . .	130
4.4 Monitoring menu . . . . .	135
4.4.1 System overview . . . . .	136
4.4.2 IBM Easy Tier Reports . . . . .	138
4.4.3 Events . . . . .	139
4.4.4 Performance . . . . .	140
4.4.5 Background Tasks . . . . .	141
4.5 Pools . . . . .	142
4.6 Volumes . . . . .	142
4.7 Hosts . . . . .	143
4.8 Copy Services . . . . .	143
4.9 Access . . . . .	144
4.9.1 Ownership groups . . . . .	144
4.9.2 Users by groups . . . . .	150
4.9.3 Audit log . . . . .	154
4.10 Settings . . . . .	156
4.10.1 Notifications menu . . . . .	156
4.10.2 Network . . . . .	159
4.10.3 Using the management GUI . . . . .	162
4.10.4 Security menu . . . . .	168
4.10.5 System menus . . . . .	170
4.10.6 Support menu . . . . .	183
4.10.7 GUI Preferences menu . . . . .	185
4.11 Additional frequent tasks in the GUI . . . . .	187
4.11.1 Renaming components . . . . .	187
4.11.2 Changing the system topology . . . . .	191

4.11.3	Restarting the GUI service .....	197
<b>Chapter 5.</b>	<b>Storage pools .....</b>	<b>199</b>
5.1	Working with storage pools .....	200
5.1.1	Creating storage pools .....	202
5.1.2	Managed disks in a storage pool .....	205
5.1.3	Actions on storage pools .....	205
5.1.4	Child pools .....	212
5.1.5	Encrypted storage pools .....	217
5.2	Working with external controllers and MDisks .....	217
5.2.1	External storage controllers .....	218
5.2.2	Actions on external storage controllers .....	220
5.2.3	Working with external MDisks .....	222
5.2.4	Actions for external MDisks .....	224
5.3	Working with internal drives and arrays .....	231
5.3.1	Working with drives .....	231
5.3.2	RAID and distributed RAID .....	239
5.3.3	Creating arrays .....	243
5.3.4	Actions on arrays .....	248
<b>Chapter 6.</b>	<b>Volumes .....</b>	<b>255</b>
6.1	Introduction to volumes .....	256
6.2	Volume characteristics .....	256
6.2.1	Volume type .....	257
6.2.2	Managed mode and image mode .....	258
6.2.3	Size .....	260
6.2.4	Performance .....	261
6.2.5	Volume copies .....	262
6.2.6	I/O operations data flow .....	264
6.2.7	Storage efficiency .....	266
6.2.8	Encryption .....	271
6.2.9	Cache mode .....	271
6.2.10	I/O throttling .....	272
6.2.11	Volume protection .....	272
6.2.12	Secure data deletion .....	273
6.3	Virtual volumes .....	273
6.4	Volumes in multi-site topologies .....	274
6.4.1	HyperSwap topology .....	275
6.4.2	Stretched topology .....	276
6.5	Operations on volumes .....	276
6.5.1	Creating volumes .....	277
6.5.2	Creating custom volumes .....	285
6.5.3	Creating volumes in multi-site topologies .....	289
6.5.4	I/O throttling .....	295
6.5.5	Volume protection .....	300
6.5.6	Modifying a volume .....	300
6.5.7	Deleting a volume .....	310
6.5.8	Mapping a volume to a host .....	312
6.5.9	Migrating a volume to another storage pool .....	315
6.6	Volume operations by using the CLI .....	322
6.6.1	Displaying volume information .....	322
6.6.2	Creating a volume .....	323
6.6.3	Creating a thin-provisioned volume .....	325

6.6.4	Creating a volume in image mode	326
6.6.5	Adding a volume copy	327
6.6.6	Splitting a mirrored volume	333
6.6.7	Modifying a volume	335
6.6.8	Deleting a volume	335
6.6.9	Volume protection	336
6.6.10	Expanding a volume	337
6.6.11	HyperSwap volume modification with CLI	338
6.6.12	Mapping a volume to a host	339
6.6.13	Listing volumes mapped to the host	340
6.6.14	Listing hosts mapped to the volume	341
6.6.15	Deleting a volume to host mapping	341
6.6.16	Migrating a volume	342
6.6.17	Migrating a fully managed volume to an image mode volume	343
6.6.18	Shrinking a volume	343
6.6.19	Listing volumes using the MDisk	344
6.6.20	Listing MDisks used by the volume	344
6.6.21	Listing volumes defined in the storage pool	345
6.6.22	Listing storage pools in which a volume has its extents	345
6.6.23	Tracing a volume from a host back to its physical disks	347
<b>Chapter 7. Hosts</b>		<b>351</b>
7.1	Host attachment overview	352
7.2	Host clusters	353
7.3	NVMe over Fibre Channel	353
7.4	N_Port ID Virtualization support	354
7.4.1	NPIV prerequisites	357
7.4.2	Enabling NPIV on a new system	357
7.4.3	Enabling NPIV on a system	359
7.5	Hosts operations by using the GUI	365
7.5.1	Creating hosts	365
7.5.2	Host clusters	378
7.5.3	Advanced host administration	382
7.5.4	Adding and deleting host ports	400
7.5.5	Host mappings overview	410
7.5.6	Listing Volumes by Host	412
7.5.7	Listing Volumes by Host Cluster	415
7.6	Performing hosts operations by using the command-line interface	415
7.6.1	Creating a host by using the CLI	415
7.6.2	Performing advanced host administration by using the CLI	418
7.6.3	Adding and deleting a host port by using the CLI	421
7.6.4	Host cluster operations	424
<b>Chapter 8. Storage migration</b>		<b>429</b>
8.1	Storage migration overview	430
8.1.1	Interoperability and compatibility	431
8.1.2	Prerequisites	431
8.2	Storage migration wizard	432
<b>Chapter 9. Advanced features for storage efficiency</b>		<b>449</b>
9.1	IBM Easy Tier	450
9.1.1	Easy Tier concepts	450
9.1.2	Implementing and tuning Easy Tier	456
9.1.3	Monitoring Easy Tier activity	460

9.2 Thin-provisioned volumes . . . . .	466
9.2.1 Concepts . . . . .	467
9.2.2 Implementation . . . . .	467
9.3 Unmap . . . . .	468
9.3.1 SCSI unmap command . . . . .	468
9.3.2 Back-end SCSI unmap . . . . .	469
9.3.3 Host SCSI unmap . . . . .	469
9.3.4 Offload IO throttle . . . . .	470
9.4 Data Reduction Pools . . . . .	471
9.4.1 Introduction to DRP . . . . .	471
9.4.2 DRP benefits . . . . .	472
9.4.3 Planning for DRPs . . . . .	473
9.4.4 Implementing DRP with compression and deduplication . . . . .	475
9.5 Compression with standard pools . . . . .	482
9.5.1 Real-time Compression concepts . . . . .	482
9.5.2 Implementing RtC . . . . .	483
9.6 Saving estimation for compression and deduplication . . . . .	484
9.6.1 Evaluate compression savings by using IBM Comprestimator . . . . .	484
9.6.2 Evaluating compression and deduplication . . . . .	485
9.7 Overprovisioning and data reduction on external storage . . . . .	486
<b>Chapter 10. Advanced Copy Services . . . . .</b>	<b>491</b>
10.1 IBM FlashCopy . . . . .	492
10.1.1 Business requirements for FlashCopy . . . . .	492
10.1.2 FlashCopy principles and terminology . . . . .	494
10.1.3 FlashCopy mapping . . . . .	494
10.1.4 Consistency groups . . . . .	495
10.1.5 Crash consistent copy and hosts considerations . . . . .	496
10.1.6 Grains and bitmap: I/O indirection . . . . .	497
10.1.7 Interaction with cache . . . . .	504
10.1.8 Background Copy Rate . . . . .	504
10.1.9 Incremental FlashCopy . . . . .	506
10.1.10 Starting FlashCopy mappings and consistency groups . . . . .	507
10.1.11 Multiple target FlashCopy . . . . .	509
10.1.12 Reverse FlashCopy . . . . .	514
10.1.13 FlashCopy and image mode Volumes . . . . .	516
10.1.14 FlashCopy mapping events . . . . .	517
10.1.15 Thin-provisioned FlashCopy . . . . .	519
10.1.16 Serialization of I/O by FlashCopy . . . . .	520
10.1.17 Event handling . . . . .	520
10.1.18 Asynchronous notifications . . . . .	521
10.1.19 Interoperation with Metro Mirror and Global Mirror . . . . .	521
10.1.20 FlashCopy attributes and limitations . . . . .	522
10.2 Managing FlashCopy by using the GUI . . . . .	523
10.2.1 FlashCopy presets . . . . .	523
10.2.2 FlashCopy window . . . . .	526
10.2.3 Creating a FlashCopy mapping . . . . .	528
10.2.4 Single-click snapshot . . . . .	538
10.2.5 Single-click clone . . . . .	540
10.2.6 Single-click backup . . . . .	542
10.2.7 Creating a FlashCopy consistency group . . . . .	543
10.2.8 Creating FlashCopy mappings in a Consistency Group . . . . .	544
10.2.9 Showing related Volumes . . . . .	547

10.2.10	Moving FlashCopy mappings across Consistency Groups . . . . .	548
10.2.11	Removing FlashCopy mappings from Consistency Groups . . . . .	549
10.2.12	Modifying a FlashCopy mapping . . . . .	551
10.2.13	Renaming FlashCopy mappings . . . . .	552
10.2.14	Deleting FlashCopy mappings . . . . .	555
10.2.15	Deleting a FlashCopy consistency group . . . . .	556
10.2.16	Starting FlashCopy mappings . . . . .	558
10.2.17	Stopping FlashCopy mappings . . . . .	559
10.2.18	Memory allocation for FlashCopy . . . . .	560
10.3	Transparent Cloud Tiering . . . . .	562
10.3.1	Considerations for using Transparent Cloud Tiering . . . . .	563
10.3.2	Transparent Cloud Tiering as backup solution and data migration . . . . .	563
10.3.3	Restoring data by using Transparent Cloud Tiering . . . . .	564
10.3.4	Transparent Cloud Tiering restrictions . . . . .	564
10.4	Implementing Transparent Cloud Tiering . . . . .	565
10.4.1	Domain Name System configuration . . . . .	565
10.4.2	Enabling Transparent Cloud Tiering . . . . .	566
10.4.3	Creating cloud snapshots . . . . .	569
10.4.4	Managing cloud snapshots . . . . .	572
10.4.5	Restoring cloud snapshots . . . . .	573
10.5	Volume mirroring and migration options . . . . .	576
10.6	Remote Copy . . . . .	578
10.6.1	IBM SAN Volume Controller and IBM Storwize system layers . . . . .	579
10.6.2	Multiple IBM Spectrum Virtualize systems replication . . . . .	580
10.6.3	Importance of write ordering . . . . .	582
10.6.4	Remote Copy intercluster communication . . . . .	584
10.6.5	Metro Mirror overview . . . . .	585
10.6.6	Synchronous Remote Copy . . . . .	586
10.6.7	Metro Mirror features . . . . .	587
10.6.8	Metro Mirror attributes . . . . .	587
10.6.9	Practical use of Metro Mirror . . . . .	588
10.6.10	Global Mirror overview . . . . .	589
10.6.11	Asynchronous Remote Copy . . . . .	589
10.6.12	Global Mirror features . . . . .	591
10.6.13	Using Global Mirror with change volumes . . . . .	593
10.6.14	Distribution of work among nodes . . . . .	595
10.6.15	Background copy performance . . . . .	595
10.6.16	Thin-provisioned background copy . . . . .	596
10.6.17	Methods of synchronization . . . . .	596
10.6.18	Practical use of Global Mirror . . . . .	597
10.6.19	IBM Spectrum Virtualize HyperSwap topology . . . . .	597
10.6.20	Consistency Protection for Global Mirror and Metro Mirror . . . . .	597
10.6.21	Valid combinations of FlashCopy, Metro Mirror, and Global Mirror . . . . .	598
10.6.22	Remote Copy configuration limits . . . . .	598
10.6.23	Remote Copy states and events . . . . .	599
10.7	Remote Copy commands . . . . .	606
10.7.1	Remote Copy process . . . . .	606
10.7.2	Listing available system partners . . . . .	607
10.7.3	Changing the system parameters . . . . .	607
10.7.4	System partnership . . . . .	608
10.7.5	Creating a Metro Mirror/Global Mirror consistency group . . . . .	609
10.7.6	Creating a Metro Mirror/Global Mirror relationship . . . . .	610
10.7.7	Changing Metro Mirror/Global Mirror relationship . . . . .	610

10.7.8	Changing Metro Mirror/Global Mirror consistency group	610
10.7.9	Starting Metro Mirror/Global Mirror relationship	610
10.7.10	Stopping Metro Mirror/Global Mirror relationship	611
10.7.11	Starting Metro Mirror/Global Mirror consistency group	611
10.7.12	Stopping Metro Mirror/Global Mirror consistency group	611
10.7.13	Deleting Metro Mirror/Global Mirror relationship	612
10.7.14	Deleting Metro Mirror/Global Mirror consistency group	612
10.7.15	Reversing Metro Mirror/Global Mirror relationship	612
10.7.16	Reversing Metro Mirror/Global Mirror consistency group	613
10.8	Native IP replication	613
10.8.1	Native IP replication technology	613
10.8.2	IP partnership limitations	615
10.8.3	IP Partnership and data compression	617
10.8.4	VLAN support	617
10.8.5	IP partnership and terminology	618
10.8.6	States of IP partnership	619
10.8.7	Remote Copy groups	620
10.8.8	Supported configurations	621
10.9	Managing Remote Copy by using the GUI	634
10.9.1	Creating Fibre Channel partnership	637
10.9.2	Creating Remote Copy relationships	638
10.9.3	Creating a consistency group	643
10.9.4	Renaming Remote Copy relationships	652
10.9.5	Renaming a Remote Copy consistency group	653
10.9.6	Moving stand-alone Remote Copy relationships to consistency group	655
10.9.7	Removing Remote Copy relationships from consistency group	657
10.9.8	Starting Remote Copy relationships	659
10.9.9	Starting a Remote Copy consistency group	660
10.9.10	Switching a relationship copy direction	660
10.9.11	Switching a consistency group direction	662
10.9.12	Stopping Remote Copy relationships	663
10.9.13	Stopping a consistency group	665
10.9.14	Deleting Remote Copy relationships	666
10.9.15	Deleting a consistency group	667
10.10	Remote Copy memory allocation	669
10.11	Troubleshooting Remote Copy	670
10.11.1	1920 error	670
10.11.2	1720 error	672
<b>Chapter 11</b>	<b>Ownership groups</b>	<b>673</b>
11.1	Ownership groups principles of operations	674
11.2	Implementing ownership groups on a new system	675
11.2.1	Creating an ownership group	675
11.2.2	Assigning users to an ownership group	676
11.2.3	Creating ownership group resources	677
11.2.4	Listing ownership group resources	679
11.2.5	Actions on ownership groups	680
11.3	Migrating objects to ownership groups	680
<b>Chapter 12</b>	<b>Encryption</b>	<b>685</b>
12.1	General types of encryption across IBM Spectrum Virtualize	686
12.1.1	Externally virtualized storage	686
12.1.2	Serial-attached SCSI internal storage	686



12.1.3	Non-Volatile Memory Express internal storage . . . . .	686
12.2	Planning for encryption . . . . .	687
12.3	Defining encryption of data at-rest . . . . .	687
12.3.1	Encryption methods . . . . .	688
12.3.2	Encrypted data . . . . .	688
12.3.3	Encryption keys . . . . .	691
12.3.4	Encryption licenses . . . . .	692
12.4	Activating encryption . . . . .	692
12.4.1	Obtaining an encryption license . . . . .	693
12.4.2	Starting the activation process during initial system setup . . . . .	693
12.4.3	Starting the activation process on a running system . . . . .	696
12.4.4	Activate the license automatically . . . . .	697
12.4.5	Manual license activation . . . . .	701
12.5	Enabling encryption . . . . .	703
12.5.1	Starting the Enable Encryption wizard . . . . .	704
12.5.2	Enabling encryption by using USB flash drives . . . . .	706
12.5.3	Enabling encryption by using key servers . . . . .	711
12.5.4	Enabling encryption by using both providers . . . . .	724
12.6	Configuring more providers . . . . .	729
12.6.1	Adding key servers as a second provider . . . . .	729
12.6.2	Adding USB flash drives as a second provider . . . . .	732
12.7	Migrating between providers . . . . .	733
12.7.1	Changing from USB flash drive provider to encryption key server . . . . .	734
12.7.2	Changing from encryption key server to USB flash drive provider . . . . .	734
12.7.3	Migrating between different key server types . . . . .	735
12.8	Recovering from a provider loss . . . . .	737
12.9	Using encryption . . . . .	737
12.9.1	Encrypted pools . . . . .	738
12.9.2	Encrypted child pools . . . . .	739
12.9.3	Encrypted arrays . . . . .	740
12.9.4	Encrypted MDisk . . . . .	741
12.9.5	Encrypted volumes . . . . .	744
12.9.6	Restrictions . . . . .	746
12.10	Rekeying an encryption-enabled system . . . . .	746
12.10.1	Rekeying by using a key server . . . . .	746
12.10.2	Rekeying by using USB flash drives . . . . .	749
12.11	Disabling encryption . . . . .	752
<b>Chapter 13. Reliability, availability, and serviceability, and monitoring and troubleshooting . . . . .</b>		<b>753</b>
13.1	Reliability, availability, and serviceability . . . . .	754
13.1.1	IBM SAN Volume Controller nodes . . . . .	754
13.1.2	Power . . . . .	759
13.2	Shutting down an IBM SAN Volume Controller cluster . . . . .	759
13.3	Removing or adding a node to or from the system . . . . .	762
13.4	Configuration backup . . . . .	765
13.4.1	Backing up by using the CLI . . . . .	765
13.4.2	Saving the backup by using the GUI . . . . .	767
13.5	Software update . . . . .	769
13.5.1	Precautions before the update . . . . .	769
13.5.2	IBM Spectrum Virtualize upgrade test utility . . . . .	770
13.5.3	Updating IBM Spectrum Virtualize V8.3.1 . . . . .	771
13.5.4	Updating IBM Spectrum Virtualize with a hot spare node . . . . .	779

13.5.5	Updating the IBM SAN Volume Controller system manually . . . . .	780
13.6	Health checker feature . . . . .	781
13.7	Troubleshooting and fix procedures . . . . .	783
13.7.1	Managing event log . . . . .	784
13.7.2	Running a fix procedure . . . . .	786
13.7.3	Resolving alerts in a timely manner . . . . .	788
13.7.4	Event log details . . . . .	788
13.8	Monitoring . . . . .	790
13.8.1	The Call Home function and email notification . . . . .	790
13.8.2	Disabling and enabling notifications . . . . .	797
13.8.3	Remote Support Assistance . . . . .	797
13.8.4	SNMP configuration . . . . .	802
13.8.5	Syslog notifications . . . . .	805
13.9	Audit log . . . . .	806
13.10	Collecting support information by using the GUI, CLI, and USB . . . . .	809
13.10.1	Collecting information by using the GUI . . . . .	809
13.10.2	Collecting logs by using the CLI . . . . .	812
13.10.3	Collecting logs by using USB . . . . .	814
13.10.4	Uploading files to the Support Center . . . . .	815
13.11	Service Assistant Tool . . . . .	816
13.12	IBM Storage Insights Monitoring . . . . .	819
13.12.1	Capacity Monitoring . . . . .	821
13.12.2	Performance monitoring . . . . .	822
13.12.3	Logging Support Tickets by using IBM Storage Insights . . . . .	824
13.12.4	Managing support tickets by using IBM Storage Insights and uploading logs . . . . .	831
	<b>Appendix A. Performance data and statistics gathering . . . . .</b>	<b>833</b>
	IBM SAN Volume Controller performance overview . . . . .	834
	Performance considerations . . . . .	834
	IBM Spectrum Virtualize performance perspectives . . . . .	835
	Performance monitoring . . . . .	836
	Collecting performance statistics . . . . .	836
	Real-time performance monitoring . . . . .	838
	Performance data collection and IBM Spectrum Control . . . . .	846
	<b>Appendix B. CLI setup . . . . .</b>	<b>849</b>
	Setting up the CLI . . . . .	850
	Basic setup on a Windows host . . . . .	851
	Basic setup on a UNIX or Linux host . . . . .	860
	<b>Appendix C. Terminology . . . . .</b>	<b>863</b>
	Commonly used terms . . . . .	864
	<b>Abbreviations and acronyms . . . . .</b>	<b>893</b>
	<b>Related publications . . . . .</b>	<b>897</b>
	IBM Redbooks . . . . .	897
	Help from IBM . . . . .	897

# Notices

This information was developed for products and services offered in the US. This material might be available from IBM in other languages. However, you may be required to own a copy of the product or product version in that language in order to access it.

IBM may not offer the products, services, or features discussed in this document in other countries. Consult your local IBM representative for information on the products and services currently available in your area. Any reference to an IBM product, program, or service is not intended to state or imply that only that IBM product, program, or service may be used. Any functionally equivalent product, program, or service that does not infringe any IBM intellectual property right may be used instead. However, it is the user's responsibility to evaluate and verify the operation of any non-IBM product, program, or service.

IBM may have patents or pending patent applications covering subject matter described in this document. The furnishing of this document does not grant you any license to these patents. You can send license inquiries, in writing, to:

*IBM Director of Licensing, IBM Corporation, North Castle Drive, MD-NC119, Armonk, NY 10504-1785, US*

INTERNATIONAL BUSINESS MACHINES CORPORATION PROVIDES THIS PUBLICATION "AS IS" WITHOUT WARRANTY OF ANY KIND, EITHER EXPRESS OR IMPLIED, INCLUDING, BUT NOT LIMITED TO, THE IMPLIED WARRANTIES OF NON-INFRINGEMENT, MERCHANTABILITY OR FITNESS FOR A PARTICULAR PURPOSE. Some jurisdictions do not allow disclaimer of express or implied warranties in certain transactions, therefore, this statement may not apply to you.

This information could include technical inaccuracies or typographical errors. Changes are periodically made to the information herein; these changes will be incorporated in new editions of the publication. IBM may make improvements and/or changes in the product(s) and/or the program(s) described in this publication at any time without notice.

Any references in this information to non-IBM websites are provided for convenience only and do not in any manner serve as an endorsement of those websites. The materials at those websites are not part of the materials for this IBM product and use of those websites is at your own risk.

IBM may use or distribute any of the information you provide in any way it believes appropriate without incurring any obligation to you.

The performance data and client examples cited are presented for illustrative purposes only. Actual performance results may vary depending on specific configurations and operating conditions.

Information concerning non-IBM products was obtained from the suppliers of those products, their published announcements or other publicly available sources. IBM has not tested those products and cannot confirm the accuracy of performance, compatibility or any other claims related to non-IBM products. Questions on the capabilities of non-IBM products should be addressed to the suppliers of those products.

Statements regarding IBM's future direction or intent are subject to change or withdrawal without notice, and represent goals and objectives only.

This information contains examples of data and reports used in daily business operations. To illustrate them as completely as possible, the examples include the names of individuals, companies, brands, and products. All of these names are fictitious and any similarity to actual people or business enterprises is entirely coincidental.


## COPYRIGHT LICENSE:

This information contains sample application programs in source language, which illustrate programming techniques on various operating platforms. You may copy, modify, and distribute these sample programs in any form without payment to IBM, for the purposes of developing, using, marketing or distributing application programs conforming to the application programming interface for the operating platform for which the sample programs are written. These examples have not been thoroughly tested under all conditions. IBM, therefore, cannot guarantee or imply reliability, serviceability, or function of these programs. The sample programs are provided "AS IS", without warranty of any kind. IBM shall not be liable for any damages arising out of your use of the sample programs.

# Trademarks

IBM, the IBM logo, and ibm.com are trademarks or registered trademarks of International Business Machines Corporation, registered in many jurisdictions worldwide. Other product and service names might be trademarks of IBM or other companies. A current list of IBM trademarks is available on the web at “Copyright and trademark information” at <http://www.ibm.com/legal/copytrade.shtml>

The following terms are trademarks or registered trademarks of International Business Machines Corporation, and might also be trademarks or registered trademarks in other countries.

AIX®	IBM Cloud®	Redbooks®
DB2®	IBM FlashCore®	Redbooks (logo)  ®
DS8000®	IBM FlashSystem®	Storwize®
Easy Tier®	IBM Research®	System Storage™
FICON®	IBM Security™	Tivoli®
FlashCopy®	IBM Spectrum®	XIV®
HyperSwap®	IBM Spectrum Storage™	
IBM®	PowerHA®	

The following terms are trademarks of other companies:

Performance View, are trademarks or registered trademarks of Kenexa, an IBM Company.

SoftLayer, are trademarks or registered trademarks of SoftLayer, Inc., an IBM Company.

Intel, Intel Xeon, Intel logo, Intel Inside logo, and Intel Centrino logo are trademarks or registered trademarks of Intel Corporation or its subsidiaries in the United States and other countries.

The registered trademark Linux® is used pursuant to a sublicense from the Linux Foundation, the exclusive licensee of Linus Torvalds, owner of the mark on a worldwide basis.

Microsoft, Windows, and the Windows logo are trademarks of Microsoft Corporation in the United States, other countries, or both.

Java, and all Java-based trademarks and logos are trademarks or registered trademarks of Oracle and/or its affiliates.

Red Hat, are trademarks or registered trademarks of Red Hat, Inc. or its subsidiaries in the United States and other countries.

UNIX is a registered trademark of The Open Group in the United States and other countries.

VMware, VMware vSphere, and the VMware logo are registered trademarks or trademarks of VMware, Inc. or its subsidiaries in the United States and/or other jurisdictions.

Other company, product, or service names may be trademarks or service marks of others.

# Preface

This IBM® Redbooks® publication is a detailed technical guide to the IBM System Storage™ SAN Volume Controller, which is powered by IBM Spectrum® Virtualize V8.3.1.

IBM SAN Volume Controller is a virtualization appliance solution that maps virtualized volumes that are visible to hosts and applications to physical volumes on storage devices. Each server within the storage area network (SAN) has its own set of virtual storage addresses that are mapped to physical addresses. If the physical addresses change, the server continues running by using the same virtual addresses that it had before. Therefore, volumes or storage can be added or moved while the server is still running.

The IBM virtualization technology improves the management of information at the *block* level in a network, which enables applications and servers to share storage devices on a network.

**Applicability:** This edition applies to IBM Spectrum Virtualize V8.3.1 and the associated hardware and software that are detailed within. The screen captures might differ from the generally available (GA) version because parts of this book were written with pre-GA code. On 11 February 2020, IBM announced the arrival of IBM SAN Volume Controller SA2 and SV2. This book was written specifically for prior versions of IBM SAN Volume Controller, but most of the general principles apply. If you are in any doubt as to their applicability, then you should work with your local IBM representative.

**IBM Knowledge Center:** In this book, we provide links to IBM Knowledge Center and a description of the relevant section that provides more information. Our starting point is the [IBM SAN Volume Controller documentation page](#), and the reader might have to select the product or relevant section that applies to their environment.

## Authors

This book was produced by a team of specialists from around the world working at IBM Redbooks, San Jose Center.



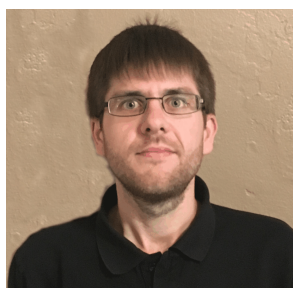
**Jack Armstrong** is a Storage Support Specialist for IBM Systems Group at Hursley, UK. He joined IBM as part of the Apprenticeship Scheme in 2012 and has 7 years of experience working with IBM Storage and providing support to thousands of customers across Europe and beyond. He also works with IBM Enhanced Technical Support Services to help clients to expand and improve their storage environments.



**Tiago Bastos** is a SAN and Storage Disk specialist at IBM Brazil. He has over 20 years in the IT arena, and is an IBM Certified Master IT Specialist. Certified for IBM Storwize®, he works on storage as a service (SaaS) implementation projects. His areas of expertise include planning, configuring, and troubleshooting IBM DS8000®, IBM FlashSystem®, IBM SAN Volume Controller, and IBM XIV®; lifecycle management; and copy services.



**Pawel Brodacki** is an Infrastructure Architect with 20 years of experience in IT who works for IBM Poland since 2003. His main focus for the last 5 years is on virtual infrastructure architecture from storage to servers to software-defined networks (SDNs). Before changing his profession to system architecture, he was an IBM Certified IT Specialist working on various infrastructure, virtualization, and disaster recovery (DR) projects. His experience includes SAN, storage, highly available (HA) systems, DR solutions, IBM System x and IBM Power Systems servers, and several types of operating systems (OSs) (Linux, IBM AIX®, and Microsoft Windows). Pawel has obtained certifications from IBM, Red Hat, and VMware. Pawel holds a master's degree in biophysics from the University of Warsaw College of Inter-Faculty Individual Studies in Mathematics and Natural Sciences.



**Markus Döllinger** is a Software Developer and Support Engineer for IBM Germany. He joined IBM as a student in 2009 and has 8 years of experience as an IBM Storage specialist. He works as an L3 support engineer and developer for the IBM Spectrum Virtualize family of products. His responsibilities include providing storage support for customers worldwide, leading the design and implementation of support tools, and developing IBM Spectrum Virtualize Software. He has an active role in the pro-active support initiative to provide real-time feedback to customers with regards to best practice recommendations based their specific storage environment. He holds a Bachelor of Science degree in Computer Science.



**Jon Herd** is an IBM Storage Technical Advisor working for the Europe, Middle East, and Africa (EMEA) Storage Competence Center (ESCC) in Mainz, Germany. He advises customers in the United Kingdom, Ireland, and Sweden about a portfolio of IBM Storage products, including IBM FlashSystem products. Jon has been with IBM for more than 45 years, and has held various technical roles, including EMEA level 2 support for mainframe servers and technical education development. He has written many IBM Redbooks publications on IBM FlashSystem products, and is an IBM Redbooks Platinum level author. He holds IBM certifications in Product Services at an expert level and Technical IT Specialist at an experienced level. He is the chair of the UKI Professions Board for Product Services. He is a certified Chartered Member of the British Computer Society (MBCS-CITP) and a Certified Member of the Institution of Engineering and Technology (MIET).





**Sergey Kubin** is a subject matter expert (SME) for IBM Storage and SAN technical support. He holds an Electronics Engineer degree from Ural Federal University in Russia and has more than 15 years of experience in IT. At IBM, he works for IBM Technology Support Services, where he provides support and guidance about IBM Spectrum Virtualize family systems for customers in Europe, the Middle East, and Russia. His expertise includes SAN, block-level, and file-level storage systems and technologies. He is IBM Certified Specialist for IBM FlashSystem Family Technical Solutions.



**Carsten Larsen** is an IBM Certified Senior IT Specialist working for the Technical Services Support organization at IBM Denmark, where he delivers consultancy services to IBM clients within the storage arena. Carsten joined IBM in 2007 when he left HP, where he worked with storage arrays and UNIX for 10 years. While working for IBM, Carsten obtained several Brocade and NetApp certifications. Carsten is the author of several IBM Redbooks publications.



**Hartmut Lonzer** is the IBM Storwize Territory Account Manager for Germany (D), Austria (A), and Switzerland (CH) (DACH). Before this position, he was OEM Alliance Manager for Lenovo in IBM Germany. He works at the IBM Germany headquarters in Ehningen.

His main focus is on the IBM FlashSystem Family and the IBM SAN Volume Controller. His experience with the IBM SAN Volume Controller and IBM FlashSystem products goes back to the beginning of these products. Hartmut has been with IBM in various technical and sales roles for 42 years.



**Jon Tate** is a Project Manager for IBM System Storage SAN Solutions at the ITSO, San Jose Center. Before joining the ITSO in 1999, he worked in the IBM Technical Support Center, providing Level 2 and 3 support for IBM mainframe storage products. Jon has 34 years of experience in storage software and management, services, and support. He is an IBM Certified IT Specialist, an IBM SAN Certified Specialist, and is Project Management Professional (PMP) certified. He is also the UK Chairman of the Storage Networking Industry Association (SNIA).

Thanks to the following for their contributions that made this book possible:

Alex Ainscow, Djihed Afifi, Christopher Bulmer, Debbie Butts, Philip Clark, Carlos Fuente, Sally Neate, Evelyn Perez, Suri Polisetti, Matt Smith, Andy Walsh  
**IBM Hursley, UK**

James Whitaker  
**IBM Manchester, UK**

Karen Brown, Mary Connell, Navin Manohar, Terry Niemeyer, Jim Olson, Andy Walls, Brent Yardley  
**IBM US**

Shrirang Bhagwat, Abhishek Jaiswal, Aakanksha Mathur, Sanjay Pathak, Sudharsan Vangal  
**IBM India**

Jorge Enrique Escalante Ramonet  
**IBM Mexico**

Special thanks to the Broadcom Inc. staff in San Jose, California for their support of this residency in terms of equipment and support in many areas:

Sangam Racherla, Brian Steffler, Marcus Thordal  
**Broadcom Inc.**

## Now you can become a published author, too!

Here's an opportunity to spotlight your skills, grow your career, and become a published author—all at the same time! Join an IBM Redbooks residency project and help write a book in your area of expertise, while honing your experience using leading-edge technologies. Your efforts will help to increase product acceptance and customer satisfaction, as you expand your network of technical contacts and relationships. Residencies run from two to six weeks in length, and you can participate either in person or as a remote resident working from your home base.

Find out more about the residency program, browse the residency index, and apply online at:

[ibm.com/redbooks/residencies.html](http://ibm.com/redbooks/residencies.html)

## Comments welcome

Your comments are important to us!

We want our books to be as helpful as possible. Send us your comments about this book or other IBM Redbooks publications in one of the following ways:

- ▶ Use the online **Contact us** review Redbooks form found at:

[ibm.com/redbooks](http://ibm.com/redbooks)

- ▶ Send your comments in an email to:

[redbooks@us.ibm.com](mailto:redbooks@us.ibm.com)

- ▶ Mail your comments to:

IBM Corporation, IBM Redbooks  
Dept. HYTD Mail Station P099  
2455 South Road  
Poughkeepsie, NY 12601-5400



## Stay connected to IBM Redbooks

- ▶ Look for us on LinkedIn:  
<http://www.linkedin.com/groups?home=&gid=2130806>
- ▶ Explore new Redbooks publications, residencies, and workshops with the IBM Redbooks weekly newsletter:  
<https://www.redbooks.ibm.com/Redbooks.nsf/subscribe?OpenForm>
- ▶ Stay current on recent Redbooks publications with RSS Feeds:  
<http://www.redbooks.ibm.com/rss.html>





# Introduction and system overview

This chapter defines the concept of *storage virtualization* and provides an overview of its application in addressing the challenges of modern storage environment. It also contains an overview of each of the products that make up the IBM SAN Volume Controller family.

This chapter includes the following topics:

- ▶ 1.1, “Storage virtualization terminology” on page 2
- ▶ 1.2, “Latest changes and enhancements” on page 5
- ▶ 1.3, “IBM SAN Volume Controller architecture” on page 7
- ▶ 1.4, “IBM SAN Volume Controller models” on page 14
- ▶ 1.5, “IBM SAN Volume Controller components” on page 20
- ▶ 1.6, “Business continuity” on page 45
- ▶ 1.7, “Management and support tools” on page 49
- ▶ 1.8, “Useful IBM SAN Volume Controller web links” on page 51

## 1.1 Storage virtualization terminology

*Storage virtualization* is a term that is used extensively throughout the storage industry. It can be applied to various technologies and underlying capabilities. In reality, most storage devices technically can claim to be virtualized in one form or another. Therefore, this chapter starts by defining the concept of storage virtualization as it is used in this book.

We describe storage virtualization in the following way:

- ▶ Storage virtualization is a technology that makes one set of resources resemble another set of resources, preferably with more wanted characteristics.
- ▶ Storage virtualization is a logical representation of resources that is not constrained by physical limitations and hides part of the complexity. It also adds or integrates new functions with services, and can be nested or applied to multiple layers of a system.

The virtualization model consists of the following layers:

- ▶ Application: The user of the storage domain.
- ▶ Storage domain:
  - File, record, and namespace virtualization and file and record subsystem
  - Block virtualization
  - Block subsystem

Applications typically read and write data as vectors of bytes or records. However, storage presents data as vectors of blocks of a constant size (512 or in the newer devices, 4096 bytes per block).

The *file, record, and namespace virtualization* and *file and record subsystem* layers convert records or files that are required by applications to vectors of blocks, which are the language of the *block virtualization* layer. The block virtualization layer maps requests of the higher layers to physical storage blocks, which are provided by *storage devices* in the *block subsystem*.

Each of the layers in the storage domain abstracts away complexities of the lower layers and hides them behind an easy-to-use, standard interface that is presented to upper layers. The resultant decoupling of logical storage space representation and its characteristics that are visible to servers (storage consumers) from underlying complexities and intricacies of storage devices is a key concept of storage virtualization.

The focus of this publication is *block-level virtualization* at the *block virtualization layer*, which is implemented by IBM as IBM Spectrum Virtualize Software that is running on IBM SAN Volume Controller, and the IBM FlashSystem and IBM Storwize families. The IBM SAN Volume Controller is implemented as a clustered appliance in the storage network layer. The IBM FlashSystem and IBM Storwize families are deployed as modular storage systems that can virtualize their internally and externally attached storage.

IBM Spectrum Virtualize uses the Small Computer System Interface (SCSI) protocol to communicate with its clients and presents storage space as SCSI logical units (LUs), which are identified by SCSI logical unit numbers (LUNs).

**Note:** Although LUs and LUNs are different entities, the term LUN in practice is often used to refer to a logical disk, that is, an LU.

Although most applications do not directly access storage but work with files or records, the operating system (OS) of a host must convert these abstractions to the language of storage, that is, vectors of storage blocks that are identified by logical block addresses (LBAs) within an LU.

Inside IBM Spectrum Virtualize, each of the externally visible LUs is internally represented by a volume, which is an amount of storage that is taken out of a storage pool. Storage pools are made of managed disks (MDisks), that is, they are LUs that are presented to the storage system by external virtualized storage or arrays that consist of internal disks. LUs that are presented to IBM Spectrum Virtualize by external storage usually correspond to RAID arrays that are configured on that storage.

The hierarchy of objects, from a file system block down to a physical block on a physical drive, is shown in Figure 1-1.

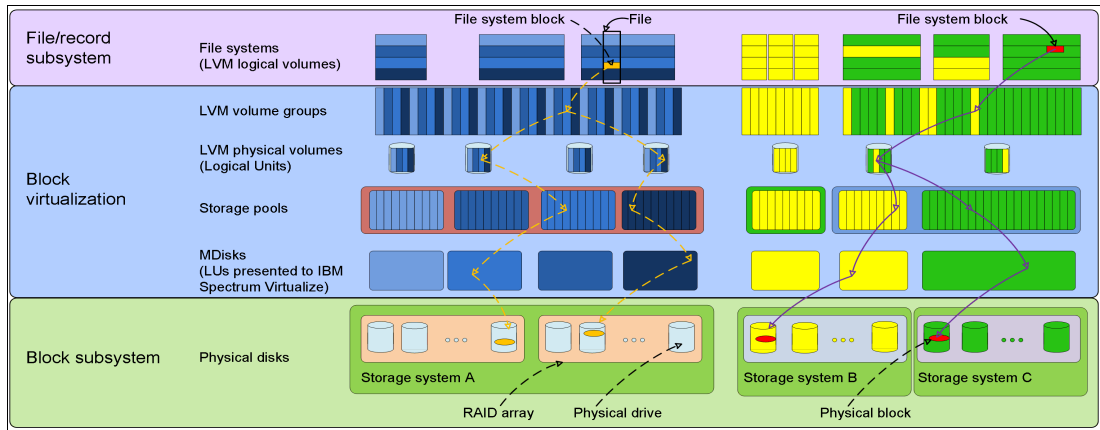


Figure 1-1 Block-level virtualization overview

With storage virtualization, you can manage the mapping between logical blocks within an LU that is presented to a host, and blocks on physical drives. This mapping can be as simple or as complicated as required by a use case. A logical block can be mapped to one physical block or for increased availability, multiple blocks that are physically stored on different physical storage systems, and in different geographical locations.

Importantly, the mapping can be dynamic: With IBM Easy Tier, IBM Spectrum Virtualize can automatically change underlying storage to which groups of blocks (extent) are mapped to better match a host’s performance requirements with the capabilities of the underlying storage systems.

IBM Spectrum Virtualize gives a storage administrator a wide range of options to modify volume characteristics: from volume resize to mirroring, creating a point-in-time (PiT) copy with IBM FlashCopy®, and migrating data across physical storage systems. Importantly, all the functions that are presented to the storage users are independent from the characteristics of the physical devices that are used to store data. This decoupling of the storage feature set from the underlying hardware and ability to present a single, uniform interface to storage users that masks underlying system complexity is a powerful argument for adopting storage virtualization with IBM Spectrum Virtualize.

Storage virtualization is implemented on many layers. Figure 1-1 on page 3 shows an example where a file system block is mirrored by the host's OS (left side of the figure) by using features of the logical volume manager (LVM) or the IBM Spectrum Virtualize system at the storage pool level (as shown on the right side of Figure 1-1 on page 3). Although the result is similar (the data block is written to two different arrays), the effort that is required for per-host configuration is disproportionately larger than for a centralized solution with organization-wide storage virtualization that is done on a dedicated system and managed from a single GUI.

IBM Spectrum Virtualize includes the following key features:

- ▶ Simplified storage management by providing a single management interface for multiple storage systems and a consistent user interface for provisioning heterogeneous storage.
- ▶ Online volume migration. IBM Spectrum Virtualize enables moving the data from one set of physical drives to another set in a way that is not apparent to the storage users and without over-straining the storage infrastructure. The migration can be done within a specific storage system (from one set of disks to another set) or across storage systems. Either way, the host that uses the storage is not aware of the operation, and no downtime for applications is needed.
- ▶ Enterprise-level Copy Services functions. Performing Copy Services functions within IBM Spectrum Virtualize removes dependencies on the capabilities and intercompatibility of the virtualized storage subsystems. Therefore, it enables the source and target copies to be on any two virtualized storage subsystems.
- ▶ Improved storage space usage because of the pooling of resources across virtualized storage systems.
- ▶ Opportunity to improve system performance as a result of volume striping across multiple virtualized arrays or controllers, and the benefits of cache that is provided by IBM Spectrum Virtualize hardware.
- ▶ Improved data security by using data-at-rest encryption.
- ▶ Data replication, including replication to cloud storage by using advanced copy services for data migration and backup solutions.
- ▶ Data reduction techniques for space efficiency, such as thin provisioning, Data Reduction Pools (DRPs), deduplication, and IBM Real-time Compression (RtC) Appliance. Today, open systems typically use less than 50% of the provisioned storage capacity. IBM Spectrum Virtualize can enable significant savings, increase the effective capacity of storage systems up to five times, and decrease the floor space, power, and cooling that are required by the storage system.

**Note:** RtC is available only for earlier generation IBM Storwize V5000 and V7000 systems. The newer IBM SAN Volume Controllers, the IBM FlashSystem 9100, 9200, 7200 and the IBM Storwize 5100 and V7000 Gen3 systems, do not support RtC. They support software compression only through DRPs.

IBM Storwize and IBM FlashSystem families are scalable solutions running on a (HA) platform that can use diverse back-end storage systems to provide all the benefits to various attached hosts.

## Summary

Storage virtualization is a fundamental technology that enables the realization of flexible and reliable storage solutions. It helps enterprises to better align IT architecture with business requirements, simplify their storage administration, and facilitate their IT departments efforts to meet business demands.

IBM Spectrum Virtualize running on IBM Storwize and IBM FlashSystem families is a mature, 10th-generation virtualization solution that uses open standards and complies with the SNIA storage model. All the products are appliance-based storage, and use in-band block virtualization engines that move the control logic (including advanced storage functions) from a multitude of individual storage devices to a centralized entity in the storage network.

IBM Spectrum Virtualize can improve the usage of your storage resources, simplify storage management, and improve the availability of business applications.

## 1.2 Latest changes and enhancements

IBM Spectrum Virtualize V8.3.1 is another step in the product line development that brings new features and enhancements. The major software changes in subsequent code releases are described in this section.

The IBM Spectrum Virtualize V8.3.1 code brought the following changes:

- ▶ Ownership groups:
  - An ownership group defines a subset of users and objects within the system. You can create ownership groups to further restrict access to specific resources that are defined in the ownership group. Only users with Security Administrator roles can configure and manage ownership groups.
  - Ownership groups restrict access for users in the ownership group to only those objects that are defined within that ownership group. An owned object can belong to one ownership group. Users in an ownership group are restricted to viewing and managing objects within their ownership group. Users that are not in an ownership group can continue to view or manage all the objects on the system based on their defined user role, including objects within ownership groups. When the user within an ownership group logs on to the management GUI or command-line interface (CLI), only resources that they have access through the ownership group are available.

The following objects can be assigned to ownership groups:

- Child pools
  - Volumes
  - Volume groups
  - Hosts
  - Host clusters
  - Host mappings
  - FlashCopy mappings
  - FlashCopy consistency groups
  - User groups
- ▶ Priority flow control (PFC).

PFC is an Ethernet protocol that supports the ability to select the priority of different types of traffic within the network. With PFC, administrators can reduce network congestion by slowing or pausing certain classes of traffic on ports, thus providing better bandwidth for more important traffic. The system supports PFC on various supported Ethernet-based protocols on three types of traffic classes: system, host attachment, and storage traffic.

- ▶ Support for Easy Tier overallocation limit for pools with IBM FlashCore® Module devices as the top tier of storage.

The system includes Easy Tier, which is a function that responds to the presence of drives in a storage pool that also contains hard disk drives (HDDs). The system automatically and nondisruptively moves frequently accessed data from HDD MDisks to flash-based storage MDisks, thus placing such data in a faster tier of storage.

The system supports these tiers:

- Storage-class memory (SCM).

SCM tier exists when the pool contains drives that use persistent memory technologies that improve the endurance and speed of current flash storage device technologies.

- Tier 0 flash.

The tier 0 flash tier exists when the pool contains high-performance flash drives.

- Tier 1 flash.

The tier 1 flash tier exists when the pool contains tier 1 flash drives. Tier 1 flash drives typically offer larger capacities, but slightly lower performance and write endurance characteristics.

- Enterprise tier.

The enterprise tier exists when the pool contains enterprise-class MDisks, which are disk drives that are optimized for performance.

- Nearline (NL) tier.

The NL tier exists when the pool contains NL-class MDisks, which are disk drives that are optimized for capacity.

- ▶ Easy Tier automatic data placement requirements, recommendations, and limitations.
- ▶ A hardware upgrade from a IBM FlashSystem 9100 Model AF7 to a Model AF8 is now available.
- ▶ Support for 32 GB Fibre Channel (FC) PCIe adapters.
- ▶ Support for expanding distributed arrays.

You can dynamically expand distributed arrays to increase the available capacity of the array or create more rebuild space. As part of the expansion, the system automatically migrates data for optimal performance for the new expanded configuration. Expansion of distributed arrays support the incremental growth of the available capacity for arrays and are compatible with other functions, such as Easy Tier and data migrations.

- ▶ Support for pool-level volume protection.

Volume protection prevents active volumes or host mappings from being deleted inadvertently if the system detects recent I/O activity. This global setting is enabled by default on new systems. You can either set this value to apply to all volumes that are configured on your system, or control whether the system-level volume protection is enabled or disabled on specific pools.

- ▶ Support for Simple Network Management Protocol (SNMP) protocol Version 3 enhanced security features:
  - SNMP is a standard protocol for managing networks and exchanging messages. The system can send SNMP messages that notify personnel about an event. You can use an SNMP manager to view the SNMP messages that the system sends. The system supports both SNMP Version 2 and Version 3.



- Some systems support setting up SNMP notifications for events. Event notifications are reported to the SNMP destinations of your choice. To specify an SNMP destination, you must provide a valid Internet Protocol (IP) address. A maximum of six SNMP destinations can be specified. For SNMP Version 2 servers, the community string is required and the default value is public. You can use the Management Information Base (MIB) file for SNMP to configure a network management program to receive SNMP messages that are sent by the system. This file can be used with SNMP messages from all versions of the software.
- ▶ Support for enhanced auditing features for syslog servers.
 

The syslog protocol is a standard protocol for forwarding log messages from a sender to a receiver on an IP network. The system can send syslog messages that notify personnel about an event. You can set up syslog event notifications with either the management GUI or the CLI.
- ▶ Enhanced password security.
 

The user must change the default password to a different password for the first login or system setup. If the password is ever reset to the default password, the user must immediately change the password to a different password.
- ▶ Improvements to the terms and definitions that relate to capacity were updated.
- ▶ Support for the new SCM technology (such as Optane and 3DXP).
- ▶ Three-site replication with limited availability at general availability (GA) and subject to RPQ initially.
 

Data is replicated from the primary site to two alternative sites. The remaining two sites are aware of the difference between themselves, which ensures that if there is a disaster at any one of the sites, the remaining two sites can establish a consistent\_synchronized Remote Copy (RC) relationship among themselves with minimal data transfer, that is, within the expected recovery point objective (RPO).
- ▶ Secure Drive Erase is the ability to erase any customer data from a Non-Volatile Memory Express (NVMe) or serial-attached SCSI (SAS) solid-state drive (SSD) drive before it is removed from either the control and expansion enclosures.

**Note:** Any references to internal drives or internal drive functions do not apply to SAN Volume Controller Models SA2 and SA2 because these models do not support internal drives. These functions might apply to externally virtualized drives and arrays, and thus be available.

## 1.3 IBM SAN Volume Controller architecture

This section explains the major concepts underlying SAN Volume Controller and presents a brief history of the product. Also, it describes the architectural overview and the terminologies that are used in a virtualized storage environment. Finally, it introduces the software and hardware components and the other functions that are available with Version 8.3.1.

All of the concepts included in this chapter are described in more detail in later chapters.

### 1.3.1 Brief history of IBM SAN Volume Controller

SAN Volume Controller (machine type (MT) 2145 / 2147) and its embedded software engine (IBM Spectrum Virtualize) are based on an IBM project that was started in the second half of 1999 at the IBM Almaden Research Center. The project was called COMmodity PARTs Storage System (COMPASS). However, most of the software was developed at IBM Hursley Labs in the UK.

One goal of this project was to create a system that was almost exclusively composed of commercial off the shelf (COTS) standard parts. As with any enterprise-level storage control system, it had to deliver a level of performance and availability that was comparable to the highly optimized storage controllers of previous generations. The idea of building a storage control system that is based on a scalable cluster of lower performance servers rather than a monolithic architecture of two nodes is still a compelling idea.

COMPASS also had to address a major challenge for the heterogeneous open systems environment: To reduce the complexity of managing storage on block devices.

The first documentation that covered this project was released to the public in 2003 in the article “The software architecture of a SAN storage control system”, which you can find at [IBM Systems Journal, Vol 42 No 2, 2003](#).

The results of the COMPASS project defined the fundamentals for the product architecture. The first release of IBM System Storage SAN Volume Controller was announced in July 2003.

Each release brought new and more powerful hardware nodes, which approximately doubled the I/O performance and throughput of its predecessors, provided new functions, and offered more interoperability with new elements in host environments, disk subsystems, and the storage area network (SAN).

The most recently (at the time of writing) released hardware nodes, IBM SAN Volume Controller 2145 / 2147-SV2 and SA2, are based on a two Intel Xeon Cascade Lake 16- and 8-core processors configurations respectively.

### 1.3.2 IBM SAN Volume Controller architectural overview

SAN Volume Controller is a SAN block aggregation virtualization appliance that is designed for attachment to various host computer systems.

The following major approaches are used today for the implementation of block-level aggregation and virtualization:

- ▶ Symmetric: In-band appliance

Virtualization splits the storage that is presented by the storage systems into smaller chunks that are known as *extents*. These extents are then concatenated by using various policies to make virtual disks (VDisks) (*volumes*). With symmetric virtualization, host systems can be isolated from the physical storage. Advanced functions, such as data migration, can run without reconfiguring the host.

With symmetric virtualization, the *virtualization engine* is the central configuration point for the SAN. The virtualization engine directly controls access to the storage and to the data that is written to the storage. As a result, locking functions that provide data integrity and advanced functions (such as cache and Copy Services) can be run in the virtualization engine itself.

Therefore, the virtualization engine is a central point of control for device and advanced function management. Symmetric virtualization enables you to build a firewall in the storage network. Only the virtualization engine can grant access through the firewall.

Symmetric virtualization can have disadvantages. The main disadvantage that is associated with symmetric virtualization is *scalability*. Scalability can cause poor performance because all input/output (I/O) must flow through the virtualization engine. To solve this problem, you can use an *n*-way cluster of virtualization engines that has failover capacity.

You can scale the extra processor power, cache memory, and adapter bandwidth to achieve the level of performance that you want. More memory and processing power are needed to run advanced services, such as Copy Services and caching. SAN Volume Controller uses symmetric virtualization. Single virtualization engines, which are known as *nodes*, are combined to create clusters. Each cluster can contain 2 - 8 nodes.

► **Asymmetric: Out-of-band or controller-based**

With asymmetric virtualization, the virtualization engine is outside the data path and performs a metadata-style service. The metadata server contains all of the mapping and the locking tables, and the storage devices contain only data. In asymmetric virtual storage networks, the data flow is separated from the control flow.

A separate network or SAN link is used for control purposes. Because the control flow is separated from the data flow, I/O operations can use the full bandwidth of the SAN. A separate network or SAN link is used for control purposes.

Asymmetric virtualization can have the following disadvantages:

- Data is at risk to increased security exposures, and the control network must be protected with a firewall.
- Metadata can become complicated when files are distributed across several devices.
- Each host that accesses the SAN must know how to access and interpret the metadata. Therefore, specific device drivers or agent software must be running on each of these hosts.
- The metadata server cannot run advanced functions such as caching or Copy Services because it only “knows” about the metadata and not about the data itself.

Figure 1-2 shows variations of the two virtualization approaches.

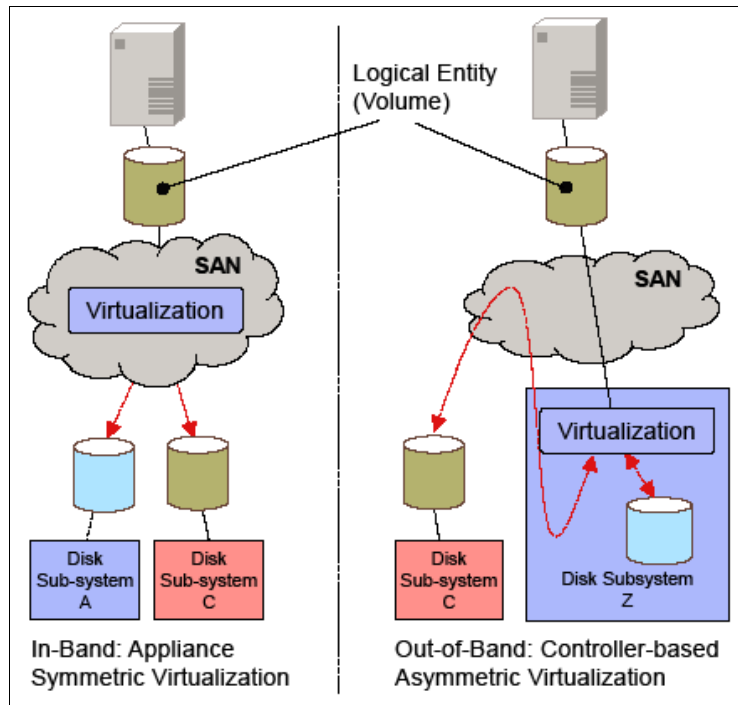


Figure 1-2 Overview of block-level virtualization architectures

Although these approaches provide essentially the same cornerstones of virtualization, interesting side-effects can occur.

The controller-based approach has high functionality, but it fails in terms of scalability or upgradeability. Because of the nature of its design, no true decoupling occurs with this approach, which becomes an issue for the lifecycle of this solution, such as with a controller. Data migration issues and questions are challenging, such as how to reconnect the servers to the new controller, and how to reconnect them online without any effect on your applications.

With this approach, you replace a controller and implicitly replace your entire virtualization solution. In addition to replacing the hardware, other actions (such as updating or repurchasing the licenses for the virtualization feature, and advanced copy functions) might be necessary.

With a SAN or fabric-based appliance solution that is based on a scale-out cluster architecture, lifecycle management tasks, such as adding or replacing new disk subsystems or migrating data between them, are simple. Servers and applications remain online, data migration occurs transparently on the virtualization platform, and licenses for virtualization and copy services require no update. There are no other costs when disk subsystems are replaced.

Only the fabric-based appliance solution provides an independent and scalable virtualization platform that can provide enterprise-class Copy Services that is open for future interfaces and protocols. By using the fabric-based appliance solution, you can choose the disk subsystems that best fit your requirements, and you are not locked into specific SAN hardware.

For these reasons, IBM chose the SAN-based appliance approach with inline block aggregation for the implementation of storage virtualization with IBM Spectrum Virtualize.

SAN Volume Controller includes the following key characteristics:

- ▶ It is highly scalable, which provides an easy growth path to two-*n* nodes (grow in a pair of nodes due to the cluster function).
- ▶ It is SAN interface-independent. It supports FC, Fibre Channel over Ethernet (FCoE), and internet Small Computer Systems Interface (iSCSI), but it is also open for future enhancements.
- ▶ It is host-independent for fixed block-based Open Systems environments.
- ▶ It is external storage RAID controller-independent, which provides a continuous and ongoing process to qualify more types of controllers.
- ▶ It can use disks that are internal disks that are attached to the nodes (flash drives) or externally direct-attached in expansion enclosures.

On the SAN storage that is provided by the disk subsystems, SAN Volume Controller offers the following services:

- ▶ Creates a single pool of storage.
- ▶ Provides LU virtualization.
- ▶ Manages logical volumes.
- ▶ Mirrors logical volumes.

SAN Volume Controller running IBM Spectrum Virtualize V8.3.1 also provides these functions:

- ▶ Large scalable cache.
- ▶ Copy Services.
- ▶ IBM FlashCopy (PiT copy) function, including thin-provisioned FlashCopy to make multiple targets affordable.
- ▶ IBM Transparent Cloud Tiering (TCT) function that enables SAN Volume Controller to interact with cloud service providers (CSPs).
- ▶ Metro Mirror (MM) (synchronous copy).
- ▶ Global Mirror (GM) (asynchronous copy).
- ▶ Data migration.
- ▶ Space management (Thin-provisioning and compression).
- ▶ Easy Tier to automatically migrate data between storage types of different performance, based on disk workload.
- ▶ Encryption of external attached storage.
- ▶ Supporting IBM HyperSwap®.
- ▶ Supporting VMware vSphere Virtual Volumes (VVOLs) and Microsoft Offloaded Data Transfer (ODX).
- ▶ Direct attachment of hosts.
- ▶ Hot spare nodes with a standby function of single or multiple nodes.

### 1.3.3 IBM Spectrum Virtualize

IBM Spectrum Virtualize is a key member of the IBM Spectrum Storage™ portfolio. It is a software-enabled storage virtualization engine that provides a single point of control for storage resources within the data centers. IBM Spectrum Virtualize is a core software engine of well-established and industry-proven IBM storage virtualization solutions, such as SAN Volume Controller, the IBM Storwize family (IBM Storwize V5000 and IBM Storwize V7000), IBM FlashSystem 5100, IBM FlashSystem 7200, IBM FlashSystem V9000, and IBM FlashSystem 9100 and 9200.

**For more information:** For more information about the IBM Spectrum Storage portfolio, see [IBM Spectrum Storage Portfolio](#).

**Naming:** With the introduction of the IBM Spectrum Storage family, the *software* that runs on SAN Volume Controller and IBM Storwize family products is called IBM Spectrum Virtualize. The name of the underlying *hardware* platform remains intact.

The objectives of IBM Spectrum Virtualize are to manage storage resources in your IT infrastructure and protect huge volumes of data that organizations use for several types of workloads. In addition, a goal is to ensure that the resources and data are used to the advantage of your business. These processes take place quickly, efficiently, and in real time, while avoiding increases in administrative costs.

IBM Spectrum Virtualize is the core software engine of the whole family of IBM Storwize products, and the contents of this book are intentionally related to the deployment considerations of SAN Volume Controller.

**Terminology note:** In this book, the term *IBM SAN Volume Controller* is used to refer to both models of the most recent products because the text applies to both.

### 1.3.4 IBM SAN Volume Controller topology

SAN-based storage can be managed by SAN Volume Controller in one or more pairs of hardware nodes. This configuration is referred to as a *clustered system*. These nodes are normally attached to the SAN fabric, with RAID controllers and host systems. The SAN fabric is zoned to enable the SAN Volume Controller to “see” the RAID storage controllers, and for the hosts to communicate with SAN Volume Controller.

Within this software release, SAN Volume Controller also supports Internet Protocol networks. This feature enables the hosts and storage controllers to communicate with SAN Volume Controller to build a storage virtualization solution.

Typically, the hosts cannot see or operate on the same physical storage (LUN) from the RAID controller that is assigned to SAN Volume Controller. If the same LUNs are not shared, storage controllers can be shared between the SAN Volume Controller and direct host access. The zoning capabilities of the SAN switch must be used to create distinct zones to ensure that this rule is enforced. SAN fabrics can include standard FC, FCoE, iSCSI over Ethernet, or possible future types.

Figure 1-3 shows a conceptual diagram of a storage system that uses SAN Volume Controller. It shows several hosts that are connected to a SAN fabric or local area network (LAN). In practical implementations that have HA requirements (most of the target clients for SAN Volume Controller), the SAN fabric cloud represents a redundant SAN. A *redundant SAN* consists of a fault-tolerant arrangement of two or more counterpart SANs, which provide alternative paths for each SAN-attached device.

Both scenarios (the use of a single network and the use of two physically separate networks) are supported for iSCSI-based and LAN-based access networks to SAN Volume Controller. Redundant paths to volumes can be provided in both scenarios. For simplicity, Figure 1-3 shows only one SAN fabric and two zones: host and storage. In a real environment, it is a best practice to use two redundant SAN fabrics. SAN Volume Controller can be connected to up to four fabrics.

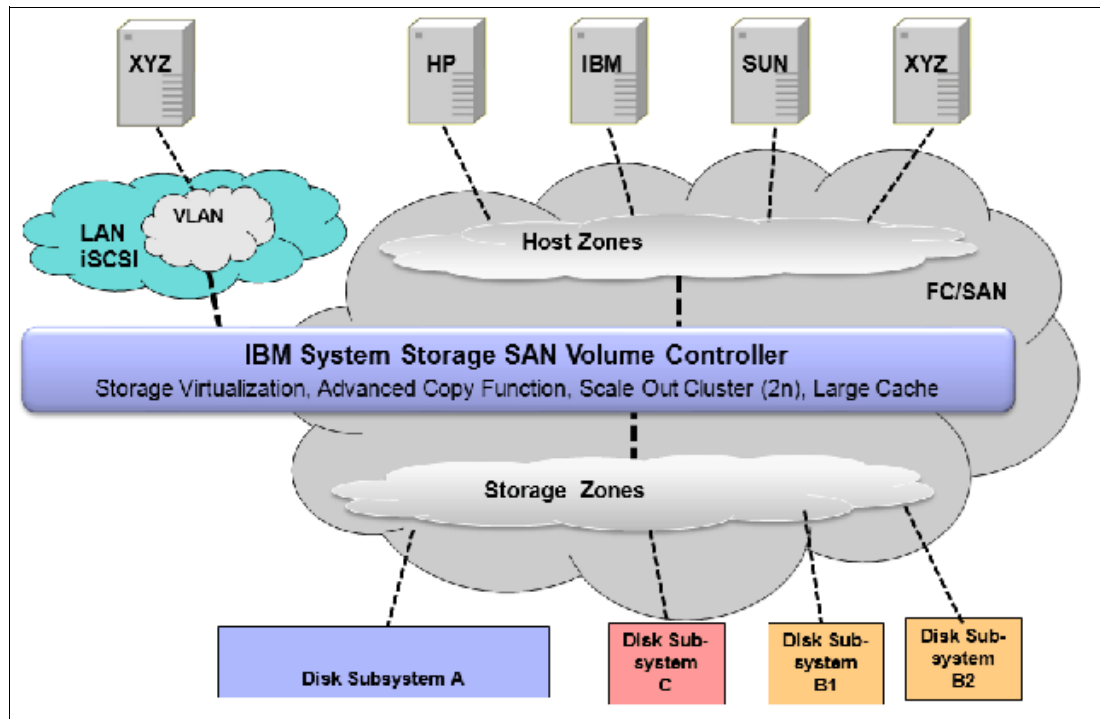


Figure 1-3 SAN Volume Controller conceptual and topology overview

A clustered system of SAN Volume Controller nodes that are connected to the same fabric presents *logical disks* or volumes to the hosts. These volumes are created from managed LUNs or MDisks that are presented by the RAID disk subsystems.

The following distinct zones are shown in the fabric:

- ▶ A host zone, in which the hosts can see and address the SAN Volume Controller nodes
- ▶ A storage zone, in which the SAN Volume Controller nodes can see and address the MDisks or LUNs that are presented by the RAID subsystems

As explained in 1.3.2, “IBM SAN Volume Controller architectural overview” on page 8, hosts are not permitted to operate on the RAID LUNs directly. All data transfer happens through the SAN Volume Controller nodes. This flow is referred to as *symmetric virtualization*.

For iSCSI-based access, the use of two networks and separating iSCSI traffic within the networks by using a dedicated virtual local area network (VLAN) path for storage traffic prevents any IP interface, switch, or target port failure from compromising the iSCSI connectivity across servers and storage controllers.

## 1.4 IBM SAN Volume Controller models

The following models of SAN Volume Controller are supported at the IBM Spectrum Virtualize V8.3.1 code level.

- ▶ IBM SAN Volume Controller Model DH8
- ▶ IBM SAN Volume Controller Model SV1
- ▶ IBM SAN Volume Controller Model SA2
- ▶ IBM SAN Volume Controller Model SV2

### 1.4.1 IBM SAN Volume Controller DH8

Figure 1-4 shows the front view of the IBM SAN Volume Controller DH8.



Figure 1-4 IBM SAN Volume Controller DH8 front view

Figure 1-5 shows the rear view of the IBM SAN Volume Controller DH8.

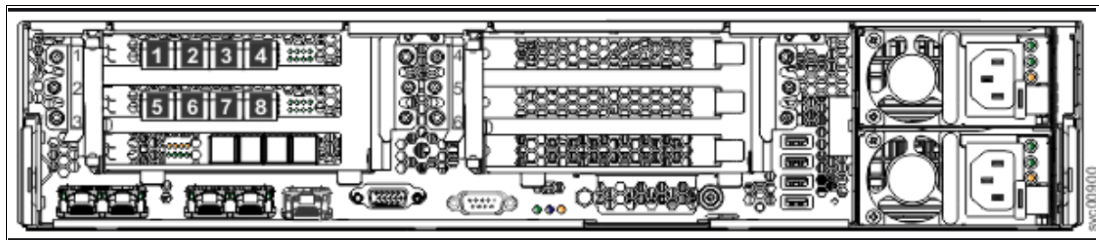


Figure 1-5 IBM SAN Volume Controller DH8 rear view

### 1.4.2 IBM SAN Volume Controller SV1

Figure 1-6 shows the front view of the IBM SAN Volume Controller SV1.



Figure 1-6 IBM SAN Volume Controller SV1 front view



Figure 1-7 shows the rear view of the IBM SAN Volume Controller SV1.

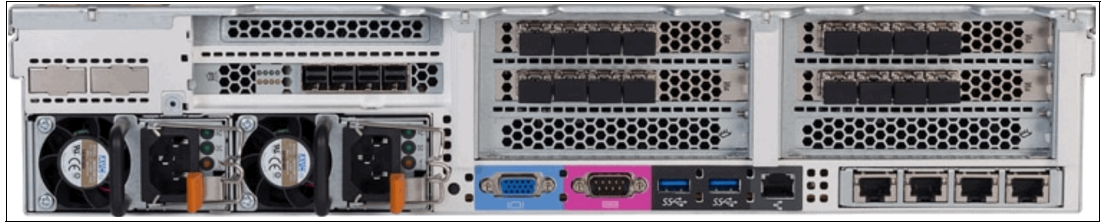


Figure 1-7 IBM SAN Volume Controller SV1 rear view

### 1.4.3 IBM SAN Volume Controller SV2 and SA2

Figure 1-8 shows the front view of the IBM SAN Volume Controller SV2 and SA2.



Figure 1-8 IBM SAN Volume Controller SV2 and SA2 front view

Figure 1-9 shows the rear view of the IBM SAN Volume Controller SV2 / SA2.

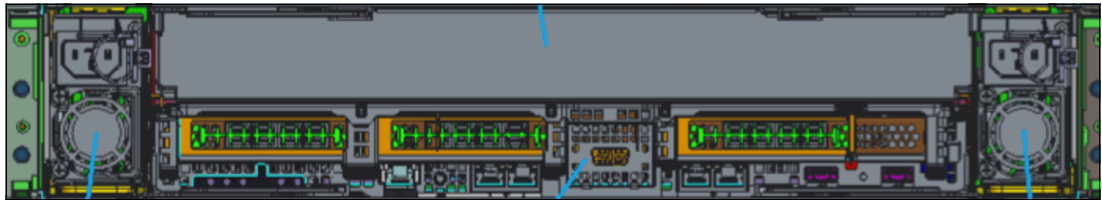


Figure 1-9 IBM SAN Volume Controller SV2 and SA2 rear view

Figure 1-10 shows a picture of internal hardware components of a SAN Volume Controller SV2 and SA2 node canister. To the left of the picture is the front of the canister where fan modules and battery backup are located, followed by two Cascade Lake CPUs and Dual Inline Memory Module (DIMM) slots and PCIe risers for adapters on the right.

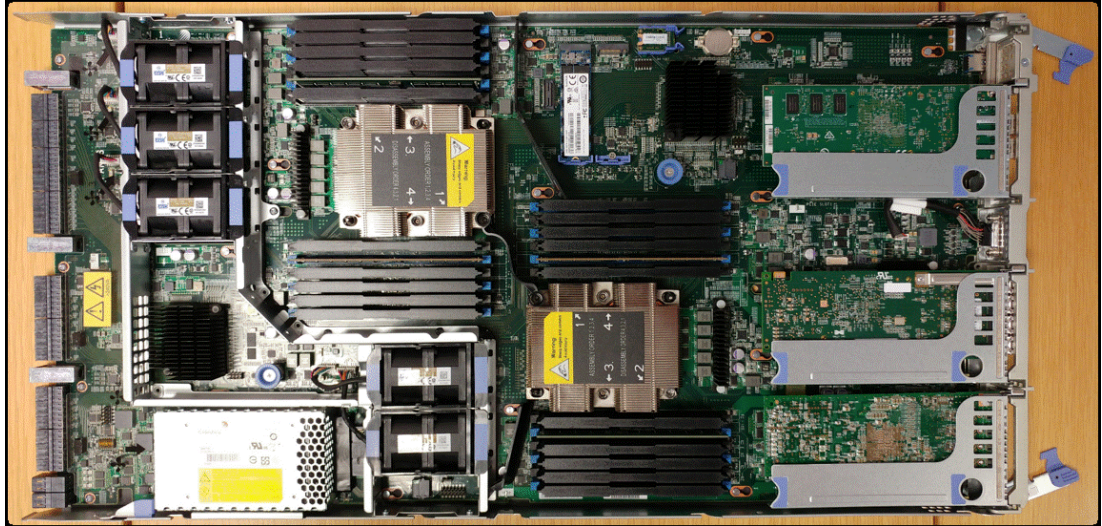


Figure 1-10 Internal hardware components

The SAN Volume Controller cluster consists of nodes that can virtualize external disk subsystems.

Figure 1-11 shows the internal architecture on the IBM SAN Volume Controller SV2 and SA2 models. You can see that the PCIe switch is still present, but has no outbound connections because these models do not support any internal drives. The PCIe switch is used for internal monitoring purposes within the SAN Volume Controller enclosure.

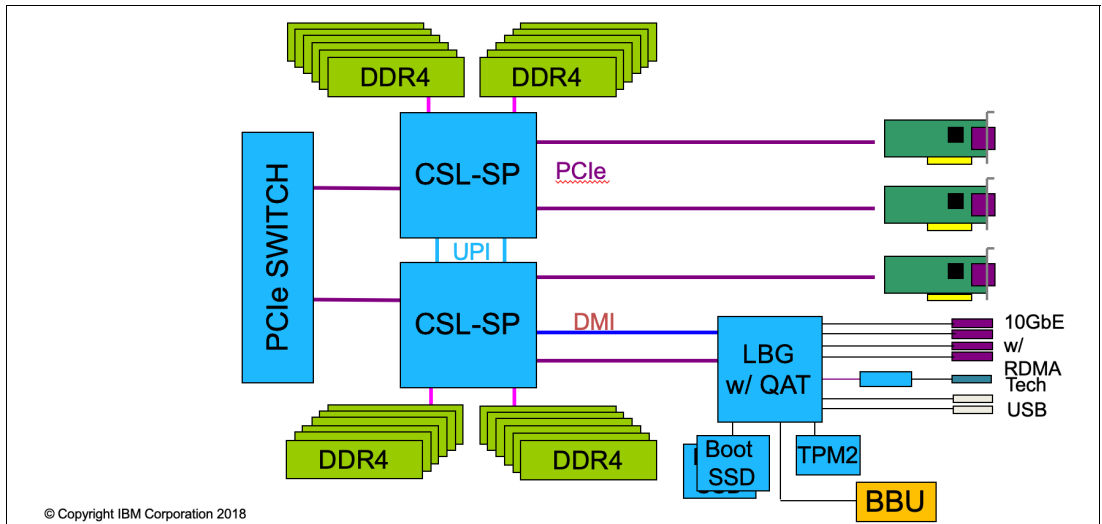


Figure 1-11 IBM SAN Volume Controller SV2 and SA2 internal architecture

**Note:** IBM SAN Volume Controller SV2 and SA2 do not support any type of expansion enclosures.

## 1.4.4 SAN Volume Controller model comparisons

All the SAN Volume Controller models are all delivered in a 2U 19-inch rack-mounted enclosure. At the time of writing, there are four models of the SAN Volume Controller, which are described in Table 1-1 and Table 1-2.

**Additional information:** For the most up-to-date information about features, benefits, and specifications of the SAN Volume Controller models, go to:

[IBM SAN Volume Controller Models](#)

The information in this book is valid at the time of writing and covers IBM Spectrum Virtualize V8.3.1. However, as SAN Volume Controller matures, expect to see new features and enhanced specifications.

Table 1-1 IBM SAN Volume Controller base models (1 of 2)

Feature	IBM SAN Volume Controller 2145-SV1 (2147-SV1) <sup>a</sup>	IBM SAN Volume Controller 2145-DH8
Processor	Two Intel Xeon E5 v4 Series, 8-cores, and 3.2 GHz	One Intel Xeon E5 v2 Series, 8-cores, and 2.6 GHz
Base cache memory	64 GB	32 GB
I/O ports and management	Three 10 Gb Ethernet ports for 10 Gb iSCSI connectivity and system management	Three 1 Gb Ethernet ports for 1 Gb iSCSI connectivity and system management
Technician port	One 1 Gb Ethernet	One 1 Gb Ethernet
Maximum host interface adapters slots	4	4
USB ports	4	4
SAS chain	2	2
Max number of dense drawers per SAS chain	4	4
Integrated battery units	2	2
Power supplies and cooling units	2	2

a. Model 2147 is identical to 2145, but with an included enterprise support option from IBM.

Table 1-2 IBM SAN Volume Controller base models (2 of 2)

Feature	IBM SAN Volume Controller 2145-SV2 (2147-SV2) <sup>a</sup>	IBM SAN Volume Controller 2145-SA2 (2147-SA2) <sup>a</sup>
Processor	Two Intel Cascade Lake 5218 Series, 16-cores, and 2.30 GHz (Gold)	Two Intel Cascade Lake 4208 Series, 8-cores, and 2.10 GHz (Silver)
Base cache memory	128 GB	128 GB
I/O ports and management	Four 10 Gb Ethernet ports for 10 Gb iSCSI connectivity and system management	Four 10 Gb Ethernet ports for 10 Gb iSCSI connectivity and system management

Feature	IBM SAN Volume Controller 2145-SV2 (2147-SV2) <sup>a</sup>	IBM SAN Volume Controller 2145-SA2 (2147-SA2) <sup>a</sup>
Technician port	One 1 Gb Ethernet	One 1 Gb Ethernet
Max host interface adapters slots	3	3
USB ports	2	2
SAS chain	N/A	N/A
Max number of dense drawers per SAS chain	N/A	N/A
Integrated battery units	1	1
Power supplies and cooling units	2	2

a. Model 2147 is identical to 2145, but with an included enterprise support option from IBM.

The following optional features are available for IBM SAN Volume Controller SV1:

- ▶ A 256 GB cache upgrade fully unlocked with code Version 8.2 or later
- ▶ A 4-port 16 Gb FC adapter for 16 Gb FC connectivity
- ▶ A 4-port 10 Gb Ethernet adapter for 10 Gb iSCSI/FCoE connectivity
- ▶ Compression accelerator card for RtC
- ▶ A 4-port 12 Gb SAS expansion enclosure attachment card

The following optional features are available for IBM SAN Volume Controller DH8:

- ▶ Extra processor with 32 GB cache upgrade
- ▶ A 4-port 16 Gb FC adapter for 16 Gb FC connectivity
- ▶ A 4-port 10 Gb Ethernet adapter for 10 Gb iSCSI/FCoE connectivity
- ▶ Compression accelerator card for RtC
- ▶ A 4-port 12 Gb SAS expansion enclosure attachment card

**Important:** IBM SAN Volume Controller models 2145 / 2147-SV1 and 2145 / 2147-DH8 can contain a 16 Gb FC or a 10 Gb Ethernet adapter, but only one 10 Gbps Ethernet adapter is supported.

The following optional features are available for IBM SAN Volume Controller SV2 and SA2:

- ▶ A 768 GB cache upgrade
- ▶ A 4-port 16 Gb FC/FC over NVMe adapter for 16 Gb FC connectivity
- ▶ A 4-port 32 Gb FC/FC over NVMe adapter for 32 Gb FC connectivity
- ▶ A 2-port 25 Gb iSCI/RDMA over Converged Ethernet (RoCE) / NVMe over Ethernet port
- ▶ A 2-port 25 Gb iSCSI/Internet Wide-area RDMA Protocol (iWARP) / NVMe over Ethernet port
- ▶ A 2-port 50/100 GbE iSCSI / RoCE / NVMe over Ethernet (not at GA) port
- ▶ A 2-port 50/100 GbE iSCSI / iWARP / NVMe over Ethernet (not at GA) port

The SV2 and SA2 systems have dual CPU sockets and three adapter slots along with four 10-GbE RJ45 ports on board.

**Note:** IBM SAN Volume Controller models SA2 and SV2 do not support FCoE.

The comparison of current and previous models of SAN Volume Controller is shown in Table 1-3. Expansion enclosures are not included in the list.

Table 1-3 Historical overview of SAN Volume Controller models

Model	Cache (GB)	FC (Gbps)	iSCSI (Gbps)	Hardware base	Announced
2145-DH8	32 - 64	8 and 16	1, optional 10	x3550 M4	06 May 2014
2145-SV1	64 - 256	16	10	Xeon E5 v4	23 August 2016
2147-SV1	64 - 256	16	10	Xeon E5 v4	23 August 2016
2145-SV2	128 - 768	16 and 32	25, 50, and 100	Intel Xeon Cascade Lake	06 March 2020
2147-SV2	128 - 768	16 and 32	25, 50, and 100	Intel Xeon Cascade Lake	06 March 2020
2145-SA2	128 - 768	16 and 32	25, 50, and 100	Intel Xeon Cascade Lake	06 March 2020
2147-SA2	128 - 768	16 and 32	25, 50, and 100	Intel Xeon Cascade Lake	06 March 2020

The IBM SAN Volume Controller DH8 and SV1 expansion enclosure consists of an enclosure and drives. Each enclosure contains two canisters that can be replaced and maintained independently. IBM SAN Volume Controller DH8 and SV1 support three types of expansion enclosure. The expansion enclosure models are 12F, 24F, and 92F dense drawers.

**Note:** IBM SAN Volume Controller SV2 and SA2 do not support any type of SAS expansion enclosures.

Expansion Enclosure model 12F features two expansion canisters and holds up to twelve 3.5-inch SAS drives in a 2U 19-inch rack mount enclosure.

Expansion Enclosure model 24F supports up to 24 internal flash drives, 2.5-inch SAS drives, or a combination of them. Expansion Enclosure 24F also features two expansion canisters in a 2U 19-inch rack-mounted enclosure.

Expansion Enclosure model 92F supports up to ninety-two 3.5-inch drives in a 5U 19-inch rack-mounted enclosure. Also, it is called dense expansion drawers, or *dense drawers*.

Figure 1-12 shows an example of an IBM SAN Volume Controller DH8 with eight expansion enclosures that are attached.

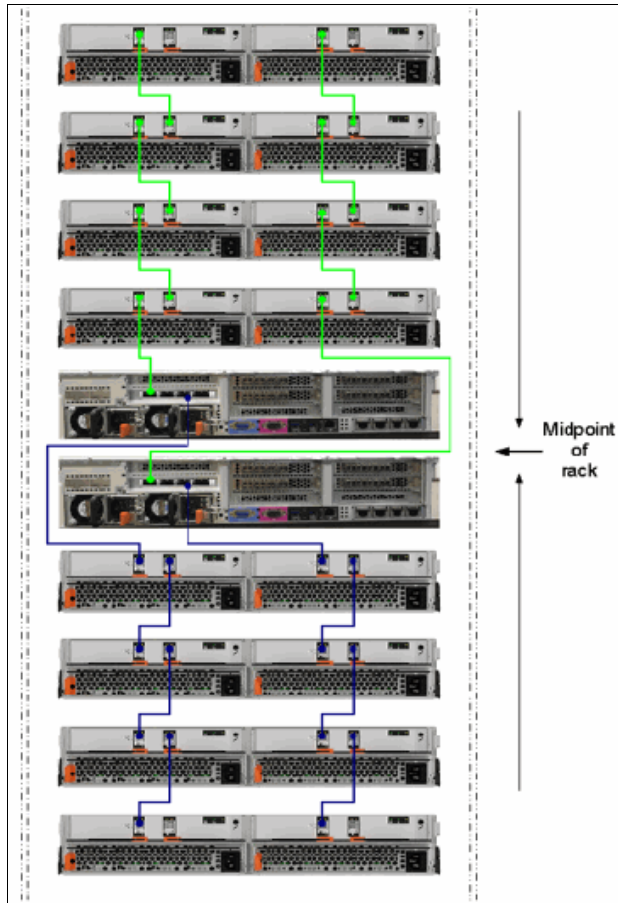


Figure 1-12 IBM SAN Volume Controller DH8 with expansion enclosures

## 1.5 IBM SAN Volume Controller components

SAN Volume Controller provides block-level aggregation and volume management for attached disk storage. In simpler terms, SAN Volume Controller manages several back-end storage controllers or locally attached disks.

SAN Volume Controller maps the physical storage within those controllers or disk arrays into logical disk images, or *volumes*, that can be seen by application servers and workstations in the SAN. It logically sits between hosts and storage arrays, presenting itself to hosts as the storage provider (*target*) and presenting itself to storage arrays as one large host (*initiator*).

The SAN is zoned such that the application servers cannot “see” the back-end physical storage. This configuration prevents any possible conflict between SAN Volume Controller and the application servers that are trying to manage the back-end storage.

In the next topics, the term *IBM SAN Volume Controller* is used to refer to both models of the SAN Volume Controller product. However, SAN Volume Controller is based on the components that are described next.

## 1.5.1 Nodes

Each SAN Volume Controller hardware unit is called a *node*. Each node is an individual server in a SAN Volume Controller clustered system on which The SAN Volume Controller software runs. The node provides the virtualization for a set of volumes, cache, and copy services functions. The SAN Volume Controller nodes are deployed in pairs (*cluster*), and one or multiple pairs constitute a *clustered system* or *system*. A system can consist of one pair and a maximum of four pairs.

One of the nodes within the system is known as the *configuration node*. The configuration node manages the configuration activity for the system. If this node fails, the system chooses a new node to become the configuration node.

Because the active nodes are installed in pairs, each node provides a failover function to its partner node if a node fails.

## 1.5.2 I/O groups

Each pair of SAN Volume Controller nodes is also referred to as an *I/O group*. A SAN Volume Controller clustered system can have 1 - 4 I/O groups.

A specific *volume* is always presented to a host server by a single I/O group of the system. The I/O group can be changed.

When a host server performs I/O to one of its volumes, all the I/Os for a specific volume are directed to one specific I/O group in the system. Under normal conditions, the I/Os for that specific volume are always processed by the same node within the I/O group. This node is referred to as the *preferred node* for this specific volume.

Both nodes of an I/O group act as the preferred node for their own specific subset of the total number of volumes that the I/O group presents to the host servers. However, both nodes also act as failover nodes for their respective partner node within the I/O group. Therefore, a node takes over the I/O workload from its partner node when required.

In a SAN Volume Controller based environment, the I/O handling for a volume can switch between the two nodes of the I/O group. So, it is a best practice that servers are connected to two different fabrics through different FC host bus adapters (HBAs) to use multipath drivers to give redundancy.

The SAN Volume Controller I/O groups are connected to the SAN so that all application servers that are accessing volumes from this I/O group have access to this group. Up to 512 host server objects can be defined per I/O group. The host server objects can access volumes that are provided by this specific I/O group.

If required, host servers can be mapped to more than one I/O group within the SAN Volume Controller system. Therefore, they can access volumes from separate I/O groups. You can move volumes between I/O groups to redistribute the load between the I/O groups. Modifying the I/O group that services the volume can be done concurrently with I/O operations if the host supports nondisruptive volume moves.

It also requires a rescan at the host level to ensure that the multipathing driver is notified that the allocation of the preferred node changed, and the ports (by which the volume is accessed) changed. This modification can be done in the situation where one pair of nodes becomes overused.



### 1.5.3 System

The system or clustered system consists of 1 - 4 I/O groups. Certain configuration limitations are then set for the individual system. For example, the maximum number of volumes that is supported per system is 10,000, or the maximum number of MDisks that is supported is ~28 PiB (32 PB) per system.

All configuration, monitoring, and service tasks are performed at the system level. Configuration settings are replicated to all nodes in the system. To facilitate these tasks, a management IP address is set for the system.

A process is provided to back up the system configuration data on to storage so that it can be restored if there is a disaster. This method does not back up application data. Only the SAN Volume Controller system configuration information is backed up.

For remote data mirroring, two or more systems must form a *partnership* before relationships between mirrored volumes are created.

For more information about the maximum configurations that apply to the system, I/O group, and nodes, search for “Configuration Limits and Restrictions for IBM System Storage IBM SAN Volume Controller” at [IBM Support home page](#).

### 1.5.4 Expansion enclosures

Expansion enclosures are rack-mounted hardware that contains several components of the system: canisters, drives, and power supplies. Enclosures can be used to extend the capacity of the system. They are supported only on IBM SAN Volume Controller 2145 or 2147-DH8 or IBM SAN Volume Controller 2145 or 2147-SV1 nodes. For other models of the system, you must use external storage systems to provide capacity for data. The term *enclosure* is also used to describe the hardware and other parts that are plugged into the enclosure.

If you have IBM SAN Volume Controller 2145 or 2147-DH8 or IBM SAN Volume Controller 2145 or 2147-SV1 controller nodes, you can add expansion enclosures to expand the available capacity of the system. (Other controller models do not support expansion enclosures.) Each system can have a maximum of four I/O groups, with two chains of expansion enclosures that are attached to each I/O group.

On each SAS chain, the systems can support up to a SAS chain weight of 10. Each 2145 / 2147-92F Expansion Enclosure adds a value of 2.5 to the SAS chain weight. Each 2145 / 2147-12F or 2145 / 2147-24F Expansion Enclosure adds a value of 1 to the SAS chain weight. For example, each of the following expansion enclosure configurations has a total SAS weight of 10:

- ▶ Four 2145 or 2147-92F enclosures per SAS chain
- ▶ Ten 2145 or 2147-12F enclosures per SAS chain
- ▶ Two 2145 or 2147-92F enclosures and five 2145 or 2147-24F enclosures per SAS chain

An expansion enclosure houses the following additional hardware: power supply units (PSUs), canisters, and drives. Enclosure objects report the connectivity of the enclosure. Each expansion enclosure is assigned a unique ID. It is possible to change the enclosure ID later.

Figure 1-13 on page 23 shows the front view of the 2145 or 2147 model 12F Expansion Enclosure.





Figure 1-13 IBM SAN Volume Controller Expansion Enclosure front view 2145 / 2147 model 12F

The drives are positioned in four columns of three horizontally mounted drive assemblies. The drive slots are numbered 1 - 12, starting at the upper left and moving left to right, top to bottom.

Figure 1-14 shows the front view of the 2145 or 2147 model 24F Expansion Enclosures.

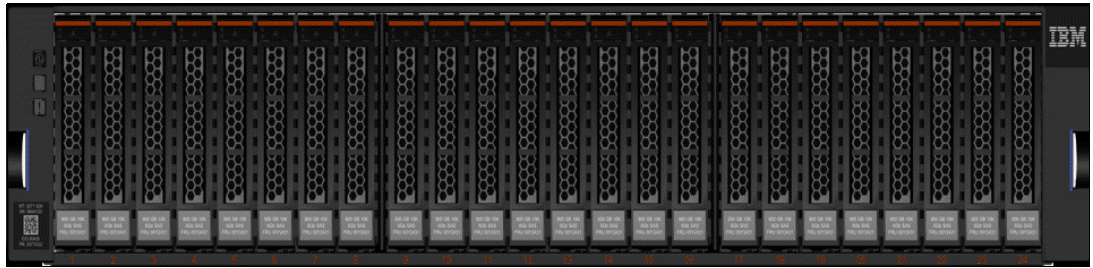


Figure 1-14 IBM SAN Volume Controller Expansion Enclosure front view 2145 / 2147 model 24F

The drives are positioned in one row of 24 vertically mounted drive assemblies. The drive slots are numbered 1 - 24, starting from the left. A vertical center drive bay molding is between slots 8 and 9 and another between slots 16 and 17.

**Note:** IBM SAN Volume Controller SV2 and SA2 do not support any type of expansion enclosures.

## 1.5.5 Dense expansion drawers

Dense expansion drawers, also known as *dense drawers*, are optional disk expansion enclosures that are 5U rack-mounted. Each chassis features two expansion canisters, two power supplies, two expander modules, and a total of four fan modules.

Each dense drawer can hold up to 92 drives that are positioned in four rows of 14 and an additional three rows of 12 mounted drives assemblies. The two Secondary Expander Modules (SEMs) are centrally located in the chassis. One SEM addresses 54 drive ports, and the other addresses 38 drive ports.

Dense drawers can support 4 TB, 6 TB, 8 TB, 10 TB, 12 TB, 14 TB, and 16 TB 3.5" NL SAS drives.

Each canister in the dense drawer chassis features two SAS ports that are numbered 1 and 2. Using SAS port1 is mandatory because the expansion enclosure must be attached to a SAN Volume Controller node or another expansion enclosure. SAS connector 2 is optional because it is used to attach to more expansion enclosures.

**Note:** IBM SAN Volume Controller SV2 and SA2 do not support any type of expansion enclosures.

Figure 1-15 shows a dense expansion drawer.



Figure 1-15 Dense expansion drawer

## 1.5.6 Flash drives

Flash drives can be used to overcome a growing problem that is known as the *memory bottleneck* or *storage bottleneck* specifically in single-layer cell (SLC) or multi-level cell (MLC) NAND flash-based disks.

### Storage bottleneck problem

The memory or storage bottleneck describes the steadily growing gap between the time that is required for a CPU to access data that is in its cache memory (typically in nanoseconds) and data that is on external storage (typically in milliseconds).

Although CPUs and cache and memory devices continually improve their performance, mechanical disks that are used as external storage generally do not improve their performance.

Figure 1-16 shows these access time differences.

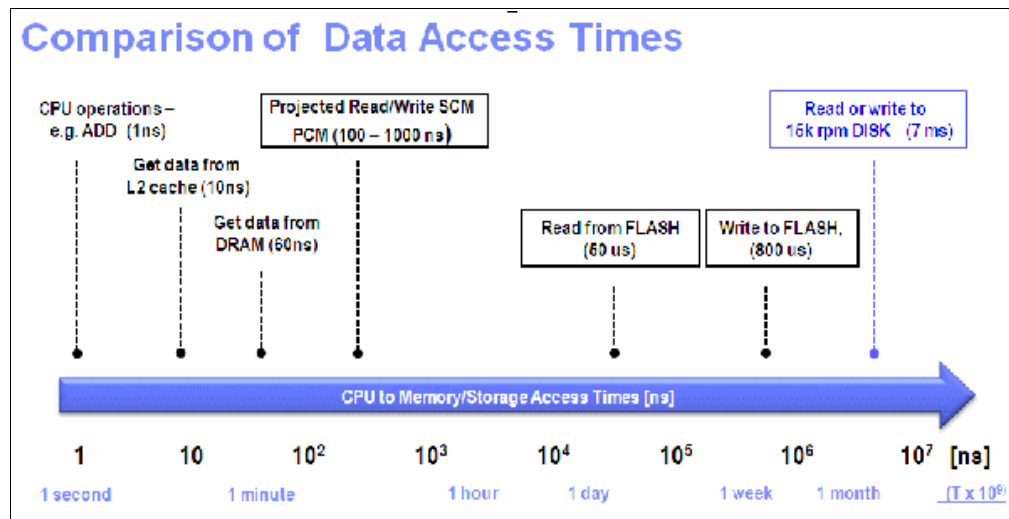


Figure 1-16 The memory or storage bottleneck

The actual times that are shown are not that important, but a noticeable difference exists between accessing data that is in cache and data that is on an external disk.

In the example that is shown in Figure 1-16, we added a second scale that gives you an idea of how long it takes to access the data in a scenario where a single CPU cycle takes 1 second. This scale shows the importance of future storage technologies closing or reducing the gap between access times for data that is stored in cache or memory versus access times for data that is stored on an external medium.

Since magnetic disks were first introduced by IBM in 1956 (the Random Access Memory Accounting System, also known as the *IBM 305 RAMAC*), they showed remarkable performance regarding capacity growth, form factor, and size reduction; price savings (cost per GB); and reliability.

However, the number of I/Os that a disk can handle and the response time that it takes to process a single I/O did not improve at the same rate, although they certainly did improve. In actual environments, you can expect from today's enterprise-class FC SAS disk up to 200 input/output operations per second (IOPS) per disk with an average response time (latency) of approximately 6 ms per I/O.

Table 1-4 shows a comparison of drive types and IOPS.

Table 1-4 Comparison of drive types to IOPS

Drive type	IOPS
NL - SAS	100
SAS 10,000 revolutions per minute (RPM)	150
SAS 15,000 RPM	250
Flash	> 500,000 read; 300,000 write

Today's spinning disks continue to advance in capacity, up to several terabytes, form factor/footprint (8.89 cm (3.5 inches), 6.35 cm (2.5 inches), and 4.57 cm (1.8 inches)), and price (cost per gigabyte), but they are not getting much faster.

The limiting factor is the number of RPM that a disk can perform (approximately 15,000). This factor defines the time that is required to access a specific data block on a rotating device. Small improvements likely will occur in the future. However, a significant step, such as doubling the RPM (if technically even possible), inevitably has an associated increase in power usage and price that is an inhibitor.

### **Flash drive solution**

Flash drives can provide a solution for this dilemma, and no rotating parts means improved robustness and lower power usage. A remarkable improvement in I/O performance and a massive reduction in the average I/O response times (latency) are the compelling reasons to use flash drives in today's storage subsystems.

Enterprise-class flash drives typically deliver 500,000 read and 300,000 write IOPS with typical latencies of 50  $\mu$ s for reads and 800  $\mu$ s for writes. Their form factors of 4.57 cm (1.8 inches) / 6.35 cm (2.5 inches) / 8.89 cm (3.5 inches) and their interfaces (FC / SAS / Serial Advanced Technology Attachment (SATA)) make them easy to integrate into existing disk shelves. The IOPS metrics improve when flash drives are consolidated in storage arrays (flash arrays). In this case, the read and write IOPS are seen in millions for specific 4 KB data blocks.

### **Flash-drive market**

The flash drive storage market is rapidly evolving. The key differentiator among today's flash drive products is not the storage medium but the logic in the disk internal controllers. The top priorities in today's controller development are optimally handling what is referred to as *wear-out leveling*, which defines the controller's capability to ensure a device's durability, and closing the remarkable gap between read and write I/O performance.

Today's flash drive technology is only the first step into the world of high-performance persistent semiconductor storage. A group of the approximately 10 most promising technologies is collectively referred to as *SCM*.

### **Read-intensive flash drives**

Generally, there are two types of SSDs in the market for enterprise storage: The MLC and SLC. The most common SSD technology is MLC. They are commonly found in consumer products, such as portable electronic devices. However, they are also present in some enterprise storage products. Enterprise class SSDs are built on mid-endurance to high-endurance MLC flash technology, mostly known as mainstream endurance SSD.

MLC SSDs use the multicell to store data and features the wear-leveling method, which is the process to evenly spread data across all memory cells on the SSD. This method helps to eliminate potential hotspots that are caused by repetitive write-erase cycles. SLC SSDs use a single cell to store 1 bit of data, and that makes them faster.

To support particular business demands, IBM Spectrum Virtualize qualified the use of read-intensive SSDs with applications where the read operations are high. The IBM Spectrum Virtualize GUI presents new attributes when managing disk drives by using the GUI and the CLI. The new function reports the *write-endurance* limits (in percentages) for each qualified read-intensive SSD that is installed in the system.

Read-intensive SSDs are available as an optional purchase product for the SAN Volume Controller and the IBM Storwize Family. For more information about read-intensive SSDs and IBM Spectrum Virtualize, see *Read Intensive Flash Drives*, REDP-5380.

## IBM FlashCore Module and NVMe drives

Figure 1-17 shows an IBM FlashCore Module (FCM) (NVMe) with a capacity of 19.2 TB that is built by using 64-layer TLC flash memory and an Everspin MRAM cache into a U.2 form factor.



Figure 1-17 FlashCore Module (NVMe)

FCM modules (NVMe) are designed for high parallelism and optimized for 3D TLC and updated FPGAs. IBM also enhanced the FCM modules by adding a read cache to reduce latency on highly compressed pages, and four-plane programming to lower the overall power during writes. FCM modules offer hardware-assisted compression up to 3:1, and are FIPS 140-2 compliant.

FCM modules carry the IBM patented Variable Stripe RAID at the FCM level, and use distributed RAID (DRAID) to protect data at the system level. VSR and DRAID together optimize RAID rebuilds by offloading rebuilds to DRAID and offers protection against FCM failures.

FCM modules can be configured to use 4.8 TB, 9.6 TB, 19.2 TB, or 38.4 TB drives.

**Note:** At the time of writing, FCM drives are only available on the IBM FlashSystem 7200, IBM FlashSystem 9100, IBM FlashSystem 9200, IBM Storwize 5100, and IBM Storwize V7000 Gen3. IBM SAN Volume Controller SV2 and SA2 do not support any internal drive types.

### Storage-class memory

SCM promises a massive improvement in performance (in IOPS), and a real density, cost, and energy efficiency compared to today's flash-drive technology. IBM Research® is actively engaged in these new technologies.

For more information about nanoscale devices, see the following website:

[Information about Nanoscale Devices](#)

For a comprehensive overview of the flash drive technology in a subset of the well-known Storage Networking Industry Association (SNIA) Technical Tutorials, see this website:

[Overview of Flash-drive Technology](#)





However, the application servers do not “see” the MDisks at all. Rather, they see several logical disks, which are known as *VDisks* or *volumes*. These disks are presented by the SAN Volume Controller I/O groups through the SAN (FC/FCoE) or LAN (iSCSI) to the servers. The MDisks are placed into storage pools where they are divided into several extents.

For more information about the total storage capacity that is manageable per system regarding the selection of extents, search for “Configuration Limits and Restrictions for IBM System Storage IBM SAN Volume Controller” at the following support website:

[IBM Support home page](#)

A volume is host-accessible storage that is provisioned out of one *storage pool*, or, if it is a mirrored volume, out of two storage pools.

The maximum size of an MDisk is 1 PiB. A SAN Volume Controller system supports up to 4096 MDisks (including internal RAID arrays). When an MDisk is presented to the SAN Volume Controller, it can be one of the following statuses:

► Unmanaged MDisk

An MDisk is reported as unmanaged when it is not a member of any storage pool. An unmanaged MDisk is not associated with any volumes and has no metadata that is stored on it. SAN Volume Controller does not write to an MDisk that is in unmanaged mode except when it attempts to change the mode of the MDisk to one of the other modes. SAN Volume Controller can see the resource, but the resource is not assigned to a storage pool.

► Managed MDisk

Managed mode MDisks are always members of a storage pool, and they contribute extents to the storage pool. Volumes (if not operated in image mode) are created from these extents. MDisks that are operating in managed mode might have metadata extents that are allocated from them and can be used as *quorum disks*. This mode is the most common and normal mode for an MDisk.

► Image mode MDisk

*Image mode* provides a direct block-for-block conversion from the MDisk to the volume by using virtualization. This mode is provided to satisfy the following major usage scenarios:

- Image mode enables the virtualization of MDisks that already contain data that was written directly and not through SAN Volume Controller. Rather, it was created by a direct-connected host.

This mode enables a client to insert SAN Volume Controller into the data path of an existing storage volume or LUN with minimal downtime. For more information about the data migration process, see Chapter 8, “Storage migration” on page 429.

Image mode enables a volume that is managed by a SAN Volume Controller to be used with the native copy services function that is provided by the underlying RAID controller. To avoid the loss of data integrity when the SAN Volume Controller is used in this way, it is important that you disable the SAN Volume Controller cache for the volume.

- The SAN Volume Controller can migrate to image mode, which enables the SAN Volume Controller to export volumes and access them directly from a host without the SAN Volume Controller in the path.

Each MDisk that is presented from an external disk controller has an online path count that is the number of nodes that has access to that MDisk. The *maximum count* is the maximum number of paths that is detected at any point by the system. The *current count* is what the system sees at this point. A current value that is less than the maximum can indicate that SAN fabric paths were lost.

SSDs that are in the IBM SAN Volume Controller 2145-CG8 or flash space, which are presented by the external flash enclosures of the IBM SAN Volume Controller 2145 / 2147-DH8 or SV1 nodes, are presented to the cluster as MDisks. To determine whether the selected MDisk is an SSD or flash, click the link on the MDisk name to display the Viewing MDisk Details window.

If the selected MDisk is an SSD or flash that is on a SAN Volume Controller, the Viewing MDisk Details window displays values for the Node ID, Node Name, and Node Location attributes. Alternatively, you can select **Work with Managed Disks** → **Disk Controller Systems** from the portfolio. On the Viewing Disk Controller window, you can match the MDisk to the disk controller system that has the corresponding values for those attributes.

## 1.5.8 Cache

The primary benefit of storage cache is to improve I/O response time. Reads and writes to a magnetic disk drive experience seek time and latency time at the drive level, which can result in 1 ms - 10 ms of response time (for an enterprise-class disk).

- ▶ The IBM SAN Volume Controller Model SV1 features 64 GB of memory with options for 256 GB of memory in a 2U 19-inch rack mount enclosure.
- ▶ The IBM SAN Volume Controller Model SA2 and SV2 features 128 GB of memory with options for 756 GB of memory in a 2U 19-inch rack mount enclosure.

The SAN Volume Controller provides a flexible cache model, and the node's memory can be used as read or write cache.

Cache is allocated in 4 kibibyte (KiB) segments. A *segment* holds part of one track. A *track* is the unit of locking and destaging granularity in the cache. The cache virtual track size is 32 KiB (eight segments). A track might be only partially populated with valid pages. The SAN Volume Controller combines writes up to a 256 KiB track size if the writes are in the same tracks before destaging. For example, if 4 KiB is written into a track, another 4 KiB is written to another location in the same track.

Therefore, the blocks that are written from the SAN Volume Controller to the disk subsystem can be any size of 512 bytes - 256 KiB. The large cache and advanced cache management algorithms enable it to improve the performance of many types of underlying disk technologies.

The SAN Volume Controller capability to manage in the background the destaging operations that are incurred by writes (in addition to still supporting full data integrity) assists with the SAN Volume Controller capability in achieving good database performance.

The cache is separated into two layers: Upper cache and lower cache. Figure 1-19 on page 31 shows the separation of the upper and lower cache.



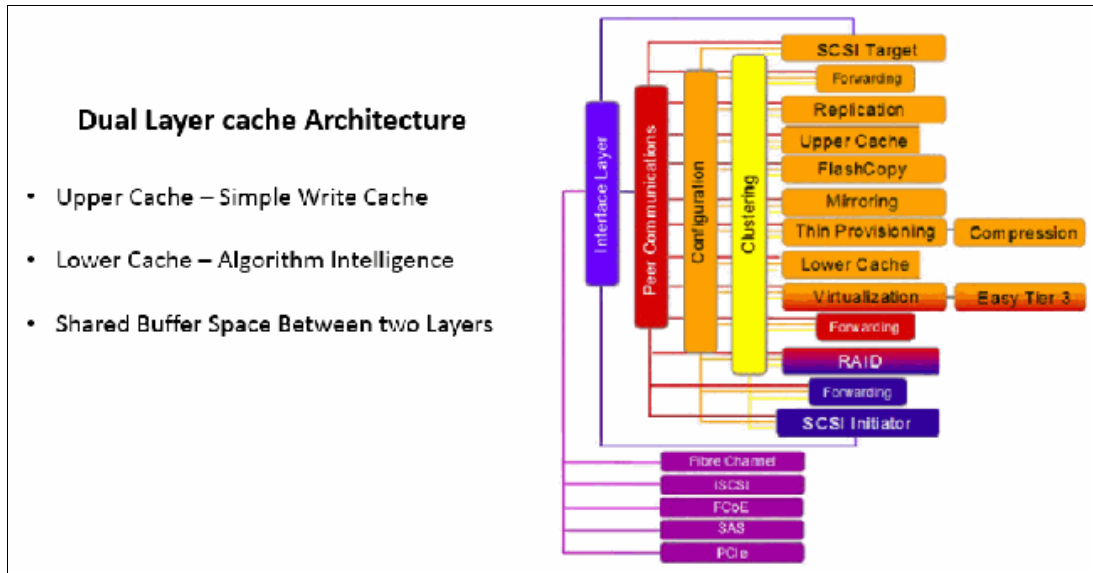


Figure 1-19 Separation of upper and lower cache

The upper cache delivers the following functions, which enable the SAN Volume Controller to streamline data write performance:

- ▶ Provides fast write response times to the host by being as high up in the I/O stack as possible.
- ▶ Provides partitioning.

The lower cache delivers the following additional functions:

- ▶ Ensures that the write cache between two nodes is in sync.
- ▶ Caches partitioning to ensure that a slow back end cannot use the entire cache.
- ▶ Uses a destaging algorithm that adapts to the amount of data and the back-end performance.
- ▶ Provides read caching and prefetching.

Combined, the two levels of cache also deliver the following functions:

- ▶ Pins data when the LUN goes offline.
- ▶ Provides enhanced statistics for IBM Tivoli® Storage Productivity Center, and maintains compatibility with an earlier version.
- ▶ Provides trace for debugging.
- ▶ Reports medium errors.
- ▶ Resynchronizes cache correctly and provides the atomic write function.
- ▶ Ensures that other partitions continue operation when one partition becomes 100% full of pinned data.
- ▶ Supports fast-write (two-way and one-way), flush-through, and write-through.
- ▶ Integrates with T3 recovery procedures.
- ▶ Supports two-way operation.
- ▶ Supports none, read-only, and read/write as user-exposed caching policies.
- ▶ Supports flush-when-idle.

- ▶ Supports expanding cache as more memory becomes available to the platform.
- ▶ Supports credit throttling to avoid I/O skew and offer fairness/balanced I/O between the two nodes of the I/O group.
- ▶ Enables switching of the preferred node without needing to move volumes between I/O groups.

Depending on the size, age, and technology level of the disk storage system, the total available cache in the SAN Volume Controller nodes can be larger, smaller, or about the same as the cache that is associated with the disk storage.

Because hits to the cache can occur in either the SAN Volume Controller or the disk controller level of the overall system, the system as a whole can take advantage of the larger amount of cache wherever the cache is. Therefore, if the storage controller level of the cache has the greater capacity, expect hits to this cache to occur in addition to hits in the SAN Volume Controller cache.

In addition, regardless of their relative capacities, both levels of cache tend to play an important role in enabling sequentially organized data to flow smoothly through the system. The SAN Volume Controller cannot increase the throughput potential of the underlying disks in all cases because this increase depends on both the underlying storage technology and the degree to which the workload exhibits *hotspots* or sensitivity to cache size or cache algorithms.

IBM SAN Volume Controller V7.3 introduced a major upgrade to the cache code, and in association with the IBM SAN Volume Controller 2145 / 2147-DH8 hardware it provided an additional cache capacity upgrade. A base SAN Volume Controller node configuration included 32 GB of cache. Adding the second processor and cache upgrade for RtC took a single node to a total of 64 GB of cache. A single I/O group with support for RtC contained 128 GB of cache, and an 8-node SAN Volume Controller system with a maximum cache configuration contained a total of 512 GB of cache.

In IBM SAN Volume Controller V8.1, these limits are enhanced with the IBM SAN Volume Controller 2145 / 2147-SV1 appliance. Before this release, the SAN Volume Controller memory manager (PLMM) could address only 64 GB of memory. In Version 8.1, the underlying PLMM is rewritten and the structure size increased. The cache size can be upgraded up to 256 GB and the whole memory can be used. However, the write cache is still assigned to a maximum of 12 GB and compression cache to a maximum of 34 GB. The remaining installed cache is used as read cache (including allocation for features like FlashCopy, GM, or MM).

**Important:** When upgrading to IBM SAN Volume Controller V8.1 where there is more than 64 GB of physical memory installed (but not used), the error message 1199 Detected hardware needs activation displays after the upgrade in the GUI event log (and error code 0x841 as a result of the `lseventlog` command in the CLI).

A different memory management feature must be activated in the SAN Volume Controller code by running a fix procedure in the GUI or by running `chnodehw <node_id> -force`. The system restarts. Do not run the command on more than one node at a time.

## 1.5.9 Quorum disk

A *quorum disk* is an MDisk or a managed drive that contains a reserved area that is used exclusively for system management. A system automatically assigns quorum disk candidates. Quorum disks are used when there is a problem in the SAN fabric or when nodes are shut down, which leaves half of the nodes remaining in the system. This type of problem causes a loss of communication between the nodes that remain in the system and those that do not remain.

The nodes are split into groups where the remaining nodes in each group can communicate with each other, but not with the other group of nodes that were formerly part of the system. In this situation, some nodes must stop operating and processing I/O requests from hosts to preserve data integrity while maintaining data access. If a group contains less than half the nodes that were active in the system, the nodes in that group stop operating and processing I/O requests from hosts.

It is possible for a system to split into two groups with each group containing half the original number of nodes in the system. A quorum disk determines which group of nodes stops operating and processing I/O requests. In this tiebreaker situation, the first group of nodes that accesses the quorum disk is marked as the owner of the quorum disk. As a result, the owner continues to operate as the system and handles all I/O requests.

If the other group of nodes cannot access the quorum disk or discover that the quorum disk is owned by another group of nodes, it stops operating as the system and does not handle I/O requests. A system can have only one active quorum disk that is used for a tiebreaker situation. However, the system uses three quorum disks to record a backup of system configuration data that is used if there is a disaster. The system automatically selects one active quorum disk from these three disks.

The other quorum disk candidates provide redundancy if the active quorum disk fails before a system is partitioned. To avoid the possibility of losing all of the quorum disk candidates with a single failure, assign quorum disk candidates on multiple storage systems.

**Quorum disk requirements:** To be considered eligible as a quorum disk, a LUN must meet the following criteria:

- ▶ It must be presented by a disk subsystem that supports SAN Volume Controller quorum disks.
- ▶ It is manually enabled as a quorum disk candidate by running the `chcontroller -allowquorum yes` command.
- ▶ It must be in managed mode (no image mode disks).
- ▶ It must have sufficient free extents to hold the system state information and the stored configuration metadata.
- ▶ It must be visible to all of the nodes in the system.

**Quorum disk placement:** If possible, the SAN Volume Controller places the quorum candidates on separate disk subsystems. However, after the quorum disk is selected, no attempt is made to ensure that the other quorum candidates are presented through separate disk subsystems.

**Important:** Quorum disk placement verification and adjustment to separate storage systems (if possible) reduce the dependency from a single storage system, and can increase the quorum disk availability.

You can list the quorum disk candidates and the active quorum disk in a system by running the `lscquorum` command.

When the set of quorum disk candidates is chosen, it is fixed. However, a new quorum disk candidate can be chosen in one of the following conditions:

- ▶ When the administrator requests that a specific MDisk becomes a quorum disk by running the `chquorum` command.
- ▶ When an MDisk that is a quorum disk is deleted from a storage pool.
- ▶ When an MDisk that is a quorum disk changes to image mode.

An offline MDisk is not replaced as a quorum disk candidate.

For disaster recovery (DR) purposes, a system must be regarded as a single entity so that the system and the quorum disk can be collocated.

Special considerations are required for the placement of the active quorum disk for a stretched or split cluster and split I/O group configurations. For more information, see [IBM Knowledge Center](#).

**Important:** Running a SAN Volume Controller system without a quorum disk can seriously affect your operation. A lack of available quorum disks for storing metadata prevents any migration operation (including a forced MDisk delete).

Mirrored volumes can be taken offline if no quorum disk is available. This behavior occurs because the synchronization status for mirrored volumes is recorded on the quorum disk.

During the normal operation of the system, the nodes communicate with each other. If a node is idle for a few seconds, a heartbeat signal is sent to ensure connectivity with the system. If a node fails for any reason, the workload that is intended for the node is taken over by another node until the failed node is restarted and readmitted into the system (which happens automatically).

If the Licensed Internal Code on a node becomes corrupted, which results in a failure, the workload is transferred to another node. The code on the failed node is repaired, and the node is readmitted into the system (which is an automatic process).

## IP quorum configuration

In a stretched configuration or HyperSwap configuration, you must use a third, independent site to house quorum devices. To use a quorum disk as the quorum device, this third site must use FC or IP connectivity together with an external storage system. In a local environment, no extra hardware or networking, such as FC or SAS-attached storage, is required beyond what is normally always provisioned within a system.

To use an IP-based quorum application as the quorum device for the third site, no FC connectivity is used. Java applications are run on hosts at the third site. However, there are strict requirements on the IP network, and some disadvantages with using IP quorum applications.

Unlike quorum disks, all IP quorum applications must be reconfigured and redeployed to hosts when certain aspects of the system configuration change. These aspects include adding or removing a node from the system, or when node service IP addresses are changed.

For stable quorum resolutions, an IP network must provide the following requirements:

- ▶ Connectivity from the hosts to the service IP addresses of all nodes. If IP quorum is configured incorrectly, the network must also deal with the possible security implications of exposing the service IP addresses because this connectivity can also be used to access the service GUI.
- ▶ Port 1260 is used by IP quorum applications to communicate from the hosts to all nodes.
- ▶ The maximum round-trip delay must not exceed 80 ms, which means 40 ms each direction.
- ▶ A minimum bandwidth of 2 MBps for node-to-quorum traffic.

Even with IP quorum applications at the third site, quorum disks at site one and site two are required because they are used to store metadata. To provide quorum resolution, run the **mkquorumapp** command to generate a Java application that is copied from the system and run on a host at a third site. The maximum number of applications that can be deployed is five. Currently, supported Java Runtime Environment (JRE) versions are IBM Java 7.1 and IBM Java 8.

### 1.5.10 Disk tier

It is likely that the MDisks (LUNs) that are presented to the SAN Volume Controller system have various performance attributes because of the type of disk or RAID array on which they are placed. The MDisks can be 15,000 disk RPM FC or SAS disks, NL SAS, SATA, or even flash drives. Therefore, a storage tier attribute is assigned to each MDisk with the default of `generic_hdd`.

### 1.5.11 Storage pool

A *storage pool* is a collection of up to 128 MDisks that provides the pool of storage from which volumes are provisioned. A single system can manage up to 1024 storage pools. The size of these pools can be changed (expanded or shrunk) at run time by adding or removing MDisks without taking the storage pool or the volumes offline. Expanding a storage pool with a single drive is not possible.

At any point, an MDisk can be a member in one storage pool only, except for image mode volumes.

Figure 1-20 shows the relationships of the SAN Volume Controller entities to each other.

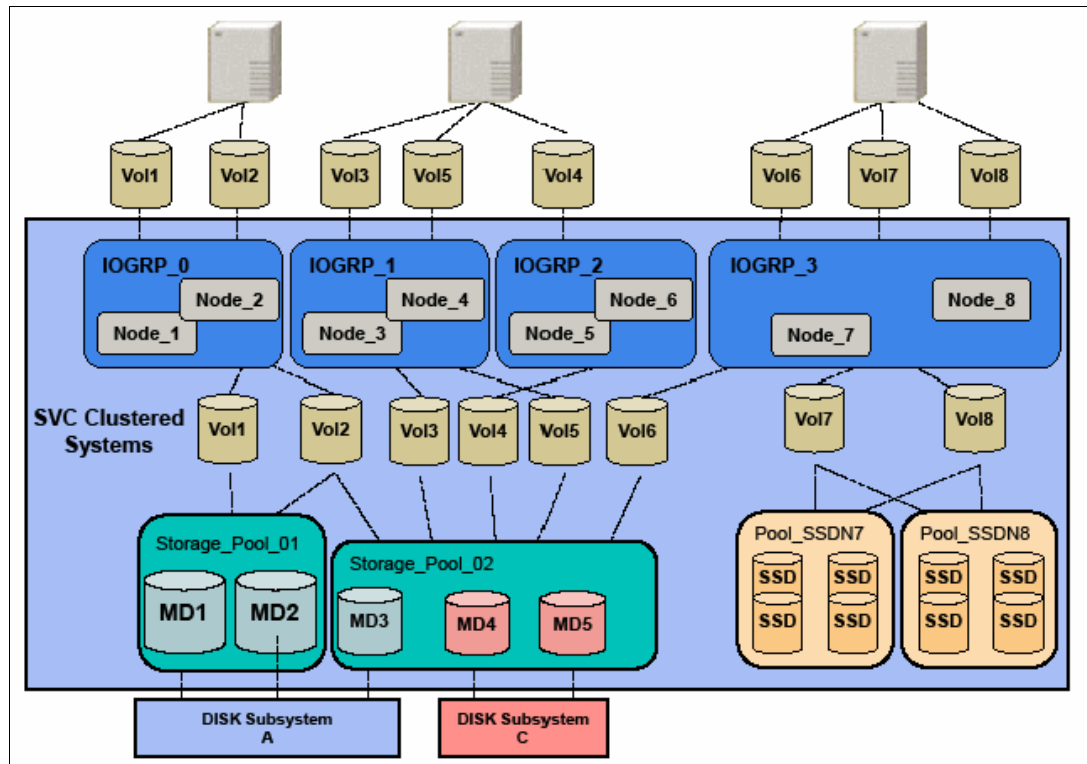


Figure 1-20 Overview of a SAN Volume Controller clustered system with an I/O group

Each MDisk in the storage pool is divided into several extents. The size of the extent is selected by the administrator when the storage pool is created and cannot be changed later. The size of the extent is 16 MiB - 8192 MiB.

It is a best practice to use the same extent size for all storage pools in a system. This approach is a prerequisite for supporting volume migration between two storage pools. If the storage pool extent sizes are not the same, you must use volume mirroring to copy volumes between pools.

The SAN Volume Controller limits the number of extents in a system to  $2^{22} \approx 4$  million. Because the number of addressable extents is limited, the total capacity of a SAN Volume Controller system depends on the extent size that is chosen by the SAN Volume Controller administrator.

## 1.5.12 Volumes

*Volumes* are logical disks that are presented to the host or application servers by the SAN Volume Controller. Hosts and application servers can see only the logical volumes that are created from combining extents from a storage pool.

There are three types of volumes in terms of extents management:

- ▶ Striped

A striped volume is allocated one extent in turn from each MDisk in the storage pool. This process continues until the space that is required for the volume is satisfied.

It is also possible to supply a list of MDisks to use.

Figure 1-21 shows how a striped volume is allocated, assuming that 10 extents are required.

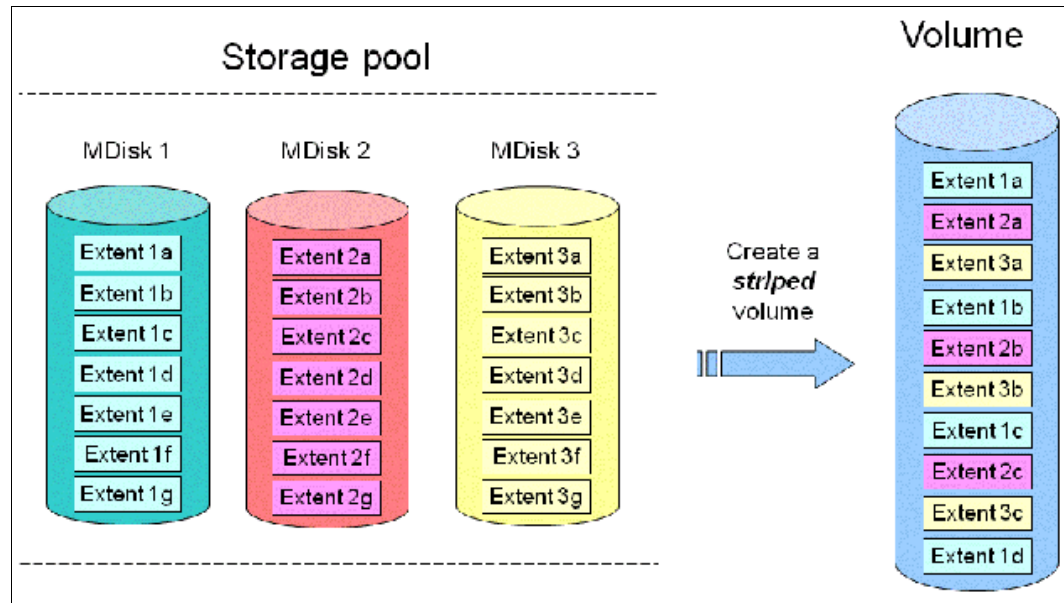


Figure 1-21 Striped volume

► Sequential

A sequential volume is where the extents are allocated sequentially from one MDisk to the next MDisk (Figure 1-22).

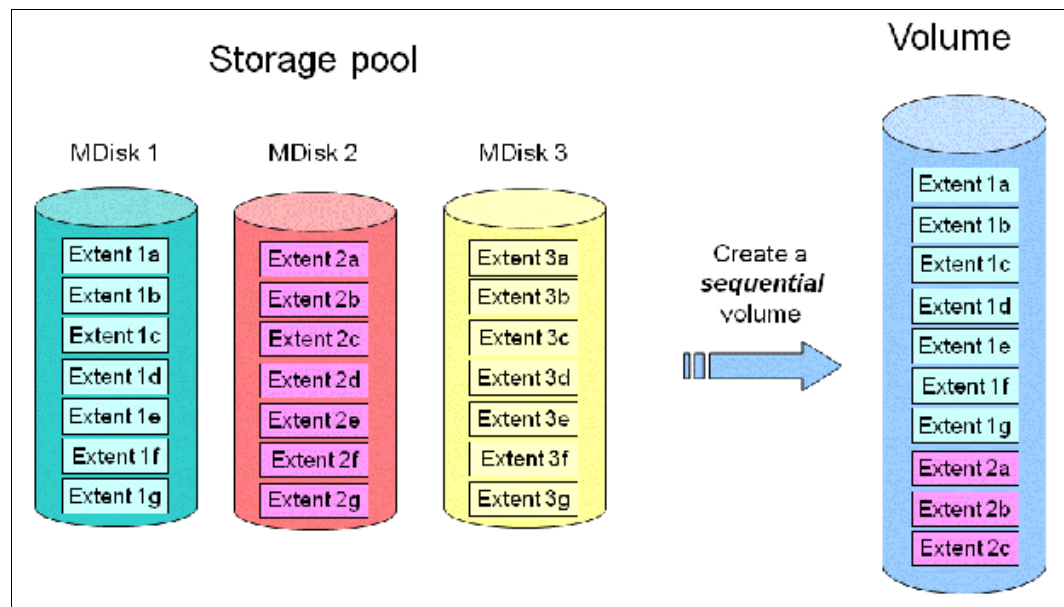


Figure 1-22 Sequential volume

► Image mode

Image mode volumes (Figure 1-23) are special volumes that have a direct relationship with one MDisk. The most common use case of image volumes is a data migration from your old (typically non-virtualized) storage to the SAN Volume Controller based virtualized infrastructure.

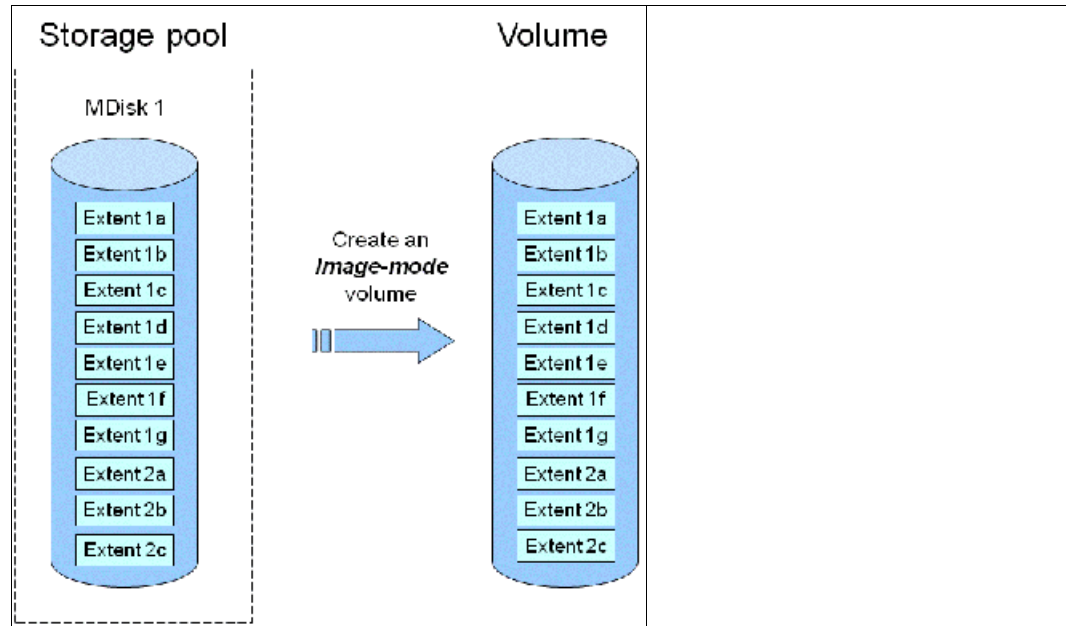


Figure 1-23 Image mode volume

When the image mode volume is created, a direct mapping is made between extents that are on the MDisk and the extents that are on the volume. The LBA  $x$  on the MDisk is the same as the LBA  $x$  on the volume, which ensures that the data on the MDisk is preserved as it is brought into the clustered system.

Some virtualization functions are not available for image mode volumes, so it is useful to migrate the volume into a new storage pool. After the migration completion, the MDisk becomes a managed MDisk.

If you add an MDisk containing any historical data to a storage pool, all data on the MDisk is lost. Ensure that you create image mode volumes from MDisks that contain data before adding MDisks to the storage pools.

### 1.5.13 Easy Tier

Easy Tier is a performance function that automatically migrates or moves extents off a volume to or from one MDisk storage tier to another MDisk storage tier. Since V7.3, the SAN Volume Controller code can support a three-tier implementation.

Easy Tier monitors the host I/O activity and latency on the extents of all volumes with the Easy Tier function that is turned on in a multitier storage pool over a 24-hour period. Then, it creates an extent migration plan that is based on this activity, and then dynamically moves high-activity or hot extents to a higher disk tier within the storage pool. It also moves extents whose activity dropped off or cooled down from the high-tier MDisks back to a lower-tiered MDisk.

Easy Tier supports the new SCM drives with a new tier that is called `tier_scm`.



**Turning on or off Easy Tier:** The Easy Tier function can be turned on or off at the storage pool level and the volume level.

The automatic load-balancing function is enabled by default on each volume, and cannot be turned off by using the GUI. This load-balancing feature is not considered as an Easy Tier function, although it uses the same principles.

The Easy Tier function can make it more appropriate to use smaller storage pool extent sizes. The usage statistics file can be offloaded from the SAN Volume Controller nodes. Then, you can use IBM Storage Tier Advisor Tool (STAT) to create a summary report. STAT is available on the web at no initial cost at the following link:

[IBM STAT tool](#)

A more detailed description of Easy Tier is provided in Chapter 9, “Advanced features for storage efficiency” on page 449.

## 1.5.14 Hosts

Volumes can be mapped to a *host* to enable access for a specific server to a set of volumes. A host within the SAN Volume Controller is a collection of HBA worldwide port names (WWPNs) or iSCSI Qualified Names (IQNs) that are defined on the specific server.

**Note:** iSCSI names are internally identified by “fake” WWPNs, which are WWPNs that are generated by the SAN Volume Controller. Volumes can be mapped to multiple hosts, for example, a volume that is accessed by multiple hosts of a server system.

iSCSI is an alternative way of attaching hosts starting with SAN Volume Controller V7.7. In addition, back-end storage can be attached by using iSCSI. This configuration is useful for migration purposes from non- Fibre Channel based environments to the new virtualized solution.

Node failover can be handled without having a multipath driver that is installed on the iSCSI server. An iSCSI-attached server can reconnect after a node failover to the original target IP address, which is now presented by the partner node. To protect the server against link failures in the network or HBA failures, using a multipath driver is mandatory.

Volumes are LUN-masked to the host’s HBA WWPNs by a process called *host mapping*. Mapping a volume to the host makes it accessible to the WWPNs or IQNs that are configured on the host object. For a SCSI over Ethernet connection, the IQN identifies the iSCSI target (destination) adapter. Host objects can have IQNs and WWPNs.

## 1.5.15 Host cluster

A *host cluster* is a host object in the SAN Volume Controller. It is a combination of two or more servers that is connected to SAN Volume Controller through an FC, FCoE, or an iSCSI connection. A host cluster object can see the same set of volumes. Therefore, volumes can be mapped to a host cluster to enable all hosts to have a common mapping.

**Note:** SAN Volume Controller models SA2 and SV2 do not support FCoE.

## 1.5.16 RAID

When planning your network, consider the type of RAID configuration. The SAN Volume Controller supports either a traditional array or a distributed array.

An array can contain 2 - 16 drives, and several arrays create the capacity for a pool. For redundancy, spare drives (*hot spares*) are allocated to assume read/write operations if any of the other drives fail. The rest of the time, the spare drives are idle and do not process requests for the system.

When an array member drive fails, the system automatically replaces the failed member with a hot spare drive and rebuilds the array to restore its redundancy. Candidate and spare drives can be manually exchanged with array members.

Distributed array configurations can contain 4 - 128 drives. Distributed arrays remove the need for separate drives that are idle until a failure occurs. Rather than allocating one or more drives as spares, the spare capacity is distributed over specific rebuild areas across all the member drives. Data can be copied faster to the rebuild area and redundancy is restored much more rapidly. Additionally, as the rebuild progresses, the performance of the pool is more uniform because all the available drives are used for every volume extent.

After the failed drive is replaced, data is copied back to the drive from the distributed spare capacity. Unlike hot spare drives, read/write requests are processed on other parts of the drive that are not being used as rebuild areas. The number of rebuild areas is based on the width of the array.

## 1.5.17 Encryption

The SAN Volume Controller provides optional encryption of data at rest, which protects against the potential exposure of sensitive user data and user metadata that is stored on discarded, lost, or stolen storage devices. Encryption of system data and system metadata is not required, so system data and metadata are not encrypted.

Planning for encryption involves purchasing a licensed function and then activating and enabling the function on the system.

To encrypt data that is stored on drives, the nodes capable of encryption must be licensed and configured to use encryption. When encryption is activated and enabled on the system, valid encryption keys must be present on the system when the system unlocks the drives or the user generates a new key.

In IBM Spectrum Virtualize V7.4, hardware encryption was introduced with the software encryption option introduced in Version 7.6. Encryption keys can either be managed by IBM Security™ Key Lifecycle Manager (SKLM) or stored on USB flash drives that are attached to a minimum of one of the nodes. Since Version 8.1, IBM Spectrum Virtualize provides a combination of SKLM and USB key repositories.

SKLM is an IBM solution to provide the infrastructure and processes to locally create, distribute, backup, and manage the lifecycle of encryption keys and certificates. Before activating and enabling encryption, you must determine the method of accessing key information during times when the system requires an encryption key to be present.

When SKLM is used as a key manager for the SAN Volume Controller encryption, you can run into a deadlock situation if the key servers are running on encrypted storage that is provided by the SAN Volume Controller. To avoid a deadlock situation, ensure that the SAN Volume Controller can communicate with an encryption server to get the unlock key after a power-on or restart scenario. Up to four SKLM servers are supported.

Data encryption is protected by the Advanced Encryption Standard (AES) algorithm that uses a 256-bit symmetric encryption key in XTS mode, as defined in the Institute of Electrical and Electronics Engineers (IEEE) 1619-2007 standard as XTS-AES-256. That data encryption key is itself protected by a 256-bit AES key wrap when stored in non-volatile form.

Because data security and encryption play significant roles in today's storage environments, see more details in Chapter 12, "Encryption" on page 685.

### 1.5.18 iSCSI

iSCSI is an alternative means of attaching hosts and external storage controllers to the SAN Volume Controller.

The iSCSI function is a software function that is provided by the IBM Spectrum Virtualize code, not hardware. In Version 7.7, IBM introduced software capabilities to enable the underlying virtualized storage to attach to SAN Volume Controller by using the iSCSI protocol.

The iSCSI protocol enables the transport of SCSI commands and data over an IP network (TCP/IP), which is based on IP routers and Ethernet switches. iSCSI is a block-level protocol that encapsulates SCSI commands. Therefore, it uses an existing IP network rather than FC infrastructure.

The major functions of iSCSI include encapsulation and the reliable delivery of CDB transactions between initiators and targets through the IP network, especially over a potentially unreliable IP network.

Every iSCSI node in the network must have an iSCSI name and address:

- ▶ An *iSCSI name* is a location-independent, permanent identifier for an iSCSI node. An iSCSI node has one iSCSI name, which stays constant for the life of the node. The terms *initiator name* and *target name* also refer to an iSCSI name.
- ▶ An *iSCSI address* specifies the iSCSI name of an iSCSI node and a location of that node. The address consists of a host name or IP address, a TCP port number (for the target), and the iSCSI name of the node. An iSCSI node can have any number of addresses, which can change at any time, particularly if they are assigned by way of Dynamic Host Configuration Protocol (DHCP). A SAN Volume Controller node represents an iSCSI node and provides statically allocated IP addresses.

### 1.5.19 IBM Real-time Compression

RtC is an attractive solution to address the increasing requirements for data storage, power, cooling, and floor space. When applied, RtC can save storage space so more data can be stored, and fewer storage enclosures are required to store a data set.

RtC provides the following benefits:

- ▶ Compression for active primary data. RtC can be used with active primary data.
- ▶ Compression for replicated/mirrored data. Remote volume copies can be compressed in addition to the volumes at the primary storage tier. This process also reduces storage requirements in MM and GM destination volumes.
- ▶ No changes to the existing environment are required. RtC is part of the storage system.
- ▶ Overall savings in operational expenses. More data is stored and fewer storage expansion enclosures are required. Reducing rack space has the following benefits:
  - Reduced power and cooling requirements. More data is stored in a system, requiring less power and cooling per gigabyte or used capacity.
  - Reduced software licensing for more functions in the system. More data is stored per enclosure, which reduces the overall spending on licensing.
- ▶ Disk space savings are immediate. The space reduction occurs when the host writes the data. This process is unlike other compression solutions in which some or all the reduction is realized only after a post-process compression batch job is run.

When compression is applied, it is a best practice to monitor the overall performance and CPU utilization. Compression can be implemented without any impact to the existing environment, and it can be used with storage processes running.

**Note:** SAN Volume Controller models SV2 and SA2 do not support RtC software compression. They support only the newer DRP software compression.

## 1.5.20 Data Reduction Pools

DRPs represent a significant enhancement to the storage pool concept. The reason is that the virtualization layer is primarily a simple layer that runs the task of lookups between virtual and physical extents.

DRP is a new type of storage pool, implementing techniques such as thin-provisioning, compression, and deduplication to reduce the amount of physical capacity that is required to store data. In addition, DRP decreases the network infrastructure that is required. Savings in storage capacity requirements translate into the reduction of the cost of storing the data.

With the storage pools, you can automatically de-allocate and reclaim the capacity of thin-provisioned volumes that contain deleted data and enable this reclaimed capacity to be reused by other volumes. Data reduction provides more performance from compressed volumes due to the implementation of the new log structured pool.

## 1.5.21 Deduplication

Data deduplication is one of the methods of reducing storage needs by eliminating redundant copies of a file. Data reduction is a way to decrease the storage disk and network infrastructure that is required, optimize the usage of existing storage disks, and improve data recovery infrastructure efficiency. Existing data or new data is standardized into chunks that are examined for redundancy. If data duplicates are detected, then pointers are shifted to reference a single copy of the chunk, and the duplicate data sets are then released.

Deduplication has several benefits, such as storing more data per physical storage system, saving energy by using fewer disk drives, and decreasing the amount of data that must be sent across a network to another storage for backup replication and for DR.

## 1.5.22 IP replication

IP replication was introduced in Version 7.2, and it enables data replication between IBM Spectrum Virtualize family members. IP replication uses the IP-based ports of the cluster nodes.

The IP replication function is transparent to servers and applications like traditional FC-based mirroring is. All remote mirroring modes (MM, GM, and GMCV) are supported.

The configuration of the system is straightforward, and IBM Storwize family systems normally “find” each other in the network and can be selected from the GUI.

IP replication includes Bridgeworks SANSlide network optimization technology, and it is available at no additional charge. Remember, remote mirror is a chargeable option, but the price does not change with IP replication. Existing remote mirror users have access to the function at no additional charge.

IP connections that are used for replication can have long latency (the time to transmit a signal from one end to the other), which can be caused by distance or by many “hops” between switches and other appliances in the network. Traditional replication solutions transmit data, wait for a response, and then transmit more data, which can result in network utilization as low as 20% (based on IBM measurements). In addition, this scenario gets worse the longer the latency.

Bridgeworks SANSlide technology, which is integrated with the IBM Storwize family, requires no separate appliances and so requires no additional cost and no configuration steps. It uses artificial intelligence (AI) technology to transmit multiple data streams in parallel, adjusting automatically to changing network environments and workloads.

SANSlide improves network bandwidth utilization up to 3x. Therefore, customers can deploy a less costly network infrastructure, or take advantage of faster data transfer to speed replication cycles, improve remote data currency, and enjoy faster recovery.

## 1.5.23 IBM Spectrum Virtualize copy services

IBM Spectrum Virtualize supports the following copy services:

- ▶ Synchronous RC (MM)
- ▶ Asynchronous RC (GM)
- ▶ FlashCopy (PiT copy)
- ▶ TCT

Copy services functions are implemented within a single SAN Volume Controller, or between multiple members of the IBM Spectrum Virtualize family.

The copy services layer sits above and operates independently of the function or characteristics of the underlying disk subsystems that are used to provide storage resources to a SAN Volume Controller.

## 1.5.24 Synchronous or asynchronous Remote Copy

The general application of RC seeks to maintain two copies of data. Often, the two copies are separated by distance, but not always. The RC can be maintained in either synchronous or asynchronous modes. IBM Spectrum Virtualize, MM, and GM are the IBM branded terms for the functions that are synchronous RC and asynchronous RC.

Synchronous RC ensures that updates are committed at both the primary and the secondary volumes before the application considers the updates complete. Therefore, the secondary volume is fully up to date if it is needed in a failover. However, the application is fully exposed to the latency and bandwidth limitations of the communication link to the secondary volume. In a truly remote situation, this extra latency can have a significant adverse effect on application performance.

Special configuration guidelines exist for SAN fabrics and IP networks that are used for data replication. Consider the distance and available bandwidth of the intersite links.

A function of GM for low bandwidth was introduced in IBM Spectrum Virtualize. It uses change volumes that are associated with the primary and secondary volumes. These volumes are used to record changes to the RC volume, the FlashCopy relationship that exists between the secondary volume and the change volume, and between the primary volume and the change volume. This function is called *GM cycling mode*.

Figure 1-24 shows an example of this function where you can see the relationship between volumes and change volumes.

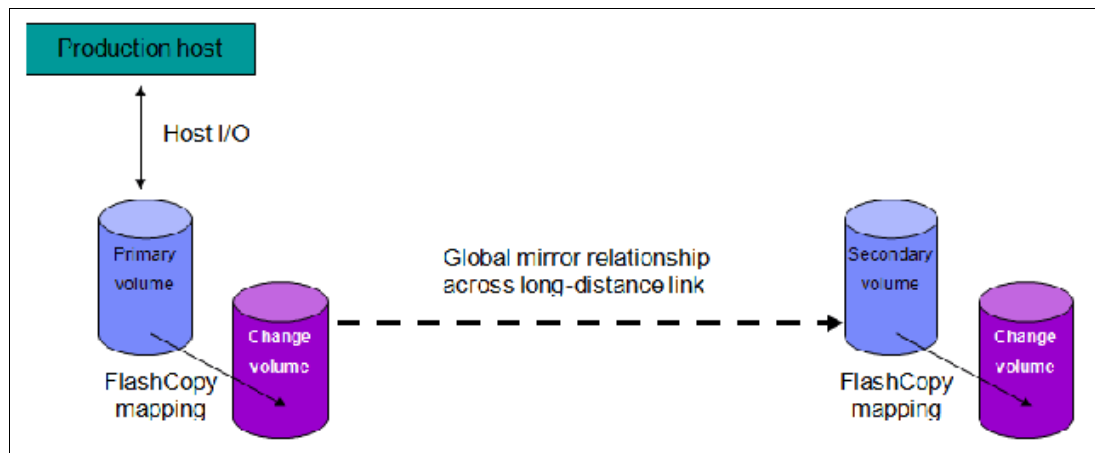


Figure 1-24 Global Mirror with change volumes

In asynchronous RC, the application acknowledges that the write is complete before the write is committed at the secondary volume. Therefore, on a failover, certain updates (data) might be missing at the secondary volume. The application must have an external mechanism for recovering the missing updates, if possible. This mechanism can involve user intervention. Recovery on the secondary site involves starting the application on this recent backup, and then rolling forward or backward to the most recent commit point.

### 1.5.25 FlashCopy and Transparent Cloud Tiering

FlashCopy and TCT are used to make a copy of a source volume on a target volume. After the copy operation starts, the original content of the target volume is lost, and the target volume has the contents of the source volume as they existed at a single PIT. Although the copy operation takes time, the resulting data at the target appears as though the copy was made instantaneously.

## FlashCopy

FlashCopy is sometimes described as an instance of a time-zero (T0) copy or a PiT copy technology.

FlashCopy can be performed on multiple source and target volumes. FlashCopy enables management operations to be coordinated so that a common single PiT is chosen for copying target volumes from their respective source volumes.

With IBM Spectrum Virtualize, multiple target volumes can undergo FlashCopy from the same source volume. This capability can be used to create images from separate PiTs for the source volume, and to create multiple images from a source volume at a common PiT. Source and target volumes can be thin-provisioned volumes.

Reverse FlashCopy enables target volumes to become restore points for the source volume without breaking the FlashCopy relationship, and without waiting for the original copy operation to complete. IBM Spectrum Virtualize supports multiple targets and multiple rollback points.

Most clients aim to integrate the FlashCopy feature for PiT copies and quick recovery of their applications and databases. An IBM solution for this goal is provided by IBM Spectrum Protect, which is described at the following website:

[Data Protection and Recovery](#)

## Transparent Cloud Tiering

IBM Spectrum Virtualize Transparent Cloud Tiering is a function that was introduced in IBM Spectrum Virtualize V7.8. TCT is an alternative solution for data protection, backup, and restore that interfaces to CSPs, such as IBM Cloud®. The TCT function helps organizations to reduce costs that are related to power and cooling when offsite data protection is required to send sensitive data out of the main site.

TCT uses IBM FlashCopy techniques that provide full and incremental snapshots of several volumes. Snapshots are encrypted and compressed before being uploaded to the cloud. Reverse operations are also supported within that function. When a set of data is transferred out to cloud, the volume snapshot is stored as object storage.

IBM Cloud Object Storage uses innovative approach and a cost-effective solution to store a large amount of unstructured data, and delivers mechanisms to provide security services, HA, and reliability.

The management GUI provides an easy-to-use initial setup, advanced security settings, and audit logs that records all backup and restore to cloud.

To learn more about IBM Cloud Object Storage, go to the following website:

[IBM Cloud Object Storage](#)

## 1.6 Business continuity

In simple terms, a *clustered system* or *system* is a collection of servers that together provide a set of resources to a client. The key point is that the client has no knowledge of the underlying physical hardware of the system. The client is isolated and protected from changes to the physical hardware. This arrangement offers many benefits including, most significantly, HA.

Resources on the clustered system act as HA versions of unclustered resources. If a node (an individual computer) in the system is unavailable or too busy to respond to a request for a resource, the request is passed transparently to another node that can process the request. The clients are “unaware” of the exact locations of the resources that they use.

The SAN Volume Controller is a collection of up to eight nodes, which are added in pairs that are known as I/O groups. These nodes are managed as a set (system), and they present a single point of control to the administrator for configuration and service activity.

The eight-node limit for a SAN Volume Controller system is a limitation that is imposed by the Licensed Internal Code, and not a limit of the underlying architecture. Larger system configurations might be available in the future.

Although the SAN Volume Controller code is based on a purpose-optimized Linux kernel, the clustered system feature is not based on Linux clustering code. The clustered system software within the SAN Volume Controller, that is, the event manager cluster framework, is based on the outcome of the COMPASS research project. It is the key element that isolates the SAN Volume Controller application from the underlying hardware nodes.

The clustered system software makes the code portable. It provides the means to keep the single instances of the SAN Volume Controller code that are running on separate systems’ nodes in sync. Therefore, restarting nodes during a code upgrade, adding nodes, removing nodes from a system, or failing nodes cannot affect SAN Volume Controller availability.

All active nodes of a system must know that they are members of the system. This knowledge is especially important in situations where it is key to have a solid mechanism to decide which nodes form the active system, such as the split-brain scenario where single nodes lose contact with other nodes. A worst case scenario is a system that splits into two separate systems.

Within a SAN Volume Controller system, the *voting set* and a quorum disk are responsible for the integrity of the system. If nodes are added to a system, they are added to the voting set. If nodes are removed, they are removed quickly from the voting set. Over time, the voting set and the nodes in the system can change so that the system migrates onto a separate set of nodes from the set on which it started.

The SAN Volume Controller clustered system implements a dynamic quorum. Following a loss of nodes, if the system can continue to operate, it adjusts the quorum requirement so that further node failure can be tolerated.

The lowest Node Unique ID in a system becomes the boss node for the group of nodes. It determines (from the quorum rules) whether the nodes can operate as the system. This node also presents the maximum two-cluster IP addresses on one or both of its nodes’ Ethernet ports to enable access for system management.

## 1.6.1 Business continuity with Stretched Clusters

Within standard implementations of the SAN Volume Controller, all the I/O group nodes are physically installed in the same location. To supply the different HA needs that customers have, the *stretched system configuration* was introduced. In this configuration, each node (from the same I/O group) on the system is physically on a different site. When implemented with mirroring technologies, such as volume mirroring or copy services, these configurations can be used to maintain access to data on the system if there are power failures or site-wide outages.



Stretched Clusters are considered HA solutions because both sites work as instances of the production environment (there is no standby location). Combined with application and infrastructure layers of redundancy, Stretched Clusters can provide enough protection for data that requires availability and resiliency.

When the SAN Volume Controller was first introduced, the maximum supported distance between nodes within an I/O group was 100 m. With the evolution of code and the introduction of new features, SAN Volume Controller V5.1 introduced support for the Stretched Cluster configuration. In this configuration, nodes within an I/O group can be separated by a distance of up to 10 km by using specific configurations.

SAN Volume Controller V6.3 began supporting Stretched Cluster configurations. In these configurations, nodes can be separated by a distance of up to 300 km with specific configurations that use FC switch inter-switch links (ISLs) between the different locations.

## 1.6.2 Business continuity with Enhanced Stretched Cluster

IBM Spectrum Virtualize V7.2 introduced the Enhanced Stretched Cluster (ESC) feature that further improved Stretched Cluster configurations. Version 7.2 introduced the *site awareness* concept for nodes and external storage, and the DR feature that enables you to manage effectively rolling disaster scenarios.

Within IBM Spectrum Virtualize V7.5, the site awareness concept is extended to hosts. This change enables more efficiency for host I/O traffic through the SAN, and an easier host path management.

IBM Spectrum Virtualize V7.6 introduces a new feature for stretched systems, the *IP quorum* application. Using an IP-based quorum application as the quorum device for the third site, no FC connectivity is required. Java applications run on hosts at the third site.

However, there are strict requirements on the IP network when using IP quorum applications. Unlike quorum disks, all IP quorum applications must be reconfigured and redeployed to hosts when certain aspects of the system configuration change.

IP quorum details can be found at [IBM Knowledge Center](#) for SAN Volume Controller by searching for the term “IP quorum”.

**Note:** Stretched Cluster and ESC features are supported only for SAN Volume Controller. They are not supported for the IBM Storwize family of products.

## 1.6.3 Business continuity with HyperSwap

The HyperSwap HA feature in the IBM Spectrum Virtualize software enables business continuity during hardware failure, power failure, connectivity failure, or disasters, such as fire or flooding. The HyperSwap feature is available on the SAN Volume Controller, IBM Storwize family, IBM Storwize V7000 Unified, and IBM FlashSystem products running IBM Spectrum Virtualize software.

The HyperSwap feature provides HA volumes that are accessible through two sites at up to 300 km apart. A fully independent copy of the data is maintained at each site. When data is written by hosts at either site, both copies are synchronously updated before the write operation is completed. The HyperSwap feature automatically optimizes itself to minimize data that is transmitted between sites and to minimize host read and write latency.

HyperSwap includes the following key features:

- ▶ Works with SAN Volume Controller and IBM Storwize, IBM FlashSystem products running IBM Spectrum Virtualize software, and V7000 Unified hardware.
- ▶ Uses intra-cluster synchronous RC (MM) capabilities along with existing change volume and access I/O group technologies.
- ▶ Makes a host's volumes accessible across two I/O groups in a clustered system by using the MM relationship in the background. They look like a single volume to the host.
- ▶ Works with the standard multipathing drivers that are available on a wide variety of host types, with no additional host support that is required to access the HA volume.

For further technical details and implementation guidelines about deploying Stretched Cluster or ESC, see *IBM Spectrum Virtualize and SAN Volume Controller Enhanced Stretched Cluster with VMware*, SG24-8211.

## 1.6.4 Automatic hot spare nodes

In previous stages of SAN Volume Controller development, the scripted *warm standby* procedure enables administrators to configure spare nodes in a cluster. In Version 8.2 the system can automatically take on the spare node to replace a failed node in a cluster or to keep the whole system under maintenance tasks, such as software upgrades. These additional nodes are called *hot spare nodes*.

Up to four nodes can be added to a single cluster, and they must match the hardware type and configuration of your active cluster nodes. For example, in a mixed node cluster you should have one of each node type. Because Version 8.3.1 is supported only on SAN Volume Controller 2145 / 2147-DH8 and 2145 / 2147-SV1, 2145 / 2147-SV2, and 2145 / 2147-SA2 nodes, this mixture is not a problem, but is something to consider. Most clients upgrade the whole cluster to a single node type. However, in addition to the node type, the hardware configurations must match, specifically the amount of memory and the number and placement of FC / compression cards must be identical.

The hot spare node essentially becomes another node in the cluster, but is not doing anything under normal conditions. Only when it is needed does it use the N\_Port ID Virtualization (NPIV) feature of the host virtual ports to take over the job of the failed node. There is approximately a minute before the cluster swaps in a node. This delay is set intentionally to avoid any thrashing when a node fails. In addition, the system must be sure that the node has definitely failed, and is not, for example, restarting.

Because you have NPIV enabled, the host should not “notice” anything during this time. The first thing that happens is the failed nodes virtual host ports fail over to the partner node. Then, when the spare swaps in, they fail over to that node. The cache flushes while only one node is in the I/O group, but when the spare swaps in you get the full cache back.

**Note:** A warm start of active node (code assert or restart) does not cause the hot spare to swap in because the restarted node becomes available within 1 minute.

The other use case for hot spare nodes is during a software upgrade. Normally, the only impact during an upgrade is slightly degraded performance. While the node that is upgrading is down, the partner in the I/O group writes through cache and handles both nodes' workload. So, to work around this limitation, the cluster uses a spare in place of the node that is upgrading. Therefore, the cache does not need to go into write-through mode.

After the upgraded node returns, it is swapped back so that you roll through the nodes as normal, but without any failover and failback at the multipathing layer. This process is handled by the NPIV ports, so the upgrades should be seamless for administrators working in large enterprise SAN Volume Controller deployments.

**Note:** After the cluster commits new code, it also automatically upgrades hot spares to match the cluster code level.

This feature is available only to SAN Volume Controller. While Storwize systems can use NPIV and get the general failover benefits, you cannot get spare canisters or split I/O group in Storwize V7000.

## 1.7 Management and support tools

The IBM Spectrum Virtualize system can be managed through the included management software that runs on the SAN Volume Controller hardware.

### 1.7.1 IBM Assist On-site and Remote Support Assistance

With the IBM Assist On-site tool, a member of the IBM Support team can view your desktop and share control of your server to provide you with a solution. This tool is a remote desktop-sharing solution that is offered through the IBM website. With it, the IBM System Services Representative (IBM SSR) can remotely view your system to troubleshoot a problem.

You can maintain a chat session with the IBM SSR so that you can monitor this activity and either understand how to fix the problem yourself or enable them to fix it for you.

To use the IBM Assist On-site tool, the master console must be able to access the internet. For more information, see [IBM remote assistance: Assist On-site](#).

When you access the website, you sign in and enter a code that the IBM SSR provides to you. This code is unique to each IBM Assist On-site session. A plug-in is downloaded on to your master console to connect you and your IBM SSR to the remote service session. The IBM Assist On-site tool contains several layers of security to protect your applications and your computers. The plug-in is removed after the next restart.

You can also use security features to restrict access by the IBM SSR. Your IBM SSR can provide you with more detailed instructions for using the tool.

The embedded part of the SAN Volume Controller V8.3.1 code is a software toolset that is called Remote Support Client. It establishes a network connection over a secured channel with Remote Support Server in the IBM network. The Remote Support Server provides predictive analysis of the SAN Volume Controller status and assists administrators with troubleshooting and fix activities. *Remote Support Assistance* is available at no extra charge, and no additional license is needed.

## 1.7.2 Event notifications

SAN Volume Controller can use SNMP traps, syslog messages, and a Call Home email to notify you and the IBM Support Center when significant events are detected. Any combination of these notification methods can be used simultaneously.

Notifications are normally sent immediately after an event is raised. Each event that SAN Volume Controller detects is assigned a notification type of Error, Warning, or Information. You can configure the SAN Volume Controller to send each type of notification to specific recipients.

### Simple Network Management Protocol traps

SNMP is a standard protocol for managing networks and exchanging messages. IBM Spectrum Virtualize can send SNMP messages that notify personnel about an event. You can use an SNMP manager to view the SNMP messages that IBM Spectrum Virtualize sends. You can use the management GUI or the CLI to configure and modify your SNMP settings.

You can use the MIB file for SNMP to configure a network management program to receive SNMP messages that are sent by the IBM Spectrum Virtualize.

### Syslog messages

The syslog protocol is a standard protocol for forwarding log messages from a sender to a receiver on an IP network. The IP network can be either Internet Protocol Version 4 (IPv4) or Internet Protocol Version 6 (IPv6).

SAN Volume Controller can send syslog messages that notify personnel about an event. The event messages can be sent in either expanded or concise format. You can use a syslog manager to view the syslog messages that SAN Volume Controller sends.

IBM Spectrum Virtualize uses UDP to transmit the syslog message. You can use the management GUI or the CLI to configure and modify your syslog settings.

### Call Home email

The Call Home feature transmits operational and error-related data to you and IBM through a Simple Mail Transfer Protocol (SMTP) server connection in the form of an event notification email. When configured, this function alerts IBM service personnel about hardware failures and potentially serious configuration or environmental issues. You can use the Call Home function if you have a maintenance contract with IBM or if the SAN Volume Controller is within the warranty period.

To send email, you must configure at least one SMTP server. You can specify as many as five more SMTP servers for backup purposes. The SMTP server must accept the relaying of email from the SAN Volume Controller clustered system IP address. Then, you can use the management GUI or the CLI to configure the email settings, including contact information and email recipients. Set the reply address to a valid email address.

Send a test email to check that all connections and infrastructure are set up correctly. You can disable the Call Home function at any time by using the management GUI or CLI.

## 1.8 Useful IBM SAN Volume Controller web links

For more information about the SAN Volume Controller-related topics, see the following websites:

- ▶ [IBM SAN Volume Controller support page](#)
- ▶ [IBM SAN Volume Controller home page](#)
- ▶ [SAN Volume Controller IBM Knowledge Center](#)





# Planning

This chapter describes steps that are required to plan the installation and configuration of an IBM System Storage SAN Volume Controller in your storage network.

This chapter is *not* intended to provide in-depth information about the described topics; it provides only general guidelines. For an enhanced analysis, see *IBM System Storage SAN Volume Controller, IBM Storwize V7000, and IBM FlashSystem 7200 Best Practices and Performance Guidelines*, SG24-7521.

**Note:** Make sure that the planned configuration is reviewed by IBM or an IBM Business Partner before implementation. Such a review can both increase the quality of the final solution and prevent configuration errors that might impact the solution delivery.

This chapter includes the following topics:

- ▶ 2.1, “General planning rules” on page 54
- ▶ 2.2, “Planning for availability” on page 55
- ▶ 2.3, “Physical installation planning” on page 56
- ▶ 2.4, “Planning for system management” on page 56
- ▶ 2.5, “Connectivity planning” on page 57
- ▶ 2.6, “Fibre Channel SAN configuration planning” on page 58
- ▶ 2.7, “IP SAN configuration planning” on page 63
- ▶ 2.8, “Planning topology” on page 67
- ▶ 2.9, “Back-end storage configuration” on page 68
- ▶ 2.10, “Internal storage configuration” on page 69
- ▶ 2.11, “Storage pool configuration” on page 69
- ▶ 2.12, “Volume configuration” on page 71
- ▶ 2.13, “Host attachment planning” on page 73
- ▶ 2.14, “Planning copy services” on page 75
- ▶ 2.15, “Data migration” on page 77
- ▶ 2.16, “Performance monitoring with IBM Storage Insights” on page 78
- ▶ 2.17, “Configuration backup procedure” on page 80

## 2.1 General planning rules

To maximize the benefit from a system, installation planning must include several important steps. These steps ensure that the system provides the best possible performance, reliability, and ease of management for your application needs.

The general rule of planning is to define your goals and then plan a solution that makes you able to reach these goals.

Consider the following points when planning a system:

- ▶ Collect and document information about application servers (hosts) that you want to attach to the system and their data:
  - Amount of data in use for each host and growth plans.
  - Data profile: Compressibility and deduplicability.
  - Host traffic profile: Percentage of reads and writes, percentage of sequential/random access patterns, and data block size.
  - Host performance requirements: Input/output operations per second (IOPS) and bandwidth.
- ▶ Perform capacity and performance sizing of a system:
  - Assess the capacity and performance capabilities of the existing back-end storage that is present in the environment and intended to be virtualized by the SAN Volume Controller.
  - Plan and install new back-end storage if it will be deployed with the SAN Volume Controller.
  - Account for future growth.
  - Verify that the capacity assessment results satisfy your performance requirements.

**Note:** Contact your IBM sales representative or IBM Business Partner to perform these calculations.

- ▶ Assess your recovery point objective (RPO) / RTO requirements and plan for high availability (HA) and Remote Copy (RC) functions. Decide whether you require a dual-site deployment with Enhanced Stretched Cluster (ESC) or HyperSwap feature, and decide whether you must implement RC and determine its type (synchronous or asynchronous). Review the extra configuration requirements that are imposed.
- ▶ Define the number of input/output (I/O) groups (control enclosures) and expansion enclosures. The number of necessary enclosures depends on the solution type, overall performance, and capacity requirements.
- ▶ Plan for host attachment interfaces, protocols, and storage area network (SAN). Consider the number of ports, bandwidth requirements, and HA.
- ▶ Perform configuration planning by defining the number of internal storage arrays and external storage arrays that will be virtualized. Define a number and the type of pools, the number of volumes, and the capacity of each of the volumes.
- ▶ Define a naming convention for the system nodes, volumes, and other storage objects.
- ▶ Plan a management Internet Protocol (IP) network and management users' authentication system.
- ▶ Plan for the physical location of the equipment in the rack.



- ▶ Verify that your planned environment is a supported configuration.

**Note:** Use [IBM System Storage Interoperation Center \(SSIC\)](#) to check compatibility.

- ▶ Verify that your planned environment does not exceed system configuration limits.

**Note:** For more information about your platform and code version, see [Configuration Limits and Restrictions](#).

- ▶ Review the planning aspects that are described in the following sections of this chapter.

## 2.2 Planning for availability

When planning the deployment of the SAN Volume Controller solution, avoid creating single points of failure (SPOFs). Plan your system availability according to the requirements of your solution. Depending on your availability needs, consider the following aspects:

- ▶ Single-site or multi-site configuration

Multi-site configurations increase solution resiliency, and can be the basis of disaster recovery (DR) solutions. Systems can be configured as a multi-site solution with sites working in active-active mode. Both synchronous and asynchronous data replication is supported by multiple inter-site link options.

- ▶ Using spare nodes.

You can purchase and configure a hot spare node to minimize the impact of hardware failures.

- ▶ Physical separation of system building blocks

A dual-rack deployment might increase the availability of your system if your back-end storage, SAN, and local area network (LAN) infrastructure also do not use a single-rack placement scheme. You can further increase system availability by ensuring that enclosures are powered from different power circuits and in different fire protection zones.

- ▶ Quorum disk placement

The SAN Volume Controller uses three managed disks (MDisks) (external storage LUs) or an IP quorum application as quorum devices for the clustered system. A best practice is to have each quorum device in a separate storage subsystem if possible. Multiple IP quorum application deployment is also recommended.

- ▶ Failure domain sizes

Failure of an MDisk takes offline the whole storage pool that contains this MDisk. To reduce the impact of an MDisk failure, consider reducing the number of back-end storage systems per storage pool and increasing the number of storage pools and reducing their size. This configuration limits the maximum performance of the pool (fewer back-end systems to share the load), increases storage management effort, can lead to less efficient storage capacity consumption, and might be subject to limitations by system configuration maximums.

- ▶ Consistency

Strive to achieve consistent availability levels of all system building blocks. For example, if the solution relies on a single switch that is placed in the same rack as one of the SAN Volume Controller nodes, investment in a dual-rack configuration for placement of the second node is not justified. Any incident affecting the rack that holds the critical switch brings down the whole system no matter where the second SAN Volume Controller node is placed.

## 2.3 Physical installation planning

You must consider several key factors when you plan the physical site of a system. The physical site must have the following characteristics:

- ▶ Sufficient rack space must exist to install controller and disk enclosures.
- ▶ The site must meet the power, cooling, and environmental requirements.

For more information about the power and environmental requirements, see [IBM Knowledge Center](#) and expand **Planning** → **Planning for hardware**.

Your system order includes a printed copy of the *Quick Installation Guide*, which also provides information about environmental and power requirements.

Create a cable connection table that follows your environment's documentation procedure to track the following connections that are required for the setup:

- ▶ Power
- ▶ Serial-attached Small Computer System Interface (SCSI) (SAS)
- ▶ Ethernet
- ▶ Fibre Channel (FC)

When planning for power, plan for a separate independent power source for each of the two redundant power supplies of a system enclosure.

When planning SAN cabling, make sure that your physical topology adheres to zoning rules and recommendations.

SAN Volume Controller physical installation and initial setup is performed by an IBM System Services Representative (IBM SSR).

## 2.4 Planning for system management

Each system's node has a *technician port*. It is a dedicated 1 gigabits per second (Gbps) Ethernet port. The initialization of a system and its basic configuration is performed by using this port. After the initialization is complete, the technician port must remain disconnected from a network and used only to service the system.

For management, each system node requires at least one Ethernet connection. The cable must be connected to port 1, which is a 10 Gbps Ethernet port (it does not negotiate speeds below 1 Gbps) on IBM SAN Volume Controller 2145-SV2 and 2145-SA2 nodes and a 1 Gbps Ethernet port on an IBM SAN Volume Controller 2145-SV1 node. For increased availability, an optional management connection may be configured over Ethernet port 2.

For configuration and management, you must allocate an IP address to each node canister, which is referred to as the *service IP address*. Both Internet Protocol Version 4 (IPv4) and Internet Protocol Version 6 (IPv6) are supported.

In addition to a service IP address on each node, each system has a *cluster management IP address*. The management IP address cannot be the same as any of the defined service IP addresses. The management IP can automatically fail over to another address if there are maintenance actions or a node failure.

For example, a system that consists of two control enclosures requires a minimum of five unique IP addresses: one for each node and one for the system as a whole.

Ethernet ports 1 and 2 are not reserved only for management. They may be also used for internet Small Computer Systems Interface (iSCSI) or IP replication traffic if they are configured to do so. However, management and service IP addresses cannot be used for host or back-end storage communication.

System management is performed by using an embedded GUI that is running on the nodes; the command-line interface (CLI) is also available. To access the management GUI, point a web browser to the system management IP address. To access the management CLI, point a Secure Shell (SSH) client to a management IP and use the default SSH protocol port (22/TCP).

By connecting to a service IP address with a browser, you can use SSH to access the *Service Assistant Interface*, which may be used for maintenance and service tasks.

When you plan your management network, note that the IP Quorum applications and Transparent Cloud Tiering (TCT) are communicating with a system through the management ports. For more information about cloud backup requirements, see 10.3, “Transparent Cloud Tiering” on page 562.

## 2.5 Connectivity planning

A SAN Volume Controller offers a wide range of connectivity options to back-end storage and hosts, such as FC technologies (“traditional” SCSI FC and Non-Volatile Memory Express (NVMe) over Fibre Channel (FC-NVMe) (also known as NVMe over Fabric (NVMe-oF))), and Fibre Channel over Ethernet (FCoE), and IP network technologies (iSCSI and iSCSI Extensions for Remote Direct Memory Access (RDMA) (iSER)). The connection options depend on the hardware configuration.

Table 2-1 lists the communication types that can be used for communicating between system nodes, hosts, and back-end storage systems. All types can be used concurrently.

Table 2-1 Communication options

Communication type	System to host	System to back-end storage	Node to node (intra-cluster)	System to system (replication)
SCSI FC	Yes	Yes	Yes	Yes
NVMe over FC	Yes	No	No	No
FCoE	Yes	Yes	Yes	Yes
iSCSI	Yes	Yes	No	No <sup>a</sup>

Communication type	System to host	System to back-end storage	Node to node (intra-cluster)	System to system (replication)
iSER	Yes	No	Yes	No <sup>a</sup>

a. Replication traffic can be sent over an IP network with native IP replication, which can be configured on both onboard 10-Gigabit Ethernet (GbE) ports and optional 25 GbE ports.

## 2.6 Fibre Channel SAN configuration planning

Each IBM SAN Volume Controller 2145-SV1 node has four card slots, which may be equipped with the following cards for Fibre Channel connectivity:

- ▶ 10 Gbps Ethernet FCoE adapter
- ▶ 16 Gbps FC adapter
- ▶ Two-port 32 Gbps FC adapter (A maximum of two adapters of this type are supported.)

Each IBM SAN Volume Controller 2145-SV2 and 2145-SA2 node has three adapter slots, which can be used for the following FC cards:

- ▶ Four-port 16 Gbps FC card
- ▶ Four-port 32 Gbps FC card

16 Gbps and 32 Gbps cards can be used for both SCSI FC and FC-NVMe attachment.

IBM SAN Volume Controller 2145-SV2 and 2145-SA2 nodes do not support FCoE attachment.

### 2.6.1 Physical topology

The switch configuration for a fabric must comply with the switch manufacturer's configuration rules, which can impose restrictions. For example, a switch manufacturer might limit the number of supported switches or ports in a SAN fabric. Operating outside of the switch manufacturer's rules is not supported.

In an environment where you have a fabric with mixed port speeds (8 Gb, 16 Gb, and 32 Gb), the best practice is to connect the system to the switch operating at the highest speed.

The connections between the system's enclosures (node-to-node traffic) and between a system and the virtualized back-end storage require the best available bandwidth. For optimal performance and reliability, ensure that paths between the system nodes and storage systems do not cross inter-switch links (ISLs). If you use ISLs on these paths, make sure that sufficient bandwidth is available. SAN monitoring is required to identify faulty ISLs.

No more than three ISL hops are permitted among nodes that are in the same system but in different I/O groups. If your configuration requires more than three ISL hops for nodes that are in the same system but in different I/O groups, contact your IBM Support Center.

Direct connection of the system FC ports to host systems or between nodes in the system without using an FC switch is supported. For more information, see [IBM Knowledge Center](#) and expand **Planning** → **Planning your network and storage network** → **Planning for a direct-attached configuration**.

## 2.6.2 Zoning

A SAN fabric must have four distinct zone classes:

- Inter-node zones** For communication between nodes in the same system
- Storage zones** For communication between the system and back-end storage
- Host zones** For communication between the system and hosts
- Inter-system zones** For remote replication

Figure 2-1 shows the system zoning classes.

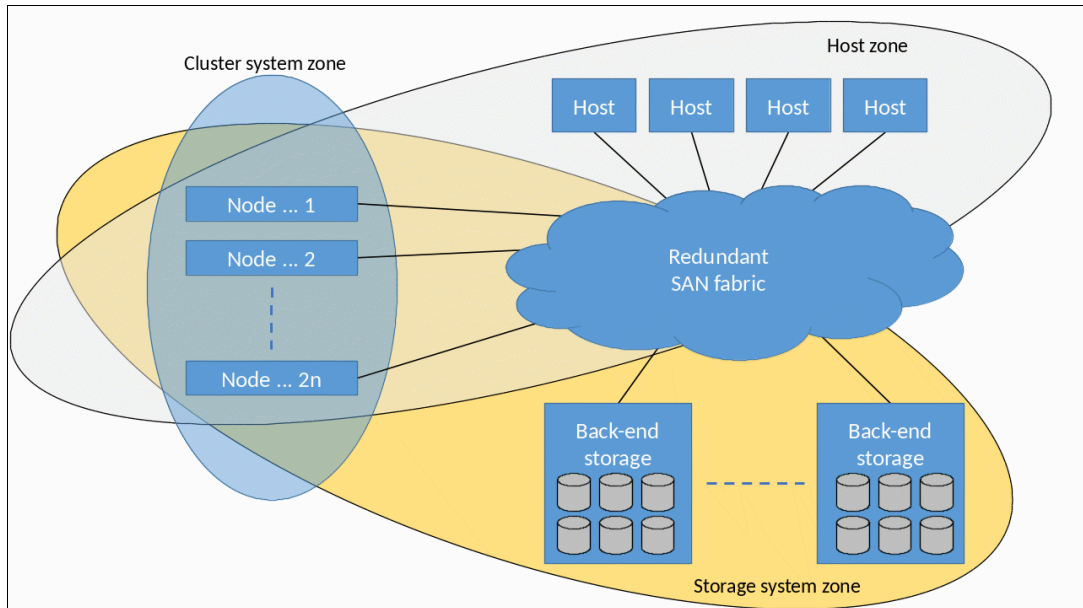


Figure 2-1 System zoning

The fundamental rules of system zoning are described in the rest of this section. However, you should review the latest zoning guidelines and requirements when designing zoning for the planned solution by going to [IBM Knowledge Center](#) and expanding **Configuring** → **Configuration details** → **SAN configuration and zoning rules summary**.

## 2.6.3 N\_Port ID Virtualization

N\_Port ID Virtualization (NPIV) is a method for virtualizing a physical FC port that is used for host I/O. By default, all new systems work in NPIV mode (the Target Port Mode attribute is set to Enabled).

NPIV mode creates a virtual worldwide port name (WWPN) for every system physical FC port. This WWPN is available only for host connection. During node maintenance, restart, or failure, the virtual WWPN from that node is transferred to the same port of the other node in the I/O group.

For more information about NPIV mode and how it works, see 7.4, “N\_Port ID Virtualization support” on page 354.

Ensure that FC switches enable each physically connected system port with the ability to create four extra NPIV ports.

When performing zoning configuration, virtual WWPNs are used only for host communication, that is, “system to host” zones must include virtual WWPNs, and internode, intersystem, and back-end storage zones must use the WWPNs of physical ports. Ensure that equivalent ports (with the same port ID) are on the same fabric and in the same zone.

For more information about other host zoning requirements, see [IBM Knowledge Center](#) and expand **Configuring** → **Configuration details** → **Zoning details** → **Zoning requirements for N\_Port ID Virtualization**.

## 2.6.4 Inter-node zone

The purpose of intracluster or inter-node zones is to enable traffic between all node canisters within the clustered system. This traffic consists of heartbeats, cache synchronization, and other data that nodes must exchange to maintain a healthy cluster state.

A pair of nodes in an I/O group performs write-cache synchronization over the SAN. All delays in this SAN path directly impact system performance.

You may create up to two inter-node zones per fabric. In each of them, place a single port per node that is designated for intracluster traffic. Each node in the system must have at least two ports with paths to all other nodes in the system. A system node cannot have more than 16 paths to another node in the same system.

Mixed port speeds are not possible for intracluster communication. All node ports within a clustered system must be running at the same speed.

## 2.6.5 Back-end storage zones

Create a separate zone for each back-end storage subsystem that is virtualized. Switch zones that contain back-end storage system ports must not have more than 40 ports. A configuration that exceeds 40 ports is not supported.

All nodes in a system must connect to the same set of back-end storage system ports on each device.

If the edge devices contain more stringent zoning requirements, follow the storage system rules to further restrict the system zoning rules.

For more information connecting back-end storage systems, see [IBM Knowledge Center](#) and expand **Configuring** → **Configuration details** → **External storage system configuration details (Fibre Channel)** and **Configuring** → **Configuring and servicing storage systems** → **External storage system configuration with Fibre Channel connections**.

## 2.6.6 Host zones

A host must be zoned to an I/O group to access volumes that are presented by this I/O group.

The preferred zoning policy is *single initiator zoning*. To implement it, create a separate zone for each host bus adapter (HBA) port, and place exactly one port from each node in each I/O group that the host accesses in this zone. For deployments with more than 64 hosts that are defined in the system, this host zoning scheme is mandatory.

For smaller installations, you may have up to 40 FC ports (including both host HBA ports and the system's virtual WWPNs) in a host zone if the zone contains similar HBAs and operating systems (OSs). A valid zone can be 32 host ports plus eight system ports.

FC-NVMe applies more limits to the host zone configuration:

- ▶ Zone up to four host ports to detect up to four ports on a node, and zone the same or more host ports to detect an extra four ports on the second node of the I/O group.
- ▶ Zone a total maximum of 16 hosts to detect a single I/O group.

Consider the following rules for zoning hosts over either SCSI or FC-NVMe:

- ▶ For any volume, the number of paths through the SAN from the host to a system must not exceed eight. For most configurations, four paths to an I/O group are sufficient.

Except by zoning, you can use a *port mask* to control the number of host paths. For more information, see Chapter 7, "Hosts" on page 351.

- ▶ Balance the host load across the system's ports. For example, zone the first host with ports 1 and 3 of each node in I/O group, zone the second host with ports 2 and 4, and so on. To obtain the best overall performance of the system, the load of each port should be equal. Assuming that a similar load is generated by each host, you can achieve this balance by zoning approximately the same number of host ports to each port.
- ▶ Spread the load across all system ports. Use all ports that are available on your machine.
- ▶ Balance the host load across HBA ports. If the host has more than one HBA port per fabric, zone each host port with a separate group of system ports.

All paths must be managed by the multipath driver on the host side. Make sure that the multipath driver on each server can handle the number of paths that is required to access all volumes that are mapped to the host.

## 2.6.7 Zoning considerations for Metro Mirror and Global Mirror

The SAN configurations that use inter-cluster Metro Mirror (MM) and Global Mirror (GM) relationships have the following extra switch zoning requirements:

- ▶ If two ISLs are connecting the sites, split the ports from each node between the ISLs, that is, exactly one port from each node must be zoned across each ISL.
- ▶ Local clustered system zoning continues to follow the standard requirement for all ports on all nodes in a clustered system to be zoned to one another.
- ▶ Review the latest requirements and recommendations at [IBM Knowledge Center](#) by expanding **Configuring** → **Configuration details** → **Zoning details** → **Zoning constraints for Metro Mirror and Global Mirror**.

When designing zoning for a geographically dispersed solution, consider the effect of the cross-site links on the performance of the local system.

Using mixed port speeds for intercluster communication can lead to port congestion, which can negatively affect the performance and resiliency of the SAN. Therefore, it is not supported.

**Note:** If you limit the number of ports that are used for remote replication to two ports on each node, you can limit the effect of a severe and abrupt overload of the intercluster link on system operations.

If all node ports (N\_Ports) are zoned for intercluster communication and the intercluster link becomes severely and abruptly overloaded, the local FC fabric can become congested so that no FC ports on the local system can perform local intracluster communication, which can result in cluster consistency disruption.

For more information about how to avoid such situations, see 2.6.8, “Port designation recommendations” on page 62.

For more information about zoning best practices, see *IBM System Storage SAN Volume Controller, IBM Storwize V7000, and IBM FlashSystem 7200 Best Practices and Performance Guidelines*, SG24-7521.

## 2.6.8 Port designation recommendations

If you have enough available FC ports on the system, designate different types of traffic to different ports. This configuration provides a level of protection against malfunctioning devices and workload spikes that might otherwise impact the system.

Intra-cluster communication must be protected because it is used for heartbeat and metadata exchange between all nodes of all I/O groups of the cluster. Nodes in one I/O group use it to synchronize write cache.

Isolating remote replication traffic to dedicated ports is beneficial because it ensures that any problems that affect the cluster-to-cluster interconnect do not affect all ports on the local cluster.

To isolate both node-to-node and system-to-system traffic, use the port designations that are shown in Figure 2-2.

Card / Port	4 ports	8 ports	12 ports	16 port
Card 1 Port 1	Host/Storage/Inter-node	Host/Storage	Host/Storage	Host/Storage
Card 1 Port 2	Host/Storage/Inter-node	Host/Storage	Host/Storage	Host/Storage
Card 1 Port 3	Host/Storage/Replication*	Inter-node	Inter-node	Host/Storage
Card 1 Port 4	Host/Storage/Replication*	Inter-node	Inter-node	Host/Storage
Card 2 Port 1		Host/Storage	Host/Storage	Intracluster
Card 2 Port 2		Host/Storage	Host/Storage	Intracluster
Card 2 Port 3		Host/Storage/Replication*	Host/Storage/Replication*	Replication or Host *
Card 2 Port 4		Host/Storage/Replication*	Host/Storage/Replication*	Host/Storage
Card 3 Port 1			Host/Storage	Host/Storage
Card 3 Port 2			Host/Storage	Replication or Host *
Card 3 Port 3			Host/Storage	Intracluster
Card 3 Port 4			Host/Storage	Intracluster
Card 4 Port 1				Host/Storage
Card 4 Port 2				Host/Storage
Card 4 Port 3				Host/Storage
Card 4 Port 4				Host/Storage
localfcportmask	0011	00001100	00000001100	0000110000110000
remotefcportmask	1100	11000000	000011000000	0000001001000000
* Use for host/storage in case no replication is in place. ** Do not use the same port for replication and inter-node traffic. *** For HyperSwap, dedicate ports for inter-node traffic				

Figure 2-2 Port masking configuration



To achieve traffic isolation, use a combination of SAN zoning and *local and partner port masking*. For more information about how to send port masks, see Chapter 3, “Initial configuration” on page 81.

Alternative port mappings that spread traffic across HBAs might allow adapters to come back online after a failure. However, they do not prevent a node from going offline temporarily to restart and attempt to isolate the failed adapter and then rejoin the cluster. Also, the mean time between failures (MTBF) of the adapter is not significantly shorter than MTBF of the non-redundant node components. The approach that is presented here accounts for all these considerations with the idea that increased complexity can lead to migration challenges in the future, so a simpler approach is better.

## 2.7 IP SAN configuration planning

IBM SAN Volume Controller 2145-SV1 nodes have three onboard 1 Gbps Ethernet interfaces that can be expanded up to four 2-port or 4-port 10Gbps-FC adapters. IBM SAN Volume Controller 2145-SV2 and 2145-SA2 nodes are equipped with four onboard 10 Gbps Ethernet network interface ports. They can operate at link speeds of 1 Gbps and 10 Gbps. Any of those ports can be used for host I/O with iSCSI, external storage virtualization with iSCSI, and for Native IP Replication. Also, ports 1 and 2 may be used for managing the system.

IBM SAN Volume Controller SV1, SV2, and AE2 nodes may also be configured with 2-port 25 Gbps Remote Direct Memory Access (RDMA)-capable Ethernet adapters. The maximum number of adapters depends on the system hardware type. Adapters can auto-negotiate link speeds of 1 - 25 Gbps. All their ports may be used for host I/O with iSCSI or iSER, external storage virtualization with iSCSI, node-to-node traffic, and for IP replication.

You can set virtual local area network (VLAN) settings to separate network traffic for Ethernet transport. The system supports VLAN configurations for the system, host attachment, storage virtualization, and IP replication traffic. VLANs can be used with priority flow control (PFC) (IEEE 802.1Qbb).

All ports may be configured with an IPv4 address, an IPv6 address, or both. Each application of a port needs a separate IP. For example, port 1 of every node can be used for management, iSCSI, and IP replication, but three unique IP addresses are required.

If node Ethernet ports are connected to different isolated networks, then a different subnet must be used for each network.

### 2.7.1 iSCSI and iSER protocols

The iSCSI protocol is a block-level access protocol that encapsulates SCSI commands into TCP/IP packets. Therefore, iSCSI uses an IP network rather than requiring the FC infrastructure.

The iSER is a network protocol that extends iSCSI to use RDMA. RDMA is provided by either the internet Wide Area RDMA Protocol (iWARP) or RDMA over Converged Ethernet (RoCE). It permits data to be transferred directly into and out of SCSI buffers, providing faster connection and processing time than traditional iSCSI connections.

iSER requires optional 25 Gbps RDMA-capable Ethernet cards. RDMA links work only between RoCE ports or between iWARP ports: from a RoCE node canister port to a RoCE port on a host, or from an iWARP node canister port to an iWARP port on a host. So, there are two types of 25 Gbps adapters that are available for a system, and they cannot be interchanged without a similar RDMA type change on the host side.

Either iSCSI or iSER works for standard iSCSI communications, that is, ones that do not use RDMA.

The 25 Gbps adapters come with SFP28 fitted, which can be used to connect to switches that use OM3 optical cables.

For more information about the Ethernet switches and adapters that are supported by iSER adapters, see [SSIC](#).

**Note:** At the time of writing, connecting a 10 Gb switch to a 25 Gb interface is supported only through a SCORE request. For more information, contact your IBM representative.

## 2.7.2 Priority flow control

PFC is an Ethernet protocol that you can use to select the priority of different types of traffic within the network. With PFC, administrators can reduce network congestion by slowing or pausing certain classes of traffic on ports, thus providing better bandwidth for more important traffic. The system supports PFC on various supported Ethernet-based protocols on three types of traffic classes: system (node-to-node), host attachment, and back-end storage traffic.

You can configure a priority tag for each of these traffic classes. The priority tag can be any value 0 - 7. You can set identical or different priority tag values to all these traffic classes. You can also set bandwidth limits to ensure quality of service (QoS) for these traffic classes by using the Enhanced Transmission Selection (ETS) setting on the network.

To use PFC and ETS, ensure that the following tasks are completed:

- ▶ Configure a VLAN on the system to use PFC capabilities for the configured IP version.
- ▶ Ensure that the same VLAN settings are configured on the all entities, including all switches between the communicating end points.
- ▶ On the switch, enable Data Center Bridging Exchange (DCBx). DCBx enables switch and adapter ports to exchange parameters that describe traffic classes and PFC capabilities. For these steps, check your switch documentation for details.
- ▶ For each supported traffic class, configure the same priority tag on the switch. For example, if you plan to have a priority tag setting of 3 for storage traffic, ensure that the priority is also set to 3 on the switch for that traffic type.
- ▶ If you are planning on using the same port for different types of traffic, ensure that ETS settings are configured on the network.

For more information, see [IBM Knowledge Center](#) and expand **Configuring** → **Configuring priority flow control**.

### 2.7.3 RDMA clustering

An IBM SAN Volume Controller may use 25 Gbps cards for node-to-node traffic. A dual-site HyperSwap configuration can also use the cards for an inter-site link.

A minimum of two dedicated RDMA-capable ports are required for node-to-node RDMA communications to ensure best performance and reliability. These ports must be configured for inter-node traffic only and cannot be used for host attachment, virtualization of Ethernet-attached external storage, or IP replication traffic.

The following limitations apply to a configuration of ports that are used for RDMA-clustering:

- ▶ Only IPv4 addresses are supported.
- ▶ Only the default value of 1500 is supported for the maximum transmission unit (MTU).
- ▶ Port masking is not supported on RDMA-capable Ethernet ports. Due to this limitation, do not exceed the maximum of four ports for node-to-node communications.
- ▶ Node-to-node communications that use RDMA-capable Ethernet ports are not supported in a network configuration that contain more than two hops in the fabric of switches.

For more information, see [IBM Knowledge Center](#) and expand **Configuring** → **Configuration details** → **Configuration details for using RDMA-capable Ethernet ports for node-to-node communications**.

**Note:** Before you configure a system that uses RDMA-capable Ethernet ports for node-to-node communications in a standard, ESC, or HyperSwap topology system, contact your IBM representative.

### 2.7.4 iSCSI back-end storage attachment

A SAN Volume Controller supports the virtualization of external storage systems that are attached through iSCSI. Onboard 10 Gbps Ethernet ports or optional 25 Gbps Ethernet ports may be used. The 25 GbE network interface controllers (NICs) work in plain iSCSI mode without using any RDMA capabilities.

Consider the following items when planning for iSCSI virtualization:

- ▶ Direct attachment between the system and external storage systems is not supported, and requires Ethernet switches between the system and the external storage.
- ▶ To avoid a SPOF, a dual-switch configuration is recommended. For full redundancy, a minimum of two paths between each initiator node and target node must be configured with each path going through a separate switch.
- ▶ Extra paths can be configured to increase throughput if both initiator and target nodes support more ports.

All planning and implementation aspects of external storage virtualization with iSCSI are described in detail in *iSCSI Implementation and Best Practices on IBM Storwize Storage Systems*, SG24-8327.

## 2.7.5 IP network host attachment

You can attach the system to iSCSI or iSER hosts by using the Ethernet ports of the systems.

For each Ethernet port on a node, a maximum of one IPv4 address and one IPv6 address can be designated for iSCSI or iSER I/O. You can configure the internet Storage Name Service (iSNS) to facilitate a scalable configuration and management of iSCSI storage devices.

The same ports can be used for iSCSI and iSER host attachment concurrently, but a single host can establish either an iSCSI or iSER session, but not both.

iSCSI or iSER hosts connect to the system through the node-port IP addresses, which are assigned to the Ethernet ports of the node. If the node fails, the address becomes unavailable and the host loses communication with the system through that node. To allow hosts to maintain access to data, the node-port IP addresses for the failed node are transferred to the partner node in the I/O group. The partner node handles requests for both its own node-port IP addresses and also for node-port IP addresses on the failed node. This process is known as *node-port IP failover*. In addition to node-port IP addresses, the iSCSI name and iSCSI alias for the failed node are also transferred to the partner node. After the failed node recovers, the node-port IP address and the iSCSI name and alias are returned to the original node.

**Note:** The cluster name and node name form parts of the iSCSI name. Changing any of them might require reconfiguration of all iSCSI hosts that communicate with the system.

iSER supports only one-way authentication through the Challenge Handshake Authentication Protocol (CHAP). iSCSI supports two types of CHAP authentication: one-way authentication (iSCSI target authenticating iSCSI initiators) and two-way (mutual) authentication (iSCSI target authenticating iSCSI initiators, and vice versa).

For more information about iSCSI host attachment, see *iSCSI Implementation and Best Practices on IBM Storwize Storage Systems*, SG24-8327.

Make sure that iSCSI initiators, host iSER adapters, and Ethernet switches that are attached to the system are supported by using [SSIC](#).

## 2.7.6 Native IP replication

Two systems can be linked over native IP links that are connected directly or by Ethernet switches to perform RC functions. RC over native IP provides a less expensive alternative to using FC configurations.

IP replication is supported on both onboard 10G bps Ethernet ports and optional 25 Gbps Ethernet ports. However, when configured over 25 Gbps ports, it does not use RDMA capabilities, and it does not provide a performance improvement compared to 10 Gbps ports.

As a best practice, use a different port for iSCSI host I/O and IP partnership traffic. Also, use a different VLAN ID for iSCSI host I/O and IP partnership traffic.

Specific intersite link requirements must be met when you are planning to use IP partnership for RC. These requirements are described at [IBM Knowledge Center](#) by expanding **Configuring** → **Configuring IP partnerships** → **Intersite link planning**. Also, see Chapter 10, “Advanced Copy Services” on page 491.

## 2.7.7 Firewall planning

After you have your IP network planned, set up the appropriate firewall rules for each data flow.

For a list of mandatory and optional network flows that are required for operating, see [IBM Knowledge Center](#) and expand **Planning** → **Planning for hardware** → **Physical installation planning** → **IP address allocation and usage**.

## 2.8 Planning topology

SAN Volume Controller supports two dual-site topologies: ESC and HyperSwap. A list of key differences between them is shown in Table 2-2.

Table 2-2 Differences between HyperSwap and Stretched Cluster

Item	Stretched Cluster	HyperSwap
Minimum number of I/O groups that are required.	1	2
Independent copies of data that are maintained.	2	2 (Four if volume mirroring to two pools in each site is configured.)
Cache that is retained if only one site is online.	No	Yes
Stale consistent data is retained during resynchronization for DR.	No	Yes
Ability to use MM, GM, or GM together with an HA solution.	Yes	No
Maximum HA volume count.	5000	1250
Licensing.	Included in base product	Requires a Remote Mirroring license.

The topologies differ in how the nodes are distributed across the sites:

- ▶ For each I/O group in the system, the “stretched” topology has one node on one site, and one node on the other site. The topology works with any number of I/O groups of 1 - 4.
- ▶ The “HyperSwap” topology places both nodes of an I/O group at the same site. Therefore, to get a volume resiliently stored at both sites, at least two I/O groups are required.

The ESC topology uses fewer system resources, which enables more HA volumes to be configured. However, during a disaster that makes one site unavailable, the SAN Volume Controller system cache on the nodes of the surviving site are disabled. The HyperSwap topology uses extra system resources to support a full independent cache on each site, enabling full performance even if one site is lost.

For more information, see the following publications:

- ▶ *IBM Storwize V7000, Spectrum Virtualize, HyperSwap, and VMware Implementation*, SG24-8317
- ▶ *IBM SAN Volume Controller Stretched Cluster with PowerVM and PowerHA*, SG24-8142
- ▶ *IBM Spectrum Virtualize and SAN Volume Controller Enhanced Stretched Cluster with VMware*, SG24-8211

## 2.9 Back-end storage configuration

External back-end storage systems (also known as *controllers*) provide their logical volumes (LUs), which are detected by a system as managed disks (MDisks) and can be used in storage pools to provision their capacity to system's hosts.

The back-end storage subsystem configuration must be planned for all external storage systems that are attached. Apply the following general guidelines:

- ▶ Most of the supported FC-attached storage controllers must be connected through an FC SAN switch. However, a limited number of systems (including IBM FlashSystem 900 and a member of the Storwize family) can be direct-attached by using FC.
- ▶ For migration purposes, all back-end storage systems are supported by FC direct attachment.
- ▶ Connect all back-end storage ports to the SAN switch up to a maximum of 16 and zone them to all of the system to maximize bandwidth. The system is designed to handle many paths to the back-end storage.
- ▶ The cluster can be connected to a maximum of 1024 worldwide node names (WWNNs). The general practice is that:
  - EMC DMX/SYMM, all HDS, and SUN/HP HDS clones use one WWNN per port. Each port appears as a separate controller to the system.
  - IBM, EMC CLARiiON, and HP use one WWNN per subsystem. Each port appears as a part of a subsystem with multiple ports, with up to a maximum of 16 ports (WWPNs) per WWNN.

However, if you plan for a configuration that might be limited by the WWNN maximum, verify the WWNN versus WWPN policy with the back-end storage vendor.

- ▶ When defining a controller configuration, avoid hybrid configurations and automated tiering solutions. Create LUs for provisioning to the system from a homogeneous disk arrays or solid-state drive (SSD) arrays.
- ▶ Do not provision all available drives on the back-end storage capacity as a single LU. A best practice is to create one LU for eight HDDs or SSDs for the back-end system.
- ▶ If your back-end storage system is not supported by the round-robin path policy, ensure that the number of MDisks per storage pool is a multiple of the number of storage ports that are available. This approach ensures sufficient bandwidth for the storage controller, and an even balance across storage controller ports.
- ▶ The SAN Volume Controller must have exclusive access to every LU that is provisioned to it from a back-end controller. Any specific LU cannot be presented to more than one system. Presenting the same back-end LU to a system and a host is not allowed.
- ▶ Data reduction (compression and deduplication) on the back-end controller is supported only with a limited set of IBM Storage systems (for example, FlashSystem 900 with FlashCore Module (FCM) modules).

In general, configure back-end controllers as though they are used as stand-alone systems. However, there might be specific requirements or limitations as to the features that are usable in the specific back-end storage system. For more information about the requirements that are specific to your back-end controller, see [IBM Knowledge Center](#) and expand **Configuring** → **Configuring and servicing storage systems**.

The system's large cache and advanced cache management algorithms also allow it to improve the performance of many types of underlying disk technologies. Because hits to the cache can occur in the upper (the system itself) and the lower (back-end controller) level of the overall solution, the solution as a whole can use the larger amount of cache wherever it is. Therefore, the system's cache also provides more performance benefits for back-end storage systems with extensive cache banks.

However, the system cannot increase the throughput potential of the underlying disks in all cases. The performance benefits depend on the underlying back-end storage technology and the workload characteristics, including the degree to which the workload exhibits hotspots or sensitivity to cache size or cache algorithms.

## 2.10 Internal storage configuration

SAN Volume Controller may be equipped with an optional SAS-attached disk expansion enclosure. For general-purpose storage pools with various I/O applications, follow the storage configuration wizard recommendations in the GUI. For specific applications with known I/O patterns, use the CLI to create arrays that suit your needs.

An array-level recommendation for all types of internal storage is distributed redundant array of independent disks (DRAID) 6 (DRAID 6). DRAID 6 outperforms other available redundant array of independent disks (RAID) levels in most applications while providing fault tolerance and high rebuild speeds.

For more information about internal storage configuration, see *IBM System Storage SAN Volume Controller, IBM Storwize V7000, and IBM FlashSystem 7200 Best Practices and Performance Guidelines*, SG24-7521.

## 2.11 Storage pool configuration

The storage pool is at the center of the many-to-many relationship between the internal drive arrays or externally virtualized logical unit numbers (LUNs), which are represented as *MDisks*, and the volumes. It acts as a container of physical disk capacity from which chunks of MDisk space, which is known as *extents*, are allocated to form volumes that are presented to hosts.

The system supports two types of pools: *Data Reduction Pools (DRP)* and *standard pools*. The type is configured when a pool is created and it cannot be changed later. The type of the pool determines the set of features that is available on the system.

Features that can be implemented only with standard pools are:

- ▶ Child pools
- ▶ VMware vSphere integration with VMware vSphere Virtual Volumes (VVOLs)
- ▶ Multi-tenancy with ownership groups
- ▶ (SAN Volume Controller 2145-SV1 only) Real-time Compression (Random Access Compression Engine (RACE))

Features that can be implemented only with DRPs are:

- ▶ Automatic capacity reclamation with SCSI UNMAP (this feature returns capacity that is marked as no longer used by a host back to storage pool)
- ▶ DRP compression (in-flight data compression)
- ▶ DRP deduplication

In addition to providing data reduction options, DRP amplifies the I/O and CPU workload, which should be accounted for during performance sizing and planning.

Another base storage pool parameter is the extent size. There are two implications of a storage pool extent size:

- ▶ The maximum volume, MDisks, and managed storage capacity depend on the extent size. The bigger the extent that is defined for the specific pool, the larger is the maximum size of this pool, the maximum MDisk size in the pool, and the maximum size of a volume that is created in the pool.
- ▶ The volume sizes must be a multiple of the extent size of the pool in which the volume is defined. Therefore, the smaller the extent size, the better control that you have over the volume size.

The system supports extent sizes 16 - 8192 mebibytes (MiB). The extent size is a property of the storage pool and it is set when the storage pool is created.

**Note:** The base pool parameters, pool type, and extent size are set during pool creation and cannot be changed later. If you must change extent size or pool type, all volumes must be migrated from a storage pool and then the pool itself must be deleted and re-created.

For more information about the relationship between a system's maximum configuration and extent size, see to [Configuration Limits and Restrictions](#) and look for your platform and code version.

When planning pools, the encryption is defined on a pool level and the encryption setting cannot be changed after a pool is created. If you create an unencrypted pool, there is no way to encrypt it later. Your only option is to delete it and re-create it as encrypted.

When planning storage pool layout, consider the following aspects:

- ▶ Pool reliability, availability, and serviceability (RAS):
  - The storage pool is a failure domain. If one array or external MDisk is unavailable, the pool and all volumes in it go offline.
  - The number and size of storage pools affects system availability. Using a larger number of smaller pools reduces the failure domain if one of the pools goes offline. However, increasing the number of storage pools affects the storage use efficiency, and the number is subject to the configuration maximum limit.
  - You cannot migrate volumes between storage pools with different types or extent sizes. However, you can use volume mirroring to create copies between storage pools.
- ▶ Pool performance:
  - Do not mix same-tier arrays or MDisks with different performance characteristics in one pool. This technique is the only way to ensure consistent performance characteristics of volumes that are created from the pool.

Arrays with different tiers in one pool may be used because their performance differences become beneficial when you use the IBM Easy Tier function.



- Create multiple storage pools if you must isolate specific workloads to separate storage.
- Ensure that performance sizing was done for selected pool type and feature set.

### 2.11.1 The storage pool and cache relationship

The system uses cache partitioning to limit the potential negative effects that a poorly performing storage controller can have on the clustered system. The cache partition allocation size is based on the number of configured storage pools. This design protects against an individual overloaded back-end storage system from filling the system write cache and degrading the performance of the other storage pools.

Table 2-3 lists the limits of the write-cache data that can be used by a single storage pool.

Table 2-3 Limit of the cache data

Number of storage pools	Upper limit
1	100%
2	66%
3	40%
4	30%
5 or more	25%

No single partition can occupy more than its upper limit of write cache capacity. When the maximum cache size is allocated to the pool, the system starts to limit incoming write I/Os for volumes that are created from the storage pool. The host writes are limited to the destage rate on a one-out-one-in basis.

Only writes that target the affected storage pool are limited. The read I/O requests for the throttled pool continue to be serviced normally. However, because the system is offloading cache data at the maximum rate that the back-end storage can sustain, read response times are expected to be affected.

All I/O that is destined for other (non-throttled) storage pools continues as normal.

## 2.12 Volume configuration

When planning a volume, consider the required performance, availability, and capacity. Every volume is assigned to an I/O group that defines which pair of system nodes services I/O requests to the volume.

**Note:** No fixed relationship exists between I/O groups and storage pools.

When a host sends I/O to a volume, it can access the volume with either of the nodes in the I/O group but each volume has a *preferred node*. Many of the multipathing driver implementations that the system supports use this information to direct I/O to the preferred node. The other node in the I/O group is used only if the preferred node is not accessible.

During volume creation, the system selects the node in the I/O group that has the fewest volumes to be the preferred node. After the preferred node is chosen, it can be changed manually, if required.

Strive to distribute volumes evenly across available I/O groups and nodes within the system.

For more information about volume types, see Chapter 6, “Volumes” on page 255.

### 2.12.1 Planning for image mode volumes

Use image mode volumes to present to hosts data that is written to the back-end storage before it was virtualized. An image mode volume directly corresponds to the MDisk from which it is created.

Image mode volumes are a useful tool in storage migration and during system implementation to a working environment.

### 2.12.2 Planning for fully allocated volumes

A fully allocated volume presents to mapped hosts the same capacity that the volume uses in the storage pool. No data reduction is performed on a pool level. However, if a fully allocated volume is provisioned from a pool with data reducing storage, the data is still reduced on a drive level.

Fully allocated volumes provide the best performance because they do not cause I/O amplification, and they require less CPU time compared to other volume types.

### 2.12.3 Planning for thin-provisioned volumes

A thin-provisioned volume presents a different capacity to mapped hosts than the capacity that the volume uses in the storage pool. Space is not allocated on a thin-provisioned volume if an incoming host write operation contains all zeros.

Using the thin-provisioned volume feature that is called *zero detect*, you can reclaim unused allocated disk space (zeros) when you convert a fully allocated volume to a thin-provisioned volume by using volume mirroring.

DRPs enhance capacity efficiency for thin-provisioned volumes by monitoring the host's capacity usage. When the host indicates that the capacity is no longer needed, the capacity is released and can be reclaimed by the DRP to be redistributed automatically. Standard pools cannot reclaim capacity.

**Note:** Avoid using thin-provisioned volumes on a data reducing back end.

### 2.12.4 Planning for compressed volumes

With compressed volumes, data is compressed as it is written to disk, which saves more space. When data is read to hosts, the data is decompressed.

Compression is available through data reduction support as part of the system. If you want volumes to use compression as part of data reduction support, compressed volumes must belong to DRPs. Additionally, systems with SNA Volume Controller 2145-SV1 nodes support RACE compression for volumes on standard pools.

Before implementing compressed volumes, perform data analysis to know your average compression ratio and ensure that performance sizing was done for compression.

**Note:** If you use compressed volumes over FCM drives, the compression ratio on a drive level must be assumed to be 1:1 to avoid array overprovisioning and running out of space.

## 2.12.5 Planning for deduplicated volumes

Deduplication can be configured for volumes that use different capacity saving methods, such as thin provisioning. Deduplicated volumes must be created in DRPs for added capacity savings. Deduplication is a type of data reduction that eliminates duplicate copies of data. Deduplication of user data occurs within a DRP and only between volumes or volume copies that are marked as deduplicated.

With deduplication, the system identifies unique chunks of data that is called *signatures* to determine whether new data is written to the storage. Deduplication is a hash-based solution, which means chunks of data are compared to their signatures rather than to the data itself. If the signature of the new data matches an existing signature that is stored on the system, then the new data is replaced with a reference. The reference points to the stored data instead of writing the data to storage. This process saves the capacity of the back-end storage by not writing new data to storage, and it might improve the performance of read operations to data that has an existing signature.

The same data pattern can occur many times, and deduplication decreases the amount of data that must be stored on the system. A part of every hash-based deduplication solution is a repository that supports looking up matches for incoming data. The system contains a database that maps the signature of the data to the volume and its virtual address. If an incoming write operation does not have a signature that is stored in the database, then a duplicate is not detected and the incoming data is stored on back-end storage.

To maximize the space that is available for the database, the system distributes this repository between all nodes in the I/O groups that contain deduplicated volumes. Each node carries a distinct portion of the records that are stored in the database. If nodes are removed or added to the system, the database is redistributed between the nodes to ensure full use of the available memory.

Before implementing deduplication, perform data analysis to estimate deduplication savings and make sure that system performance sizing was done for deduplication.

## 2.13 Host attachment planning

The system supports the attachment of a various host hardware types running different OSs with FC SAN or IP SAN. For a list of instructions that is specific to your host setup, see [IBM Knowledge Center](#) and expand **Configuring** → **Host attachment**.

### 2.13.1 Queue depth

Typically, hosts issue subsequent I/O requests to storage systems without waiting for the completion of previous ones. The number of outstanding requests is called *queue depth*. Sending multiple I/O requests in parallel (asynchronous I/O) provides significant performance benefits compared to sending them one-by-one (synchronous I/O). However, if the number of queued requests exceeds the maximum that is supported by the storage controller, you experience performance degradation.

For more information about how to calculate correct host queue depth for your environment, see [IBM Knowledge Center](#) and expand **Configuring** → **Host attachment**.

### 2.13.2 Microsoft Offloaded Data Transfer

If your Windows hosts are configured to use Microsoft Offloaded Data Transfer (ODX) to offload the copy workload to the storage controller, consider the benefits of this technology against the extra load on the storage controllers. The benefits and effects of enabling ODX are especially prominent in Microsoft Hyper-V environments with ODX enabled.

### 2.13.3 SAN boot support

The system supports SAN boot or startup for selected configurations of hosts running AIX, Microsoft Windows, and other OSs. To check whether your configuration is supported for SAN boot, see the [SSIC](#).

### 2.13.4 Planning for large deployments

Each I/O group can have up to 512 host objects defined. This limit is the same whether hosts are attached by using FC, iSCSI, or a combination of both. To allow more than 512 hosts to access the storage, you must divide them into groups of 512 hosts or less and map each group to a single I/O group only. With this approach, you can configure up to 2048 host objects on a system with four I/O groups (eight nodes).

For best performance, split each host group into two sets. For each set, configure the preferred access node for volumes that are presented to the host set to one of the I/O group nodes. This approach helps to evenly distribute load between the I/O group nodes.

**Note:** A volume can be mapped only to a host that is associated with the I/O group to which the volume belongs.

### 2.13.5 Planning for SCSI Unmap

UNMAP is a set of SCSI primitives that hosts use to indicate to a SCSI target that space that is allocated to a range of blocks on a target storage volume is no longer required. With this command, the storage controller takes measures and optimizes the system so that the space can be reused for other purposes.

The system supports end-to-end UNMAP compatibility, which means that a command that is issued by a host is processed and sent to the back-end storage device.

UNMAP processing can be controlled with two separate settings:

- ▶ One setting enables hosts to send **UNMAP** commands,
- ▶ The other setting controls whether the system sends **UNMAP** commands to back-end storage (drives and external controllers).

By default, back-end UNMAP is on, and it is a best practice to keep it turned on for most use cases. Host UNMAP is off by default.

Turn on host UNMAP if you have data reducing storage that is virtualized by the SAN Volume Controller, for example, a FlashSystem 9100 system with self-compressing drives, to benefit from end-to-end UNMAP support.

Consider turning off host UNMAP if you do not use DRPs and use HDDs as a back end.

However, if the system has a virtualized slow back-end controller (for example, back-end storage running nearline (NL) SAS drives for the lowest tier), UNMAP requests that are sent from the host as large I/O chunks might overload the back end. Thorough planning is required if you plan to virtualize a mixture of slow back-end controllers and data reducing (compressing) back-end controllers.

## 2.14 Planning copy services

SAN Volume Controller offers a set of copy services, such as IBM FlashCopy (snapshots) and RC, in synchronous and asynchronous modes. For more information about copy services, see Chapter 10, “Advanced Copy Services” on page 491.

### 2.14.1 FlashCopy guidelines

With the IBM FlashCopy function of IBM Spectrum Virtualize, you can perform a point-in-time (PiT) copy of one or more volumes. The FlashCopy function creates a PiT or time-zero (T0) copy of data that is stored on a source volume to a target volume by using a Copy on Write (CoW) and Copy on Demand mechanism.

While the FlashCopy operation is performed, the source volume is stopped briefly to initialize the FlashCopy bitmap, and then I/O can resume. Although several FlashCopy options require the data to be copied from the source to the target in the background, which can take time to complete, the resulting data on the target volume is presented so that the copy appears to complete immediately.

The FlashCopy function operates at the block level below the host OS and cache, so those levels must be flushed by the OS for a FlashCopy copy to be consistent.

When you use the FlashCopy function, observe the following guidelines:

- ▶ Both the FlashCopy source and target volumes should use the same preferred node.
- ▶ If possible, keep the FlashCopy source and target volumes on separate storage pools.

For more information about planning for the FlashCopy function, see *IBM System Storage SAN Volume Controller, IBM Storwize V7000, and IBM FlashSystem 7200 Best Practices and Performance Guidelines*, SG24-7521.

## 2.14.2 Planning for Metro Mirror and Global Mirror

MM is a copy service that provides a continuous, synchronous mirror of one volume to a second volume. The systems can be up to 300 kilometers (186.4 miles) apart. Because the mirror is updated synchronously, no data is lost if the primary system becomes unavailable. MM is typically used for DR purposes, where it is important to avoid any data loss.

GM is a copy service that is similar to MM, but copies data asynchronously. You do not have to wait for the write to the secondary system to complete. For long distances, performance is improved compared to MM. However, if a failure occurs, you might lose data.

GM uses one of two methods to replicate data. Multicycling GM is designed to replicate data while adjusting for bandwidth constraints. It is appropriate for environments where it is acceptable to lose a few minutes of data if a failure occurs.

For environments with higher bandwidth, non-cycling GM can be used so that less than a second of data is lost if a failure occurs. GM also works well when sites are more than 300 kilometers (186.4 miles) apart.

When copy services are used, all components in the SAN must sustain the workload that is generated by application hosts and the data replication workload. Otherwise, the system can automatically stop copy services relationships to protect your application hosts from increased response times.

While planning RC services, consider the following aspects:

- ▶ Copy services topology

One or more clusters can participate in a copy services relationship. One typical and simple use case is DR, where one site is active and another performs only a DR function. In such a case, the solution topology is simple, with one cluster per site and uniform replication direction for all volumes. However, multiple other topologies are possible that you can use to design a solution that optimally fits your set of requirements.

- ▶ GM versus MM

Decide which type of copy services you are going use. This decision should be requirement-driven. With MM, you prevent any data loss during a system failure, but it has more stringent requirements, especially regarding intercluster link bandwidth and latency, and remote site storage performance. Also, MM incurs a performance penalty because writes are not confirmed to host until a data reception confirmation is received from the remote site.

With GM, you can relax constraints on the system requirements at the cost of using asynchronous replication, which enables the remote site to lag behind the local site. The choice of the replication type has major effects on all other aspects of the copy services planning.

Using GM and MM between the same two clustered systems is supported. Also, the RC type may be changed from one to another one.

For native IP replication, use the RC mode of Multicycling GM (or Global Mirror with Change Volumes (GMCV).

- ▶ Intercluster link

The local and remote clusters can be connected by an FC or IP network. Each of the technologies has its own requirements concerning supported distance, link speeds, bandwidth, and vulnerability to frame or packet loss.

When planning the intercluster link, consider the peak performance that is required. This consideration is especially important for MM configurations.

The bandwidth between sites must be sized to meet the peak workload requirements. When planning the inter-site link, consider the initial sync and any future resync workloads. It might be worthwhile to secure more link bandwidth for the initial data synchronization.

If the link between the sites is configured with redundancy so that they can tolerate single failures, you must size the link so that the bandwidth and latency requirements are met even during single failure conditions.

When planning the inter-site link, note whether it is dedicated to the inter-cluster traffic or is going to be used to carry any other data. Sharing the link with other traffic might affect the link's ability to provide the required bandwidth for data replication.

- ▶ Volumes and consistency groups

Determine whether volumes can be replicated independently. Some applications use multiple volumes and require that the order of writes to these volumes is preserved in the remote site. Notable examples of such applications are databases.

If an application requires that the write order is preserved for the set of volumes that it uses, create a consistency group for these volumes.

## 2.15 Data migration

Data migration is an important part of an implementation, so you must prepare a detailed data migration plan. You might need to migrate your data for one of the following reasons:

- ▶ Redistribute a workload within a clustered system across back-end storage subsystems.
- ▶ Move a workload on to newly installed storage.
- ▶ Move a workload off old or failing storage ahead of decommissioning it.
- ▶ Move a workload to rebalance a changed load pattern.
- ▶ Migrate data from an older disk subsystem.
- ▶ Migrate data from one disk subsystem to another one.

Because multiple data migration methods are available, choose the method that best fits your environment, OS platform, type of data, and the application's service-level agreement (SLA).

Data migration methods can be divided into three classes:

- ▶ Based on the host OS, for example, by using the system's logical volume manager (LVM)
- ▶ Based on specialized data migration software
- ▶ Based on the system data migration features

For more information about system data migration tools, see Chapter 8, "Storage migration" on page 429 and 10.5, "Volume mirroring and migration options" on page 576.

With data migration, apply the following guidelines:

- ▶ Choose the data migration method that best fits your OS platform, type of data, and SLA.
- ▶ Choose where you want to place your data after migration in terms of the storage tier, pools, and back-end storage.
- ▶ Check whether enough free space is available in the target storage pool.
- ▶ To minimize downtime during the migration, plan ahead of time all of the required changes, including zoning, host definition, and volume mappings.
- ▶ Prepare a detailed operation plan so that you do not overlook anything at data migration time. Especially for a large or critical data migration, have the plan peer-reviewed and formally accepted by an appropriate technical design authority within your organization.

- ▶ Perform and verify a backup before you start any data migration.
- ▶ You might want to use the system as a data mover to migrate data from a non-virtualized storage subsystem to another non-virtualized storage subsystem. In this case, you might have to add checks that relate to the specific storage subsystem that you want to migrate.

Be careful when you are using slower disk subsystems for the secondary volumes for high-performance primary volumes because the system's cache might not be able to buffer all the writes. Flushing cache writes to slower back-end storage might impact performance of your hosts.

- ▶ Consider storage performance. The migration workload might be much higher than expected during normal operations of the system. If there is already application data on the system to which you are migrating, the application performance might suffer if the system is overloaded. Consider using host or volume level throttles when performing migration on a production environment.

## 2.16 Performance monitoring with IBM Storage Insights

IBM Storage Insights is integral to monitoring and ensuring the continued availability of the system.

Available at no additional charge, the cloud-based IBM Storage Insights product provides a single dashboard that provides a clear view of all your IBM block storage. You can make better decisions by seeing trends in performance and capacity.

With storage health information, you can focus on areas needing attention and when IBM support is needed, IBM Storage Insights simplifies uploading logs, speeds resolution with online configuration data, and provides an overview of open tickets all in one place.

IBM Storage Insights provides a unified view of IBM systems. By using it, you can see all of your IBM storage inventory as a live event feed so that you know what is going on with your storage.

IBM Storage Insights provides advanced customer service with an event filter that provides the following functions:

- ▶ The ability for you and support to view support tickets and open and close them, and to track trends.
- ▶ With the auto log collection capability, you can collect the logs and send them to IBM before IBM Support starts looking into the problem. This feature can reduce the time to solve the case by as much as 50%.

Figure 2-3 on page 79 shows the architecture of the IBM Storage Insights application, the supported products, and the three main teams who can benefit from the use of the tool.



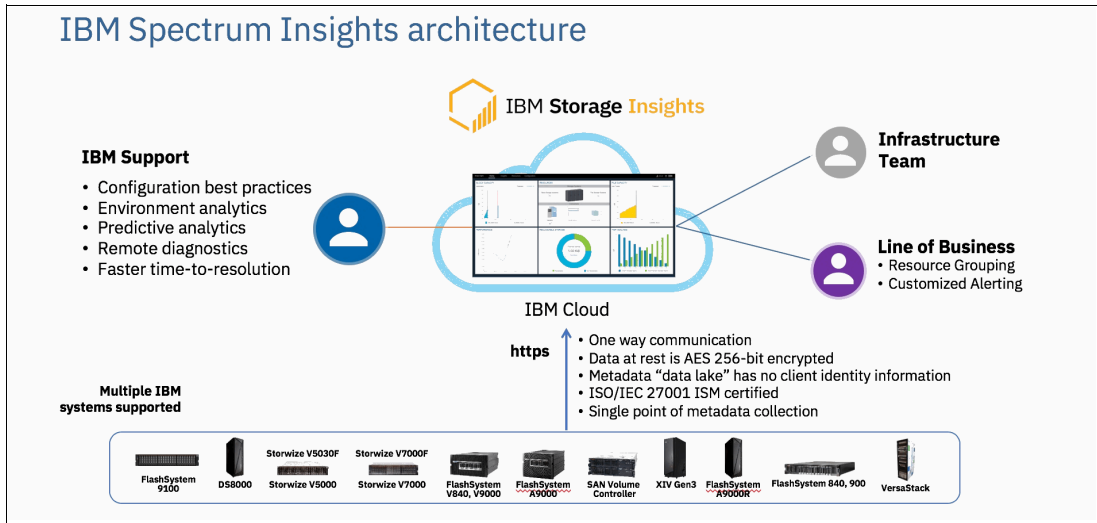


Figure 2-3 IBM Storage Insights Architecture

IBM Storage Insights provides a lightweight data collector that is deployed on a Linux, Windows, or AIX server or a guest in a virtual machine (VM) (for example, a VMware guest).

The data collector streams performance, capacity, asset, and configuration metadata to your IBM Cloud instance.

The metadata flows in one direction, that is, from your data center to IBM Cloud over HTTPS. In the IBM Cloud, your metadata is protected by physical, organizational, access, and security controls. IBM Storage Insights is ISO/IEC 27001 Information Security Management certified.

Figure 2-4 shows the data flow from systems to the IBM Storage Insights cloud.

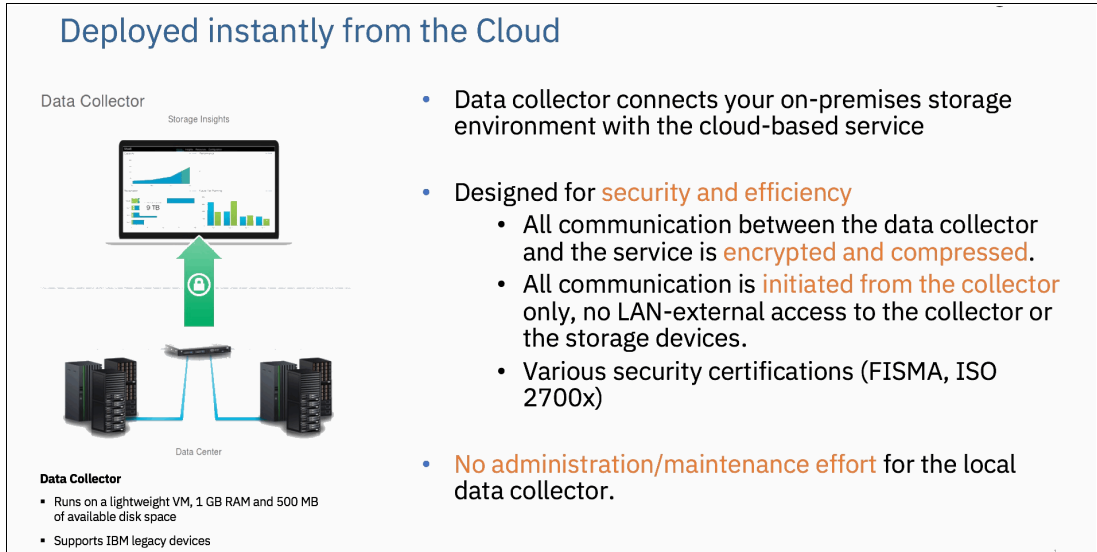


Figure 2-4 Data flow from the storage systems to the IBM Storage Insights cloud

Metadata about the configuration and operations of storage resources is collected, such as:

- ▶ Name, model, firmware, and type of storage system
- ▶ Inventory and configuration metadata for the storage system's resources, such as volumes, pools, disks, and ports

- ▶ Capacity values, such as capacity, unassigned space, used space, and the compression ratio.
- ▶ Performance metrics, such as read and write data rates, I/O rates, and response times.

The application data that is stored on the storage systems cannot be accessed by the data collector.

Access to the metadata that is collected is restricted to the following users:

- ▶ The customer who owns the dashboard.
- ▶ The administrators who are authorized to access the dashboard, such as the customer's operations team.
- ▶ The IBM Cloud team that is responsible for the day-to-day operation and maintenance of IBM Cloud instances.
- ▶ IBM Support for investigating and closing service tickets.

For more information about IBM Storage Insights and to sign up and register for the free service, see the following resources:

- ▶ [Fact Sheet](#)
- ▶ [Demonstration](#)
- ▶ [Security Guide](#)
- ▶ [Registration](#)

For more information, see 13.12, "IBM Storage Insights Monitoring" on page 819.

## 2.17 Configuration backup procedure

Save the configuration before and after any major configuration changes on the system. Saving the configuration is a crucial part of management, and various methods can be applied to back up your system configuration. A best practice is to implement an automatic configuration backup by using the configuration backup command. Make sure that you save the configuration to a host system that does not depend on the storage that is provisioned from a system whose configuration is backed up.

For more information, see 13.4, "Configuration backup" on page 765.



# Initial configuration

This chapter describes the initial configuration of the IBM SAN Volume Controller systems, provides step-by-step instructions about how to do the initial setup, and defines the base settings of the system, which are done during the implementation phase before volumes are created and provisioned.

This chapter includes the following topics:

- ▶ 3.1, “Prerequisites” on page 82
- ▶ 3.2, “System initialization” on page 83
- ▶ 3.3, “System setup” on page 87
- ▶ 3.4, “Base configuration” on page 96
- ▶ 3.5, “Configuring management access” on page 108

## 3.1 Prerequisites

**Note:** SAN Volume Controller is installed by an IBM System Services Representative (IBM SSR). You must provide all the necessary information to the IBM SSR by filling out the planning worksheets, which can be found at [IBM Knowledge Center](#) and expanding **Planning** → **Planning for hardware** → **Physical configuration planning of a system** → **Requirements and guidelines for completing the hardware location chart**.

After the IBM SSR completes their task, continue the setup by following the instructions in 3.3, “System setup” on page 87.

Before initializing and setting up the SVC, ensure that the following prerequisites are met:

- ▶ The physical components fulfill all the requirements and are correctly installed, including:
  - The control enclosures are physically installed in the racks.
  - The Ethernet and Fibre Channel (FC) cables are connected.
  - The expansion enclosures, if available, are physically installed and attached to the control enclosures that will use them.
  - The SAN Volume Controller nodes and optional expansion enclosures are powered on.
- ▶ The web browser that is used for managing the system is supported by the management GUI. For the list of supported browsers, see [IBM Knowledge Center](#).
- ▶ You have the required information, including:
  - The IPv4 (or IPv6) addresses that are assigned for the system’s management interfaces:
    - The unique cluster IP address, which is the address that is used for the management of the system.
    - Unique service IP addresses, which are used to access node service interfaces. You need one address for each node (two per control enclosure).
    - The IP subnet mask for each subnet that is used.
    - The IP gateway for each subnet that is used.
  - The licenses that might be required to use particular functions (depending on the system type):
    - Remote Copy (RC)
    - External Virtualization
    - IBM FlashCopy
    - Compression
    - Encryption
  - Information that is used by a system when performing Call Home functions, such as:
    - The company name and system installation address.
    - The name, email address, and phone number of the storage administrator whom IBM can contact if necessary.
  - (optional) The Network Time Protocol (NTP) server IP address.

- (optional) The Simple Mail Transfer Protocol (SMTP) server IP address, which is necessary only if you want to enable Call Home or want to be notified about system events through email.
- (optional) The IP addresses for Remote Support Proxy Servers, which are required only if you want to use them with the Remote Support Assistance feature.

## 3.2 System initialization

This section provides step-by-step instructions about how to create the SVC cluster.

To start the initialization procedure, connect a desktop PC or a notebook to the technician port. The *technician port* is a dedicated 1 Gb Ethernet port at the rear of each of the nodes in the control enclosure. It can be used only to initialize or service the system. It cannot be connected to an Ethernet switch because it supports only a direct connection, and it remains disconnected after the initial setup is done.

The location of a technician port on the SAN Volume Controller 2145-SV1 node is shown in Figure 3-1.



Figure 3-1 Location of the technician port on the SAN Volume Controller 2145-SV1 node

The location of a technician port on the SAN Volume Controller 2145-SV2 and 2145-SA2 nodes is shown on Figure 3-2.

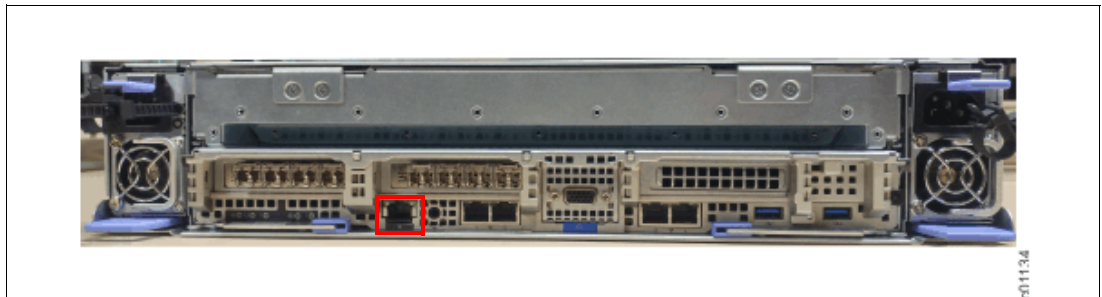


Figure 3-2 Location of the technician port on the SAN Volume Controller 2145-SV2 and 2145-SA2 nodes

The technician port runs an IPv4 DHCP server and it can assign an address to any device that is connected to this port. Ensure that your PC or notebook Ethernet adapter is configured to use a DHCP client if you want the IP to be assigned automatically. If you prefer not to use DHCP, you can set a static IP on the Ethernet port from the 192.168.0.0/24 subnet, for example, 192.168.0.2 with the netmask 255.255.255.0.

The default IP address of a technician port on a node canister is 192.168.0.1. Do not use this IP address for your PC or notebook.

**Note:** Ensure that the technician port is not connected to the organization's network. No Ethernet switches or hubs are supported on this port.

### 3.2.1 System initialization process

Before the SAN Volume Controller is initialized, each node of a new system remains in the *candidate* state and cannot process I/O. During initialization, a *cluster* is created, which at that moment consists only of one node. All other nodes except the first one must not be initialized, and they are added to the cluster by using a cluster management interface (GUI or CLI) after first one is set up.

You must specify IPv4 or an IPv6 system management addresses, which are assigned to Ethernet port 1 on each node and used to access the management GUI and CLI. After the system is initialized, you can specify other IP addresses.

**Note:** Do not perform the system initialization procedure on more than one node in a system. During initialization, candidate nodes are added automatically. You may also use the management GUI or CLI to add nodes to the system.

To do the initialization of a new system, complete the following steps:

1. Connect your PC or notebook to a technician port of any canister of the control enclosure. Ensure that you obtained a valid IPv4 address with DHCP.
2. Open a supported web browser and go to `http://install`. The browser is automatically redirected to the System Initialization wizard. You can also use the IP address `http://192.168.0.1` if you are not automatically redirected.

**Note:** During the system initialization, you are prompted to accept untrusted certificates because the system certificates are self-signed. If you are directly connected to the service interface, there is no doubt about the identity of the certificate issuer, so you can safely accept the certificates.

If the system is not in a state that allows initialization, you are redirected to the Service Assistant interface. Use the displayed error codes to troubleshoot the problem.

3. Welcome dialog box opens, as shown in Figure 3-3 on page 85. Click **Next** to start the procedure.

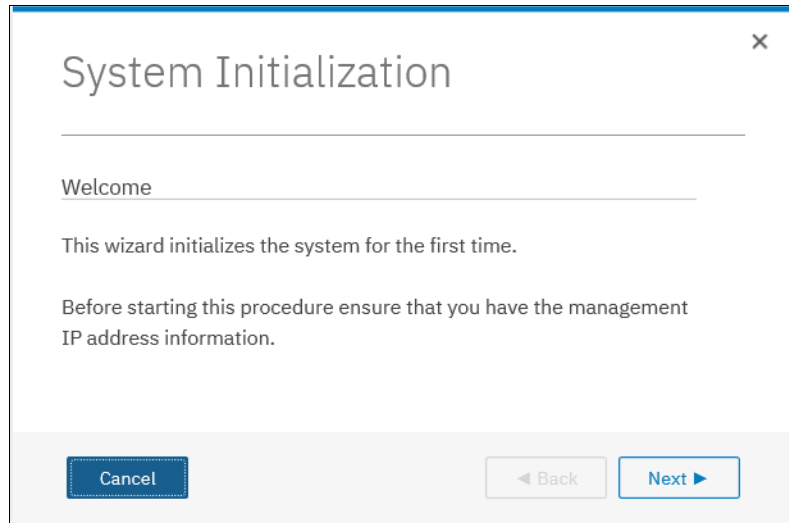


Figure 3-3 System Initialization: Welcome dialog box

4. A window opens in which two options are presented, as shown in Figure 3-4. Select the first option and click **Next**.

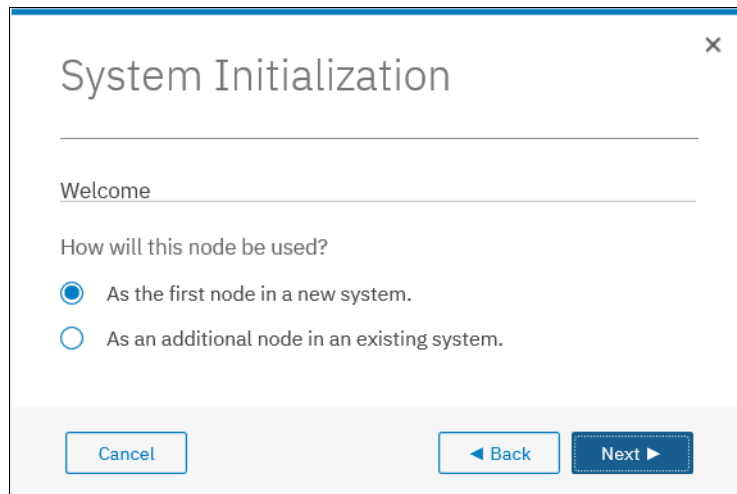
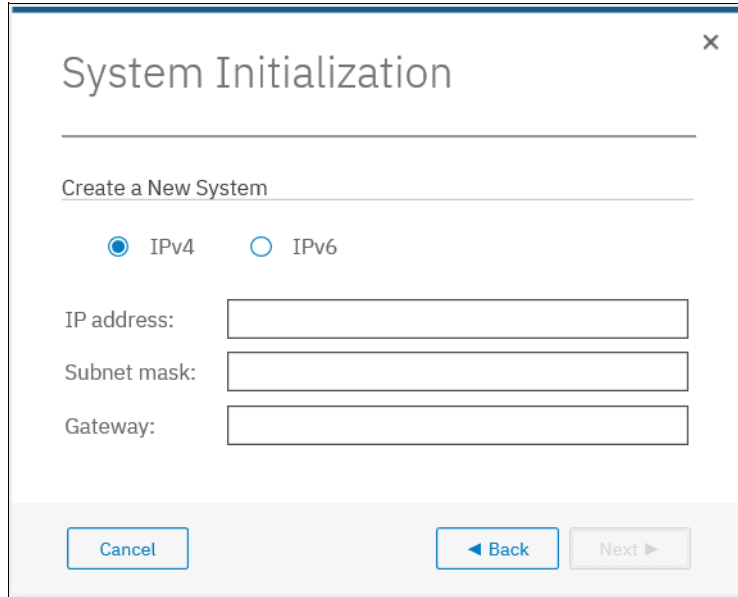


Figure 3-4 System initialization: Create a system or expand the existing one

If you select **As an additional node in an existing system**, you are prompted to disconnect from the technician port and use the GUI of an existing system to add new nodes.

5. Enter the management IP address information for the new system, as shown in Figure 3-5. Set the IP address, network mask, and gateway.

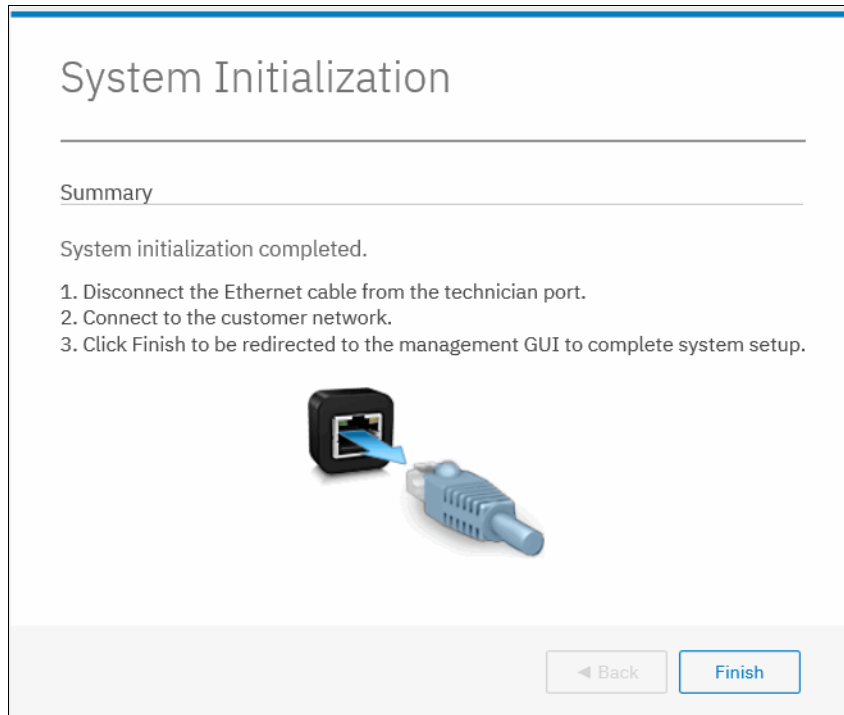


The screenshot shows a window titled "System Initialization" with a close button (X) in the top right corner. Below the title is a horizontal line, followed by the text "Create a New System". There are two radio buttons: "IPv4" (selected) and "IPv6". Below these are three input fields labeled "IP address:", "Subnet mask:", and "Gateway:". At the bottom of the window are three buttons: "Cancel", "Back" (with a left arrow), and "Next" (with a right arrow).

Figure 3-5 System Initialization: Management IP

As you click **Next**, a command to create a cluster runs.

6. A window with restart timer opens. When the timeout is reached, you can click **Next** to see the final initialization window, as shown in Figure 3-6. Follow the instructions, and browser is redirected to the management IP address to access the system GUI after you click **Finish**.



The screenshot shows a window titled "System Initialization" with a close button (X) in the top right corner. Below the title is a horizontal line, followed by the text "Summary". Below this is the text "System initialization completed." followed by a numbered list:  
1. Disconnect the Ethernet cable from the technician port.  
2. Connect to the customer network.  
3. Click Finish to be redirected to the management GUI to complete system setup.  
Below the list is an illustration of a black Ethernet port on a device with a blue Ethernet cable plugged into it. At the bottom of the window are two buttons: "Back" (with a left arrow) and "Finish".

Figure 3-6 System initialization: Complete



If you cannot connect to a network that has access to the management IP, you can continue the system setup from any other workstation that can reach it.

## 3.3 System setup

This section provides instructions about how to define the basic settings of the system by using the system setup wizard.

### 3.3.1 System setup wizard

After the initialization is complete and you are redirected to a management GUI from your PC or notebook, or you browse to the management IP address of a freshly initialized system from another workstation, you must complete the system setup wizard to define the basic settings of the system.

**Note:** Experienced users can disable the system setup wizard and complete the configuration manually. However, this method is *not recommended* for most use cases.

To disable the system setup wizard on a new system, run the following command:

```
chsystem -easysetup no
```

**Note:** During the setup wizard, you are prompted to change the default superuser password. If the wizard is bypassed, the system blocks the configuration functions until it is changed. All attempts at configuration return the following error:

```
CMMVC9473E The command failed because the superuser password must be changed before the system can be configured
```

**Note:** All configuration settings that are done by using the system setup wizard can be changed later by using the system GUI or CLI.

The first time that you connect to the management GUI, you are prompted to accept untrusted certificates because the system certificates are self-signed.

If your company policy requests certificates that are signed by a trusted certificate authority (CA), you can install them after you complete the system setup. For more information about how to perform this task, see 3.5.1, “Configuring secure communications” on page 108.

To complete the system setup wizard, complete the following steps:

1. Log in to system GUI. Until the wizard is complete, you may use only *superuser* account, as shown in Figure 3-7. Click **Sign in**.

**Note:** The default password for the superuser account is `passwd0rd` (with the number zero and not the capital letter O). The default password must be changed by using the system setup wizard or after the first CLI login. The new password cannot be set to the default one.

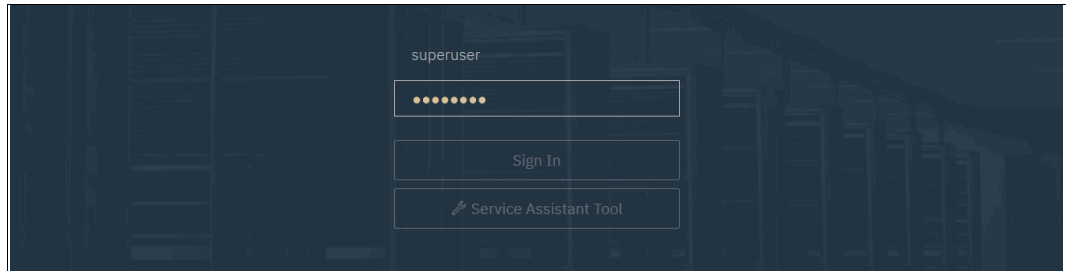


Figure 3-7 System setup: Logging in for the first time.

2. The welcome window opens, as shown in Figure 3-8. Verify the prerequisites and click **Next**.

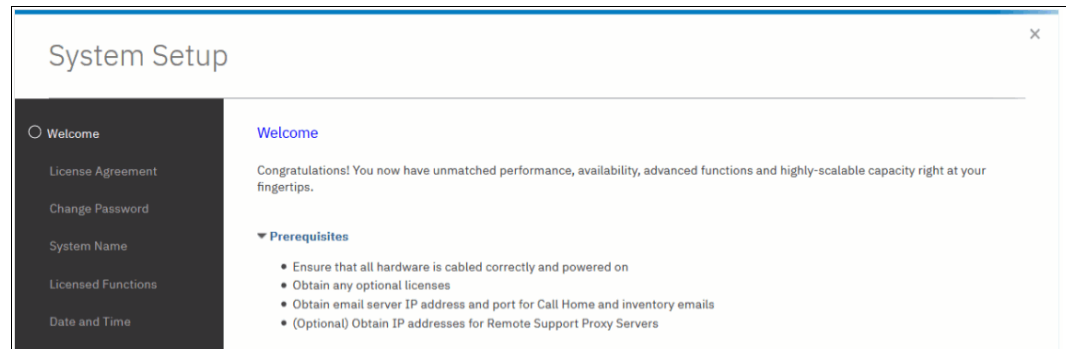


Figure 3-8 System setup: Welcome

3. The system automatically adds nodes that are connected and available in the *Candidate* status to the cluster. This task might take few minutes.
4. Carefully read the license agreement, select **I agree with the terms in the license agreement** if you want to continue the setup, as shown in Figure 3-9 on page 89, and click **Next**.

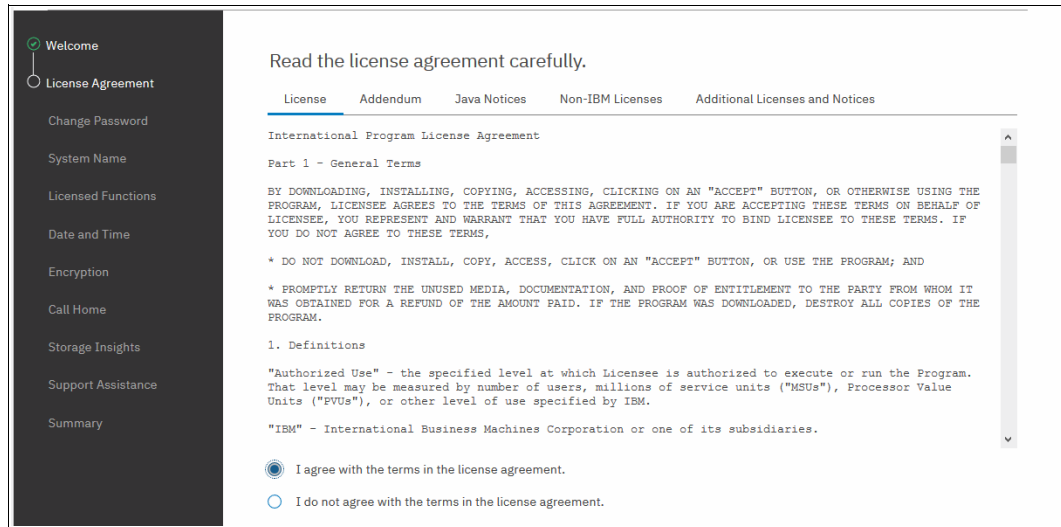


Figure 3-9 System setup: License agreement

5. Enter a new password for the superuser, as shown in Figure 3-10. A valid password is 6 - 64 characters long and it cannot begin or end with a space. Also, the password cannot be set to match the default password. Click **Apply** and then **Next**.

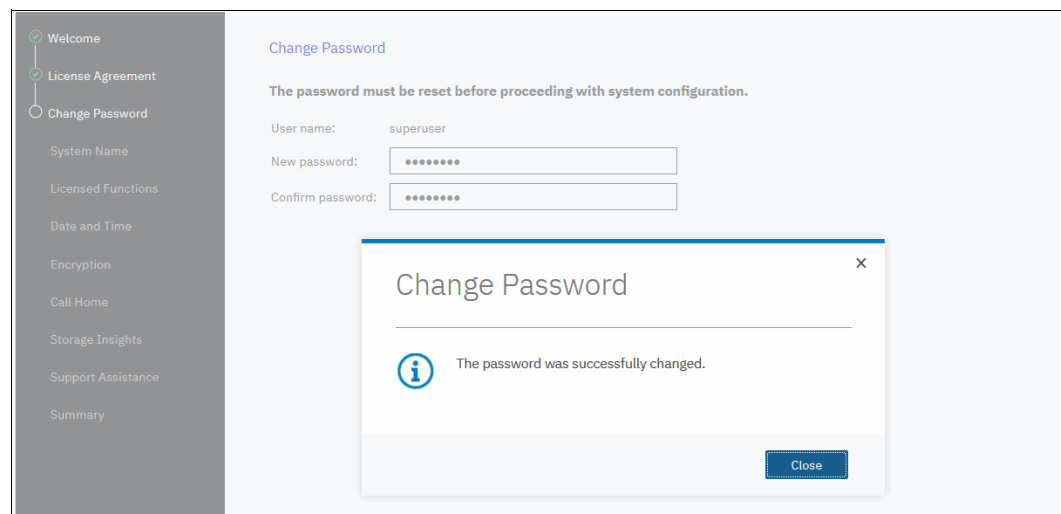


Figure 3-10 System setup: Changing the password for the superuser

**Note:** All configuration changes that are done with the system setup wizard are applied immediately, including the password change.

- Enter the name that you want to give the new system, as shown on Figure 3-11. Click **Apply** and then **Next**.

Avoid using an underscore ( \_ ) in a system name. While permitted here, it is not allowed in domain name server (DNS) short names and fully qualified domain names (FQDNs), so such naming might cause confusion and access issues. The following characters can be used: A - Z, a - z, 0 - 9, and - (hyphen).

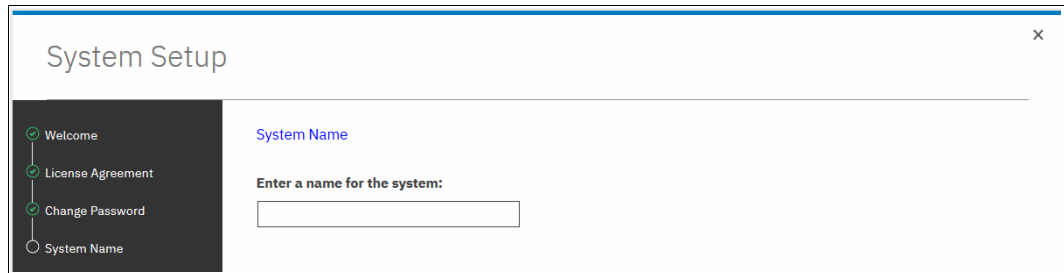


Figure 3-11 System setup: Setting the system name

- Enter the number of licensed SCUs or licensed capacity for each function, as shown in Figure 3-12.

SAN Volume Controller uses differential and capacity-based licensing. For External Virtualization and compression, differential licensing offers different pricing rates for different types of storage, and it is based on the number of Storage Capacity Units (SCUs) that are purchased. For other licensed functions, the system supports capacity-based licensing.

Make sure that the numbers you enter here match the numbers in your license authorization papers.

When done, click **Apply** and then **Next**.

**Note:** Encryption uses a key-based licensing scheme, and it is activated later in the wizard.

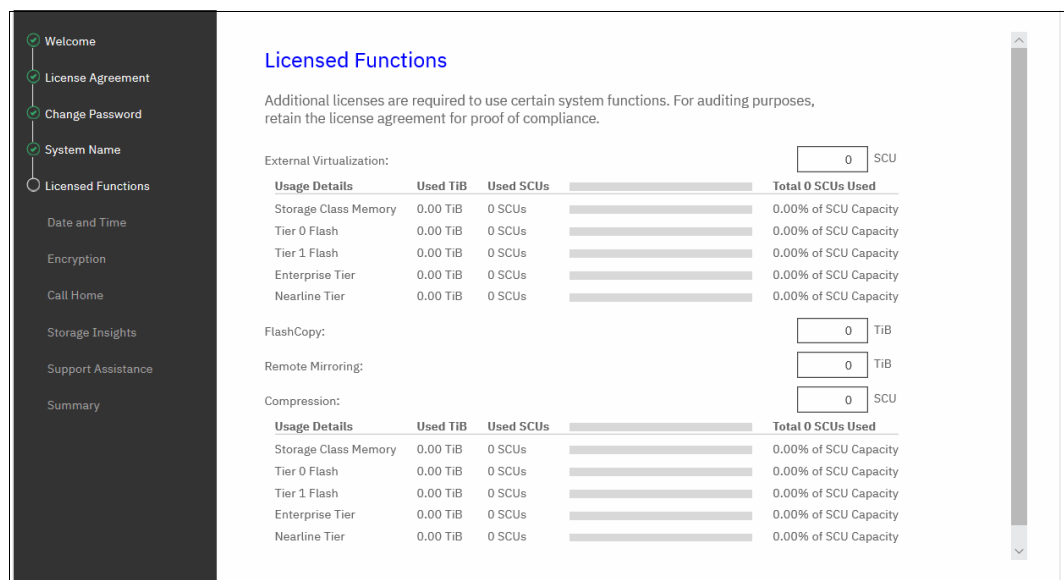


Figure 3-12 System setup: Setting the system licenses

- Enter the date and time settings. In the example that is shown in Figure 3-13, the date and time are set by using an NTP server. Generally, use an NTP server so that all of your storage area network (SAN) and storage devices have a common time stamp. This practice facilitates troubleshooting and prevents time stamp-related errors if you use a key server as an encryption key provider.

If you choose to manually enter these settings, you are prompted to input the date, time, and time zone, or you can take those settings from your web browser. You cannot use a 24-hour clock system here, but you can switch to it later by using the system GUI.

When the data is set, click **Apply** and then **Next**.

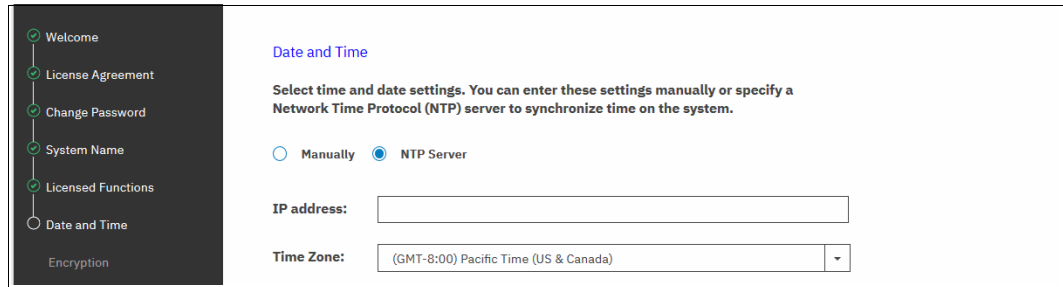


Figure 3-13 System setup: Setting the date and time

- Select whether the encryption feature was purchased for this system, as shown in Figure 3-14.

If encryption is not planned at this moment, select **No** and click **Next**. You can enable this feature later, as described in Chapter 12, “Encryption” on page 685.

If you purchased the encryption feature, you are prompted to activate your license manually or automatically. The encryption license is key-based and required for each control enclosure.

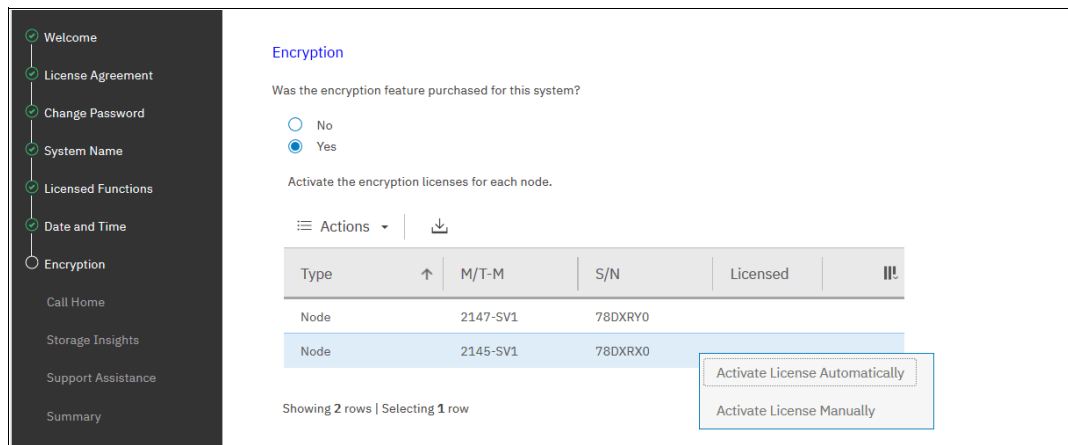


Figure 3-14 System setup: Encryption activation

You can use automatic activation if the PC or notebook that you use to connect to the GUI and run the system setup wizard has internet access. If no internet connection is available, use manual activation and follow the instructions. For more information, see Chapter 12, “Encryption” on page 685.

After the encryption license is activated, you see a green check mark next to the node serial number. After all the nodes show that the encryption is licensed, click **Next**.

10. Set up the Call Home functions, as shown in Figure 3-15. With Call Home enabled, IBM automatically opens problem reports and contacts you to verify whether replacement parts are required.

**Note:** It is a best practice to configure Call Home and keep it enabled if your system is under warranty or if you have a hardware maintenance agreement.

An IBM SSR configures Call Home when performing cluster initialization. You need to only check whether all the entered data is correct.

The system supports two methods of sending Call Home notifications to IBM:

- Cloud Call Home
- Call Home with email notifications.

Cloud Call Home is the default and preferred option for a system to report event notifications to IBM Support. With this method, the system uses RESTful APIs to connect to an IBM centralized file repository that contains troubleshooting information that is gathered from customers. This method requires no extra configuration.

The system may also be configured to use email notifications for this purpose. If this method is selected, you are prompted to enter the SMTP server IP address.

If both methods are enabled, cloud Call Home is used, and the email notifications method is kept as a backup.

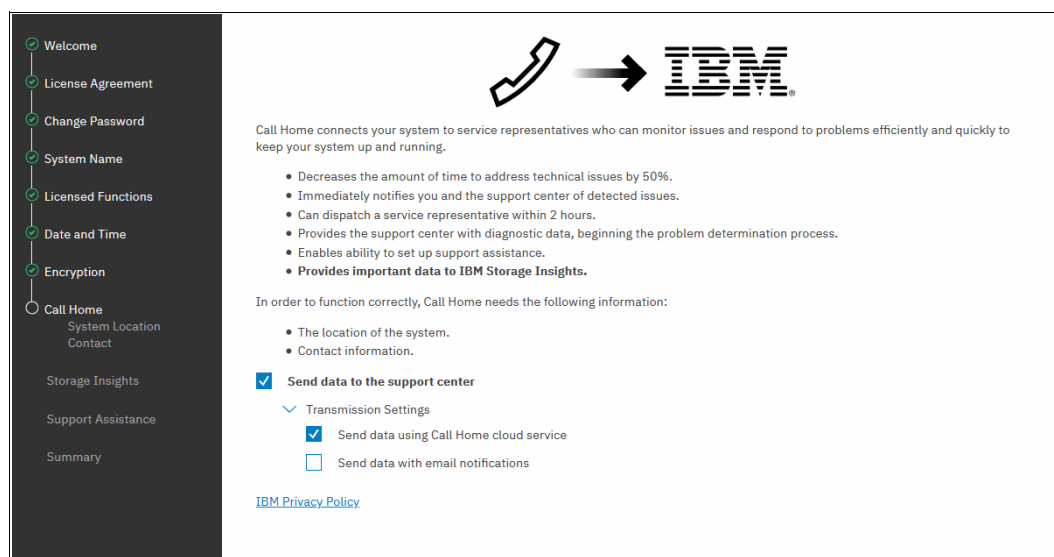


Figure 3-15 System setup: Call Home methods

For more information about setting up Call Home, including Cloud Call Home, see Chapter 13, “Reliability, availability, and serviceability, and monitoring and troubleshooting” on page 753.

If either of these methods is selected, the system location and contact information must be entered. This information is used by IBM to provide technical support. All fields in the form must be populated. In this step, the system also verifies that it can contact the Cloud Call Home servers, as shown in Figure 3-16 on page 93.

Welcome  
 License Agreement  
 Change Password  
 System Name  
 Licensed Functions  
 Date and Time  
 Encryption  
 Call Home  
 System Location  
 Contact

✓ Connection to the support center was successful!

### System Location

Service parts should be shipped to the same physical location as the system.

Company name:

System address:

City:

Figure 3-16 System setup: System location

After clicking **Next**, you can provide business-to-business contact information that IBM Support uses to contact a person who manages this machine if it is necessary, as shown in Figure 3-17.

Welcome  
 License Agreement  
 Change Password  
 System Name  
 Licensed Functions  
 Date and Time  
 Encryption  
 Call Home  
 System Location  
 Contact

### Contact

The support center contacts this person to resolve issues on the system.

Enter business-to-business contact information. To comply with privacy regulations, personal contact information for individuals with your organization is not recommended.

Name:

Email:

Phone (primary):

Figure 3-17 System setup: Contact information

If the **Email notifications** option was selected, you are prompted to enter the details for the email servers to be used for Call Home. Figure 3-18 shows an example. You can click **Ping** to verify that the email server is reachable over the network. Click **Apply** and then **Next**.

Welcome  
 License Agreement  
 Change Password  
 System Name  
 Licensed Functions  
 Date and Time

### Email Servers

Call home and event notifications are routed through this email server.

Server IP:  Port:

Figure 3-18 System setup: Email servers

11. ISAN Volume Controller systems may be used with IBM Storage Insights, which is an IBM cloud storage monitoring and management tool. During this setup phase, the system tries to contact the IBM Storage Insights web service. If it is available, you are prompted to sign up, as shown in Figure 3-19.

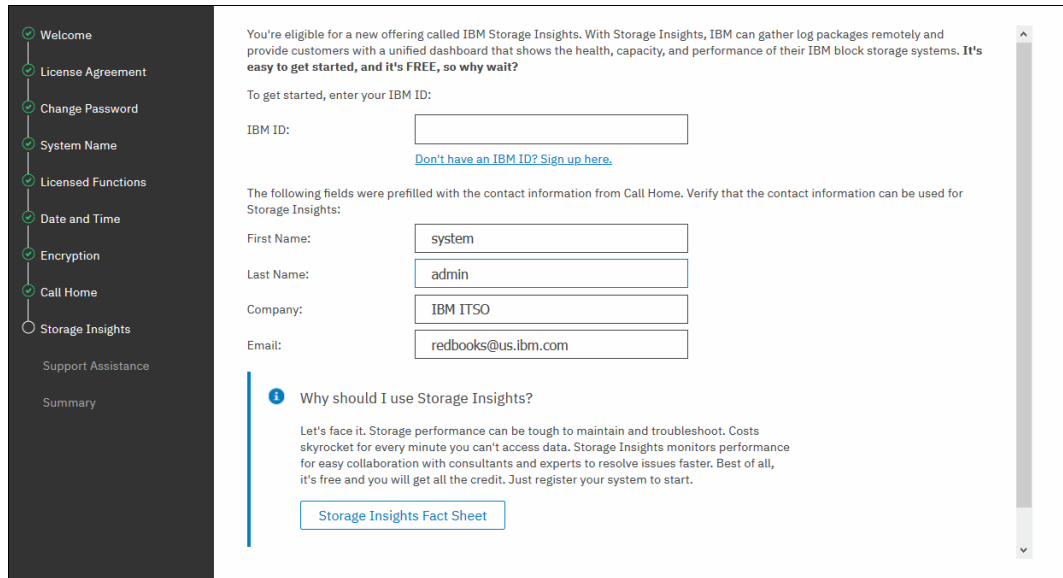


Figure 3-19 System setup: IBM Storage Insights

If a connection cannot be established, you are prompted to add the system that you are currently working on to the IBM Storage Insights setup manually, as shown in Figure 3-20.

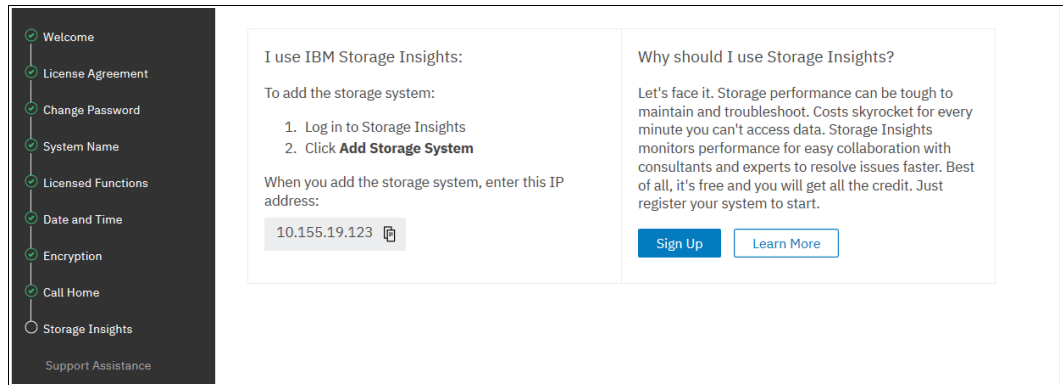


Figure 3-20 System setup: IBM Storage Insights

For more information about IBM Storage Insights, see Chapter 13, “Reliability, availability, and serviceability, and monitoring and troubleshooting” on page 753.

12. After you click **Next**, if you enabled at least one Call Home method, the Support Assistance configuration window opens, as shown in Figure 3-21 on page 95. The Support Assistance function requires Call Home, so if it is disabled, Support Assistance cannot be used.

With the Support Assistance feature, you allow IBM Support to perform maintenance tasks on your system while an IBM SSR is onsite. The IBM SSR can log in locally with your permission and a special user ID and password so that a superuser password does not need to be shared with the IBM SSR.



You can also enable Support Assistance with remote support to allow IBM Support personnel to log in remotely to the machine with your permission through a secure tunnel over the internet.

For more information about the Support Assistance feature, see Chapter 13, “Reliability, availability, and serviceability, and monitoring and troubleshooting” on page 753.

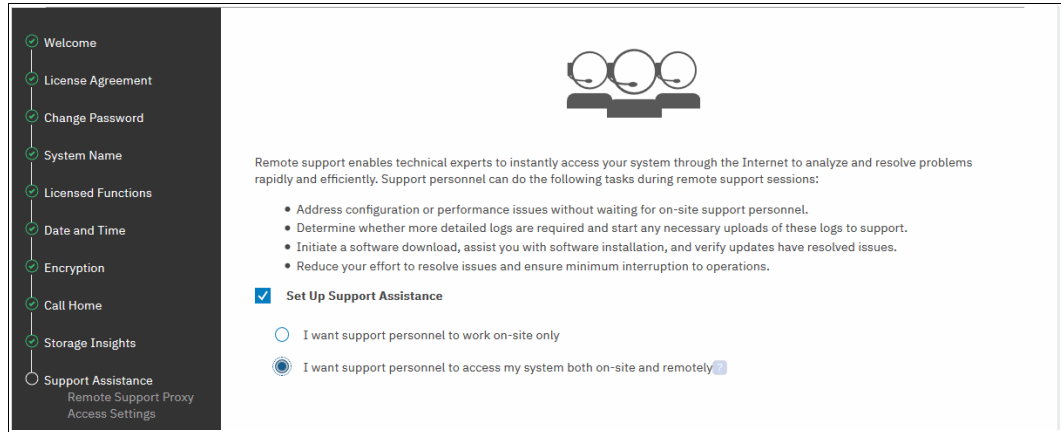


Figure 3-21 System setup: Support Assistance

If you allow remote support, you are given the IP addresses and ports of the remote support centers and an opportunity to provide proxy server details (if required) to allow the connectivity, as shown in Figure 3-22. Also, you can allow remote connectivity at any time or only after obtaining permission from the storage administrator.

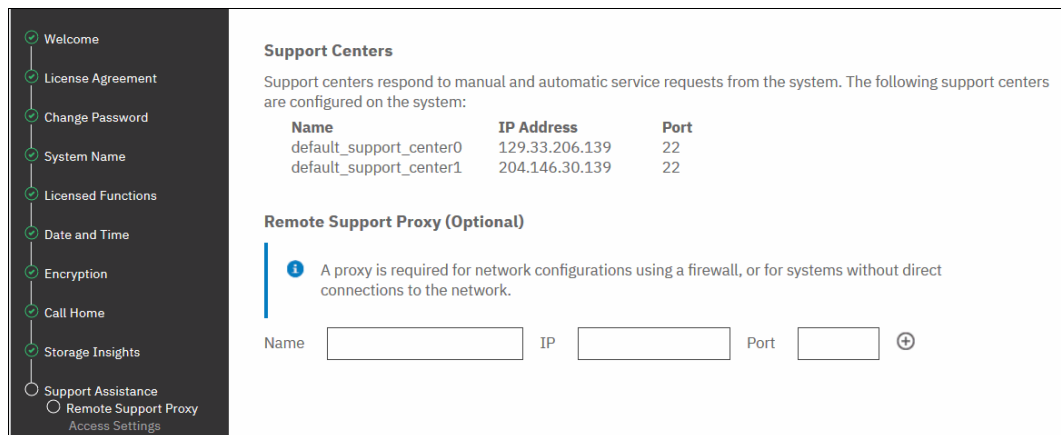


Figure 3-22 System setup: Support Centers

13. On the Summary page, the settings that were set by the system setup wizard are shown. If corrections are needed, you may return to a previous step by clicking **Back**. Otherwise, click **Finish** to be redirected to a system GUI.

After the wizard completes, your system consists only of the nodes that were available as candidates: powered off, connected to SAN and zoned together. If you have nodes, you must add them to complete system setup. For more information about how to add a node, see 3.4.2, “Adding a node or hot spare node” on page 99.

If you have no more nodes to add to this system, the system setup process is complete. All the mandatory steps of the initial configuration are done. If required, you can configure other global functions, such as system topology, user authentication, or local port masking, before configuring the volumes and provisioning them to hosts.

## 3.4 Base configuration

Tasks that are listed in this section are used to define global system configuration settings. Often, they are performed during system setup. However, they also can be performed any time later, such as when the system is expanded or the system environment is reconfigured.

### 3.4.1 Configuring Remote Direct Memory Access clustering

Up to four control enclosures may be joined in an IBM HyperSwap or a standard topology cluster. This subsection describes the configuration steps that must be performed if a system is designed for IP-based Remote Direct Memory Access (RDMA) node-to-node traffic. For FC SAN clustering, no special configuration is required on the system, but the SAN must be set up as described in Chapter 2, “Planning” on page 53.

#### Prerequisites

Before RDMA clustering is configured, ensure that the following prerequisites are met:

- ▶ 25 gigabits per second (Gbps) RDMA-capable Ethernet cards are installed in each node.
- ▶ RDMA-capable adapters in all nodes use the same technology, such as RDMA over Converged Ethernet (RoCE) or internet Wide Area RDMA Protocol (iWARP).
- ▶ RDMA-capable adapters are installed in the same slots across all the nodes of the system.
- ▶ Ethernet cables between each node are connected correctly.
- ▶ The network configuration does not contain more than two hops in the fabric of switches. The router must *not* be placed between nodes that use RDMA-capable Ethernet ports for node-to-node communication.
- ▶ The negotiated speeds on the local and remote adapters are the same.
- ▶ The local and remote port virtual local area network (VLAN) identifiers are the same. All the ports that are used for node-to node communication must be assigned to one VLAN ID, and ports that are used for host attachment must have a different VLAN ID. If you plan to use VLAN to create this separation, you must configure VLAN support on the all the Ethernet switches in your network before you define the RDMA-capable Ethernet ports on nodes in the system. On each switch in your network, set the VLAN to Trunk mode and specify the VLAN ID for the RDMA-ports that will be in the same VLAN.
- ▶ A minimum of two dedicated RDMA-capable Ethernet ports are required for node-to-node communications to ensure best performance and reliability. These ports must be configured for inter-node traffic only and must not be used for host attachment, virtualization of Ethernet-attached external storage, or IP replication traffic.
- ▶ A maximum of four RDMA-capable Ethernet ports per node are allowed for node-to-node communications.

## Configuration process

To enable RDMA clustering, IP addresses must be configured on each port of each node that is used for node-to-node communication. Complete the following steps:

1. Connect to a Service Assistant of a node by going to [https://<node\\_service\\_IP>/service](https://<node_service_IP>/service) and clicking **Change Node IP**, as shown in Figure 3-23.

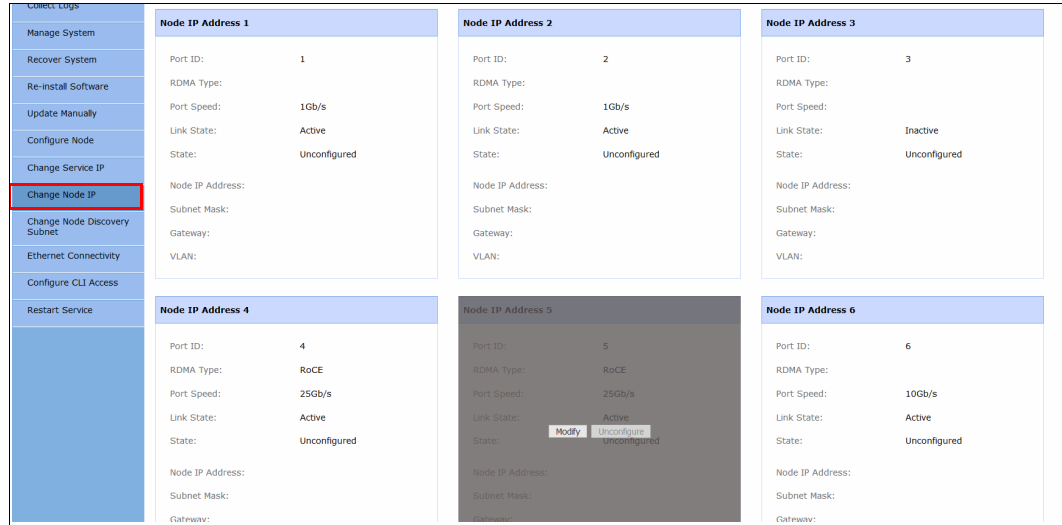


Figure 3-23 Node IP address setup for RDMA clustering

Figure 3-23 shows that ports 1 - 3 do not show any RDMA type, so they cannot be used for node-to-node traffic. Ports 4 and 5 show RDMA type RoCE, so they can be used.

2. Hover your cursor over a tile with a port and click **Modify** to set the IP address, netmask, gateway address, and VLAN ID for a port. The IP address for each port must be unique and cannot be used anywhere else on the system. The VLAN ID for ports that are used for node-to-node traffic must be the same on all nodes. When the required information is entered, click **Save** and verify that the operation completed successfully, as shown in Figure 3-24. Repeat this step for all ports that you intend to use for node-to-node traffic, with a minimum of two and a maximum of four ports per node.

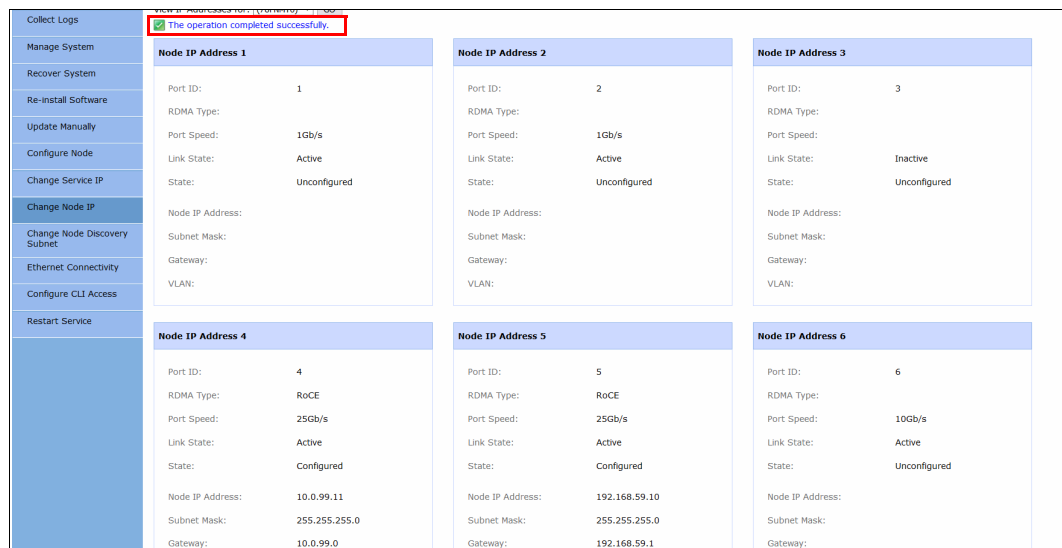


Figure 3-24 Node IP addresses configured

To list and change the node IP configuration by using the CLI, use the **sainfo lsnodeip** and **satask chnodeip** commands, as shown in Example 3-1.

*Example 3-1 Setting IP addresses for node-to-node connectivity*

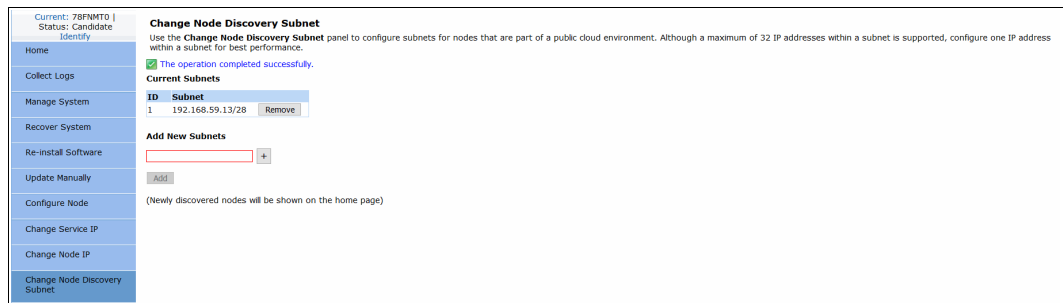
```

IBM_2145::superuser>sainfo lsnodeip
port_id    rdma_type port_speed vlan link_state state      node_IP_address
1          1Gb/s      active  unconfigured
2          1Gb/s      active  unconfigured
3          1Gb/s      active  unconfigured
4          RoCE      25Gb/s active  unconfigured
5          RoCE      25Gb/s active  unconfigured
IBM_2145::superuser>satask chnodeip -ip 10.0.99.12 -gw 10.0.99.1 -mask 255.255.255.0 -port_id 4
IBM_2145::superuser>satask chnodeip -ip 192.168.59.11 -gw 192.168.59.1 -mask 255.255.255.0
-port_id 5
IBM_2145::superuser>sainfo lsnodeip
port_id    rdma_type port_speed vlan link_state state      node_IP_address
1          1Gb/s      active  unconfigured
2          1Gb/s      active  unconfigured
3          1Gb/s      active  unconfigured
4          RoCE      25Gb/s active  configured 10.0.99.12
5          RoCE      25Gb/s active  configured 192.168.59.11

```

- Some environments might not include a stretched layer 2 subnet. In such scenarios, a layer 3 network such as in standard topologies or long-distance RDMA node-to-node HyperSwap configurations is applicable. To support the layer 3 Ethernet network, use the unicast discovery method for RDMA node-to-node communication. This method relies on unicast-based fabric discovery rather than multicast discovery.

To configure unicast discovery, see the information about the **satask addnodediscoverysubnet**, **satask rmnodediscoverysubnet**, or **sainfo lsnodediscoverysubnet** commands in [IBM Knowledge Center](#). You can also configure discovery subnets by using the Service Assistant interface menu option **Change Node Discovery Subnet**, as shown in Figure 3-25.



*Figure 3-25 Setting the node discovery subnet*

- After the IP addresses are configured on all nodes in a system, run the **sainfo lsnodeipconnectivity** command or use the Service Assistant GUI menu **Ethernet Connectivity** to verify that the partner nodes are visible on the IP network, as shown in Figure 3-26 on page 99. If necessary, troubleshoot connection problems by running the **ping** and **sainfo traceroute** commands.

Cluster ID	Status	Error Data	Port	ID	Type	IP Address	VLAN	WWNN
000002032C411602	Connected:RoCE		Local	4	RoCE	192.168.59.10		
			Remote	4	RoCE	192.168.59.11		500507680C008B01
000002032C411602	Connected:RoCE		Local	5	RoCE	10.0.99.11		
			Remote	5	RoCE	10.0.99.12		500507680C008B01

Figure 3-26 Node-to-node Ethernet connectivity

When all the nodes that are joined to the cluster are connected, the enclosure may be added to the cluster.

### 3.4.2 Adding a node or hot spare node

This procedure is the same whether you are configuring the system for the first time or expanding it later. The same process is used to add a node to an I/O group, or a hot spare node.

Before beginning this process, ensure that the new control enclosure is correctly installed and cabled to the existing system. For FC node-to-node communication, verify that correct the SAN zoning is set. For node-to-node communication over RDMA-capable Ethernet ports, ensure that the IP addresses are configured and a connection between nodes can be established.

To add a node to the system, complete the following steps:

1. In the GUI, select **Monitoring** → **System**. When a new enclosure is detected by a system, the **Add Node** button appears on the System - Overview window next to System Actions, as shown in Figure 3-27.



Figure 3-27 Add Node button

**Note:** If the **Add Node** button does not appear, review the installation instructions to verify that the new node is connected and set up correctly.

2. Click **Add Node**, and a form that you can use to assign nodes to I/O groups opens, as shown in Figure 3-28. To light the Identify light-emitting diode (LED) on a node, click the LED icon that is next to a node name. When the required node (or nodes) is assigned, click **Next**.

Figure 3-28 Assigning a node

3. Review the summary in the next window and click **Finish** to add the node or nodes to the system.

**Note:** When a node is added, the software version running on it is upgraded or rolled back to match the cluster software version. This process can take up to 30 minutes or more, and the node is added only after this process completes.

4. After the node is successfully added to the system, a success message appears. Click **Close** to return to the System Overview window and check that new node is visible and available for management.

To perform the same procedure by using a CLI, complete the following steps. For more information about the detailed syntax for each command, go to [IBM Knowledge Center](#).

1. When adding nodes, check for unpopulated I/O groups by running **lsiogrp**. Each complete I/O group has two nodes. Example 3-2 shows that only `io_grp0` has nodes, so a new control enclosure can be added to `io_grp1`.

*Example 3-2 Listing I/O groups*

```
IBM_2145:ITS0-SVC:superuser>lsiogrp
id name                node_count vdisk_count host_count site_id site_name
0  io_grp0              2          0          0          0
```

1	io_grp1	0	0	0
2	io_grp2	0	0	0
3	io_grp3	0	0	0
4	recovery_io_grp	0	0	0

- To list nodes that are available to add to the I/O group, run the `lscnodecandidate` command, as shown in Example 3-3.

*Example 3-3 Listing the candidate nodes*

```
IBM_2145:ITS0-SVC:superuser>lscnodecandidate
id          panel_name UPS_serial_number UPS_unique_id hardware
500507680C000416 KD8P1BP                               0000000000000000 DH8
```

- Add a node by running the `addnode` command, as shown in Example 3-4. The command triggers only the process, which starts in background and can take up to 30 minutes or more.

*Example 3-4 Adding a node as a spare*

```
IBM_2145:ITS0-SVC:superuser>addnode -panelname KD8P1BP -spare
```

Example 3-5 shows same command, but used to add a node to an I/O group `io_grp1`.

*Example 3-5 Adding a node to an I/O group*

```
IBM_2145:ITS0-SVC:superuser>addnode -panelname KD8P1BP -name node3 -iogrp 1
```

### 3.4.3 Changing the system topology

SAN Volume Controller supports two multi-site topologies” *HyperSwap* and *Enhanced Stretched Cluster* (ESC). For each I/O group in the system, the “stretched” topology has one node on one site and one node on the other site. The HyperSwap topology places both nodes of an I/O group at the same site. Both topologies enable full configuration of the highly available (HA) volumes through a single point of configuration.

If your solution is designed for ESC or HyperSwap, use the guidance in this section to configure your topology for either solution.

For a list of requirements for a HyperSwap or ESC configuration, see [IBM Knowledge Center](#) and expand **Configuring** → **Configuration details** → **Stretched system configuration details** and **HyperSwap system configuration details**.

To change the system topology, complete the following steps:

- In the GUI, click **Monitoring** → **System** to open the System - Overview window. Click **System Actions** and select **Modify System Topology**, as shown in Figure 3-29.



*Figure 3-29 Starting the Modify System Topology wizard*

2. The Modify Topology wizard welcome window opens. Click **Next**. You are prompted to change the default site names, as shown in Figure 3-30. The site names can indicate, for example, building locations for each site, or other descriptive information.

Assign Site Names

Enter the names:

Site 1:

Site 2:

Site 3 (quorum):

Figure 3-30 Assigning site names

3. Select **HyperSwap System** or **Stretched System** for the topology, and assign I/O groups to the sites. Click the marked icons in the center of the window to swap site assignments, as shown in Figure 3-31. Click **Next**.

Modify System Topology

Assign Nodes

Topology:

I/O Group 0:

node2	
node1	

I/O Group 0

PaloAlto

Figure 3-31 Specifying the system topology

4. If any host objects or back-end storage controllers are configured, you must assign a site for each of them. Right-click the object and click **Modify Site**, as shown in Figure 3-32 on page 103.



### Assign External Storage Systems to Sites

All external storage must have a site that is assigned before you configure a Stretched system.

At least one storage system must be assigned to the quorum site.

Actions ▼

Name	Status	Site	IQN Count	Model	
FS9100	✓ Online	PaloAlto		IBM 2145	<input type="button" value="Modify Site"/>
V7k_5_n2	✓ Online	MountainView		IBM 2145	
V7k_4_n2	✓ Online	Almaden (Quorum)		IBM 2145	
V7k_4_n1	✓ Online	Almaden (Quorum)		IBM 2145	
V7k_5_n1	✓ Online	MountainView		IBM 2145	
V7k_5_n2	✓ Online	MountainView		IBM 2145	

Showing 6 storage systems | Selecting 1 storage system

All external storage must be assigned to a site to ensure the I/O is routed correctly.

The quorum site determines which site can process I/O if communication between sites is lost or one site becomes unavailable.

Figure 3-32 Assigning controllers to the sites

- If you are configuring HyperSwap topology, set the maximum background copy operations bandwidth between the sites. *Background copy* is the initial synchronization and any subsequent resynchronization traffic for HyperSwap volumes. Use this setting to limit the impact of volume synchronization to host operations. You may also set it higher during the initial setup (when there are no host operations on the volumes yet), and set it lower when the system is in production.

As shown in Figure 3-33, you must specify the total bandwidth between the sites in megabits per second (Mbps) and what percentage of this bandwidth that can be used for background copying. Click **Next**.

An ESC topology system does not require this setting.

### Set Bandwidth Between Sites

Bandwidth between sites:  Mbps

Portion bandwidth for background copies:  %

Total background copy rate: 19 MiB/second

*To ensure that host and storage system operations are not affected by the background copy operations, you can assign a portion of the bandwidth for background copies.*

Figure 3-33 Setting the bandwidth between the sites for a HyperSwap topology

- Review the summary and click **Finish**. The wizard starts implementing changes to migrate the system to the topology.

When you later add a host or back-end storage controller objects, the GUI prompts you to set an object site during the creation process.

### 3.4.4 Configuring quorum disks or applications

Quorum devices are required for a system to hold a copy of important system configuration data. A managed disk (MDisk) from FC-attached external back-end storage or a special application that is connected over an IP network may work as a quorum device.

One of these items is selected for the *active quorum* role, which is used to resolve failure scenarios where half the nodes on the system become unavailable or a link between enclosures is disrupted. The active quorum determines which nodes can continue processing host operations and to avoid a “split brain” condition, which happens when both halves of the system continue I/O processing independently of each other.

For systems with a standard topology, quorum devices are automatically assigned from a managed MDisk. No special configuration actions are required. Optionally, an IP quorum device can be configured to provide extra redundancy.

For ESC and HyperSwap topology systems, an active quorum device must be on a third, independent site. Due to the costs that are associated with deploying a separate FC-attached storage device on a third site, an IP-based quorum device may be used for this purpose.

A stretched or HyperSwap system can be configured without a quorum device at a third site. If there is no third site, then quorum must be configured to select a site to always win a tie-breaker. If there is a loss of connectivity between the sites, then the site that is configured as the winner continues operating and processing I/O requests and the other site stops until the fault is fixed. If there is a site outage at the winning site, then the system stops processing I/O requests until this site is recovered or the manual quorum override procedure is used.

#### Creating and installing an IP quorum application

To create and install an IP quorum application, complete the following steps:

1. Select **System** → **Settings** → **IP Quorum** to download the IP quorum application, as shown in Figure 3-34. If you are using IPv6 for management IP addresses, the **Download IPv6 Application** button is available and the IPv4 option is disabled.

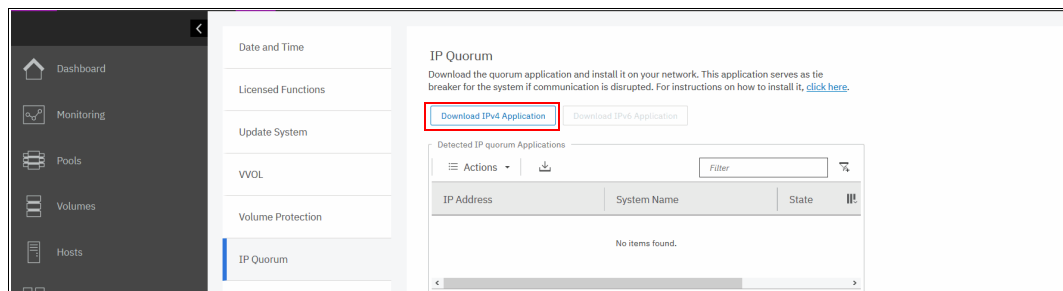


Figure 3-34 Download IPv4 quorum button

2. After you click **Download...**, a window opens, as shown in Figure 3-35 on page 105. It provides an option to create an IP application that is used for tie-breaking only, or an application that can be used as a tie-breaker and to store recovery metadata.

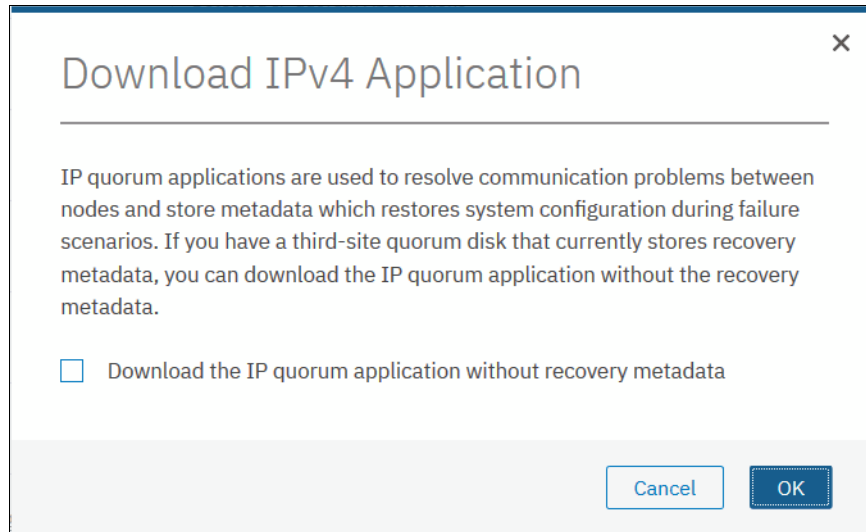


Figure 3-35 Download IP quorum application window

An application that does not store recovery metadata requires less channel bandwidth for a link between the system and the quorum application, which may be a decision-making factor for using a multi-site HyperSwap system.

For a full list of IP quorum app requirements, see [IBM Knowledge Center](#) and expand **Configuring** → **Configuration details** → **Configuring quorum** → **IP quorum application configuration**.

3. After you click **OK**, the `ip_quorum.jar` file is created. Save the file and transfer it to a supported AIX, Linux, or Windows host that can establish an IP connection to the service IP address of each system node. Move it to a separate directory and start the app, as shown in Example 3-6.

*Example 3-6 Starting the IP Quorum app on the Windows operating system*

```
C:\IPQuorum>java -jar ip_quorum.jar
=== IP quorum ===
Name set to null.
Successfully parsed the configuration, found 4 nodes.
....
```

**Note:** Add the IP quorum application to the list of auto-started applications at each start or restart or configure your operating system (OS) to run it as an auto-started service in the background.

The IP quorum log file and recovery metadata are stored in the same directory with the `ip_quorum.jar` file.

4. Check that the IP quorum application is successfully connected and running by verifying its Online status by selecting **System** → **Settings** → **IP Quorum**, as shown in Figure 3-36.

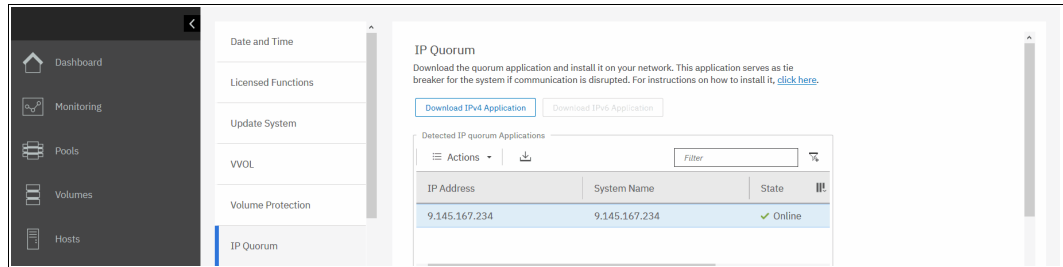


Figure 3-36 IP quorum application that is deployed and connected

## Configuring IP quorum mode

On a standard topology system, only the Standard quorum mode is supported. No additional configuration is required. On an ESC or HyperSwap topology, you may configure different tie-breaker scenarios (a tie occurs when exactly half of the nodes that were previously a member of the system are present):

- ▶ If the quorum mode is set to **Standard**, both sites have an equal chance to continue working after the tie-breaker.
- ▶ If the quorum mode is set **Preferred**, during a disruption, the system delays processing tie-breaker operations on non-preferred sites, leaving more time for the preferred site to win. If during an extended period a preferred site cannot contact the IP quorum app (for example, if it is destroyed), a non-preferred site continues working.
- ▶ If the quorum mode is set to **Winner**, the selected site always is the tie-breaker winner. If the winner site is destroyed, the remaining site may continue operating only after manual intervention.

The Preferred and Winner quorum modes are supported only with an IP quorum. For a FC-attached active quorum MDisk, only Standard mode is possible.

To set a quorum mode, select **System** → **Settings** → **IP Quorum** and click **Quorum Setting**. The Quorum Setting window opens, as shown in Figure 3-37.

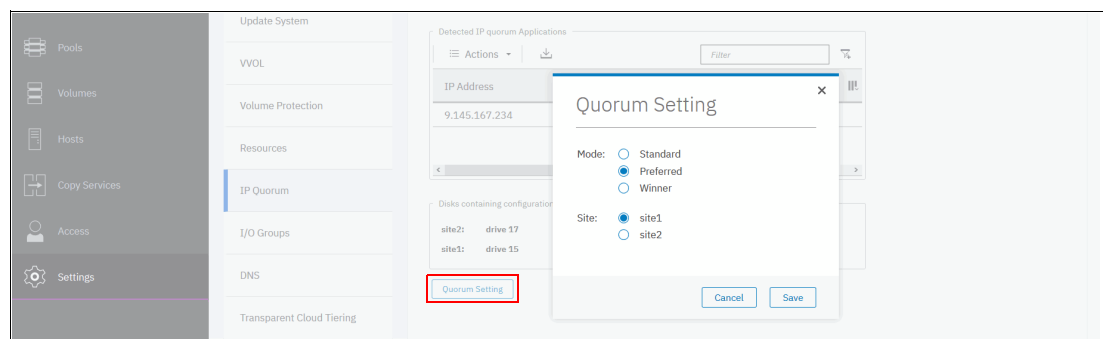


Figure 3-37 Changing the quorum mode

## 3.4.5 Configuring the local Fibre Channel port masking

With FC port masking, you control the usage of FC ports. By applying a mask, you restrict node-to-node communication or FC RC traffic on selected ports.

To decide whether your system must have port masks configured, see 2.6.8, “Port designation recommendations” on page 62.

To set the FC port mask by using the GUI, complete the following steps:

1. Select **System** → **Network** → **Fibre Channel Ports**. In a displayed list of FC ports, the ports are grouped by a system port ID. Each port is configured identically across all nodes in the system. You can click the arrow next to the port ID to expand a list and see which node ports (N\_Port) belong to the selected system port ID and their worldwide port names (WWPNs).
2. Right-click a system port ID that you want to change and select **Modify Connection**, as shown in Figure 3-38.

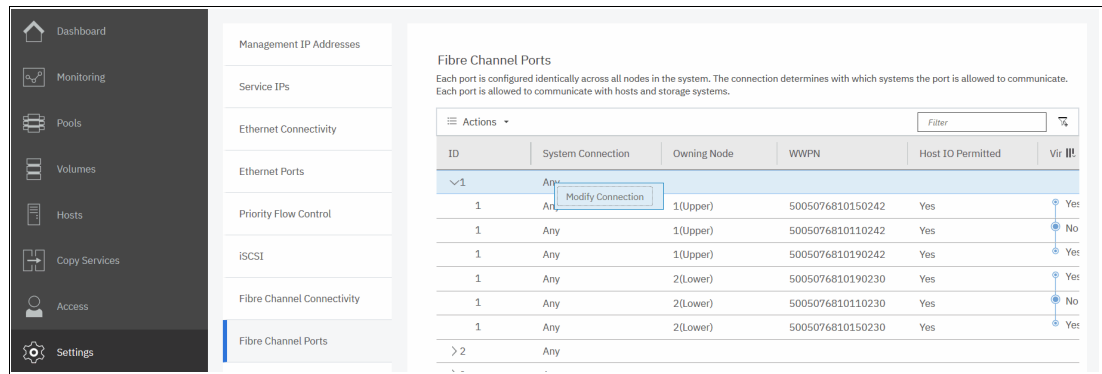


Figure 3-38 Applying a port mask by using a GUI

By default, all system ports can send and receive traffic of any kind:

- ▶ Host traffic
- ▶ Traffic to virtualized back-end storage systems
- ▶ Local system traffic (node-to-node)
- ▶ Partner system (remote replication) traffic

The first two types are always allowed, and you may control them only with SAN zoning. The other two types can be blocked by port masking. In the Modify Connection dialog box, as shown in Figure 3-39, you can choose which type of traffic that a port can send.

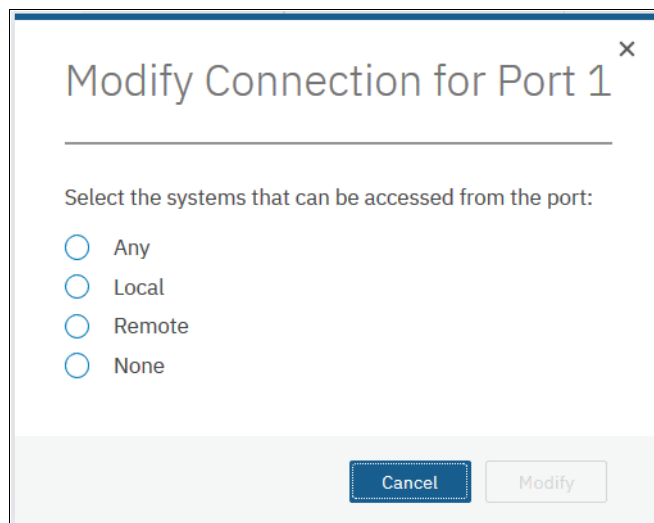


Figure 3-39 Modify Connection dialog box



*Signed* SSL certificates are issued by a trusted CA. A browser maintains a list of trusted CAs that are identified by their *root* certificate. The root certificate must be included in this list in order for the signed certificate to be trusted.

To see the details of your system certificate, select **Settings** → **Security** and click **Secure Communications**, as shown in Figure 3-40, or run the `lsystemcert` command.

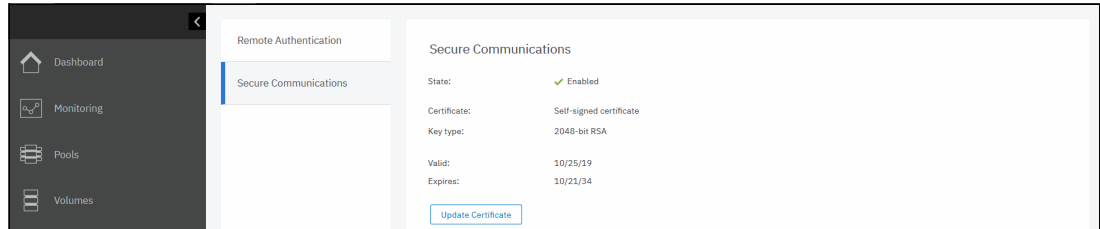


Figure 3-40 Accessing the Secure Communications window

Based on the security requirements for your system, you can create either a new self-signed certificate or install a signed certificate that is created by a third-party CA.

### Generating a self-signed certificate

If a self-signed certificate is expired or its key type does not comply with your company's security policy, you can regenerate it. To renew a self-signed certificate, complete the following steps:

1. Select **Update Certificate** on the Secure Communications window, as shown in Figure 3-40.
2. Select **Self-signed certificate** and enter the details for the new certificate. The “Key type” and “Validity days” are the only mandatory fields.

**Note:** Before re-creating a self-signed certificate, ensure that your browser supports the type of keys that you are going to use for a certificate. See your organization's security policy to ensure what key type is required.

3. Click **Update**.

You are prompted to confirm the action. Click **Yes** to proceed. Close the browser, wait approximately 2 minutes, and reconnect to the management GUI.

To regenerate an SSL certificate by using a CLI, run the `chsystemcert` command, as shown in Example 3-9. Valid values for `-keytype` are `rsa2048`, `ecdsa384`, or `ecdsa521`.

#### Example 3-9 Regenerating a self-signed certificate

```
IBM_2145:ITS0-SVC:superuser>chsystemcert -mkselfsigned -keytype ecdsa521 -validity 365
```

## Configuring a signed certificate

If your company's security policy requests certificates to be signed by a trusted authority, complete the following steps to configure a signed certificate:

1. Select **Update Certificate** in the Secure Communications window.
2. Select **Signed certificate** and enter the details for the new certificate signing request, as shown in Figure 3-41. All fields are mandatory except for the Subject Alternative Name. For the "Country" field, use a two-letter country code. Click **Generate Request**.

The screenshot shows a dialog box titled "Update Certificate" with a close button (X) in the top right corner. The dialog is divided into two main sections: "Certificate type" and "Certificate Signing Request".

**Certificate type:** There are two radio buttons. The first is "Self-signed certificate" (unselected). The second is "Signed certificate" (selected).

**Certificate Signing Request:** This section contains several input fields and a button:

- Key type:** A dropdown menu showing "2048-bit RSA".
- Country:** A text input field containing "US".
- State:** A text input field containing "CA".
- City:** A text input field containing "San Jose".
- Organization:** A text input field containing "IBM".
- Organization unit:** A text input field containing "ITSO".
- Common name:** A text input field containing "9.155.123.198".
- Subject Alternative Name:** A text input field with a help icon. The text inside reads: "Use the suggested format below:  
IP:123.45.67.91  
URI:http://www.mydomain.com  
DNS:cluster.mydomain.com  
email:support@mydomain.com".
- Email address:** A text input field containing "redbooks@us.ibm.com".
- Generate Request:** A blue button.

**Signed Certificate:** This section is partially visible at the bottom of the dialog, showing a "Signed certificate:" label and a dropdown menu with "Upload from Certificate Authority" selected.

At the bottom of the dialog, there are two buttons: "Cancel" and "Update".

Figure 3-41 Generating a certificate request

3. When prompted, save the `certificate.csr` file that contains the certificate signing request.

Until the signed certificate is installed, the Secure Communications window shows that an outstanding certificate request exists.



**Attention:** If you must update a field in the certificate request, generate a new request and submit it for signing by the proper CA. However, this process invalidates the previous certificate request and prevents the installation of the signed certificate that is associated with the original request.

4. Submit the request to the CA to receive a signed certificate. Notify the CA that you need a certificate (or certificate chain) in base64-encoded Privacy Enhanced Mail (PEM) format.
5. When you receive the signed certificate, select **Update Certificate** in the Secure Communications window again.
6. Select **Signed Certificate** and click the folder icon next to the **Signed Certificate** input field of the Update Certificate window, as shown in Figure 3-41 on page 110. Click **Update**.
7. You are prompted to confirm the action. Click **Yes** to proceed. After your certificate is installed, the GUI session disconnects. Close the browser window and wait approximately 2 minutes before reconnecting to the management GUI.
8. Reconnect to the GUI and select **Settings** → **Security** → **Secure Communications**. The window that opens should show that you are using a signed certificate, as shown in Figure 3-42.

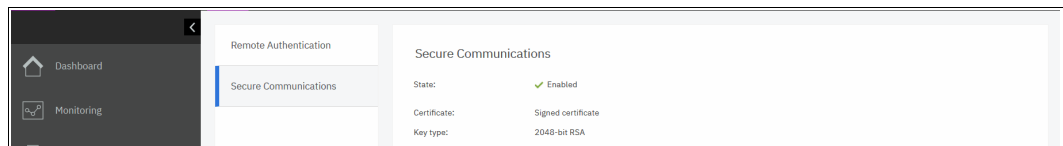


Figure 3-42 Signed certificate installed

## 3.5.2 Configuring user authentication

There are two methods of user authentication to control access to the GUI and to the CLI:

- ▶ *Local user authentication* is performed within the SAN Volume Controller system. GUI users authenticate with user name and password. CLI users must provide a user name and a Secure Shell (SSH) public key or a password.
- ▶ *Remote user authentication* allows users to authenticate to the system by using credentials that are stored on an external authentication service. With this feature, you use user credentials and user groups that are defined on the remote service to simplify user management and access, enforce password policies more efficiently, and separate user management from storage management.

Locally administered users can coexist with remote authentication.

### User roles and groups

User groups are used to determine what tasks the user is authorized to perform. Each user group is associated with a single role. Roles apply to both local and remote users on the system and are based on the user group to which the user belongs. A local user can belong only to a single group, so the role of a local user is defined by the single group to which that user belongs.

For a list of user roles and their tasks, and a description of a pre-configured user group, see [IBM Knowledge Center](#) and expand **Product overview** → **Technical overview** → **User roles**.

## Superuser account

Every system has a default user that is called the *superuser*. It cannot be deleted or modified, except for changing the password and SSH key. The superuser is a *local* user and cannot be authenticated remotely. The superuser has a *SecurityAdmin* user role, which has the most privileges within the system.

**Note:** The superuser is the only user that may log in to the Service Assistant interface. It is also the only user that may run **sa info** and **sa task** commands through the CLI.

The password for superuser is set during the system setup. The superuser password can be reset to its default value of `passwd` by using a procedure that is described in [IBM Knowledge Center](#) by expanding **Troubleshooting** → **Resolving a problem** → **Procedure: Resetting the superuser password**.

**Note:** The superuser password reset procedure uses system internal USB ports. Systems with SAN Volume Controller 2145-SV2 and 2145-SA2 nodes may be configured to disable those ports. If the USB ports are disabled and there are no users with the *SecurityAdmin* role and a known password, the superuser password cannot be reset without replacing the system hardware and deleting the system configuration.

## Local authentication

A *local user* is a user whose account is managed entirely on the system. A local user belongs to one user group only, and it must have a password, an SSH public key, or both. Each user has a name, which must be unique across all users in one system.

User names can contain up to 256 printable American Standard Code for Information Interchange (ASCII) characters. Forbidden characters are the single quotation mark ('), colon (:), percent symbol (%), asterisk (\*), comma (,), and double quotation marks ("). A user name cannot begin or end with a blank space.

Passwords for local users can be up to 64 printable ASCII characters, but cannot begin or end with a space.

When connecting to the CLI, encryption key authentication is attempted first with the user name and password combination available as a fallback. The SSH key authentication method is available for CLI and file transfer access only. For GUI access, only the password is used.

To add a user that is authenticated without a password by using only an SSH key, select **Access** → **Users by Group**, click **Add user**, and then click **Browse** to select the SSH public key for that user, as shown in Figure 3-43 on page 113. The Password field may be left blank. The system accepts public keys that are generated by PuTTY (SSH2), OpenSSH, and RFC 4716-compliant keys that are generated by other clients.

Figure 3-43 Creating a user authenticated by a SSH key

If local authentication is used, user accounts must be created for each system. If you want access for a user on multiple systems, you must define the user in each system.

### Remote authentication

A *remote user* is authenticated by using identity information that is accessible by using the Lightweight Directory Access Protocol (LDAP). The LDAP server must be available for the users to log in to the system. Remote users have their groups defined by the remote authentication service.

Users that are authenticated by an LDAP server can log in to the management GUI and the CLI. These users do not need to be configured locally for CLI access, and they do not need an SSH key that is configured to log in by using the CLI.

If multiple LDAP servers are available, you can configure more than one LDAP server to improve resiliency. Authentication requests are processed by those LDAP servers that are marked as preferred unless the connection fails or a user is not found. Requests are distributed across all preferred servers for load balancing in a round-robin fashion.

**Note:** All LDAP servers that are configured within the same system must be of the same type.

If users that are part of a group on the LDAP server are to be authenticated remotely, a user group with an identical name must exist on the system. The user group name is *case-sensitive*. The user group must also be enabled for remote authentication on the system.

A user who is authenticated remotely is granted permissions according to the role that is assigned to the user group of which the user is a member.

To configure remote authentication by using LDAP, start by enabling remote authentication by completing the following steps:

1. Select **Settings** → **Security**, click **Remote Authentication**, and then click **Configure Remote Authentication**, as shown in Figure 3-44.

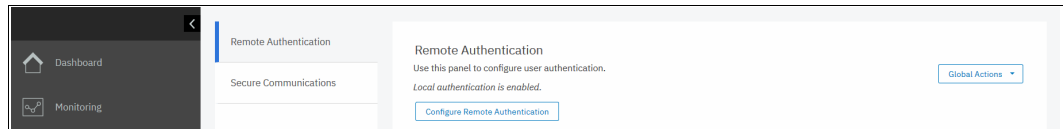


Figure 3-44 Configuring remote authentication

2. Enter the LDAP settings. These settings are not server-specific. They are applied to all LDAP servers that are configured in the system. Extra optional settings are available by clicking **Advanced Settings**, as shown in Figure 3-45.

The dialog box is titled 'Configure Remote Authentication' and has a close button (X) in the top right. It is divided into several sections: 'LDAP Type' with radio buttons for 'IBM Tivoli Directory Server', 'Microsoft Active Directory' (selected), and 'Other'; 'Security' with radio buttons for 'LDAP with StartTLS' (selected), 'LDAPS', and 'LDAP with no security'; 'Service Credentials (Optional)' with input fields for 'User Name' and 'Password'; and 'Advanced Settings' with input fields for 'User Attribute' (sAMAccountName), 'Group Attribute' (memberOf), and 'Audit Log Attribute' (userPrincipalName). At the bottom are 'Cancel', 'Back', and 'Next' buttons.

Figure 3-45 Configure Remote Authentication settings

The following settings are available:

- LDAP type:
  - **IBM Security Directory Server** (for IBM Security Directory Server).
  - **Microsoft Active Directory** (AD).
  - **Other** (other LDAP v3-capable directory servers, for example, OpenLDAP).
- Security
  - **LDAP with StartTLS**: Select this option to use the StartTLS extension (RFC 2830). It works by establishing a non-encrypted connection with an LDAP server on a standard LDAP port (389), and then performing a TLS handshake over an existing connection.
  - **LDAPS**: Select to use LDAP over SSL and establish secure connections by using port 636.
  - **None**: Select to transport data in clear text format without encryption.
- Service Credentials: Sets a user name and password for administrative binding (the credentials of a user that has the authority to query the LDAP directory). Leave it empty if your LDAP server is configured to support anonymous bind.  
For AD, a user name must be in User Principal Name (UPN) format.
- Advanced settings

Speak to the administrator of the LDAP server to ensure that these fields are completed correctly:

- **User Attribute**

This LDAP attribute is used to determine the user name of remote users. The attribute must exist in your LDAP schema and must be unique for each of your users.

This advanced setting defaults to `sAMAccountName` for AD and to `uid` for **IBM Security Directory Server** and **Other**.

- **Group Attribute**

This LDAP attribute is used to determine the user group memberships of remote users. The attribute must contain either the distinguished name of a group or a colon-separated list of group names.

This advanced setting defaults to `memberOf` for AD and **Other**, and to `ibm-allGroups` for **IBM Security Directory Server**. For **Other** LDAP type implementations, you might need to configure the `memberOf` overlay if it is not in place.

- **Audit Log Attribute**

This LDAP is an attribute that is used to determine the identity of remote users. When an LDAP user performs an audited action, this identity is recorded in the audit log. This advanced setting defaults to `userPrincipalName` for AD and to `uid` for IBM Security Directory Server and the **Other** type.

3. Enter the server settings for one or more LDAP servers, as shown in Figure 3-46.

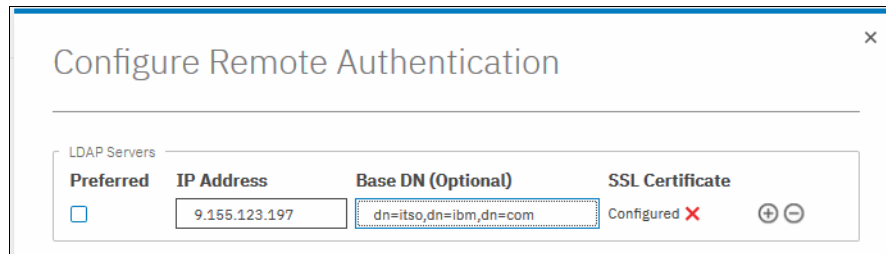


Figure 3-46 Configure Remote Authentication: Creating an LDAP server

The following settings are available:

– **Preferred**

One or more configured LDAP servers may be marked as **Preferred**. Requests are distributed among these servers, and use only non-preferred servers if all the preferred servers failed.

– **IP Address**

The IP address of the server.

– **Base DN**

The distinguished name to use as a starting point for searching for users on the server (for example, dc=itso,dc=ibm,dc=com).

– **SSL Certificate**

The SSL certificate that is used to securely connect to the LDAP server. This certificate is required only if you chose to use SSL or Transport Layer Security as a security method earlier.

Click **Browse** to select a server certificate. The system accepts certificates in base-64 encoded PEM format. To get a certificate in PEM format from your AD server, select **Base-64 Encoded X.509 (.CER)** in the MS Windows Certificate Export wizard, as shown in Figure 3-47 on page 117.

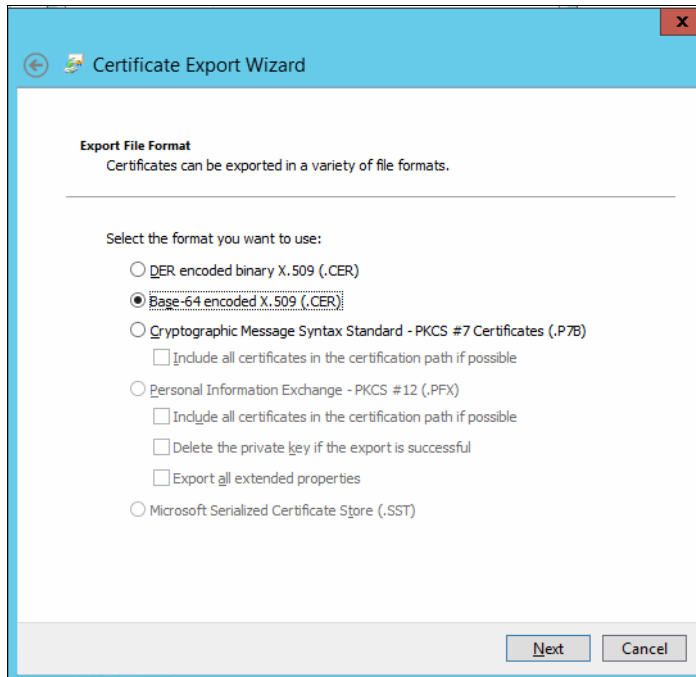


Figure 3-47 Exporting the AD certificate

**Note:** If your organization is using a tiered CA hierarchy, a server certificate that is exported for use on a system must include all the certificates in a chain. To accomplish this task, export the certificate in MS Windows in .P7B format and use third-party tools (OpenSSL) to convert it to PEM format. Otherwise, the exported certificate will not contain all certificates in the certification path.

If you set a certificate and you want to remove it, click the red cross next to **Configured**.

- Click the plus and minus signs to add or remove LDAP server records. You may define up to six servers.

Click **Finish** to save the settings.

4. To verify that LDAP is enabled, select **Settings** → **Security** → **Remote Authentication**, as shown in Figure 3-48. You may also test the server connection by selecting **Global Actions** → **Test LDAP connections** and verifying that all servers return “CMMVC07051 The LDAP task completed successfully”.

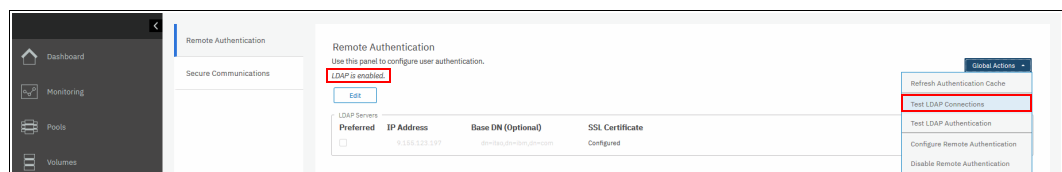


Figure 3-48 Verifying that LDAP is enabled

You can use the **Global Actions** menu to disable remote authentication and switch to local authentication only.

After remote authentication is enabled, the remote user groups must be configured. You can use the default built-in user groups for remote authentication. However, the name of the default user groups cannot be changed. If the LDAP server contains a group that you want to use and you do not want to create this group on the storage system, the name of the group must be changed on the server side to match the default name. Any user group, whether default or self-defined, must be enabled for remote authentication before LDAP authentication can be used for that group.

To create a user group with remote authentication enabled, complete the following steps:

1. Select **Access** → **Users by Group** and click **Create User Group**. Enter the name for the new group, select the **LDAP** check box, and choose a role for the users in the group, as shown in Figure 3-49.

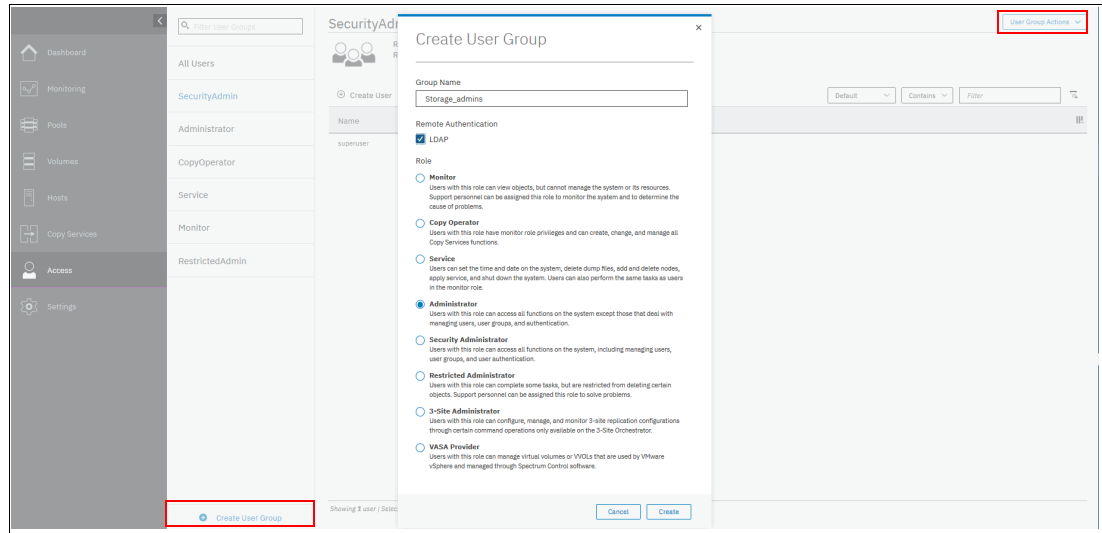


Figure 3-49 Creating a user group with remote authentication enabled

To enable LDAP for one of the existing groups, select it in the list, select **User Group Actions** → **Properties** in the upper right corner, and select the LDAP check box.

2. When you have at least one user group that is enabled for remote authentication, verify that you set up your user group on the LDAP server correctly by checking whether the following conditions are true:
  - The name of the user group on the LDAP server matches the one that you modified or created on the storage system.

**Note:** The user group name is case-sensitive.

- Each user that you want to authenticate remotely is a member of the LDAP user group that is configured for the system role.
3. To test the user authentication, select **Settings** → **Security** → **Remote Authentication**, and then select **Global Actions** → **Test LDAP Authentication** (for an example, see Figure 3-48 on page 117). Enter the user credentials of a user that is defined on the LDAP server and click **Test**. A successful test returns the message “CMMVC70751 The LDAP task completed successfully”.

A user can log in with their short name (that is, without the domain component) or with the fully qualified user name in the UPN format (user@domain).





# IBM Spectrum Virtualize GUI

This chapter describes an overview of the IBM Spectrum Virtualize GUI. The management GUI is a tool that is enabled and provided by IBM Spectrum Virtualize that helps you to monitor, manage, and configure your system.

This chapter explains the basic view and the configuration procedures that are required to get your IBM SAN Volume Controller environment running as quickly as possible by using the GUI.

This chapter explains the basic view and the configuration procedures that are required to get your system environment running as quickly as possible by using the GUI. This chapter does not describe advanced troubleshooting or problem determination and some of the complex operations (compression and encryption).

Throughout this chapter, all GUI menu items are introduced in a systematic, logical order as they appear in the GUI. However, topics that are described more in detail in other chapters of the book are only referred to here. For example, storage pools (Chapter 5, “Storage pools” on page 199), volumes (Chapter 6, “Volumes” on page 255), hosts (Chapter 7, “Hosts” on page 351), and Copy Services (Chapter 10, “Advanced Copy Services” on page 491) are described in separate chapters.

This chapter includes the following topics:

- ▶ 4.1, “Normal operations by using the GUI” on page 120
- ▶ 4.2, “Introduction to the GUI” on page 124
- ▶ 4.3, “System - Overview window” on page 129
- ▶ 4.4, “Monitoring menu” on page 135
- ▶ 4.5, “Pools” on page 142
- ▶ 4.6, “Volumes” on page 142
- ▶ 4.7, “Hosts” on page 143
- ▶ 4.8, “Copy Services” on page 143
- ▶ 4.9, “Access” on page 144
- ▶ 4.10, “Settings” on page 156
- ▶ 4.11, “Additional frequent tasks in the GUI” on page 187

## 4.1 Normal operations by using the GUI

This section describes the graphical icons and the indicators that you use to manage IBM Spectrum Virtualize. For the example in this book, we configure the system in a standard topology.

Multiple users can be logged in to the GUI. However, no locking mechanism exists, so be aware that if two users change the same object simultaneously, the last action that is entered from the GUI is the action that takes effect.

**Important:** Data entries that are made through the GUI are case-sensitive.

You must enable Java Script in your browser. For Mozilla Firefox, JavaScript is enabled by default and requires no other configuration steps. For more information, see [IBM Knowledge Center](#).

### 4.1.1 Accessing the GUI

To access the SAN Volume Controller GUI, enter the Internet Protocol (IP) address that was set during the initial setup process into your web browser. You can connect from any workstation that can communicate with the system. The login window opens (see Figure 4-1).

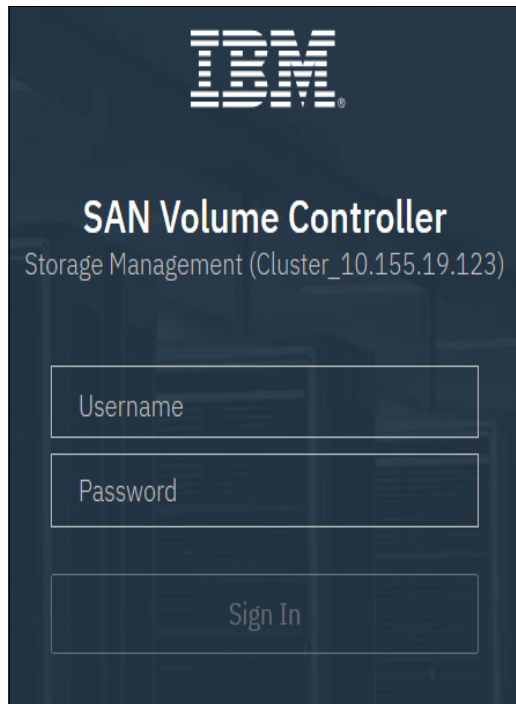


Figure 4-1 Login window of SAN Volume Controller

**Note:** If you log in to the GUI by using the configuration node, you receive another option: Service Assistant Tool (SAT). Clicking this option takes you to the service assistant instead of the cluster GUI, as shown in Figure 4-2.

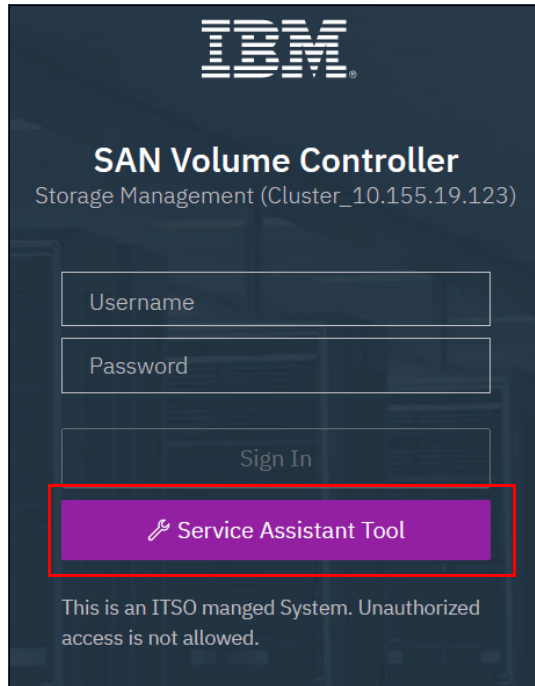


Figure 4-2 Login window of SAN Volume Controller when connected to the configuration node

It is a best practice for each user to have their own unique account. The default user accounts should be disabled for use or their passwords changed and kept secured for emergency purposes only. This approach helps to identify personnel working on the systems and track all important changes that are done by them. The *superuser* account should be used for initial configuration only.

After a successful login, the Version 8.3 Welcome window opens and displays the new system dashboard (see Figure 4-3).

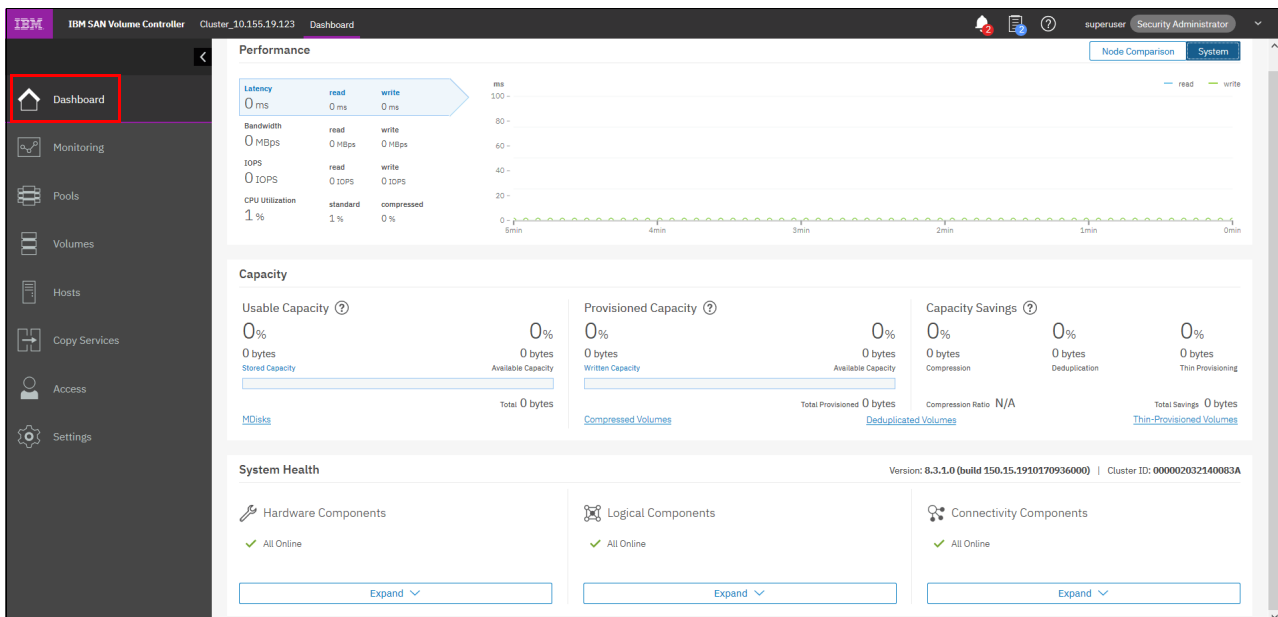


Figure 4-3 Welcome page with new dashboard

The Dashboard is divided into three sections:

- ▶ Performance
- ▶ This section provides important information about latency, bandwidth, input/output operations per second (IOPS), and CPU utilization. All this information can be viewed at either the SAN Volume Controller system level or node level. A “Node comparison” view shows the differences in characteristics of each node (see Figure 4-4). The performance graph is updated with new data every 5 seconds.

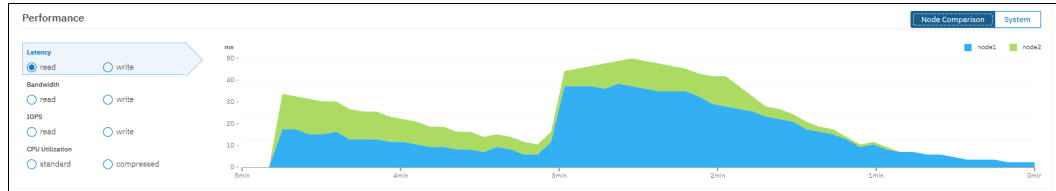


Figure 4-4 Performance statistics

▶ Capacity

This section shows the current utilization of attached storage and its usage. Apart from the usable capacity, it also shows provisioned capacity and capacity savings. You can select the **Compressed Volumes**, **Deduplicated Volumes**, or **Thin Provisioned Volumes** options to display a complete list of the options in the Volumes tab.

New with Version 8.2 is the “Overprovisioned External Storage” section (highlighted in the red box in Figure 4-5), which appears only when attached storage systems that are overprovisioned are present.

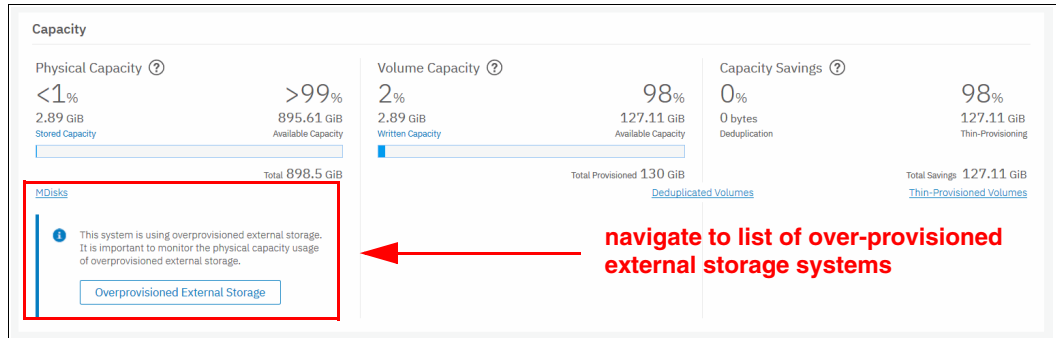


Figure 4-5 Capacity overview

Selecting this option provides a list of Overprovisioned External Systems that you can then click to see a list of related managed disks (MDisks) and pools, as shown in Figure 4-6 on page 123.

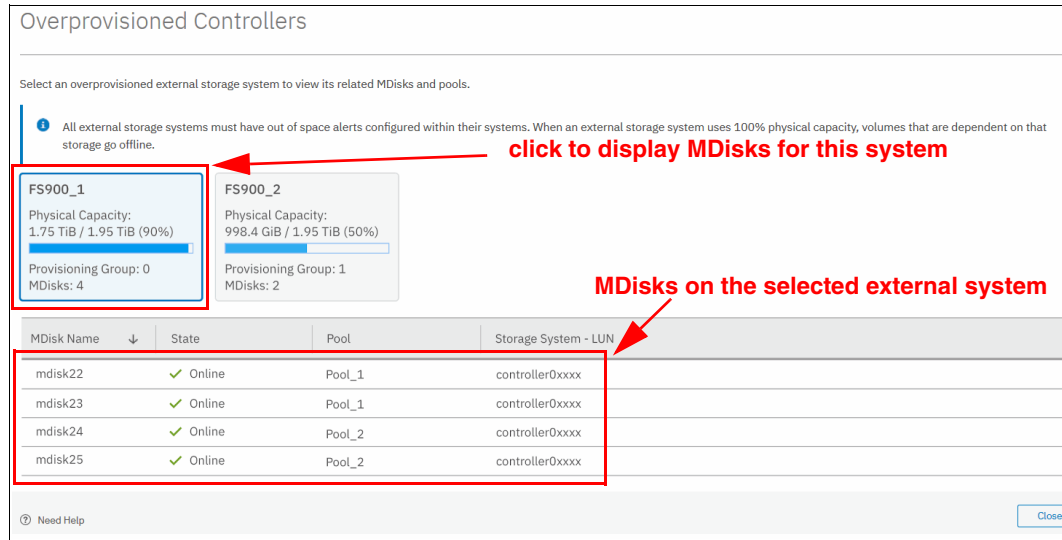


Figure 4-6 List showing overprovisioned external storage

You also see a warning when assigning MDisks to pools if the MDisk is on an overprovisioned external storage controller.

► System Health

This section indicates the status of all critical system components, which are grouped in three categories: Hardware, logical, and connectivity components, as shown in Figure 4-7. When you click **Expand**, each component is listed as a subgroup. You can then go directly to the section of GUI where the component that you are interested in is managed.

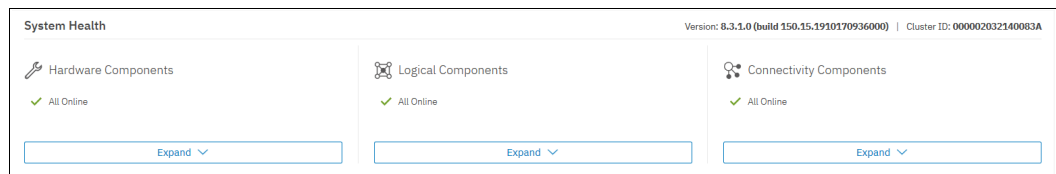


Figure 4-7 System health overview

Figure 4-8 shows the expanded view.

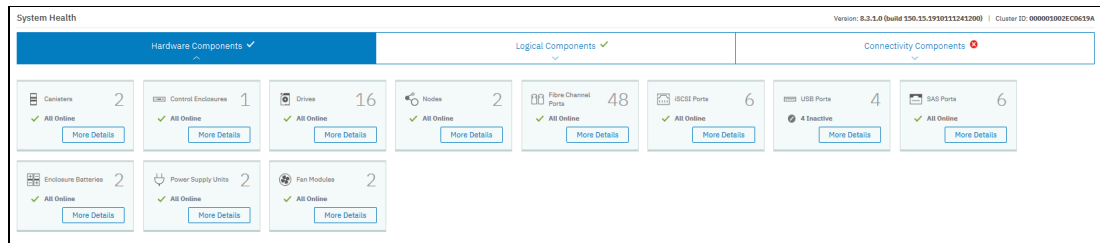


Figure 4-8 Expanded System health view

The dashboard in Version 8.3.1 displays as a welcome page instead of the system pane as in previous versions. This system overview was moved to the **Monitoring** → **System** menu.

Although the Dashboard pane provides key information about system behavior, the System menu is a preferred starting point to obtain the necessary details about your system components.

## 4.2 Introduction to the GUI

As shown in Figure 4-9, the former SAN Volume Controller GUI System pane was relocated to **Monitoring** → **System**.

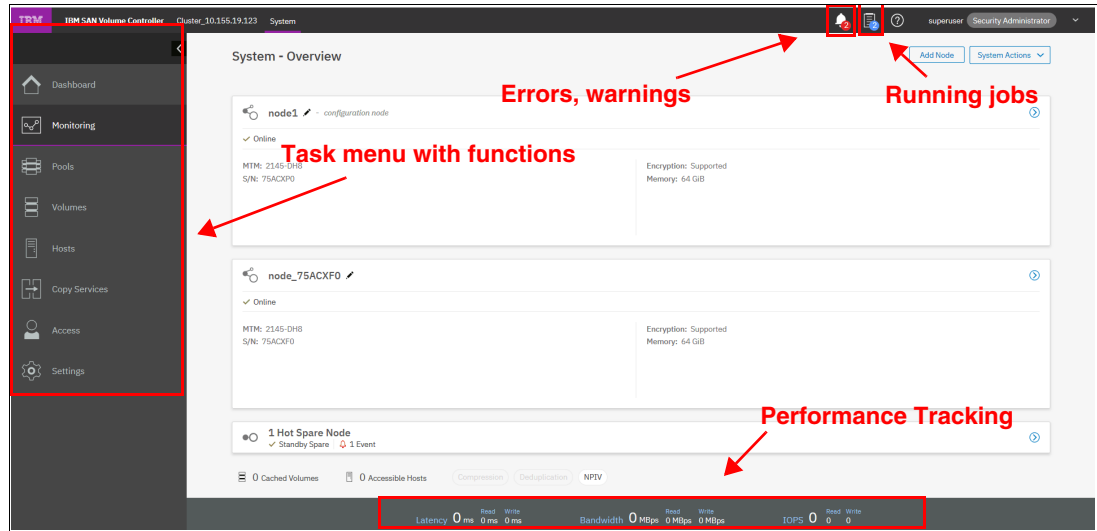


Figure 4-9 IBM SAN Volume Controller System window

### 4.2.1 Task menu

The IBM Spectrum Virtualize GUI task menu is always available on the left side of the GUI window. To browse by using this menu, click the action and choose a task that you want to display, as shown in Figure 4-10.

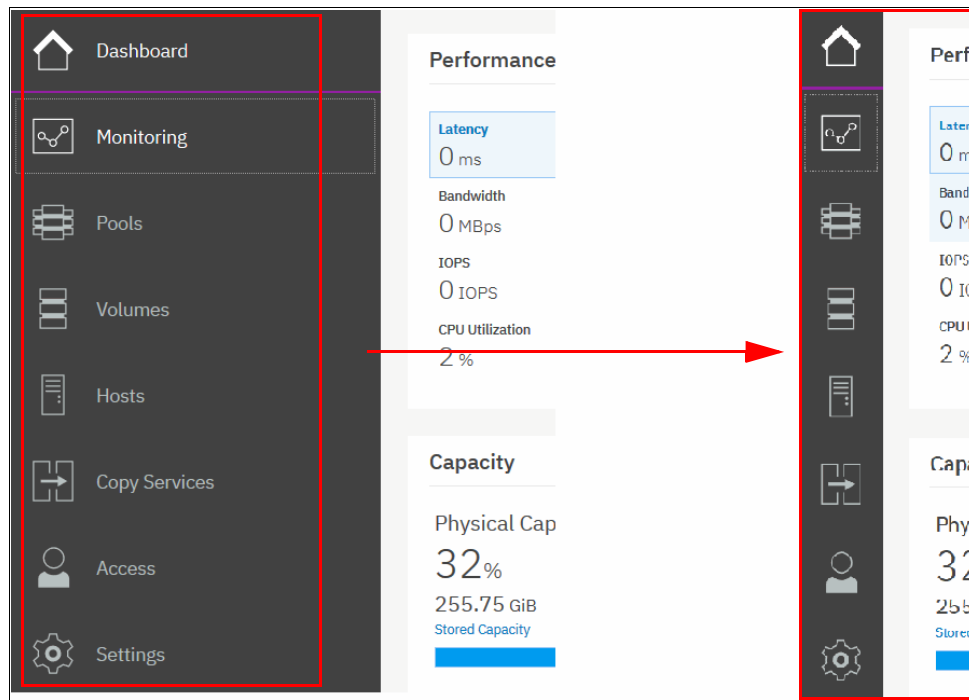


Figure 4-10 The task menu on the left side of the GUI

By reducing the horizontal size of your browser window, the wide task menu shrinks to the icons only.

## 4.2.2 Suggested tasks

After the initial configuration process is complete, IBM Spectrum Virtualize shows the information about suggested tasks that notify the administrator that several key functions are not yet configured. If necessary, this indicator can be closed and these tasks can be performed at any time. Figure 4-11 shows the suggested tasks in the System pane.

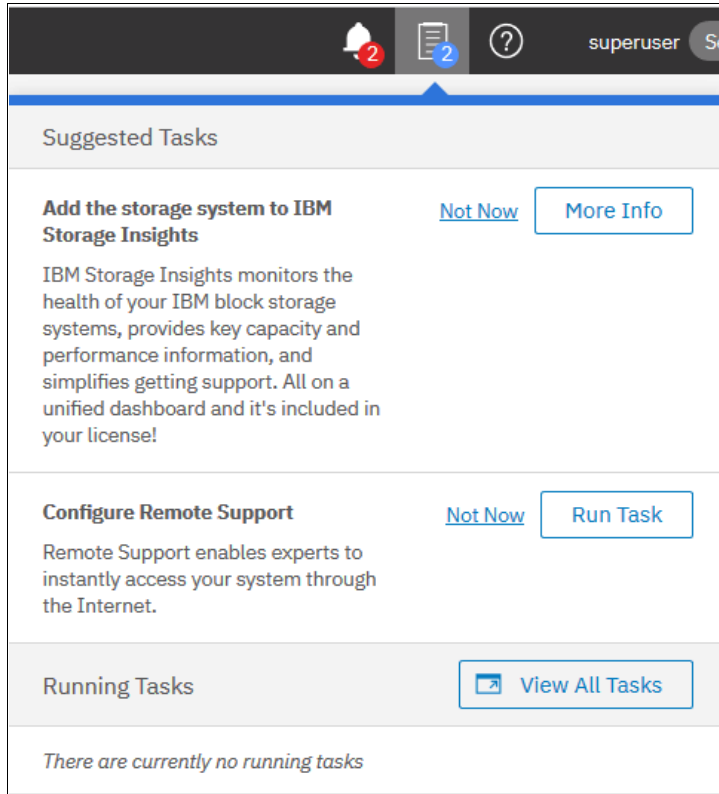


Figure 4-11 Suggested tasks

In this case, the GUI has two suggested tasks that help with the general administration of the system: You can directly perform the tasks from this window, or cancel them and run the procedure later. Other suggested tasks that typically appear after the initial system configuration are to create a volume and configure a storage pool.

The dynamic IBM Spectrum Virtualize menu contains the following panes:

- ▶ Dashboard
- ▶ Monitoring
- ▶ Pools
- ▶ Volumes
- ▶ Hosts
- ▶ Copy Services
- ▶ Access
- ▶ Settings

## 4.2.3 Notification icons and help

Three notification icons are in the upper navigation area of the GUI (see Figure 4-12). The left icon indicates warning and error alerts that were recorded in the event log. The middle icon shows running jobs and suggested tasks. The third rightmost icon offers a help menu with content that is associated with the current tasks and the currently opened GUI menu.



Figure 4-12 Notification area

### Alerts indication

The left icon in the notification area informs administrators about important alerts in the systems. Click the icon to list warning messages in yellow and errors in red (see Figure 4-13).

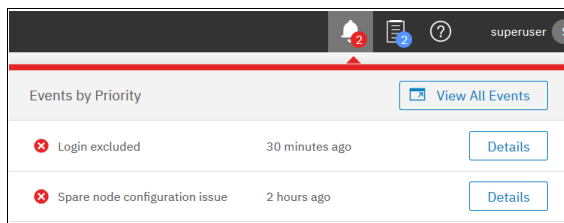


Figure 4-13 System alerts

You can go directly to the Events menu by clicking the **View All Events** option or see each event message separately by clicking the **Details** icon of the specific message, analyzing the content, and eventually running the suggested fix procedure, as shown in Figure 4-14.

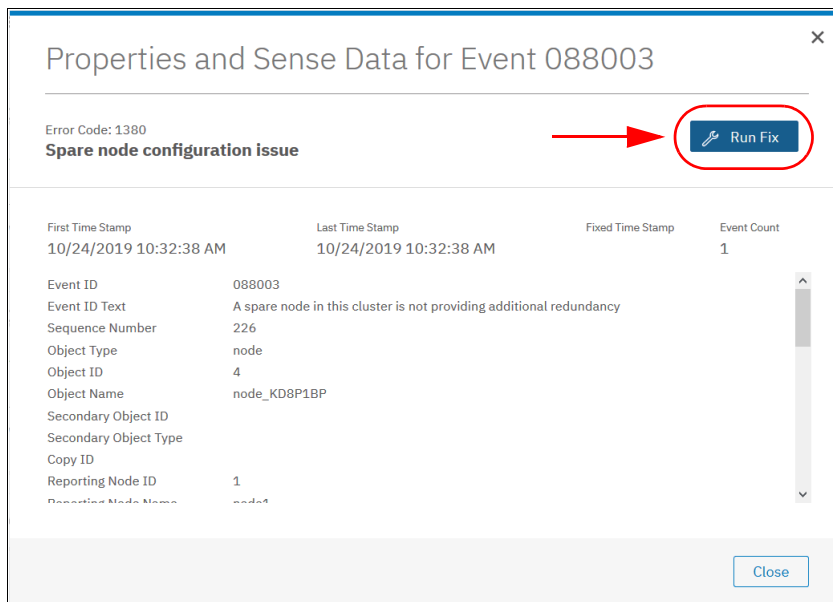


Figure 4-14 Spare node configuration issue

### Running jobs and suggested tasks

The middle icon in the notification area provides an overview of currently running tasks that are triggered by administrator. It also includes the suggested tasks that recommend that users perform specific configuration actions.



In the example that is shown in Figure 4-15, we have not yet defined remote support in the system. Therefore, the system suggests that we do so and offers us direct access to the associated **Remote Support** menu. Click **Run Task** to define remote support according to the procedure that is explained in Chapter 13, “Reliability, availability, and serviceability, and monitoring and troubleshooting” on page 753. If you do not want to define remote support now, click **Not Now** and the suggestion message disappears.

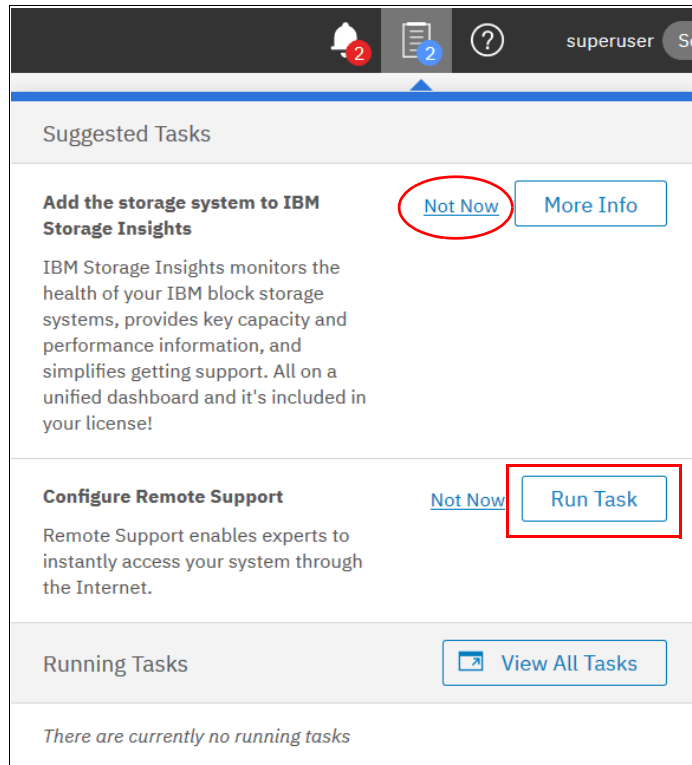


Figure 4-15 Suggested Tasks

Similarly, you can analyze the details of running tasks (all of them together in one window or of a single task). Click **View** to open the Array Initialization job, as shown in Figure 4-16.

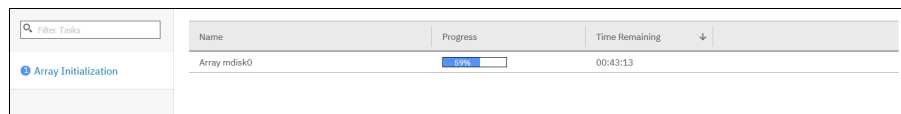


Figure 4-16 Details of running task

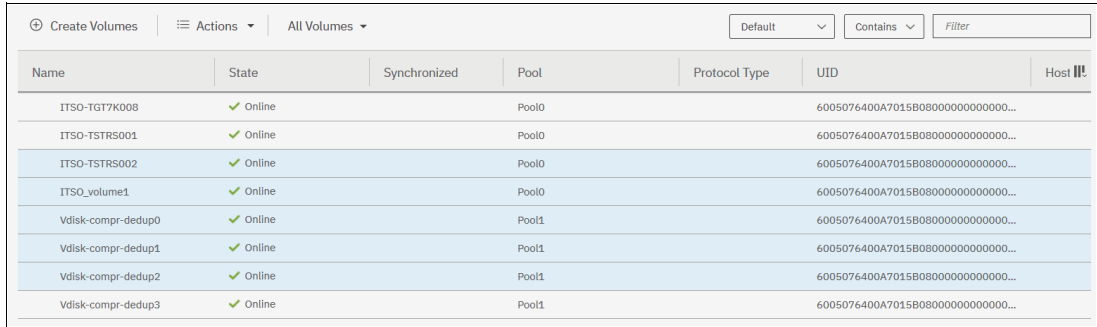
The following information can be displayed as part of the running tasks:

- ▶ Volume migration
- ▶ MDisk removal
- ▶ Image mode migration
- ▶ Extent migration
- ▶ IBM FlashCopy
- ▶ Metro Mirror (MM) and Global Mirror (GM)
- ▶ Volume formatting
- ▶ Space-efficient copy repair
- ▶ Volume copy verification and synchronization
- ▶ Estimated time for the task completion

## Making selections

Recent updates to the GUI brought improved selection making. You can now select multiple items more easily. Go to a wanted window, press and hold the Shift or Ctrl key, and make your selection.

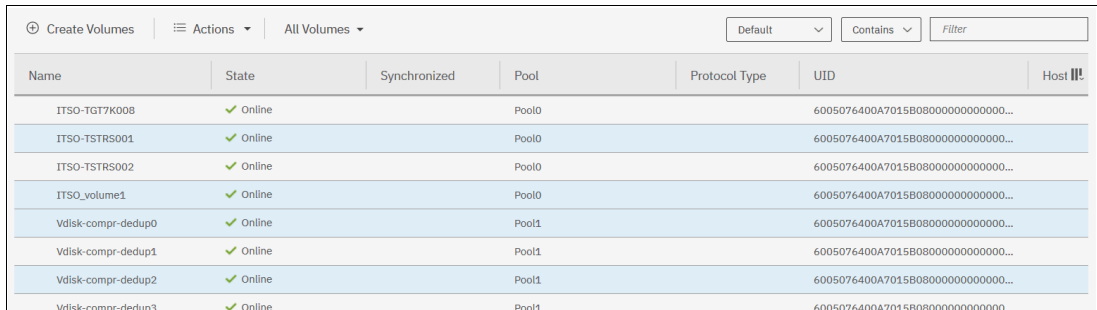
Pressing and holding the Shift key, select the first item in your list that you want, and then select the last item. All items between the two that you choose are also selected, as shown in Figure 4-17.



Name	State	Synchronized	Pool	Protocol Type	UID	Host
ITSO-TG17K008	Online		Pool0		6005076400A7015B08000000000000...	
ITSO-TSTRS001	Online		Pool0		6005076400A7015B08000000000000...	
ITSO-TSTRS002	Online		Pool0		6005076400A7015B08000000000000...	
ITSO_volume1	Online		Pool0		6005076400A7015B08000000000000...	
Vdisk-compr-dedup0	Online		Pool1		6005076400A7015B08000000000000...	
Vdisk-compr-dedup1	Online		Pool1		6005076400A7015B08000000000000...	
Vdisk-compr-dedup2	Online		Pool1		6005076400A7015B08000000000000...	
Vdisk-compr-dedup3	Online		Pool1		6005076400A7015B08000000000000...	

Figure 4-17 Selecting items by using the Shift key

Pressing and holding the Ctrl key, select any items from the entire list. You can select items that do not appear in sequential order, as shown in Figure 4-18.



Name	State	Synchronized	Pool	Protocol Type	UID	Host
ITSO-TG17K008	Online		Pool0		6005076400A7015B08000000000000...	
ITSO-TSTRS001	Online		Pool0		6005076400A7015B08000000000000...	
ITSO-TSTRS002	Online		Pool0		6005076400A7015B08000000000000...	
ITSO_volume1	Online		Pool0		6005076400A7015B08000000000000...	
Vdisk-compr-dedup0	Online		Pool1		6005076400A7015B08000000000000...	
Vdisk-compr-dedup1	Online		Pool1		6005076400A7015B08000000000000...	
Vdisk-compr-dedup2	Online		Pool1		6005076400A7015B08000000000000...	
Vdisk-compr-dedup3	Online		Pool1		6005076400A7015B08000000000000...	

Figure 4-18 Selecting items by using the Ctrl key

You can also select items by using the built-in filtering function. For more information, see 4.3.1, “Content-based organization” on page 130.

## Help

If you need help, you can select the (?) button, as shown in Figure 4-19 on page 129. You see two options:

- ▶ The first option opens a new tab with plain text information about the pane you are on and its contents.
- ▶ The second option shows the same information in IBM Knowledge Center. This option requires an internet connection, but the first option does not because the information is stored locally on the SAN Volume Controller.

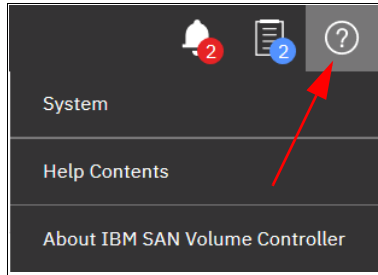


Figure 4-19 Access help menu

For example, in the System pane, you can open help that is related to the system in general, as shown in Figure 4-20.

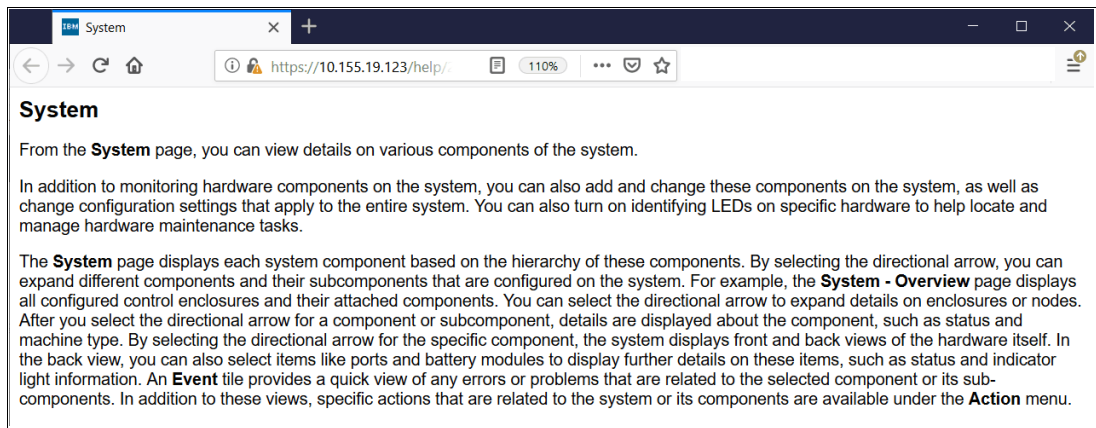


Figure 4-20 Example of System help content

## 4.3 System - Overview window

The System - Overview window is shown in Figure 4-21.

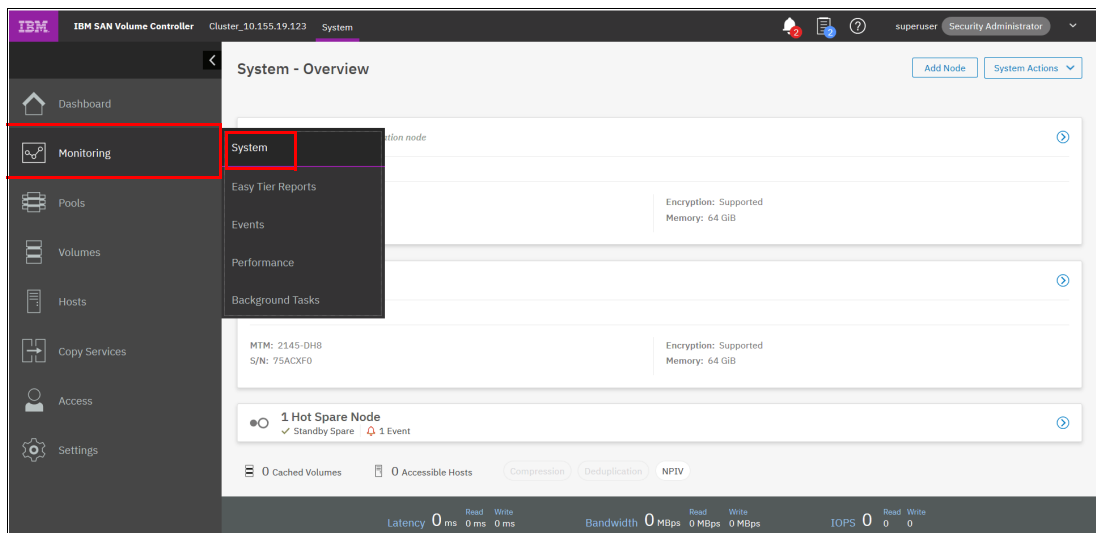


Figure 4-21 The System - Overview window

The next section describes the structure of the window and how to go to various system components to manage them more efficiently and quickly.

### 4.3.1 Content-based organization

The following sections describe several view options within the GUI in which you can filter (to minimize the amount of data that is shown on the window), sort, and reorganize the content of the window.

#### Table filtering

On most pages, a Filter box is available at the upper right side of the window. Use this option if the list of object entries is too long and you want to search for something specific.

To use search filtering, complete the following steps:

1. In the **Filter** box that is shown in Figure 4-22, enter a search term by which you want to filter. You can also use the drop-down menus to modify what the system searches for. For example, if you want an exact match to your filter, select **=** instead of **Contains**. The first drop-down list limits your filter to search through a specific column only, for example, Name and State.

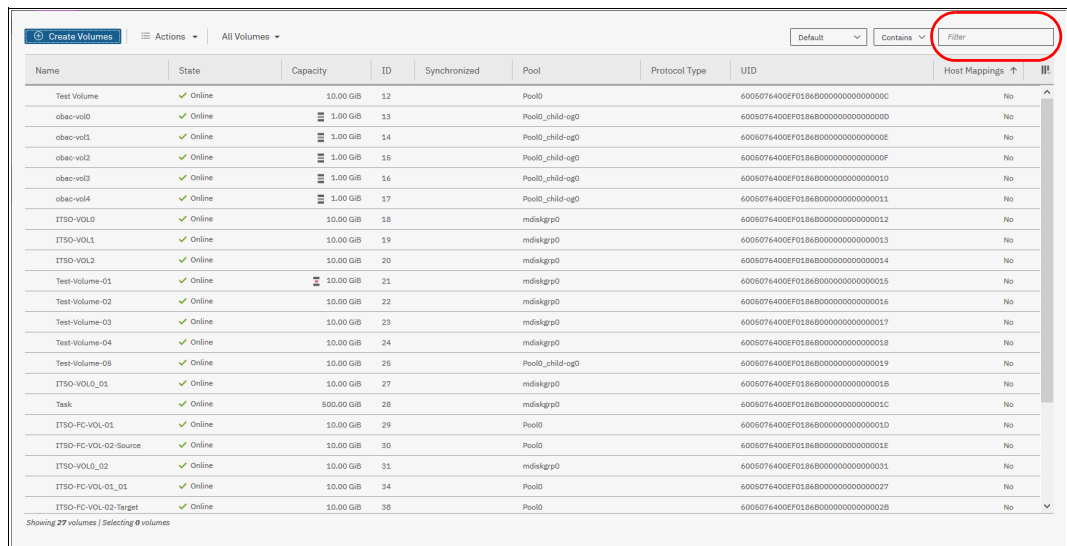


Figure 4-22 Filter search box

2. Enter the text string that you want to filter and press Enter.

By using this function, you can filter your table that is based on column names. In our example, a volume list is displayed that contains the names that include ITSO somewhere in the name. ITSO is highlighted in amber, as are any columns that contain this information, as shown in Figure 4-23 on page 131. The search option is not case-sensitive.

Create Volumes | Actions | All Volumes | Default | Contains | **ITSO**

Name	State	Capacity	ID	Synchronized	Pool	Protocol Type	UID	Host Mappings
ITSB-VOL0	✓ Online	10.00 GiB	18		mdiskgrp0		6005076400EF01868000000000000012	No
ITSB-VOL1	✓ Online	10.00 GiB	19		mdiskgrp0		6005076400EF01868000000000000013	No
ITSB-VOL2	✓ Online	10.00 GiB	20		mdiskgrp0		6005076400EF01868000000000000014	No
ITSB-VOL0_01	✓ Online	10.00 GiB	27		mdiskgrp0		6005076400EF01868000000000000018	No
ITSB-FC-VOL-01	✓ Online	10.00 GiB	29		Pool0		6005076400EF0186800000000000001D	No
ITSB-FC-VOL-02-Source	✓ Online	10.00 GiB	30		Pool0		6005076400EF0186800000000000001E	No
ITSB-VOL0_02	✓ Online	10.00 GiB	31		mdiskgrp0		6005076400EF01868000000000000031	No
ITSB-FC-VOL-01_01	✓ Online	10.00 GiB	34		Pool0		6005076400EF01868000000000000027	No
ITSB-FC-VOL-02-Target	✓ Online	10.00 GiB	38		Pool0		6005076400EF0186800000000000002B	No
ITSB-FC-VOL-01_02	✓ Online	10.00 GiB	39		Pool0		6005076400EF0186800000000000002C	No
ITSB-FC-VOL-01_03	✓ Online	10.00 GiB	40		Pool0		6005076400EF0186800000000000002D	No
ITSB-FC-VOL-01_04	✓ Online	10.00 GiB	41		Pool0		6005076400EF0186800000000000002E	No
ITSB-FC-VOL-01_05	✓ Online	10.00 GiB	42		Pool0		6005076400EF0186800000000000002F	No

Figure 4-23 Show filtered rows

- Remove this filtered view by clicking the **X** icon that displays in the Filter box or by deleting what you searched for and pressing Enter.

**Filtering:** This filtering option is available in most menu options of the GUI.

## Table information

In the table view, you can add or remove the information in the tables on most pages.

For example, on the Volumes pane, complete the following steps to add a column to the table:

1. Right-click any column headers of the table or select the icon in the left corner of the table header. A list of all of the available columns displays, as shown in Figure 4-24.

The screenshot displays the 'Volumes' pane in IBM Spectrum Virtualize. At the top, there are navigation options: 'Create Volumes', 'Actions', and 'All Volumes'. Below this is a table with columns: Name, State, Capacity, ID, and Synchronized. The 'Name' column header is highlighted with a red box, and a red arrow points to it with the text 'right click'. A context menu is open over the table, listing various columns that can be added or removed. The 'Real Capacity' column is currently selected in the menu. The table contains 30 rows of volume information, all with a state of 'Online'. At the bottom of the table, it says 'Showing 30 volumes / Selecting 0 volumes'.

Name	State	Capacity	ID	Synchronized
Test-Volume-01	Online			
Task	Online			
Linux Test	Online			
ITSO-FC-VOL-01_04	Online			
Test Volume	Online			
ITSO-VOL0	Online			
ITSO-VOL1	Online			
ITSO-VOL2	Online			
Test-Volume-02	Online			
Test-Volume-03	Online			
Test-Volume-04	Online			
Test-Volume-05	Online			
Test-Volume-06	Online			
ITSO-VOL0_01	Online			
ITSO-FC-VOL-01	Online			
ITSO-FC-VOL-02-Source	Online			
ITSO-VOL0_02	Online			
ITSO-FC-VOL-01_02_01	Online			
ITSO-FC-VOL-01_01	Online			
Redbook_Test_1	Online			
ITSO-FC-VOL-02-Target	Online			

Showing 30 volumes / Selecting 0 volumes

- Name
- Real Capacity
- State
- Capacity
- ID
- Ownership Group
- Synchronized
- Pool
- Volume Group
- Last Snapshot
- Protocol Type
- UID
- Host Mappings
- Preferred Node ID
- Used Capacity
- Cache State
- Compression Savings
- Capacity Savings
- Estimated Compression Savings
- Estimated Compression Savings %
- Estimated Thin Savings
- Estimated Thin Savings %

Figure 4-24 Add or remove details in a table

- Select the column that you want to add or remove from this table. In our example, we added the volume Real Capacity column, as shown in Figure 4-25.

<span>⊕ Create Volumes</span>   <span>☰ Actions</span> ▾   <span>All Volumes</span> ▾					
Name	↓	Real Capacity	State	Capacity	ID
obac-vol4		☰ 36.48 MiB	✓ Online	☰ 1.00 GiB	17
obac-vol3		☰ 36.48 MiB	✓ Online	☰ 1.00 GiB	16
obac-vol2		☰ 36.48 MiB	✓ Online	☰ 1.00 GiB	15
obac-vol1		☰ 36.48 MiB	✓ Online	☰ 1.00 GiB	14
obac-vol0		☰ 51.73 MiB	✓ Online	☰ 1.00 GiB	13
Test-Volume-06		10.00 GiB	✓ Online	10.00 GiB	26
Test-Volume-05		10.00 GiB	✓ Online	10.00 GiB	25
Test-Volume-04		10.00 GiB	✓ Online	10.00 GiB	24
Test-Volume-03		10.00 GiB	✓ Online	10.00 GiB	23
Test-Volume-02		10.00 GiB	✓ Online	10.00 GiB	22
Test-Volume-01		Not Applicable	✓ Online	☰ 10.00 GiB	21

Figure 4-25 Table with an added Real Capacity column

- You can repeat this process several times to create custom tables to meet your requirements.

- Return to the default table view by selecting **Restore Default View** (the last entry) in the column selection menu, as shown in Figure 4-26.

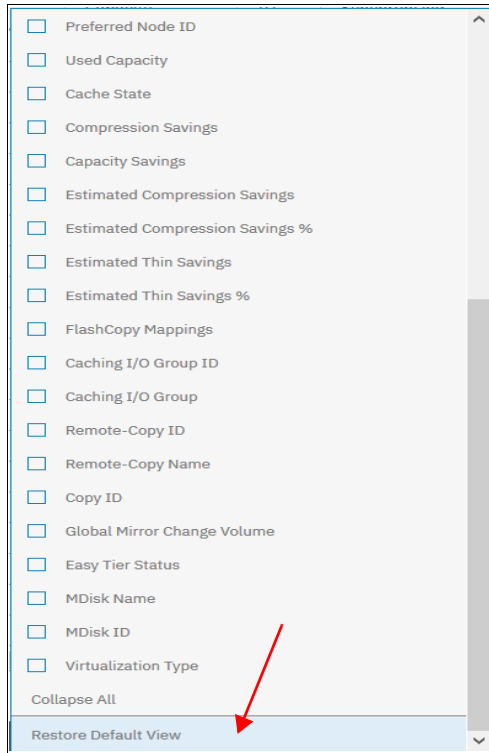


Figure 4-26 Table Restore Default View

**Sorting:** By clicking a column, you can sort a table based on that column in ascending or descending order.

### Shifting columns in tables

You can move columns by clicking and moving the column right or left, as shown in Figure 4-27. In this example, we attempt to move the Capacity column before the State column.

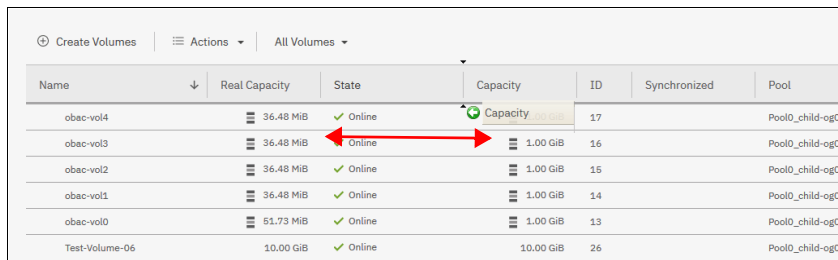


Figure 4-27 Reorganizing table columns



## 4.4 Monitoring menu

Click the **Monitoring** icon in left pane to open the **Monitoring** menu (Figure 4-28). The **Monitoring** menu offers these navigation options:

- ▶ **System:** This option opens an overview of the system. It shows all control enclosures and groups them into I/O groups if more than one control enclosure is present. Useful information about the nodes is displayed, including node status, number of events against each node, and key node information, such as cache size and serial numbers. For more information, see 4.4.1, “System overview” on page 136.
- ▶ **Events:** This option tracks all informational, warning, and error messages that occurred in the system. You can apply various filters to sort the messages according to your needs or export the messages to an external comma-separated value (CSV) file. For more information, see 4.4.3, “Events” on page 139.
- ▶ **Performance:** This option reports the general system statistics that relate to the processor (CPU) utilization, host and internal interfaces, volumes, and MDisks. With this option, you can switch between megabytes per second (MBps) or IOPS. For more information, see Figure 4.4.4 on page 140.
- ▶ **Background Tasks:** The option shows the progress of all tasks running in the background as listed in “Running jobs and suggested tasks” on page 126.

The following sections describe each option of the **Monitoring** menu (Figure 4-28).

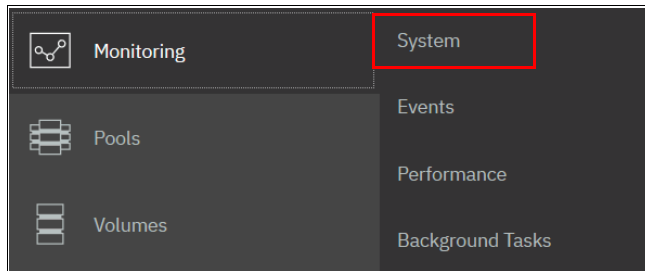


Figure 4-28 Monitoring menu

## 4.4.1 System overview

The **System** option on the **Monitoring** menu provides a general overview of the SAN Volume Controller and key information. If you have more than one I/O group, this view is presented by I/O group, with nodes displayed within their own I/O group section, as shown in Figure 4-29.

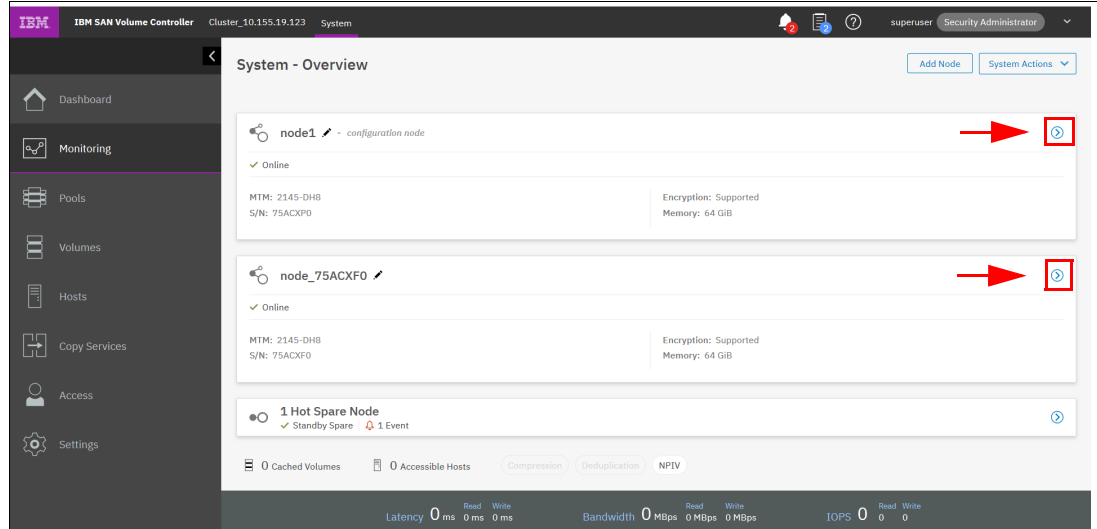


Figure 4-29 System overview window showing both nodes with the spare node

You can select a component to view more details about it by using the arrow that is highlighted in the red box in Figure 4-29. Clicking this arrow shows detailed technical attributes, such as ports that are in use, memory, serial number, node name, encryption status, and node status (online or offline). Selecting an individual component on the image of your SAN Volume Controller displays its details in the Component Details area on the right side, as shown in Figure 4-30.

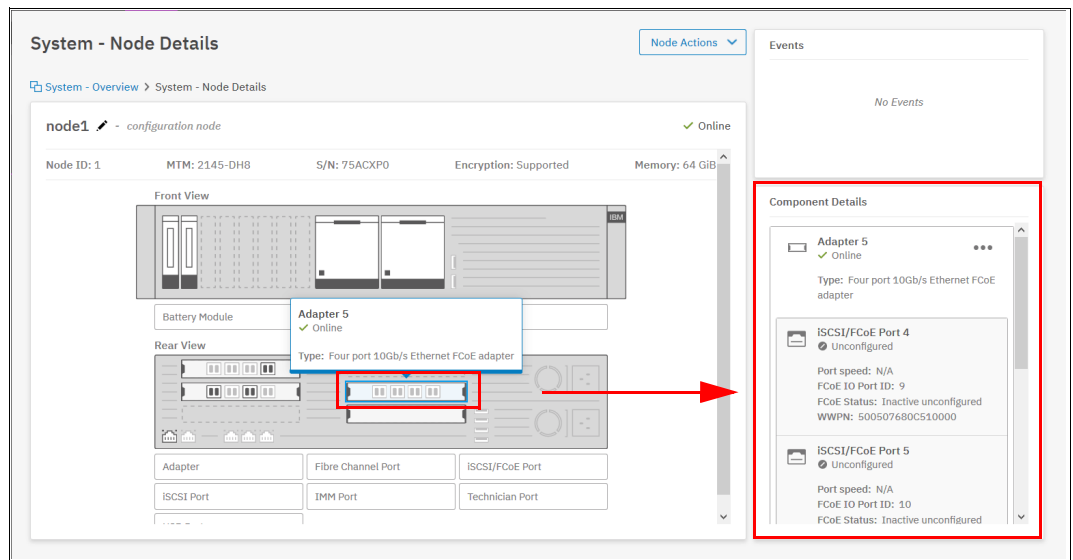


Figure 4-30 Component details

In an environment with multiple IBM Storage System clusters, you can easily direct the onsite personnel or technician to the correct device by enabling the identification LED on the front pane by completing the following steps:

1. Click **Turn Identify On** in the menu that is shown in Figure 4-31.

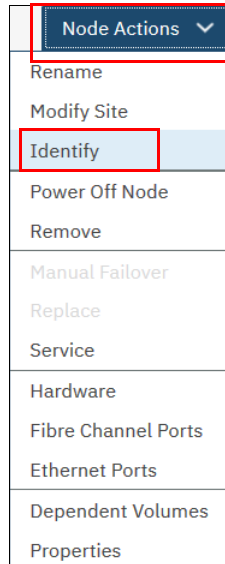


Figure 4-31 Turning on the identification LED

2. Wait for confirmation from the technician that the device in the data center was correctly identified.
3. After the confirmation, click **Turn LED Off** (Figure 4-32).

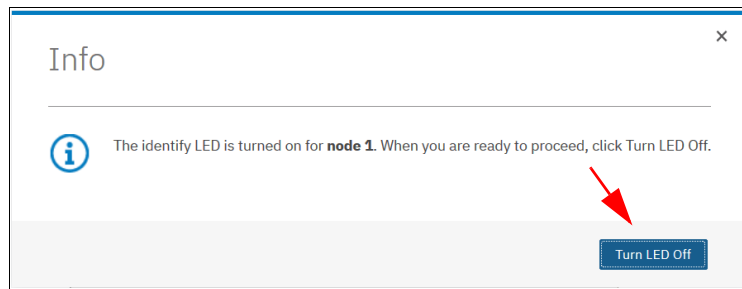


Figure 4-32 Turning off the identification LED

Alternatively, you can use the SAN Volume Controller command-line interface (CLI) to get the same results. Type the following commands in this sequence:

1. Type `svctask chnode -identify yes 1` (or `chnode -identify yes 1`).
2. Type `svctask chnode -identify no 1` (or `chnode -identify no 1`).

## 4.4.2 IBM Easy Tier Reports

The management GUI supports monitoring Easy Tier data movement in graphical reports to help you understand what is happening with the performance of your storage device. Charts for data movement, tier composition, and workload skew comparison can be viewed as web-generated HTML files in a browser, or can be downloaded as CSV files.

Data is collected by the IBM Storage Tier Advisor Tool (IBM STAT) tool in 5-minute increments. When data that is displayed in increments that are larger than 5 minutes (for example, 1 hour), the data that is displayed for that 1 hour is the sum of all the data points that were received for that 1-hour time span.

To view Easy Tier data and reports in the management GUI, select one of the following paths:

- ▶ From the management GUI, select **Monitoring** → **Easy Tier Reports**.
- ▶ From the management GUI, select **Pools** → **View Easy Tier Reports**.

### Data Movement statistics

The Data Movement chart displays the migration actions that are triggered by Easy Tier.

### Tier composition statistics

The Tier Composition chart displays the distributed workload between the top tier, middle tier, and bottom tier. Each tier is composed of one or more tier types.

### Workload Skew Comparison

The Workload Skew Comparison chart displays the percentage of IO workload compared to the total capacity.

See Figure 4-33 shows you how to export the Easy Tier Reports.

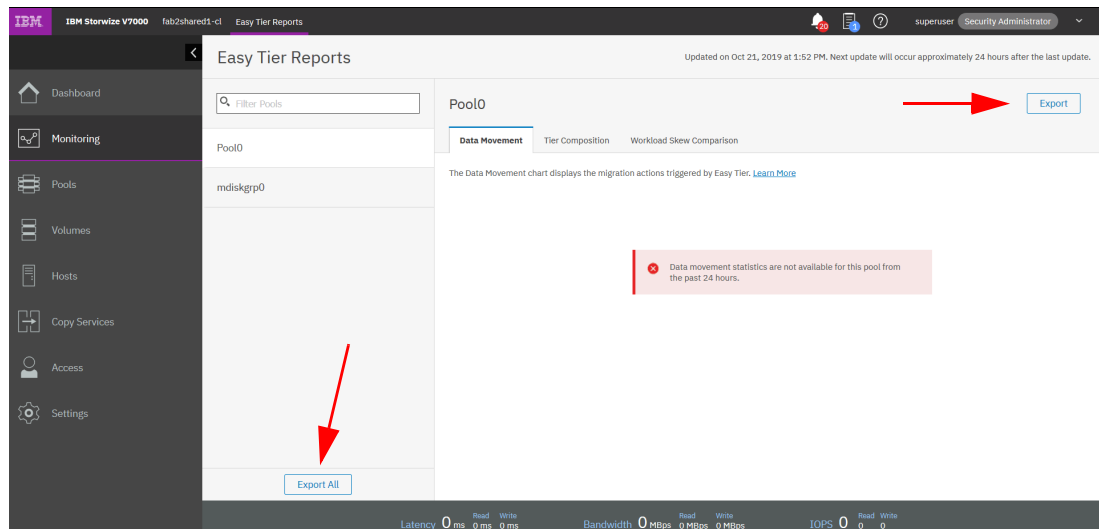


Figure 4-33 Easy Tier Reports

You can export your Easy Tier statistics to a CSV file for further analysis. For more information about Easy Tier Reports, see Chapter 9, “Advanced features for storage efficiency” on page 449.

### 4.4.3 Events

The **Events** option, which is available in the **Monitoring** menu, tracks all informational, warning, and error messages that occur in the system. You can apply various filters to sort them, or export them to an external CSV file. A CSV file can be created from the information that is shown here. Figure 4-34 provides an example of records in the system Event log.

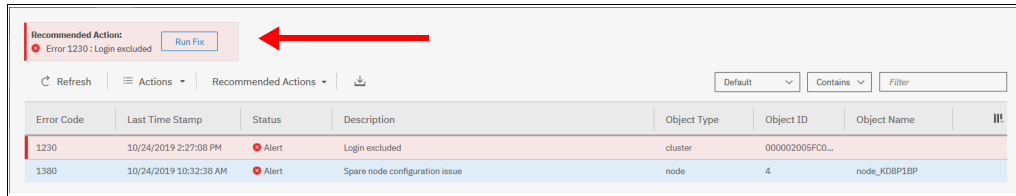


Figure 4-34 Event log list

For the error messages with the highest internal priority, perform corrective actions by running fix procedures. Click **Run Fix** (see Figure 4-34), and the fix procedure wizard opens, as shown in Figure 4-35.

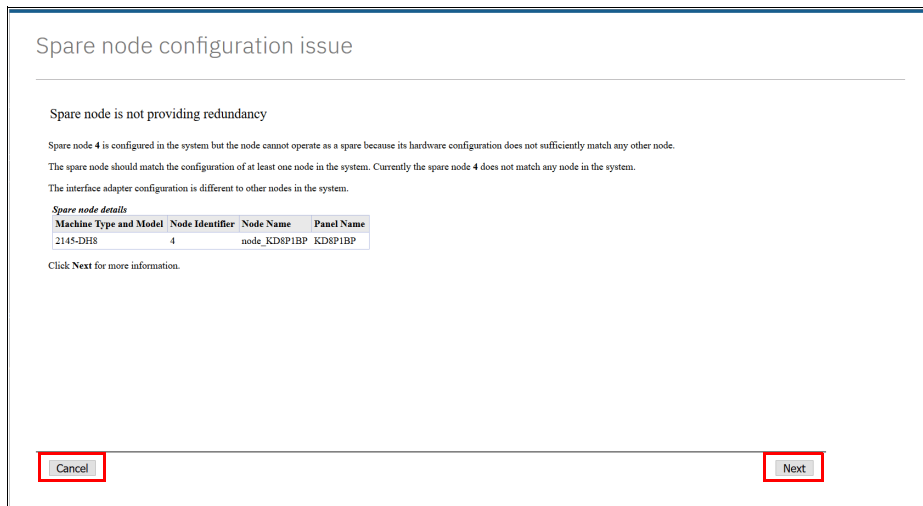


Figure 4-35 Performing a fix procedure

The wizard guides you through the troubleshooting and fixing process from a hardware or software perspective. If you determine that the problem cannot be fixed without a technician's intervention, you can cancel the procedure execution at any time.

For more information about fix procedures, see Chapter 13, "Reliability, availability, and serviceability, and monitoring and troubleshooting" on page 753.

## 4.4.4 Performance

The Performance pane reports the general system statistics that relate to processor (CPU) utilization, host and internal interfaces, volumes, and MDisks. You can switch between MBps or IOPS, and drill down in the statistics to the node level. This capability might be useful when you compare the performance of each node in the system if problems exist after a node failover occurs (see Figure 4-36).

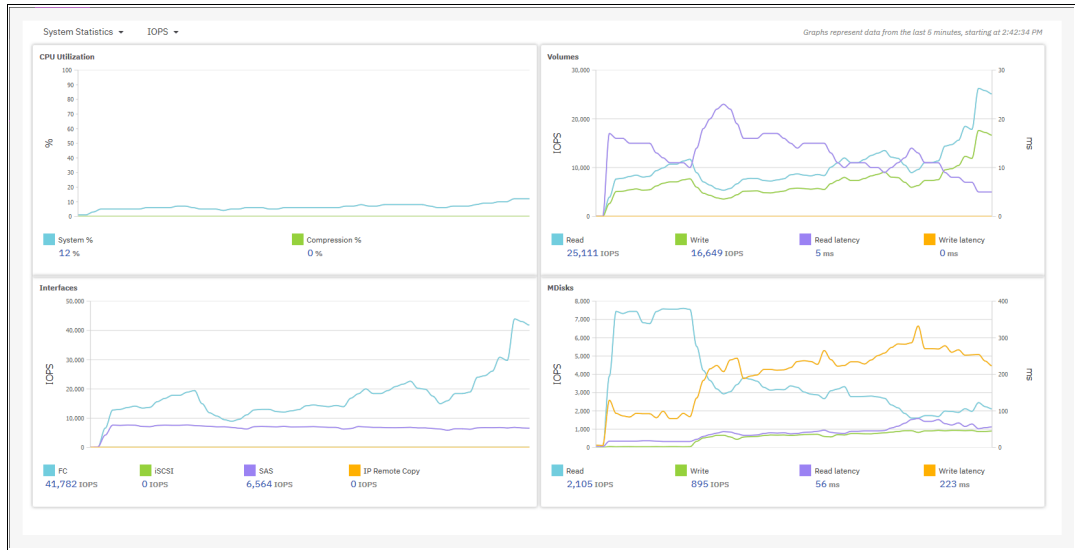


Figure 4-36 Performance statistics of the SAN Volume Controller

The performance statistics in the GUI show, by default, the latest 5 minutes of data. To see details of each sample, click the graph and select the time stamp, as shown in Figure 4-37.

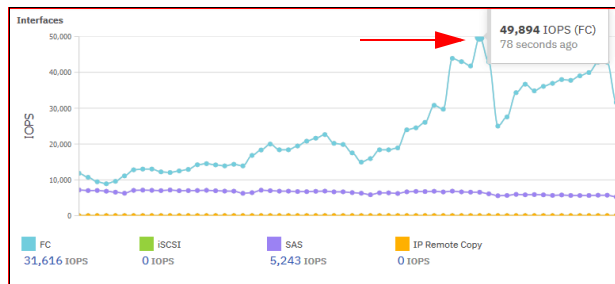


Figure 4-37 Sample details

The charts that are shown in Figure 4-37 represent 5 minutes of the data stream. For in-depth storage monitoring and performance statistics with historical data about your SAN Volume Controller system, use IBM Spectrum Control (enabled by the former IBM Tivoli Storage Productivity Center for Disk and IBM Virtual Storage Center) or IBM Storage Insights.

You can also obtain a no-charge unsupported version of the Quick Performance Overview (**qperf**) for the SAN Volume Controller and Storwize systems from [this website](#).

## 4.4.5 Background Tasks

Use the Background Tasks pane to view and manage current tasks that are running on the system (see Figure 4-38).

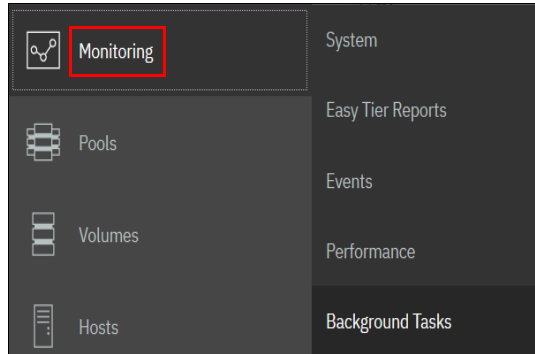
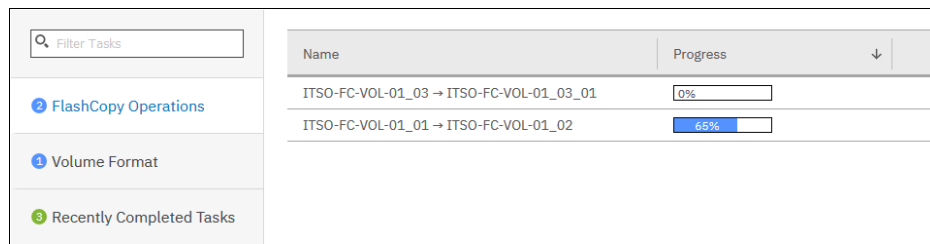


Figure 4-38 Selecting Background Tasks

This menu provides an overview of currently running tasks that are triggered by the administrator. In contrast to the Running jobs and Suggested tasks indication in the middle of top pane, it does not list the suggested tasks that administrators should consider performing. The overview provides more details than the indicator, as shown in Figure 4-39.



Name	Progress
ITSO-FC-VOL-01_03 → ITSO-FC-VOL-01_03_01	0%
ITSO-FC-VOL-01_01 → ITSO-FC-VOL-01_02	65%

Figure 4-39 List of running tasks

You can switch between each type (group) of operation, but you cannot show them all in one list (see Figure 4-40).



Name	Progress	Time Remaining
ITSO-FC-VOL-01_02_01	13%	01:14:10

Figure 4-40 Switching between types of background tasks

## 4.5 Pools

The **Pools** menu option is used to configure and manage storage pools, internal, and external storage, MDisks, and to migrate old attached storage to the system.

The **Pools** menu contains the following items accessible from GUI (see Figure 4-41):

- ▶ Pools
- ▶ Volumes by Pool
- ▶ Internal Storage
- ▶ External Storage
- ▶ MDisks by Pool
- ▶ System Migration

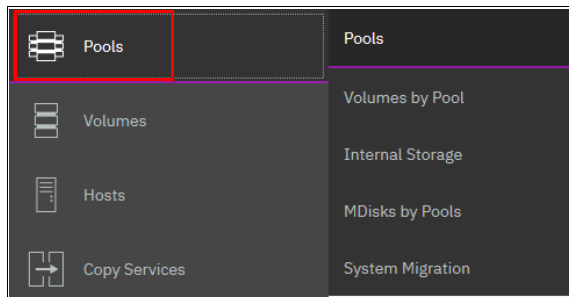


Figure 4-41 Pools menu

For more information about storage pool configuration and management, see Chapter 6, “Storage pools” on page 197.

## 4.6 Volumes

A *volume* is a logical disk that the system presents to attached hosts. By using GUI operations, you can create different types of volumes depending on the type of topology that is configured on your system.

The **Volumes** menu contains the following items, as shown in Figure 4-42 on page 142:

- ▶ Volumes
- ▶ Volumes by Pool
- ▶ Volumes by Host
- ▶ Volumes by Host Cluster
- ▶ Cloud Volumes

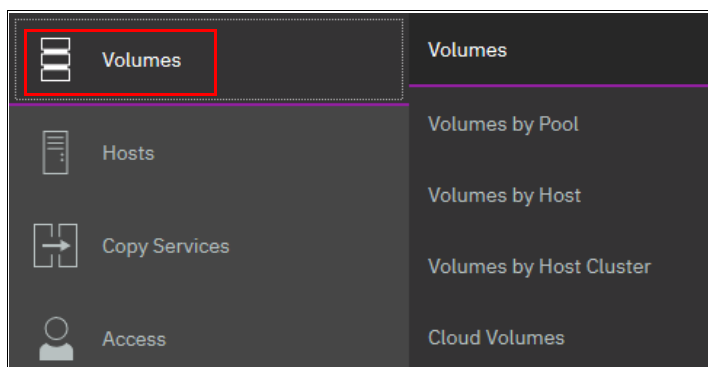


Figure 4-42 Volumes menu



For more information about these tasks and configuration and management process guidance, see Chapter 6, “Volumes” on page 255.

## 4.7 Hosts

A host system is a computer that is connected to the system through a Fibre Channel (FC) interface or an IP network. It is a logical object that represents a list of worldwide port names (WWPNs) that identify the interfaces that the host uses to communicate with the SAN Volume Controller. Both FC and serial-attached Small Computer System Interface (SCSI) (SAS) connections use WWPNs to identify the host interfaces to the systems.

The **Hosts** menu consists of the following choices, as shown in Figure 4-43:

- ▶ Hosts
- ▶ Host Clusters
- ▶ Ports by Host
- ▶ Mappings
- ▶ Volumes by Host
- ▶ Volumes by Host Cluster

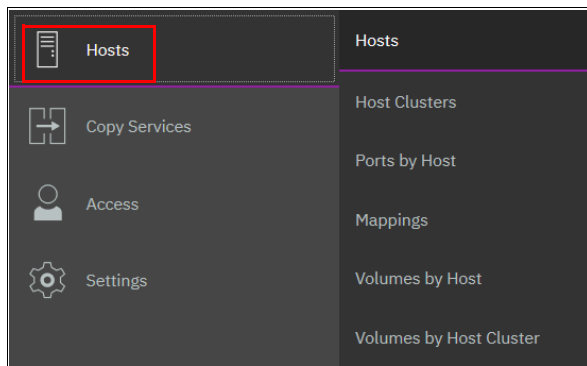


Figure 4-43 Hosts menu

For more information about configuration and management of hosts by using the GUI, see Chapter 7, “Hosts” on page 351.

## 4.8 Copy Services

The IBM Spectrum Virtualize Copy Services and Volumes Copy operations are based on the FlashCopy function. In its basic mode, the function creates copies of content on a source volume to a target volume. Any data on the target volume is lost and is replaced by the copied data.

More advanced functions allow FlashCopy operations to occur on multiple source and target volumes. Management operations are coordinated to provide a common, single point-in-time (PiT) for copying target volumes from their respective source volumes. This technique creates a consistent copy of data that spans multiple volumes.

The SAN Volume Controller Copy Services menu offers the following operations in the GUI, as shown in Figure 4-44:

- ▶ FlashCopy
- ▶ Consistency Groups

- ▶ FlashCopy Mappings
- ▶ Remote Copy (RC)
- ▶ Partnerships

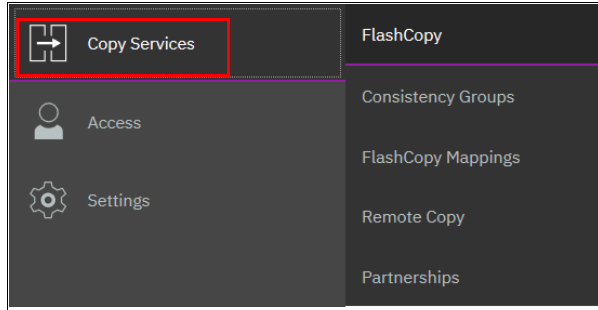


Figure 4-44 Copy Services in GUI

Because Copy Services is one of the most important features for resiliency solutions, see Chapter 10, “Advanced Copy Services” on page 491.

## 4.9 Access

The **Access** menu in the GUI maintains who can log in to the system, defines the access rights to the user, and tracks what was done by each privileged user to the system. It is logically split into three categories:

- ▶ Ownership groups
- ▶ Users by group
- ▶ Audit log

In this section, we explain how to create, modify, or remove a user, and how to see records in the audit log.

The **Access** menu is available from the left pane, as shown in Figure 4-45.

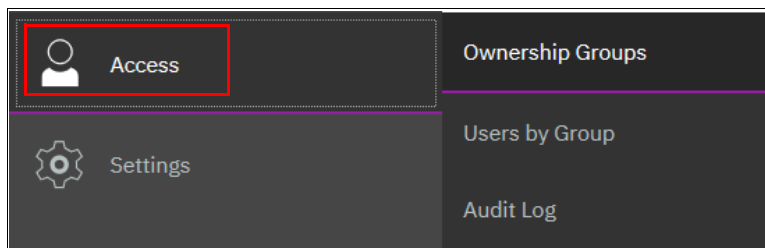


Figure 4-45 Access menu

### 4.9.1 Ownership groups

An *ownership group* defines a subset of users and objects within the system. You can create ownership groups to further restrict access to specific resources that are defined in the ownership group. Only users with Administrator or Security Administrator roles can configure and manage ownership groups. Ownership groups restrict access to only those objects that are defined within that ownership group. An owned object can belong to one ownership group. An *owner* is a user with an ownership group that can view and manipulate objects within that group.

The system supports several resources that you assign to ownership groups:

- ▶ Child pools
- ▶ Volumes
- ▶ Volume groups
- ▶ Hosts
- ▶ Host clusters
- ▶ Host mappings
- ▶ FlashCopy mappings
- ▶ FlashCopy consistency groups

When a user group is assigned to an ownership group, the users in that user group retain their role but are restricted to only those resources within the same ownership group. User groups can define the access to operations on the system, and the ownership group can further limit access to individual resources. For example, you can configure a user group with the Copy Operator role, which limits access of the user to Copy Services functions, such as FlashCopy and RC operations. Access to individual resources, such as a specific FlashCopy consistency group, can be further restricted by adding it to an ownership group. When the user logs on to the management GUI, only resources that they have access to through the ownership group are displayed. Additionally, only events and commands that are related to the ownership group in which a user belongs are viewable by those users.

### **Inheriting ownership**

Depending on the type of resource, ownership can be defined explicitly or ownership can be inherited from the user, user group, or from other parent resources. Objects inherit their ownership group from other objects whenever possible:

- ▶ Volumes inherit the ownership group from the child pool that provides capacity for the volumes.
- ▶ FlashCopy mappings inherit the ownership group from the volumes that are configured in the mapping.
- ▶ Hosts inherit the ownership group from the host cluster they belong to, if applicable.
- ▶ Host mappings inherit the ownership group from both the host and the volume to which the host is mapped.

These objects cannot be explicitly moved to a different ownership group without creating inconsistent ownership.

Ownership groups are also inherited from the user. Objects that are created by an owner inherit the ownership group of the owner. If the owner is in more than one ownership group (only possible for remote users), then the owner must choose the group when the object is created.

### **Child pools**

The following rules apply to child pools that are defined in ownership groups:

- ▶ Child pools can be assigned to an ownership group when you create a pool or change a pool.
- ▶ Users who assign the child pool to the ownership group cannot be defined within that ownership group.
- ▶ Resources that are within the child pool inherit the ownership group that is assigned for the child pool.

## **Host clusters**

The following rules apply to host clusters that are defined in ownership groups:

- ▶ If the user who is creating the host cluster is defined in only one ownership group, the host cluster inherits the ownership group of that user.
- ▶ If the user is defined in an ownership group but is also defined in multiple user groups, the host cluster inherits the ownership group. The system uses the lowest role that the user has from the user group. For example, if a user is defined in two user groups with the roles of Monitor and Copy Operator, the host cluster inherits the Monitor role.
- ▶ Only users not within an ownership group can assign ownership groups when a host cluster is created or changed.

## **Hosts that are not part of a host cluster**

The following rules apply to a host that are not part of a host cluster that is defined in ownership groups:

- ▶ If the user who is creating the host is in only one ownership group, the host cluster inherits the ownership group of that user.
- ▶ If the user is defined in an ownership group but is also defined in multiple user groups, the host inherits the ownership group. The system uses the lowest role that the user has from the user group. For example, if a user is defined in two user groups with the roles of Monitor and Copy Operator, the host inherits the Monitor role.
- ▶ Only users not within an ownership group can assign ownership groups when you create a new host or change an existing host.

## Volume groups

Volume groups can be created to manage multiple volumes that are used with Transparent Cloud Tiering (TCT) support. The following rules apply to volume groups that are defined in ownership groups:

- ▶ If the user that is creating the volume group is defined in only one ownership group, the volume group inherits the ownership group of that user.
- ▶ If the user is defined in an ownership group but is also defined in multiple user groups, the volume group inherits the ownership group. The system uses the lowest role that the user has from the user group. For example, if a user is defined in two user groups with the roles of Monitor and Copy Operator, the host inherits the Monitor role.
- ▶ Only users not within an ownership group can assign ownership groups when you create a new volume group or change an existing volume group.
- ▶ Volumes can be added to a volume group if both the volume and the volume group are within the same ownership group or if both are not in an ownership group. There are situations where a volume group and its volumes can belong to different ownership groups. Volume ownership can be inherited from the ownership group or from one or more child pools.
- ▶ The ownership of a volume group does not affect the ownership of the volumes it contains. If a volume group and its volumes are owned by different ownership groups, then the owner of the child pool that contains the volumes can change the volume directly. For example, the owner of the child pool can change the name of a volume within it. The owner of the volume group can change the volume group itself and indirectly change the volume, such as deleting a volume from the volume group. Neither the ownership group of the child pools or the owner of the volume group can directly manipulate the resources that are not defined in their ownership group.

## FlashCopy consistency groups

FlashCopy consistency groups can be created to manage multiple FlashCopy mappings. The following rules apply to FlashCopy consistency groups that are defined in ownership groups:

- ▶ If the user that is creating the FlashCopy consistency group is in only one ownership group, the FlashCopy consistency group inherits the ownership group of that user.
- ▶ If the user is defined in an ownership group but is also defined in multiple user groups, the FlashCopy consistency group inherits the ownership group. The system uses the lowest role that the user has from the user group.
- ▶ Only users not within an ownership group can assign ownership groups when a FlashCopy consistency is created or changed.
- ▶ FlashCopy mappings can be added to a consistency group if the volumes in the mapping and the consistency group are within the same ownership group. You can also add a FlashCopy mapping to a consistency group if it and all of its dependent resources are not in an ownership group.
- ▶ There are situations where a FlashCopy consistency group and its resources can belong to different ownership groups.
- ▶ As with volume groups and volumes, the ownership of the consistency group has no impact on the ownership of the mappings it contains.

## User groups

The following rules apply to user groups that are defined in ownership groups:

- ▶ If the user that is creating the user group is in only one ownership group, the user group inherits the ownership group of that user.
- ▶ If the user is with multiple user groups, the user group inherits the ownership group of the user group with the lowest role.
- ▶ Only users not within an ownership group can assign an ownership group when a user group is created or changed.

These resources inherit ownership from the parent resource. A user cannot change the ownership group of the resource, but can change the ownership group of the parent object.

## Hosts that are a part of a host cluster

The following rules apply to hosts that are defined in ownership groups:

- ▶ The host inherits the ownership group of the host cluster to which it belongs.
- ▶ If a host is removed from a host cluster within an ownership group, the host inherits the ownership group of the host cluster to which it used to belong.
- ▶ If a host is removed from a host cluster that is not within an ownership group, the host inherits no ownership groups.
- ▶ Hosts can be added to a host cluster if the host and host cluster have the same ownership group.
- ▶ Changing the ownership group of a host cluster automatically changes the ownership group of all the hosts inside the host cluster.

## Volumes

The following rules apply to volumes that are defined in ownership groups:

- ▶ The volume inherits the ownership group of the child pools that provide capacity for the volume and its copies.
- ▶ If the child pool that provides capacity for the volume or its copies is defined in different ownership groups, then the volume cannot be created in an ownership group.
- ▶ With volume groups, the volume group and its volumes can belong to different ownership groups. However, the ownership of a volume group does not impact the ownership of the volumes that it contains.

## Users

The following rules apply to users that are defined in ownership groups:

- ▶ A user inherits the ownership group of the user group to which it belongs.
- ▶ Users that use Lightweight Directory Access Protocol (LDAP) for remote authentication can belong to multiple user groups and multiple ownership groups.

## Volume-to-host mappings

The following rules apply to volume-to-host mappings that are defined in ownership groups:

- ▶ Volume-to-host mappings inherit the ownership group of the host or host cluster and volume in the mapping.
- ▶ If host or host cluster and volume are within different ownership groups, then the mapping cannot be assigned an ownership group.

## FlashCopy mappings

The following rules apply to FlashCopy mappings that are defined in ownership groups:

- ▶ FlashCopy mappings inherit the ownership group of both volumes that are defined in the mapping.
- ▶ If the volumes are within different ownership groups, then the mapping cannot be assigned to an ownership group.
- ▶ Like with FlashCopy consistency groups, it is possible for a consistency group and its mappings to belong to different ownership groups. However, the ownership of the consistency group has no impact on the ownership of the mappings that it contains.

## Configuring ownership groups

You can configure ownership groups to manage access to resources on the system. An ownership group defines a subset of users and objects within the system. You can create ownership groups to further restrict access to specific resources that are defined in the ownership group. Only users with Administrator or Security Administrator roles can configure and manage ownership groups.

## Migrating to ownership groups

If you updated your system to a software level that supports ownership groups, you must reconfigure certain resources if you want to configure ownership groups. An ownership group defines a subset of users and objects within the system. You can create ownership groups to further restrict access to specific resources that are defined in the ownership group. Only users with Administrator or Security Administrator roles can configure and manage ownership groups.

Figure 4-46 shows an example of an ownership group.

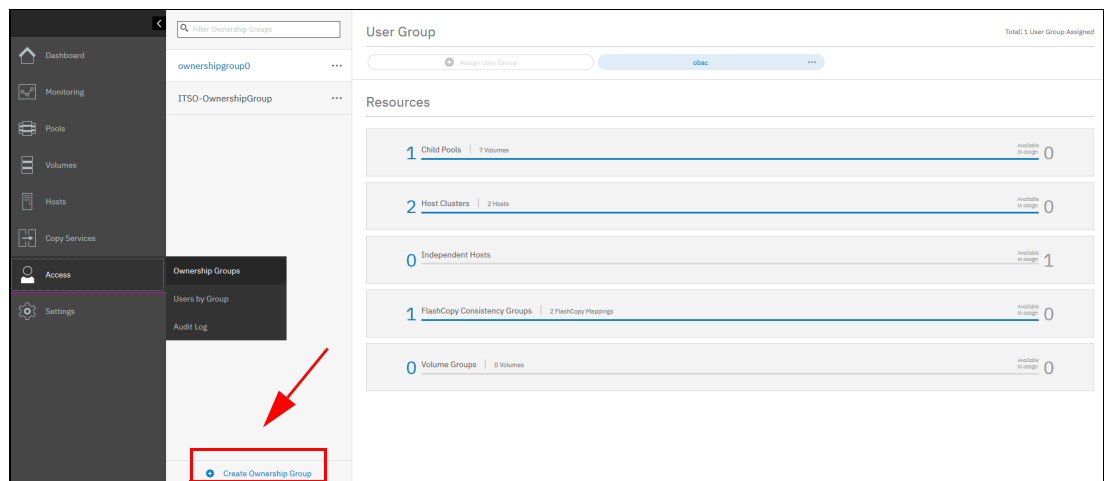


Figure 4-46 Ownership by Groups

To create an ownership group, select **Create Ownership Group**, as shown in Figure 4-46.

## 4.9.2 Users by groups

You can create local users who can access the system. These user types are defined based on the administrative privileges that they have on the system.

Local users must provide a password, Secure Shell (SSH) key, or both. Local users are authenticated through the authentication methods that are configured on the system. If the local user needs access to the management GUI, a password is needed for the user. If the user requires access to the CLI through SSH, a password or a valid SSH key file is necessary.

Local users must be part of a user group that is defined on the system. User groups define roles that authorize the users within that group to a specific set of operations on the system.

To define your user group in the SAN Volume Controller, select **Access** → **Users by Group**, as shown in Figure 4-47.

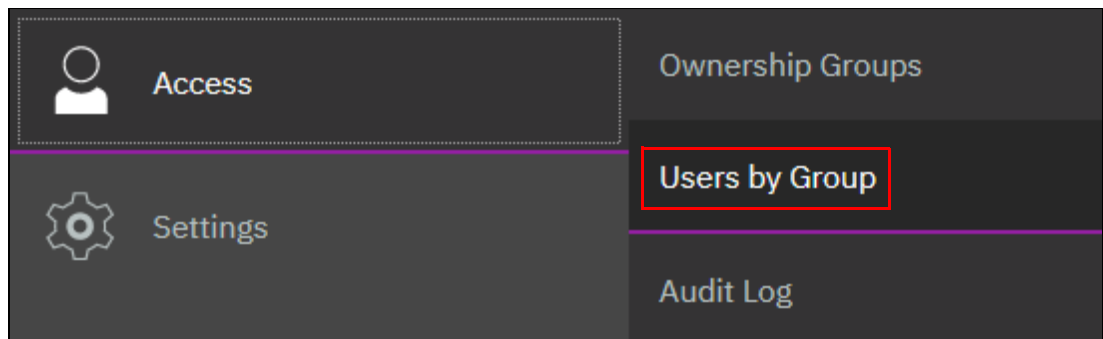


Figure 4-47 Access Users by Group

Select **Create User Group**, as shown in Figure 4-48 on page 151.



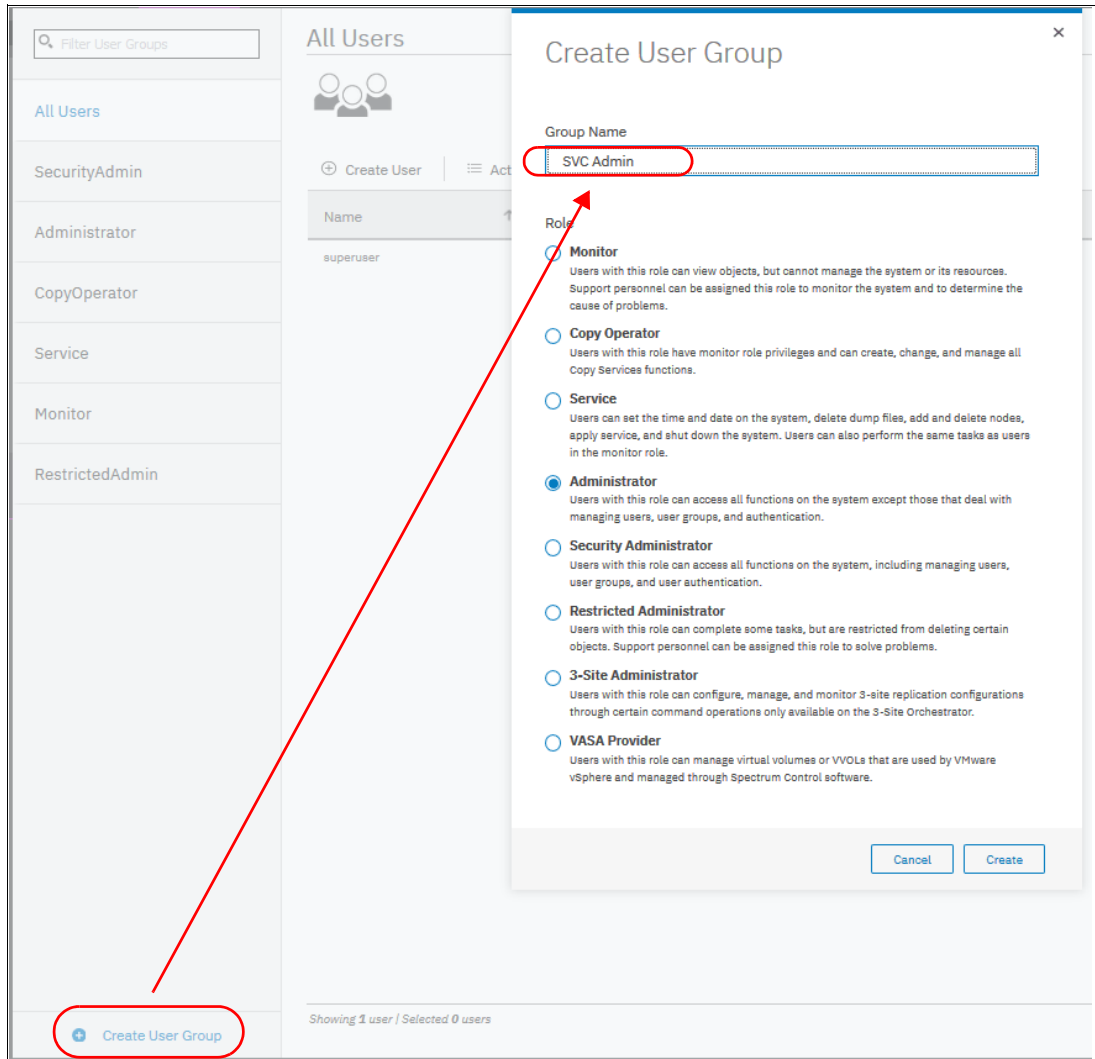


Figure 4-48 Creating a user group in SAN Volume Controller

Figure 4-49 shows the newly created User Group.

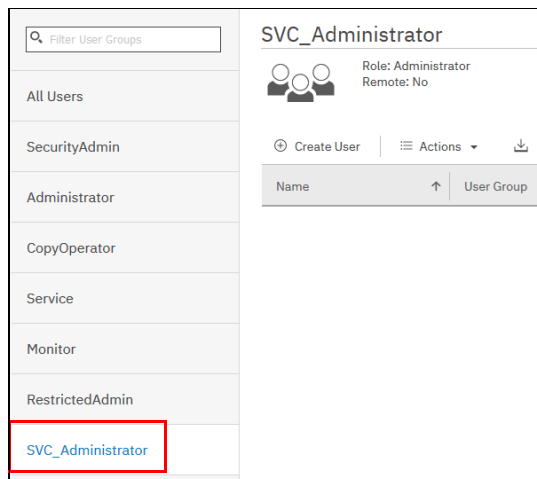


Figure 4-49 User Group

The following privileged user group roles exist in IBM Spectrum Virtualize:

- ▶ **Monitor**

These users can access all system viewing actions. Monitor role users cannot change the state of the system or the resources that the system manages. Monitor role users can access all information-related GUI functions and commands, back up configuration data, and change their own passwords.
- ▶ **Copy Operator**

These users can start and stop all existing FlashCopy, MM, and GM relationships. Copy Operator role users can run the system commands that Administrator role users can run that deal with FlashCopy, MM, and GM relationships.
- ▶ **Service**

These users can set the time and date on the system, delete dump files, add and delete nodes, apply service, and shut down the system. Users can also complete the same tasks as users in the monitor role.
- ▶ **Administrator**

These users can manage all functions of the system except for those functions that manage users, user groups, and authentication. Administrator role users can run the system commands that the Security Administrator role users can run from the CLI, except for commands that deal with users, user groups, and authentication.
- ▶ **Security Administrator**

These users can manage all functions of the system, including managing users, user groups, user authentication, and configuring encryption. Security Administrator role users can run any system commands from the CLI. However, they cannot run the **sa info** and **sa task** commands from the CLI. Only the superuser ID can run those commands.
- ▶ **Restricted Administrator**

These users can perform the same tasks and run most of the same commands as Administrator role users. However, users with the Restricted Administrator role are not authorized to run the **rmvdisk**, **rmvdiskhostmap**, **rmhost**, or **rmmdiskgrp** commands. Support personnel can be assigned this role to help resolve errors and fix problems.
- ▶ **3-Site Administrator**

These users can configure, manage, and monitor 3-site replication configurations through certain command operations that are available only on the 3-Site Orchestrator. Before you can work with 3-Site Orchestrator, a user profile must be created.
- ▶ **vSphere APIs for Storage Awareness (VASA) Provider**

These users can manage VMware vSphere Virtual Volumes (VVOLS).

## Registering a user

After you define your group (in our example, SVC\_Administrator with Administrator privileges), you can register a user (SVC\_Cluster\_Admin) within this group by clicking **Create User** and selecting **SVC Administrator** (see Figure 4-50 on page 153).

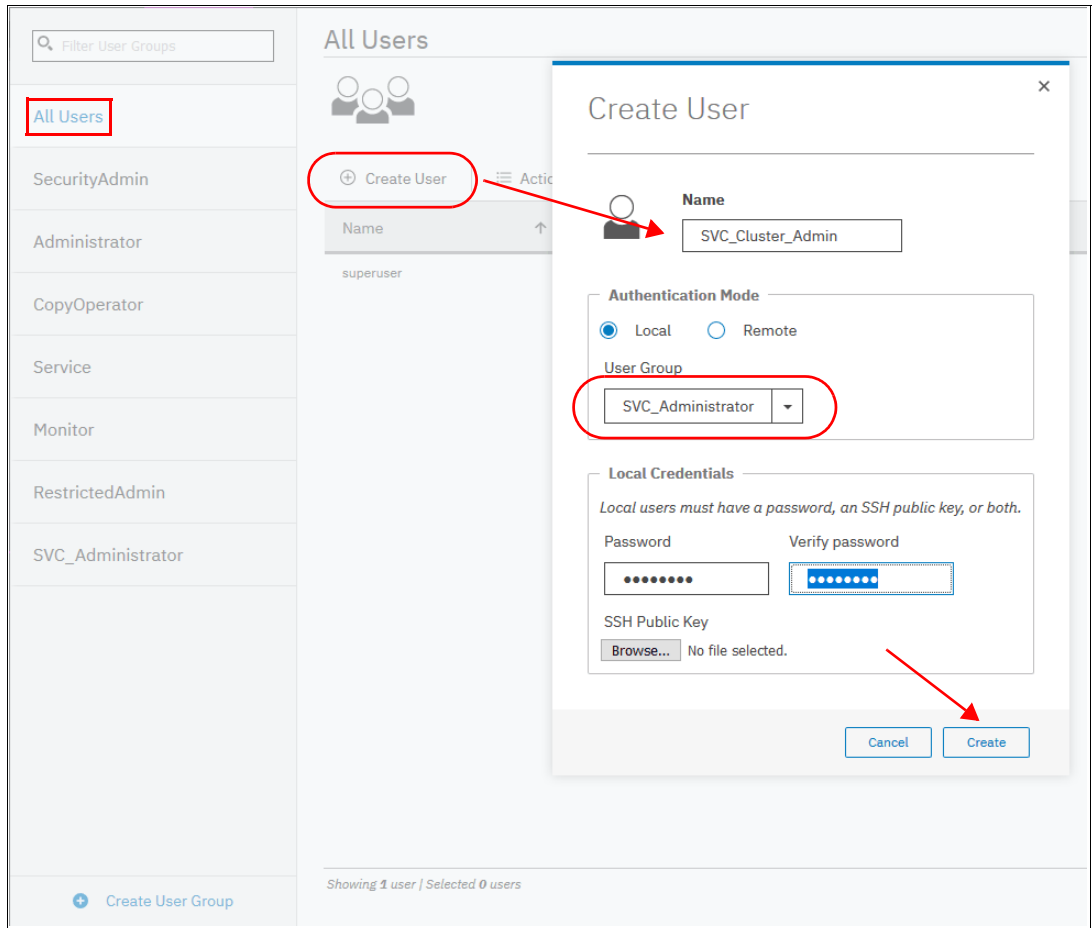


Figure 4-50 Registering a user account

## Deleting a user

To remove a user account, right-click the user in the **All Users** list and select **Delete**, as shown in Figure 4-51.

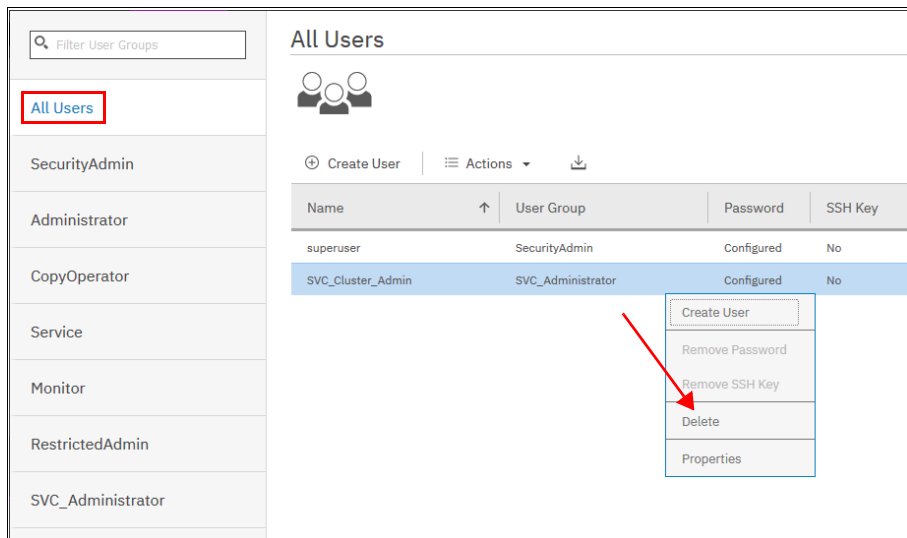


Figure 4-51 Deleting a user account

**Attention:** When you click **Delete**, the user account is directly deleted SAN Volume Controller. No other confirmation request is presented.

## Setting a new password

To set a new password for the user, right-click the user (or click **Actions**) and select **Properties**. In this window, you can either assign the user to a different group or reset their password (see Figure 4-52).

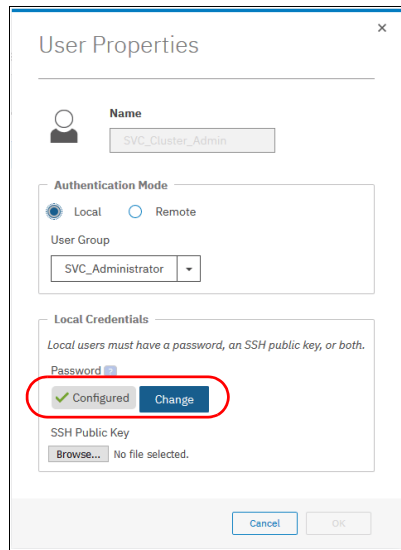


Figure 4-52 Setting a new password

## 4.9.3 Audit log

An *audit log* documents actions that are submitted through the management GUI or the CLI. You can use the audit log to monitor user activity on your system.

The audit log entries provide the following information:

- ▶ Time and date when the action or command was submitted.
- ▶ Name of the user who completed the action or command.
- ▶ IP address of the system where the action or command was submitted.
- ▶ Name of source and target node on which the command was submitted.
- ▶ Parameters that were submitted with the command, excluding confidential information.
- ▶ Results of the command or action that completed successfully.
- ▶ Sequence number and the object identifier that is associated with the command or action.

An example of the audit log is shown in Figure 4-53 on page 155.

Date and Time	User Name	IP Address	Command	Object ID	Source Node	Target Node
10/25/2019 9:59:24 AM	superuser	10.230.206.77	svctask mkuser -gui -name SVC_Cluster_Admin -password ##...	1		
10/25/2019 9:53:23 AM	superuser	10.230.206.77	svctask rmuser -gui 1			
10/25/2019 9:50:07 AM	superuser	10.230.206.77	svctask mkuser -gui -name SVC_Cluster_Admin -password ##...	1		
10/25/2019 9:31:05 AM	superuser	10.230.206.77	svctask mkusergrp -gui -name SVC_Administrator -role Administ...	6		
10/24/2019 6:00:14 PM	superuser	127.0.0.1	sataisk cpfiles -prefix /dumps/svc.config.cron*_75ACXFO -sourc...		75ACXFO	75ACXFO
10/24/2019 2:45:53 PM	superuser	10.230.192.10	svctask chiogrp -gui -maintenance no 3			
10/24/2019 2:45:53 PM	superuser	10.230.192.10	svctask chiogrp -gui -maintenance no 2			
10/24/2019 2:45:52 PM	superuser	10.230.192.10	svctask chiogrp -gui -maintenance no 0			
10/24/2019 2:45:52 PM	superuser	10.230.192.10	svctask chiogrp -gui -maintenance no 1			
10/24/2019 2:25:05 PM	superuser	10.230.192.10	svctask chnode -gui -identify no 1			
10/24/2019 2:23:34 PM	superuser	10.230.192.10	svctask chnode -gui -identify yes 1			
10/24/2019 10:32:39 AM	superuser	10.230.206.35	svctask cheventlog -fix 118 -gui			
10/24/2019 10:32:39 AM	superuser	10.230.206.35	svctask cheventlog -fix 214 -gui			
10/24/2019 10:32:39 AM	superuser	10.230.206.35	svctask cheventlog -fix 220 -gui			
10/24/2019 10:32:39 AM	superuser	10.230.206.35	svctask cheventlog -fix 223 -gui			
10/24/2019 10:26:56 AM	superuser	10.230.206.35	svctask chuser -password ### superuser			
10/24/2019 10:25:13 AM	superuser	10.230.206.35	svctask addnode -gui -name node_KDBP1BP -panelname KDBP1...			
10/24/2019 10:24:59 AM	superuser	10.230.206.35	svctask rmode -gui 3			
10/24/2019 10:09:40 AM	superuser	10.230.206.35	svctask addnode -gui -iogrp 0 -name node_75ACXFO -panelnam...			
10/24/2019 10:09:40 AM	superuser	10.230.206.35	svctask addnode -gui -name node_KDBP1OG -panelname KDBP1...			

Figure 4-53 Audit log

The following commands are not documented in the audit log:

- ▶ **dumpconfig**
- ▶ **cpdumps**
- ▶ **finder**
- ▶ **dumperrlog**

The following items are also not documented in the audit log:

- ▶ Commands that fail are not logged.
- ▶ A result code of 0 (success) or 1 (success in progress) is not logged.
- ▶ Result object ID of node type (for the **addnode** command) is not logged.
- ▶ Views are not logged.

**Important:** Failed commands are not recorded in the audit log. Commands that are triggered by IBM Support personnel are recorded with the flag **Challenge** because they use challenge-response authentication.

## 4.10 Settings

Use the Settings pane to configure system options for notifications, security, IP addresses, and preferences that are related to display options in the management GUI (see Figure 4-54).

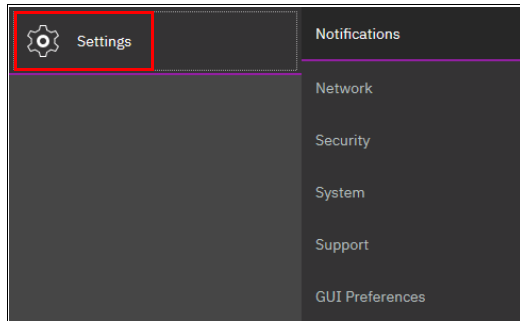


Figure 4-54 Settings menu

The following options are available for configuration from the **Settings** menu:

- ▶ **Notifications:** The system can use Simple Network Management Protocol (SNMP) traps, syslog messages, and Call Home emails to notify you and IBM Support Center when significant events are detected. Any combination of these notification methods can be used simultaneously.
- ▶ **Network:** Use the Network pane to manage the management IP addresses for the system, service IP addresses for the nodes, and internet Small Computer Systems Interface (iSCSI) and FC configurations. The system must support FC or Fibre Channel over Ethernet (FCoE) connections to your storage area network (SAN).
- ▶ **Security:** Use the Security pane to configure and manage remote authentication services.
- ▶ **System:** Use the **System** menu to manage overall system configuration options, such as licenses, updates, and date and time settings.
- ▶ **Support:** Use this option to configure and manage connections, and upload support packages to the support center.
- ▶ **GUI Preferences:** Configure welcome message after login, and refresh internals and GUI logout timeouts.

These options are described in more detail in the following sections.

### 4.10.1 Notifications menu

The SAN Volume Controller can use SNMP traps, syslog messages, and Call Home email to notify you and the IBM Support Center when significant events are detected. Any combination of these notification methods can be used simultaneously.

Notifications are normally sent immediately after an event is raised. However, events can occur because of service actions that are performed. If a recommended service action is active, notifications about these events are sent only if the events are still unfixed when the service action completes.

#### SNMP notifications

SNMP is a standard protocol for managing networks and exchanging messages. The system can send SNMP messages that notify personnel about an event. You can use an SNMP manager to view the SNMP messages that are sent by the SAN Volume Controller.

To view the SNMP configuration, click the **Settings** icon and select **Notification** → **SNMP** (see Figure 4-55).

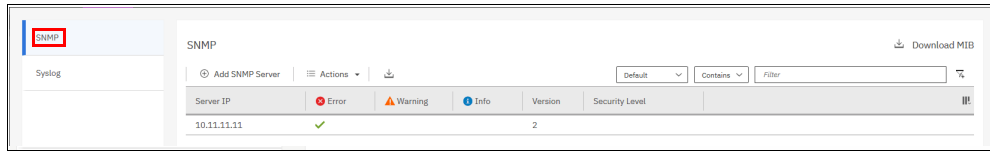


Figure 4-55 Setting SNMP server and traps

In Figure 4-55, you can view and configure an SNMP server to receive various informational, error, or warning notifications by setting the following information:

► **IP Address**

The address for the SNMP server.

► **Server Port**

The remote port number for the SNMP server. The remote port number must be a value 1 - 65535.

► **Community**

The SNMP community is the name of the group to which devices and management stations that run SNMP belong.

► **Event Notifications**

Consider the following points about event notifications:

- Select **Error** if you want the user to receive messages about problems, such as hardware failures, that must be resolved immediately.

**Important:** Browse to **Recommended Actions** to run the fix procedures on these notifications.

- Select **Warning** if you want the user to receive messages about problems and unexpected conditions. Investigate the cause immediately to determine any corrective action.

**Important:** Browse to **Recommended Actions** to run the fix procedures on these notifications.

- Select **Info** if you want the user to receive messages about expected events. No action is required for these events.

To remove an SNMP server, click the minus sign (-). To add another SNMP server, click the plus sign (+).

## Syslog notifications

The syslog protocol is a standard protocol for forwarding log messages from a sender to a receiver on an IP network. The IP network can be IPv4 or IPv6. The system can send syslog messages that notify personnel about an event. You can use the Syslog pane to view the syslog messages that are sent by the SAN Volume Controller. To view the Syslog configuration, go to the System pane and click **Settings**, and select **Notification** → **Syslog** (see Figure 4-56). A domain name server (DNS) server is required to use domain names in syslog.

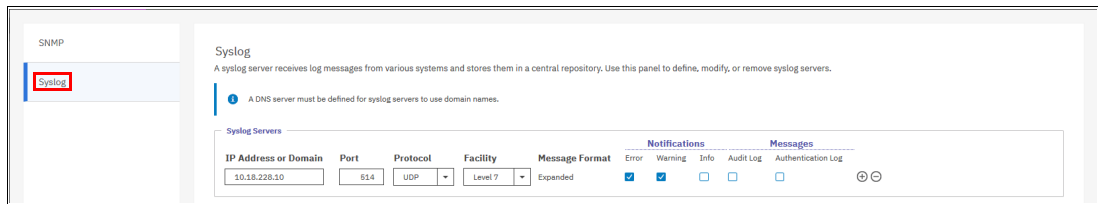


Figure 4-56 Setting the syslog messages

From this window, you can view and configure a syslog server to receive log messages from various systems and store them in a central repository by entering the following information:

- ▶ IP Address

The IP address for the syslog server.

- ▶ Port

Port number of the syslog server

- ▶ Protocol of the transmission protocol

Select **UDP** or **TCP**.

- ▶ Facility

The facility determines the format for the syslog messages. The facility can be used to determine the source of the message.

- ▶ Message Format

The message format depends on the facility. The system can transmit syslog messages in the following formats:

- The concise message format provides standard detail about the event.
- The expanded format provides more details about the event.

- ▶ Event Notifications

Consider the following points about event notifications:

- Select **Error** if you want the user to receive messages about problems, such as hardware failures, that must be resolved immediately.

**Important:** Browse to **Recommended Actions** to run the fix procedures on these notifications.

- Select **Warning** if you want the user to receive messages about problems and unexpected conditions. Investigate the cause immediately to determine whether any corrective action is necessary.

**Important:** Browse to **Recommended Actions** to run the fix procedures on these notifications.





## Management IP addresses

To view the management IP addresses of IBM Spectrum Virtualize, select **Settings** → **Network**, and click **Management IP Addresses**. The GUI shows the management IP address by pointing to the network ports, as shown in Figure 4-58.

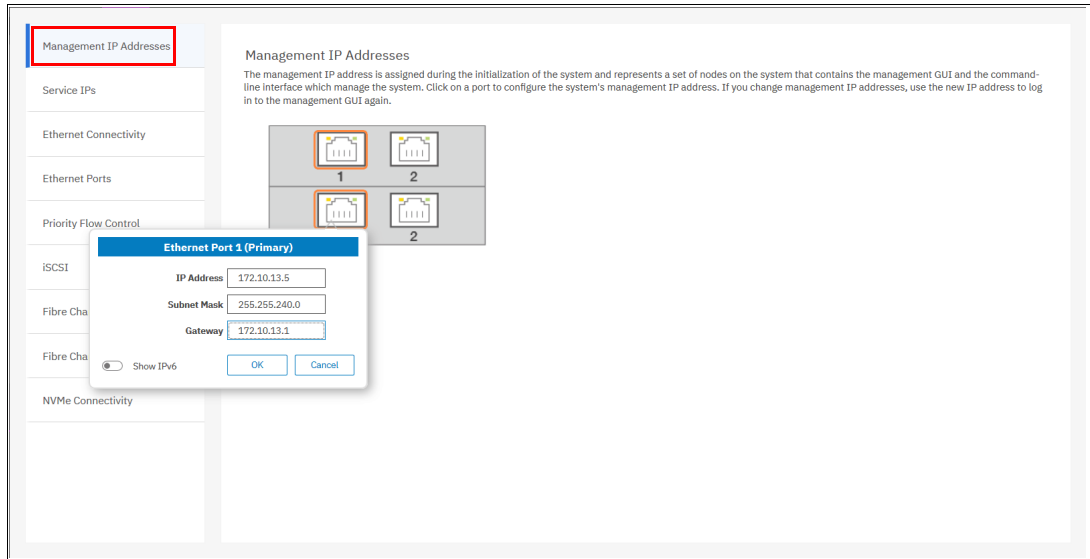


Figure 4-58 Viewing the management IP addresses

## Service IP information

To view the Service IP information of your IBM Spectrum Virtualize installation, select **Settings** → **Network**, as shown in Figure 4-57 on page 159. Click the **Service IP Address** option to view the properties, as shown in Figure 4-59.

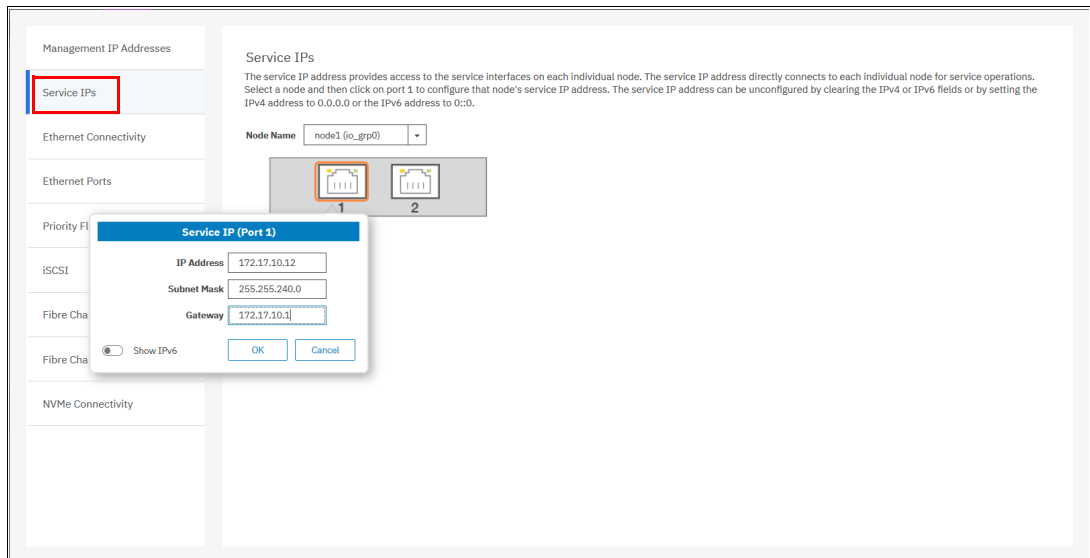


Figure 4-59 Viewing service IP addresses

The service IP address is commonly used to provide access to the network interfaces on each individual node.

Instead of reaching the management IP address, the service IP address directly connects to each individual node for service operations. You can select a node from the drop-down list and then click any of the ports that are shown in the GUI. The service IP address can be configured to support IPv4 or IPv6.

## Ethernet ports

Ethernet ports for each node are at the rear of the system and used to connect the system to hosts, external storage systems, and to other systems that are part of RC partnerships. Depending on the model of your system, supported connection types include FC, when the ports are FCoE-capable, iSCSI, and iSCSI Extensions for Remote Direct Memory Access (RDMA) (iSER). iSER connections use either the RDMA over Converged Ethernet (RoCE) protocol or the internet Wide Area RDMA Protocol (iWARP). The panel indicates whether a specific port is being used for a specific purpose and traffic.

You can modify how the port is used by selecting **Actions**. Select either **Modify Remote Copy**, **Modify iSCSI Hosts**, or **Modify Storage Ports** to change the use of the port. You can also display the login information for each host that is logged in to a selected node.

To display this information, select **Settings** → **Network** → **Ethernet Ports** and right-click the node and select **IP Login Information**. This information can be used to detect connectivity issues between the system and hosts and to improve the configuration of iSCSI host to optimize performance. Select **Ethernet Ports** for an overview from the menu, as shown in Figure 4-60. For planning, see 2.5, “Connectivity planning” on page 57.

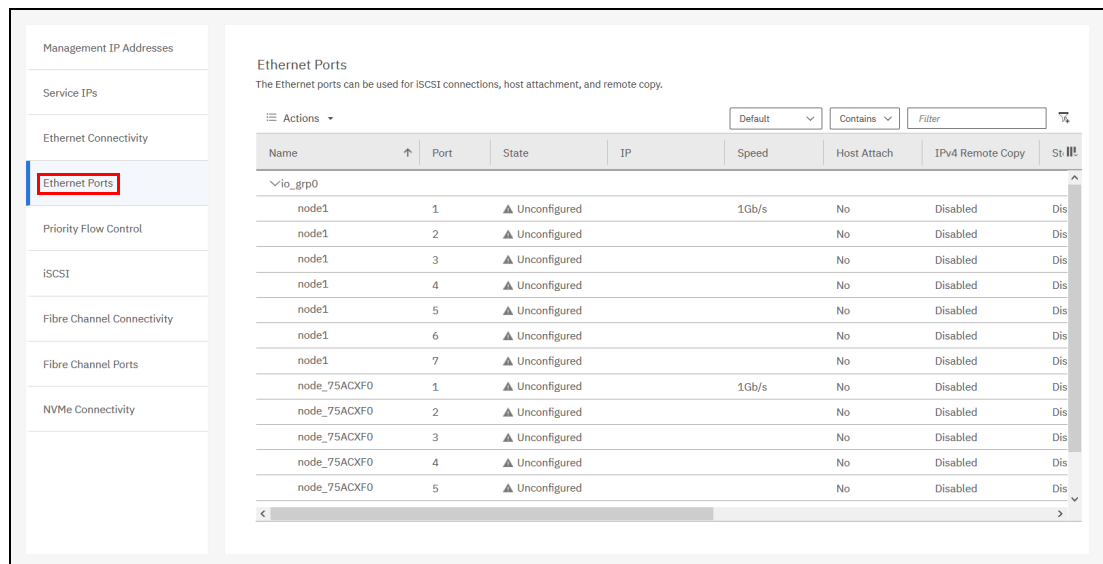


Figure 4-60 Ethernet Ports

## Priority flow control

Priority flow control (PFC) is an Ethernet protocol that you can use to select the priority of different types of traffic within the network. With PFC, administrators can reduce network congestion by slowing or pausing certain classes of traffic on ports, thus providing better bandwidth for more important traffic. The system supports PFC on various supported Ethernet-based protocols on three types of traffic classes: system, host attachment, and storage traffic. PFC requires virtual local area network (VLAN) configuration on the corresponding Ethernet Ports.

You can configure a priority tag for each of these traffic classes. The priority tag can be any value 0 - 7. You can set identical or different priority tag values to all these traffic classes. You can also set bandwidth limits to ensure quality of service (QoS) for these traffic classes by using the Enhanced Transmission Selection (ETS) setting on the network. When you plan to configure PFC, follow these guidelines and examples.

To use PFC and ETS, ensure that the following tasks are completed:

- ▶ Ensure that ports support 10 Gb or higher bandwidth to use PFC settings.
- ▶ Configure a virtual local area network (VLAN) on the system to use PFC capabilities for the configured IP version.
- ▶ Ensure that the same VLAN settings are configured on the all entities, including all switches between the communicating end points.
- ▶ Configure the QoS values (priority tag values) for host attachment, storage, or system traffic by running the `chsystemethernet` command.
- ▶ To enable priority flow for host attachment traffic on a port, make sure that the host flag is set to `yes` on the configured IP on that port.
- ▶ To enable priority flow for storage traffic on a port, make sure that storage flag is set to `yes` on the configured IP on that port.
- ▶ On the switch, enable the Data Center Bridging Exchange (DCBx). DCBx enables switch and adapter ports to exchange parameters that describe traffic classes and PFC capabilities. For these steps, check your switch documentation for details.
- ▶ For each supported traffic class, configure the same priority tag on the switch. For example, if you plan to have a priority tag setting of 3 for storage traffic, ensure that the priority is also set to 3 on the switch for that traffic type.
- ▶ If you are planning on using the same port for different types of traffic, ensure that the ETS settings are configured on the network.

### 4.10.3 Using the management GUI

To set PFC on the system, complete the following steps:

1. In the management GUI, select **Settings** → **Network** → **Priority Flow Control**, as shown in Figure 4-61 on page 163.

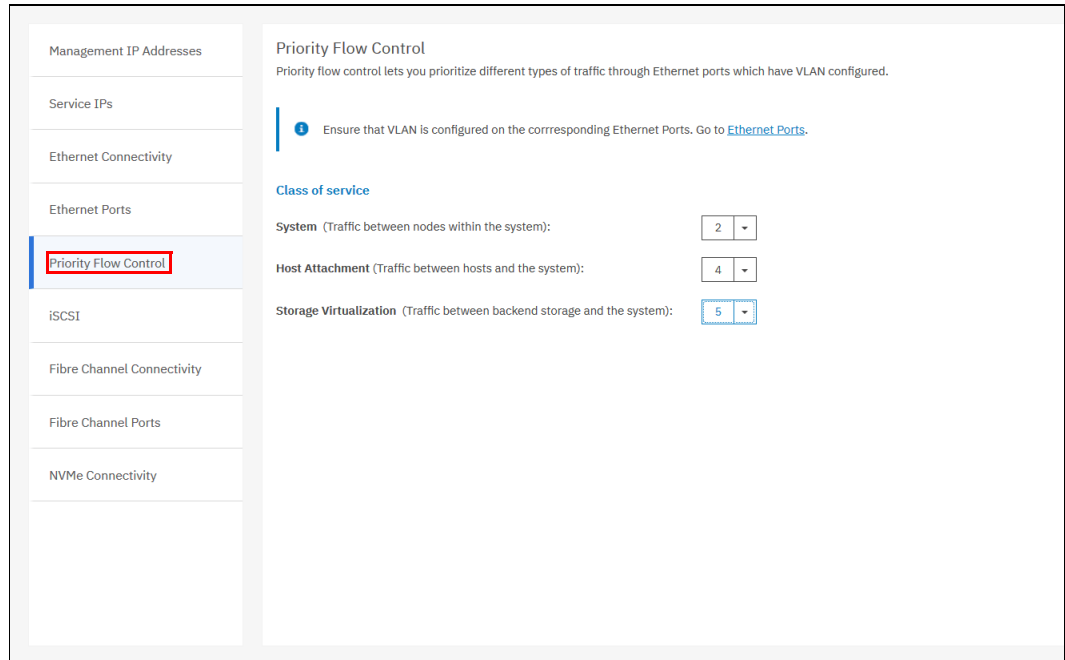


Figure 4-61 Priority flow control

2. For each of following classes of service, select the priority setting for that traffic type:

- System

Set a value 0 - 7 for the system traffic, which includes communication between nodes within the system. The system priority tag is supported on iSCSI connections and on systems that support RDMA over Ethernet connections between nodes. Ensure that you set the same priority tag on the switch to use PFC capabilities.

- Host attachment

Set the priority tag 0 - 7 for system to host traffic. The host attachment priority tag is supported on iSCSI connections and on systems that support RDMA over Ethernet connections. Ensure that you set the same priority tag on the switch to use PFC capabilities.

- Storage virtualization

Set the priority tag 0 - 7 for system to external storage traffic. The storage virtualization priority tag is supported on storage traffic over iSCSI connections. Ensure that you set the same priority tag on the switch to use PFC capabilities.

Make sure that IP is configured with VLAN.

## iSCSI information

From the iSCSI pane in the **Settings** menu, you can display and configure parameters for the system to connect to iSCSI-attached hosts, as shown in Figure 4-62.

Node Name	iSCSI Alias	iSCSI Name (IQN)
node1		iqn.1986-03.com.ibm:2145.cluster10.155.19.123.node1
node_75ACXF0		iqn.1986-03.com.ibm:2145.cluster10.155.19.123.node75acxf0
node_KDBP1BP		

Figure 4-62 iSCSI Configuration pane

The following parameters can be updated:

- **System Name**

It is important to set the system name correctly because it is part of the iSCSI Qualified Name (IQN) for the node.

**Important:** If you change the name of the system after iSCSI is configured, you might need to reconfigure the iSCSI hosts.

To change the system name, click the system name and specify the new name.

**System name:** You can use the letters A - Z and a - z, the numbers 0 - 9, and the underscore (\_) character. The name can be 1 - 63 characters.

- **iSCSI aliases (optional)**

An *iSCSI alias* is a user-defined name that identifies the node to the host. Complete the following steps to change an iSCSI alias:

- Click an iSCSI alias.
- Specify a name for it.

Each node has a unique iSCSI name that is associated with two IP addresses. After the host starts the iSCSI connection to a target node, this IQN from the target node is visible in the iSCSI configuration tool on the host.

- **Internet Storage Name Service (iSNS) and Challenge Handshake Authentication Protocol (CHAP)**

You can specify the IP address for the iSNS. Host systems use the iSNS server to manage iSCSI targets and for iSCSI discovery.

You can also enable CHAP to authenticate the system and iSCSI-attached hosts with the specified shared secret.

The CHAP secret is the authentication method that is used to restrict access for other iSCSI hosts that use the same connection. You can set the CHAP for the whole system under the system properties or for each host definition. The CHAP must be identical on the server and the system and host definition. You can create an iSCSI host definition without using CHAP.

## Fibre Channel information

As shown in Figure 4-63, you can use the Fibre Channel Connectivity pane to display the FC connectivity between nodes and other storage systems and hosts that attach through the FC network. You can filter by selecting one of the following fields:

- ▶ All nodes, storage systems, and hosts
- ▶ Systems
- ▶ Nodes
- ▶ Storage systems
- ▶ Hosts

You can view Fibre Channel Connectivity as shown in Figure 4-63.

Fibre Channel Connectivity

Display the connectivity between nodes and other storage systems and hosts that are attached through the Fibre Channel network.

View connectivity for: All nodes, storage systems, and hosts [Show Results](#)

Name	System Name	Remote WWPN	Remote ...	Local WWPN	Local Port	Local NP...	State	Node
KDBP1CG	Cluster_10.155.	500507680C21041D	AB0200	500507680C210000	5	AB0100	Active	node
		500507680C21041D	AB0200	500507680C250000	5	AB0101	Inactive	node
node_KDBP1BP	Cluster_10.155.	500507680C220416	AB0100	500507680C140000	4	AB0200	Active	node
		500507680C220416	AB0100	500507680C180000	4	AB0201	Inactive	node
		500507680C22041D	AB0100	500507680C150508	1	AB0201	Inactive	node
KDBP1CG	Cluster_10.155.	500507680C22041D	AB0100	500507680C110508	1	AB0200	Active	node
node_75ACXF0	Cluster_10.155.	500507680C220508	AB0200	500507680C230000	7	AB0100	Active	node
		500507680C220508	AB0200	500507680C270000	7	AB0101	Inactive	node
		500507680C230000	AB0100	500507680C260508	6	AB0201	Inactive	node
node1	Cluster_10.155.	500507680C230000	AB0100	500507680C220508	6	AB0200	Active	node
node_KDBP1BP	Cluster_10.155.	500507680C240416	AB0200	500507680C140508	4	AB0100	Active	node
		500507680C240416	AB0200	500507680C180508	4	AB0101	Inactive	node

You can change the WWPN notation from the actions menu.

Figure 4-63 Fibre Channel Connectivity

In the Fibre Channel Ports pane, you can use this view to display how the FC port is configured across all control node canisters in the system. This view helps, for example, to determine which other clusters and hosts the port may communicate with, and which ports are virtualized. “No” indicates that this port cannot be online on any node other than the owning node (see Figure 4-64).

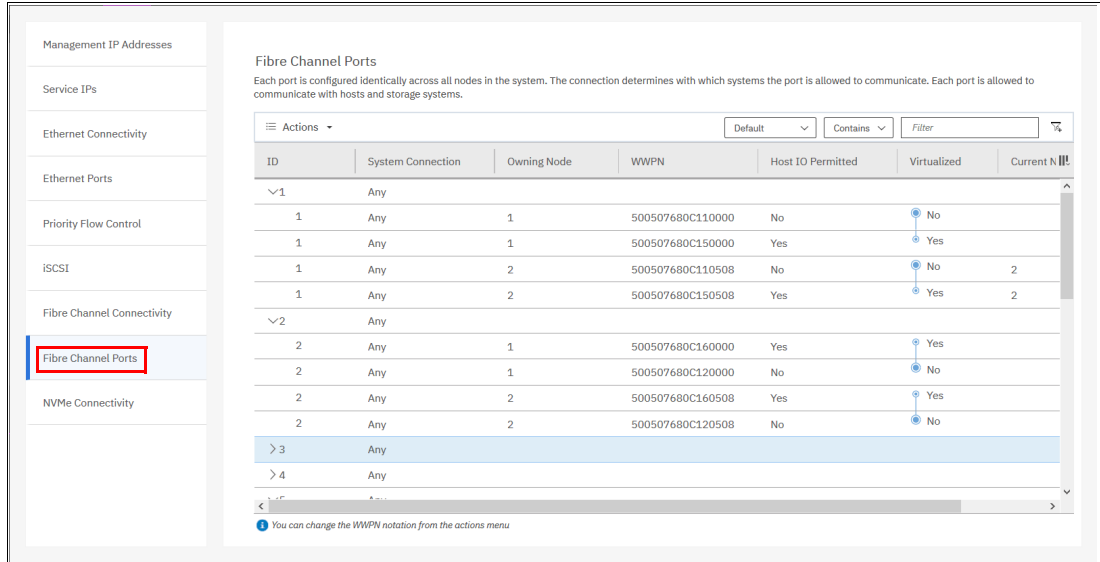


Figure 4-64 Viewing Fibre Channel Port properties

## Non-Volatile Memory Express connectivity

Use the NVMe Connectivity page to display connections between nodes and hosts that use a Non-Volatile Memory Express (NVMe) over Fibre Channel Connection (IBM FICON®). For other specific information about NVMe over Fibre Channel (FC-NVMe) such as interoperability requirements, go to [IBM Support](#) and search for “Configuration Limits and Restrictions”.

If your system supports an FC-NVMe connection between nodes and hosts, you can display details about each side of the connection. To display node details, select the node from the drop-down menu and select **Show Results**. You can also display the host details for the connection or for all hosts and nodes. Use this window to troubleshoot issues between nodes and hosts that use FC-NVMe connections.

For these connections, the Status column displays the current state of the connection. The following states for the connection are possible:

- ▶ **Active**  
Indicates that the connection between the node and host is being used.
- ▶ **Inactive**  
Indicates that the connection between the node and host is configured, but no FC-NVMe operations occurred in the last 5 minutes. Since the system sends periodic heartbeat message to keep the connection open between the node and the host, it is unusual to see an inactive state for the connection. However, it can take up to 5 minutes for the state to change from inactive to active.



If the inactive state remains beyond the 5-minute refresh interval, it can indicate a connection problem between the host and the node. If a connection problem persists between the host and the node, a reduced node login count or the status of the host indicates it is degraded, which you can view by selecting **Hosts** → **Hosts by Port** in the management GUI. Verify these values in the management GUI, and view the messages by selecting **Monitoring** → **Events**.

Figure 4-65 shows the NVMe connectivity when it is available.

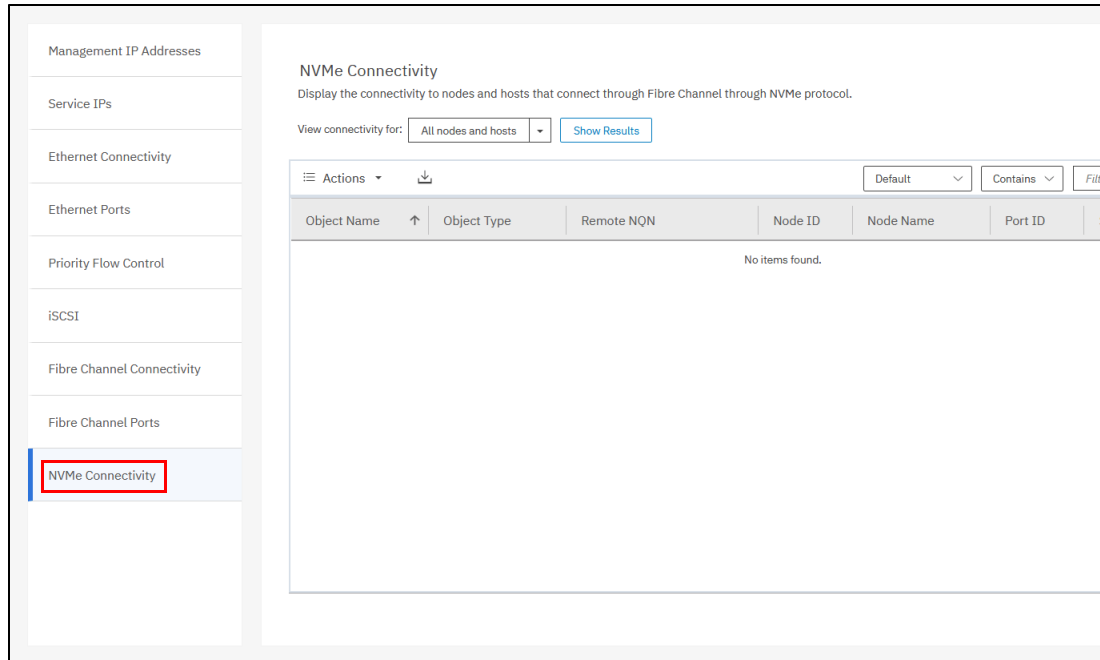


Figure 4-65 NVMe Connectivity window

Consider FC-NVMe target limits when you plan and configure the hosts. Include the following points in your plan:

- ▶ An NVMe host can connect to four NVMe controllers on each system node. The maximum per node is four with an extra four in failover.
- ▶ Zone up to four ports in a single host to detect up to four ports on a node. To allow failover and avoid outages, zone the same or extra host ports to detect an extra four ports on the second node in the I/O group.
- ▶ A single I/O group can contain up to 256 FC-NVMe I/O controllers. The maximum number of I/O controllers per node is 128 plus an extra 128 in failover. Zone a total maximum of 16 hosts to detect a single I/O group. Also, consider that a single system target port may have up to 16 NVMe I/O controllers.

When you install and configure attachments between the system and a host that runs the Linux operating system (OS), follow specific guidelines. For more information about these guidelines, see [Linux specific guidelines](#).

## 4.10.4 Security menu

Use the Security option from the **Settings** menu (as shown in Figure 4-66) to view and change security settings, authenticate users, and manage secure connections.

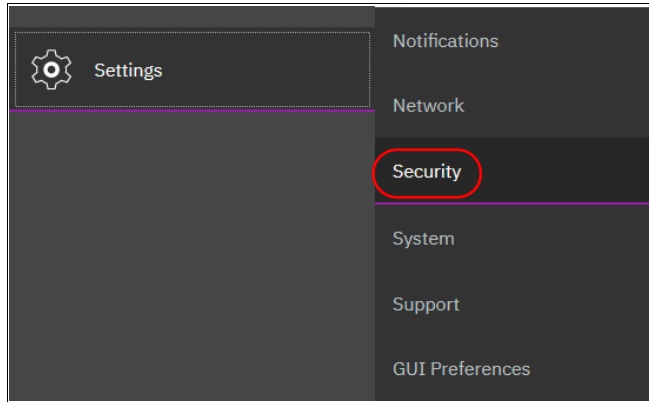


Figure 4-66 Security settings

### Remote Authentication

In the Remote Authentication pane, you can configure remote authentication with LDAP, as shown in Figure 4-67. By default, the system has local authentication that is enabled. When you configure remote authentication, you do not need to configure users on the system or assign more passwords. Instead, you can use your passwords and user groups that are defined on the remote service to simplify user management and access, enforce password policies more efficiently, and separate user management from storage management.

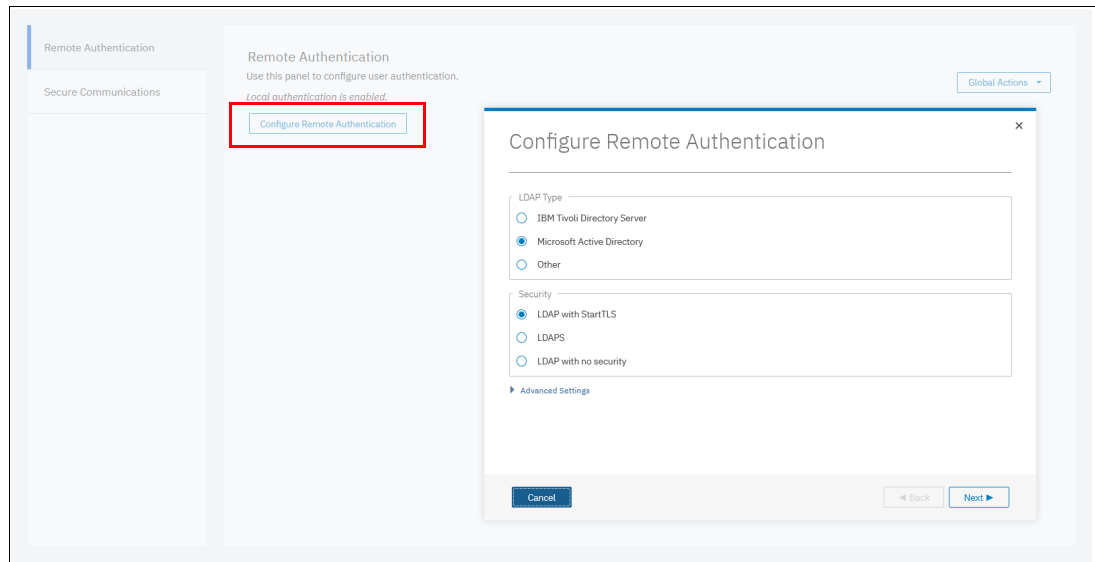


Figure 4-67 Configuring Remote Authentication

For more information about how to configure remote logon, see the following resources [IBM Knowledge Center](#).

## Encryption

As shown in Figure 4-68, you can enable or disable the encryption function on a SAN Volume Controller System. For more information, see Chapter 12, “Encryption” on page 685 for more information.

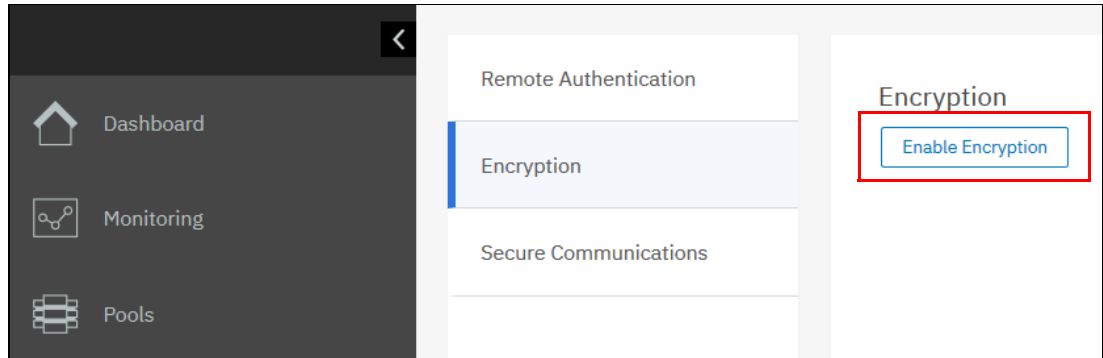


Figure 4-68 Enable Encryption

## Secure Communications

To enable or manage secure communications, select the **Secure Communications** pane, as shown in Figure 4-69. Before you create a request for either type of certificate, ensure that your current browser does not have restrictions about the type of keys that are used for certificates.

Some browsers limit the use of specific key-types for security and compatibility issues. Select **Update Certificate** to add new certificate details, including certificates that were created and signed by a third-party certificate authority (CA).

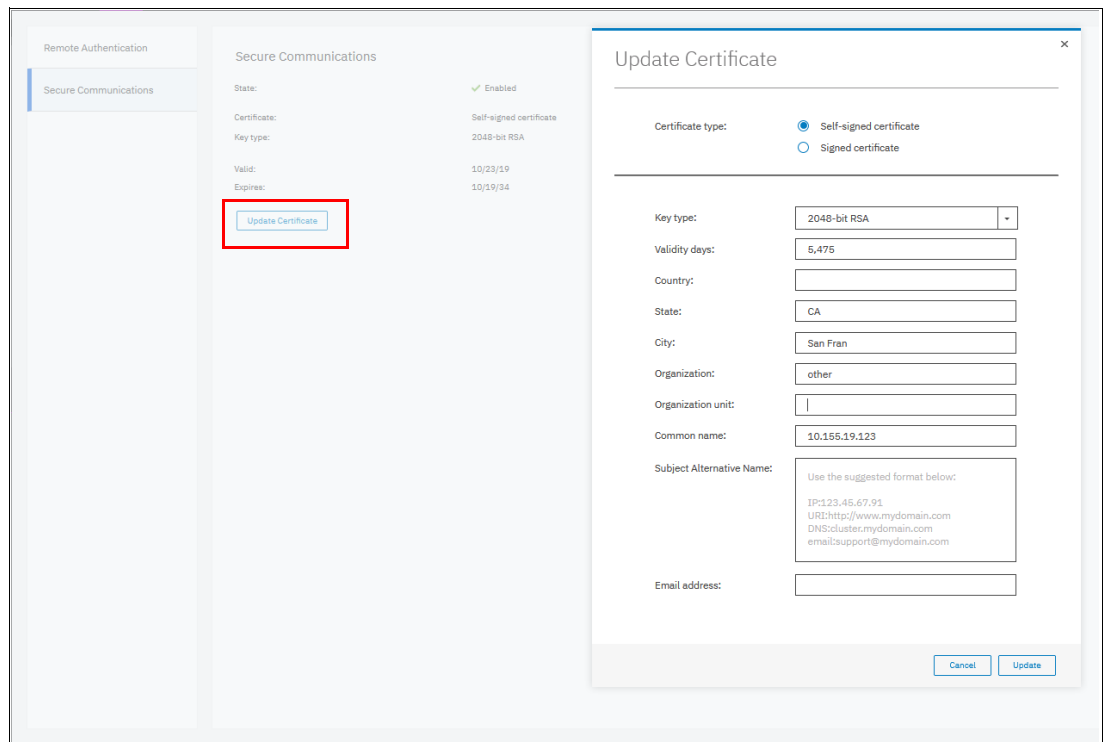


Figure 4-69 Configuring Secure Communications and Updating Certificates

## 4.10.5 System menus

Click the **System** option from the **Settings** menu (see Figure 4-70) to view and change the time and date settings, work with licensing options, download configuration settings, work with VMware VVOLs and IP Quorum, or download software upgrade packages.

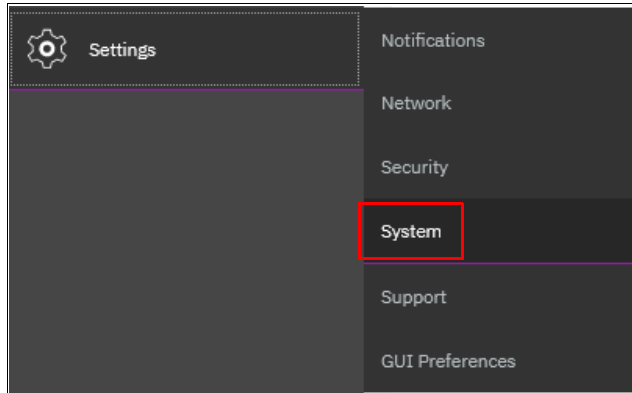


Figure 4-70 System option

### Date and time

Complete the following steps to view or configure the date and time settings:

1. From the SAN Volume Controller System pane, click **Settings** and click **System**.
2. In the left column, select **Date and Time**, as shown in Figure 4-71.

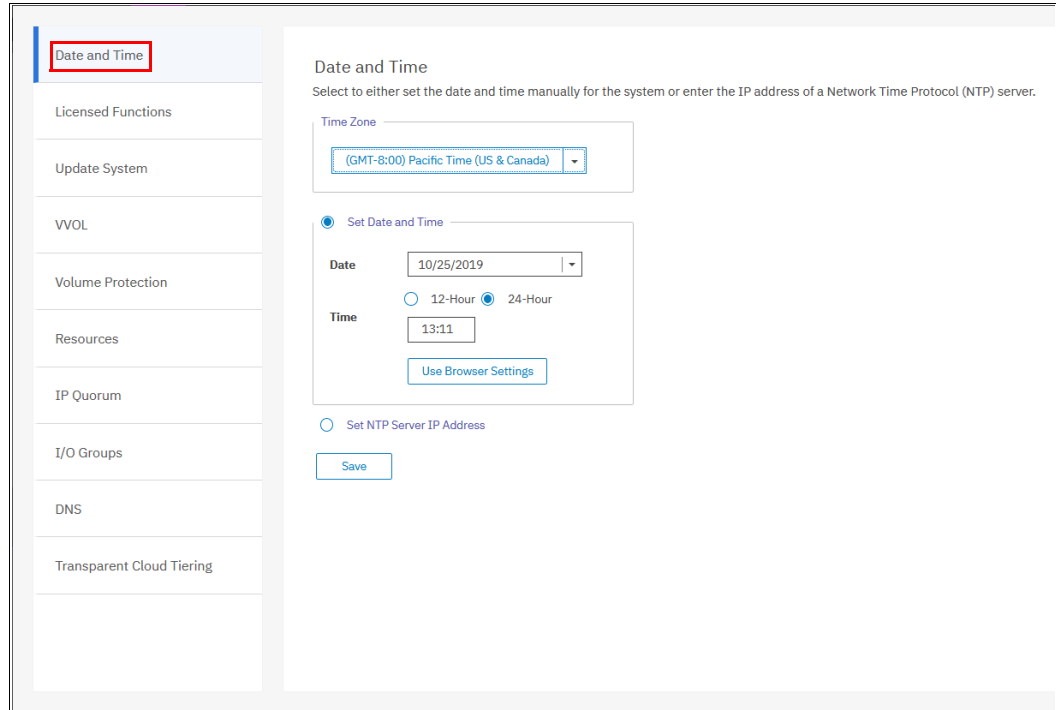


Figure 4-71 Date and Time window

3. From this pane, you can modify the following information:

- Time zone

Select a time zone for your system by using the drop-down list.

- Date and time

The following options are available:

- If you are not using a Network Time Protocol (NTP) server, select **Set Date and Time**, and then manually enter the date and time for your system, as shown in Figure 4-72. You can click **Use Browser Settings** to automatically adjust the date and time of your SAN Volume Controller system with your local workstation date and time.

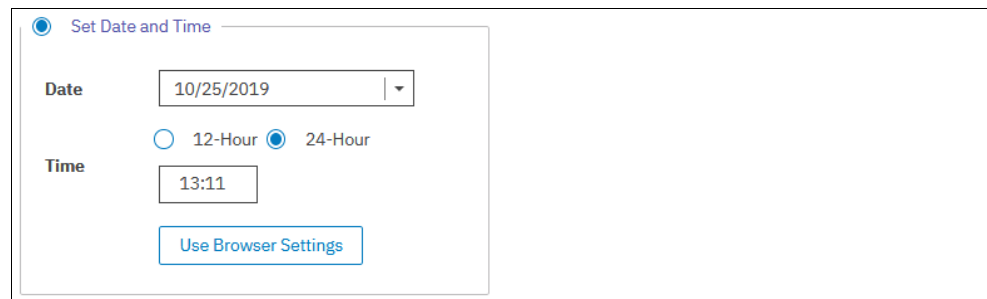


Figure 4-72 Set Date and Time window

- If you are using an NTP server, select **Set NTP Server IP Address**, and then enter the IP address of the NTP server, as shown in Figure 4-73.



Figure 4-73 Set NTP Server IP Address window

4. Click **Save**.

## Licensing

The base license that is provided with your system includes the use of its basic functions. However, the extra licenses can be purchased to expand the capabilities of your system. Administrators are responsible for purchasing extra licenses and configuring the systems within the license agreement, which includes configuring the settings of each licensed function on the system.

The SAN Volume Controller supports both differential and capacity-based licensing. For virtualization and compression functions, differential licensing charges different rates for different types of storage, which provides cost-effective management of capacity across multiple tiers of storage. Licensing for these functions are based on the number of Storage Capacity Units (SCUs) that are purchased. With other functions, like remote mirroring and FlashCopy, the license grants a specific number of terabytes for that function.

Differential licensing is granted per SCU. Each SCU corresponds to a different amount of usable capacity based on the type of storage.

Table 4-1 shows the different storage types and the associated SCU ratios.

Table 4-1 SCU ratio per storage type

Storage type	SCU ratio
Tier storage-class memory (SCM), Tier 0 and Tier 1	One SCU equates to 1 TiB of solid-state drive (SSD) or flash storage.
Enterprise Tier	One SCU equates to 1.18 TiB of Enterprise storage.
Nearline (NL) Tier	One SCU equates to 4.00 TiB of NL storage.

License settings are initially entered into a system initialization wizard. They can also be changed later.

To change the license settings or view the current license usage, select **Settings** → **System** → **Licensed Functions**, as shown in Figure 4-74.

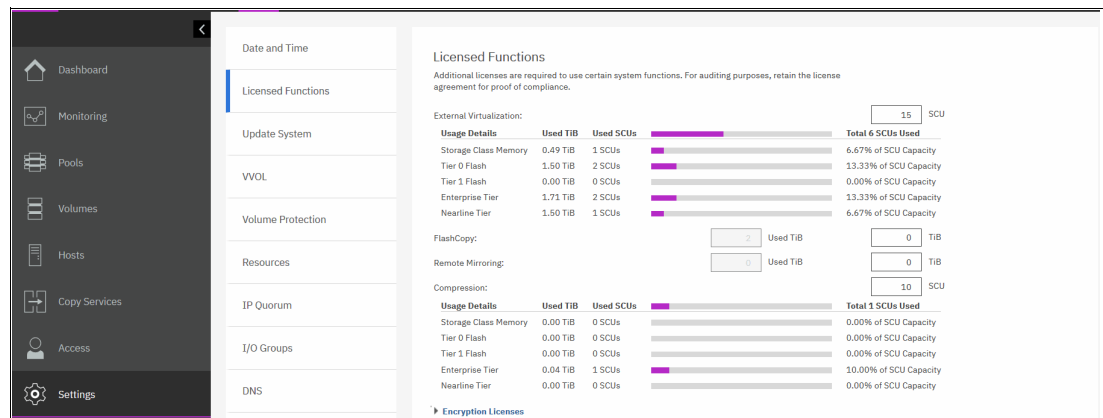


Figure 4-74 Licensed functions pane

The following elements are displayed:

► External Virtualization

You can enter the number of SCU units that are licensed for External Virtualization for your SAN Volume Controller environment.

When monitoring External Virtualization license usage, consider the following items:

- The license accounts for usable MDisk capacity. For example, one SCU is used for one 1 TB Tier0 MDisk when it is assigned to a storage pool, independently of the amount of the actual data written on the MDisk.
- SCU is used for complete and incomplete chunks of MDisk capacity. For example, if a combined capacity of all NL Tier MDisks in your system is 5 TB, two SCUs are needed: one SCU for 4 TB of NL storage, and one SCU for another 4 TB (even if 1 TB is used).

► FlashCopy Value

The FlashCopy function copies the contents of a source volume to a target volume. It is also used to create cloud snapshots of volumes in systems that have TCT enabled. FlashCopy is licensed in terabytes.

**Important:** The used capacity for FlashCopy mappings is the sum of all the volumes that are the source volumes of a FlashCopy mapping and volumes with cloud snapshots.

► Remote Mirroring Value

Enter the capacity that is available for MM and GM relationships. The remote-mirroring function configures a relationship between two volumes. This function mirrors updates that are made to one volume to another volume. The volumes can be in the same system or on two different systems. Remote mirroring is licensed in terabytes.

**Important:** The used capacity for GM and MM is the sum of the capacities of all of the volumes that are in an MM or GM relationship. Both master and auxiliary volumes are counted.

► Compression Value

With the compression function, data is compressed as it is written to disk, saving extra capacity for the system. A compression license can be purchased for a specific quantity of SCUs, which can be divided among different tiers of storage.

SCU usage for compression is calculated in the same way as for an External Virtualization license. But for a Compression License, *used capacity* (and not *usable capacity*) is calculated.

For example, if your compressed data occupies 4 TB in an NL tier, you need one SCU that is licensed for compression independently of the total compressed virtual disk (VDisk) capacity.

**Note:** Only IBM Real-Time Compression (RtC) for volumes in standard pools is accounted for on a Compression License. Compressed volumes in Data Reduction Pools (DRPs) are not accounted for by this license.

► Encryption Licenses

In addition to these enclosure-based licensed functions, the system also supports encryption through a key-based license. Key-based licensing requires an authorization code to activate encryption on the system. Only certain models of the control enclosures support encryption.

During system setup, you can select to activate the license with the authorization code. The authorization code is sent with the licensed function authorization documents that you receive after purchasing the license.

Encryption is activated on a per system basis, and an active license is required for each node that uses encryption. During system setup, the system detects the nodes that support encryption, and a license should be applied to each one. If more nodes are added and require encryption, more encryption licenses must be purchased and activated.

**Note:** To monitor license usage, you can run the `lslicense` CLI command, as described in [IBM Knowledge Center](#).

## Updating your system

For more information about the update procedure that uses the GUI, see Chapter 13, “Reliability, availability, and serviceability, and monitoring and troubleshooting” on page 753.

## VMware vSphere Virtual Volumes

IBM Spectrum Virtualize can manage VVOLs directly in cooperation with VMware. It enables VMware virtual machines (VMs) to get the assigned disk capacity directly SAN Volume Controller rather than from the ESXi data store. This technique enables storage administrators to control the appropriate usage of storage capacity, and to enable enhanced features of storage virtualization directly to the VM (such as replication, thin-provisioning, compression, and encryption).

VVOL management is enabled in the SAN Volume Controller in the System section, as shown in Figure 4-75. The NTP server must be configured before enabling VVOLs management. As a best practice, use the same NTP server for ESXi and your SAN Volume Controller.

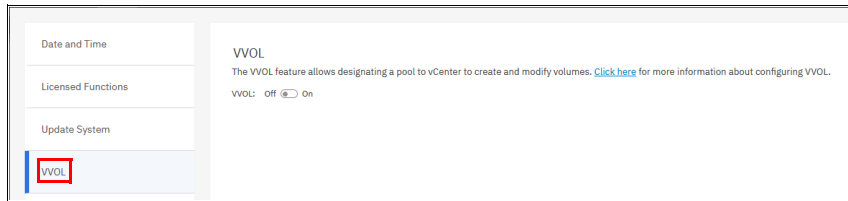


Figure 4-75 Enabling VVOLs management

**Restriction:** You cannot enable VVOLs support until the NTP server is configured in SAN Volume Controller.

For more information about VVOLs, see the following publications:

- ▶ *Quick-start Guide to Configuring VMware Virtual Volumes for Systems Powered by IBM Spectrum Virtualize*, REDP-5321.
- ▶ *Configuring VMware Virtual Volumes for Systems Powered by IBM Spectrum Virtualize*, SG24-8328.

## Volume protection

Volume protection prevents active volumes or host mappings from being deleted inadvertently if the system detects recent I/O activity.

**Note:** This global setting is enabled by default on new systems. You can either set this value to apply to all volumes that are configured on your system, or control whether the system-level volume protection is enabled or disabled on specific pools.

To prevent an active volume from being deleted unintentionally, administrators can use the system-wide setting to enable volume protection. They can also specify a period that the volume must be idle before it can be deleted. If volume protection is enabled and the period is not expired, the volume deletion fails even if you force the deletion.



**Note:** The system-wide volume protection and the pool-level protection must both be enabled for protection to be active on a pool. The pool-level protection depends on the system-level setting to ensure that protection is applied consistently for volumes within that pool. If system-level protection is enabled but pool-level protection is not enabled, any volumes in the pool can be deleted even when the setting is configured at the system level.

When you delete a volume, the system verifies whether it is a part of a host mapping, FlashCopy mapping, or remote-copy relationship. For a volume that contains these dependencies, the volume cannot be deleted unless the force option is selected. However, the force option does not delete a volume if it has recent I/O activity and volume protection is enabled. The force option overrides the volume dependencies, not the volume protection setting.

The following actions are affected by this setting:

- ▶ Deleting a volume
- ▶ Deleting a volume copy
- ▶ Deleting a host or a host cluster mapping
- ▶ Deleting a storage pool
- ▶ Deleting a host from an I/O group
- ▶ Deleting a host or host cluster
- ▶ Deleting a defined host port
- ▶ Creating a remote-copy relationship

Figure 4-76 shows the menu for volume protection.

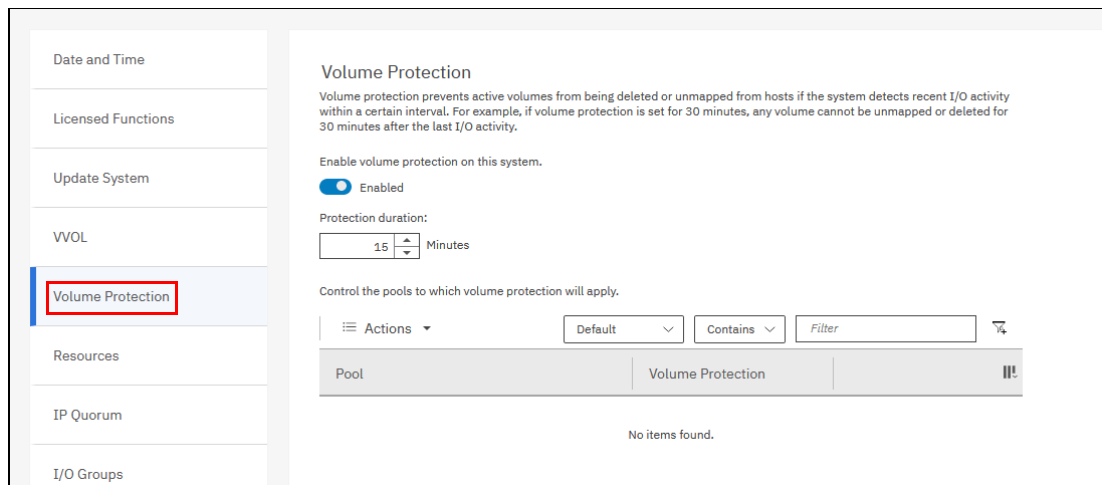


Figure 4-76 Volume protection

## Resources

Use this option to change the memory limits for the Copy Services and Redundant Array of Independent Disks (RAID) functions for an I/O group.

Copy Services features and RAID require that small amounts of volume cache be converted from cache memory into bitmap memory to enable the functions to operate. If you do not have enough bitmap space that is allocated when you try to use one of the functions, you cannot complete the configuration.

The total memory that can be dedicated to these functions is not defined by the physical memory in the system. The memory is constrained by the software functions that use the memory.

To use the Resource option, select **Settings** → **System** → **Resources**, as shown in Figure 4-77.

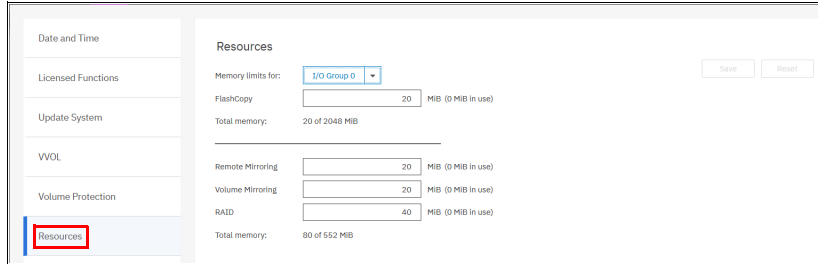


Figure 4-77 Resources allocation

Table 4-2 provides an example of the amount of memory that is required for remote mirroring functions, FlashCopy functions, and volume mirroring.

Table 4-2 Examples of allocation of bitmap memory

Function	Grain size [kibibyte (KiB)]	1 MiB of memory provides the following volume capacity for the specified I/O group
RC	256	2 TiB of total MM, GM, or HyperSwap volume capacity
FlashCopy	256	2 TiB of total FlashCopy source volume capacity
FlashCopy	64	512 GiB of total FlashCopy source volume capacity
Incremental FlashCopy	256	1 TiB of incremental FlashCopy source volume capacity
Incremental FlashCopy	64	256 GiB of incremental FlashCopy source volume capacity
Volume Mirroring	256	2 TiB of mirrored volume capacity

## IP Quorum

IBM Spectrum Virtualize also supports an IP quorum application. By using an IP-based quorum application as the quorum device for the third site, IBM FICON® is not required. Java applications run on hosts at the third site.

To install the IP quorum device, complete the following steps:

1. If your SAN Volume Controller is configured for IPv4, click **Download IPv4 Application**. If it is configured for IPv6, select **Download IPv6 Application**. In our example, IPv4 is the option, as shown in Figure 4-78.

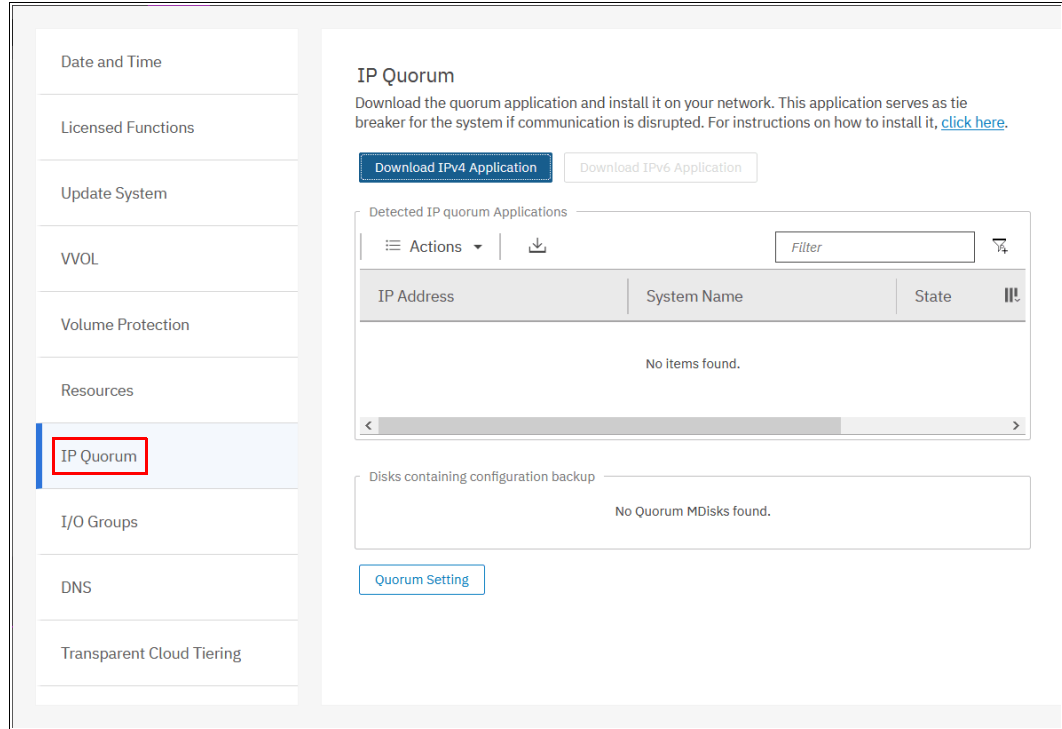


Figure 4-78 IP Quorum settings

2. After you select your correct IP configuration, IBM Spectrum Virtualize generates an IP Quorum Java application, as shown in Figure 4-79. The application can be saved and installed in a host that is to run the IP quorum application.

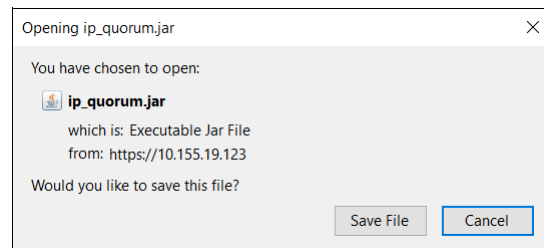


Figure 4-79 IP Quorum Java Application

3. After you download the IP quorum application, you must install the application on a separate host or server.
4. If you change the configuration by adding a node, changing a service IP address, or changing Secure Sockets Layer (SSL) certificates, you must download and install the IP quorum application again.
5. On the host, you must use the Java command line to initialize the IP quorum application. On the server or host on which you plan to run the IP quorum application, create a separate directory that is dedicated to the IP quorum application.

6. Run the **ping** command on the host server to verify that it can establish a connection with the service IP address of each node in the system.
7. Change to the folder where the application is, and run the following command:

```
java -jar ip_quorum.jar
```

**Note:** The IP quorum application always must be running.

8. To verify that the IP quorum application is installed and active, select **Settings** → **System** → **IP Quorum**. The new IP quorum application is displayed in the table of detected applications. The system automatically selects MDisks for quorum disks.

An IP quorum application can also act as the quorum device for systems that are configured with a single-site or standard topology that does not have any external storage configured. The IP quorum mode is set to Standard when the system is configured for standard topology; the quorum mode of Preferred or Winner is in effect only if the system topology is not set to Standard. To change the quorum mode for the IP quorum application, select **Quorum Setting** and set the mode to **Preferred** or **Winner**. This configuration gives the tie-breaker capability to a system, and automatically resumes I/O processing if half of the system's nodes or enclosures are inaccessible. For specific quorum settings, see Figure 4-80.

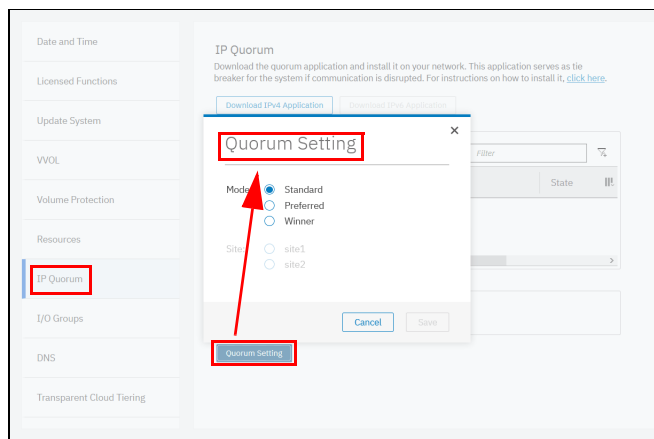


Figure 4-80 Quorum settings

## I/O groups

For ports within an I/O group, you can enable virtualization of FC ports that are used for host I/O operations. With N\_Port ID Virtualization (NPIV), the FC port consists of both a physical port and a virtual port. When port virtualization is enabled, ports do not come up until they are ready to handle I/O, which improves host behavior. In addition, path failures due to an offline node are masked from hosts.

The target port mode on the I/O group indicates the current state of port virtualization:

- ▶ **Enabled:** The I/O group contains virtual ports that are available to use.
- ▶ **Disabled:** The I/O group does not contain any virtualized ports.
- ▶ **Transitional:** The I/O group contains physical FC and virtual ports that are being used. You cannot change the target port mode directly from Enabled to Disabled states, or vice versa. The target port mode must be in a transitional state before it can be changed to Disabled or Enabled states.

The system can be in the transitional state for an indefinite period while the system configuration is changed. However, system performance can be affected because the number of paths from the system to the host doubled. To avoid increasing the number of paths substantially, use zoning or other means to temporarily remove some of the paths until the state of the target port mode is enabled.

The port virtualization settings of I/O groups are available by selecting **Settings** → **System** → **I/O Groups**, as shown in Figure 4-81.

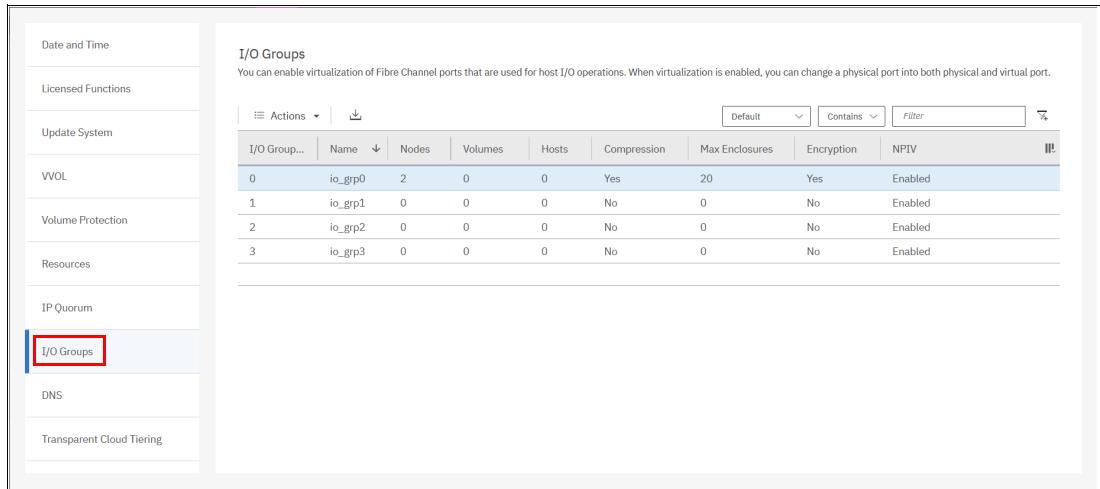


Figure 4-81 I/O groups port virtualization

You can change the status of the port by right-clicking the I/O group and selecting **Change NPIV Settings**, as shown in Figure 4-82.

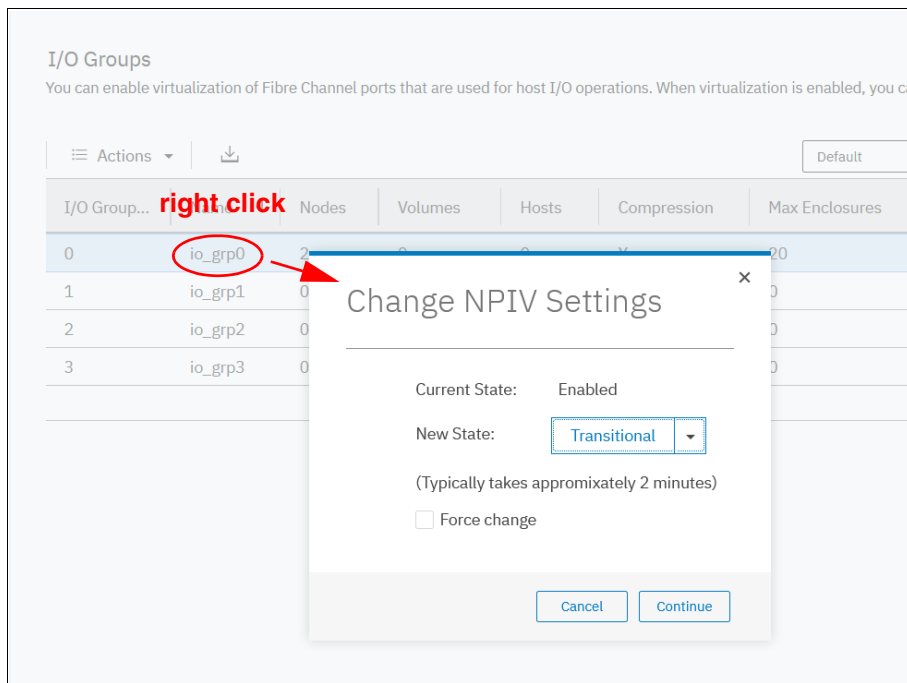


Figure 4-82 Changing the port mode

## Domain name server

IBM Spectrum Virtualize allows DNS entries to be manually set up in the SAN Volume Controller. The information about the DNS servers in the SAN Volume Controller helps the system to access the DNS servers to resolve names of the computer resources that are in the external network.

To view and configure DNS server information in IBM Spectrum Virtualize, complete the following steps:

1. In the left pane, click the **DNS** icon (see Figure 4-83) and enter the **IP address** and **Name** of each DNS server. IBM Spectrum Virtualize supports up two DNS servers for IPv4 or IPv6 (see Figure 4-84).

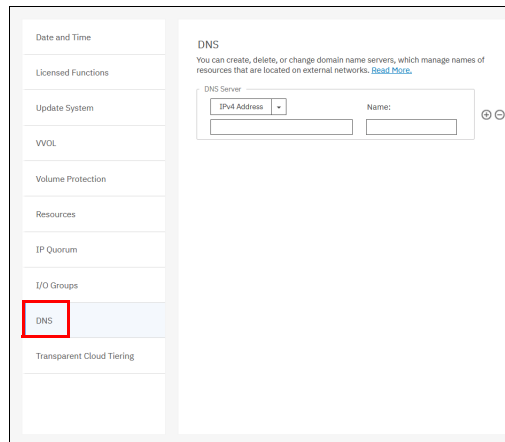


Figure 4-83 Selecting DNS

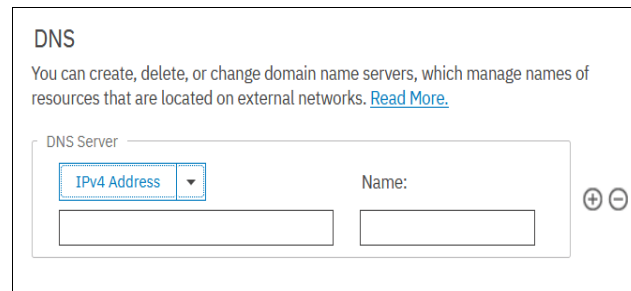


Figure 4-84 DNS information

- c. Click **Save** after you enter the DNS server information.

## Transparent Cloud Tiering

TCT is a licensed function that enables volume data to be copied and transferred to cloud storage. The system supports creating connections to cloud service providers (CSPs) to store copies of volume data in private or public cloud storage.

With TCT, administrators can move older data to cloud storage to free up capacity on the system. PiT snapshots of data can be created on the system and then copied and stored on the cloud storage. An external CSP manages the cloud storage, which reduces storage costs for the system. Before data can be copied to cloud storage, a connection to the CSP must be created from the system.

A cloud account is an object on the system that represents a connection to a CSP by using a particular set of credentials. These credentials differ depending on the type of CSP that is being specified. Most CSPs require the host name of the CSP and an associated password, and some CSPs also require certificates to authenticate users of the cloud storage.

Public clouds use certificates that are signed by well-known CAs. Private CSPs can use a self-signed certificate or a certificate that is signed by a trusted CA. These credentials are defined on the CSP and passed to the system through the administrators of the CSP. A cloud account defines whether the system can successfully communicate and authenticate with the CSP by using the account credentials.

If the system is authenticated, it can then access cloud storage to copy data to the cloud storage or restore data that is copied to cloud storage back to the system. The system supports one cloud account to a single CSP. Migration between providers is not supported.

**Note:** Before enabling TCT, consider the following requirements:

- ▶ Ensure that the DNS server is configured on your system and accessible.
- ▶ Determine whether your company's security policies require enabled encryption. If yes, ensure that the encryption licenses are properly installed and that the encryption is enabled.

Each CSP requires different configuration options. SAN Volume Controller supports the following CSPs:

- ▶ IBM Cloud (formerly known as IBM SoftLayer® Object Storage)  
The system can connect to IBM Cloud, which is a cloud computing platform that combines platform as a service (PaaS) with infrastructure as a service (IaaS).
- ▶ OpenStack Swift  
OpenStack Swift is a standard cloud computing architecture from which administrators can manage storage and networking resources in a single private cloud environment. Standard application programming interfaces (APIs) can be used to build customizable solutions for a private cloud solution.
- ▶ Amazon Simple Storage Service (Amazon S3)  
Amazon S3 provides programmers and storage administrators with flexible and secure public cloud storage. Amazon S3 is also based on Object Storage standards and provides a web-based interface to manage, back up, and restore data over the web.

To view your IBM Spectrum Virtualize cloud provider settings, from the SAN Volume Controller Settings pane, click **Settings** and select **System**. Then, select **Transparent Cloud Tiering**, as shown in Figure 4-85.

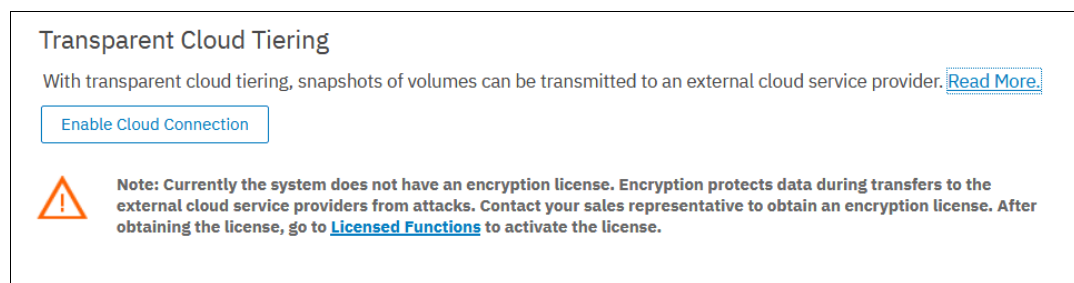


Figure 4-85 Transparent Cloud Tiering settings

By using this view, you can enable and disable features of your TCT and update the system information concerning your CSP. Use this pane to set the following options:

- ▶ CSP.
- ▶ Cloud Object Storage URL.
- ▶ The tenant or the container information that is associated to your Cloud Object Storage.
- ▶ User name of the cloud object account.
- ▶ Application programming interface (API) key
- ▶ The container prefix or location of your object.
- ▶ Encryption.
- ▶ Bandwidth.

For more information about how to configure and enable TCT, see 10.4, “Implementing Transparent Cloud Tiering” on page 565.

### **Automatic Configuration for Virtualization**

With automatic configuration, you can configure IBM FlashSystem 9100 or IBM Storwize systems as external storage to an existing SAN Volume Controller system. Automatic configuration implements optimal virtualization with a SAN Volume Controller system based on best practices. For these supported systems, the process can be completed in the management GUI during system setup or later as part of storage configuration.

The automatic configuration process is intended for new systems. If you want to virtualize this system by using a SAN Volume Controller system, no other objects, such as volumes or pools, can be configured on the system.

Before the wizard completes automatic configuration, you are prompted to complete the following prerequisite tasks:

- ▶ Add enclosures.

If you have any control or expansion enclosures to include as part of the external storage to be virtualized, you can add them. If you do not have more enclosures to add, this part can be skipped. If you do have enclosures to add but they are not automatically detected by the management GUI, verify the cabling between the system and the enclosures. For more information, see the installation information that came with the system.
- ▶ Verify zoning on the IBM SAN Volume Controller system.

In FICON, zoning is the process of grouping multiple ports to form a virtual, private storage network. Ports that are members of a zone can communicate with each other, but are isolated from ports in other zones. Before automatic configuration can be completed, verify that the SAN Volume Controller is zoned correctly as part of its SAN configuration.
- ▶ Define a host cluster to represent the IBM SAN Volume Controller  

As part of the prerequisite steps, you must create a host cluster that represents the SAN Volume Controller system. Creating a host cluster simplifies the port management between the SAN Volume Controller and the systems that are virtualized as part of the automatic configuration. In the host cluster that represents the SAN Volume Controller, a host represents a node, and each host port represents a port on a node. The management GUI displays WWPNs that are associated with node ports.

After these prerequisite steps are completed, the automatic configuration process begins. During this process, the following actions are completed automatically:

- ▶ Creates the appropriate RAID arrays based on the technology type of the drives.
- ▶ Creates a pool for each array.
- ▶ Provisions all usable capacity in each pool to volumes based on best practices.
- ▶ Maps all volumes to the SAN Volume Controller system for virtualization as MDisks.



After the process finishes, you are prompted to complete tasks on the SAN Volume Controller system to begin using this system as external storage.

Figure 4-86 shows how to enable Automatic Configuration for Virtualization.

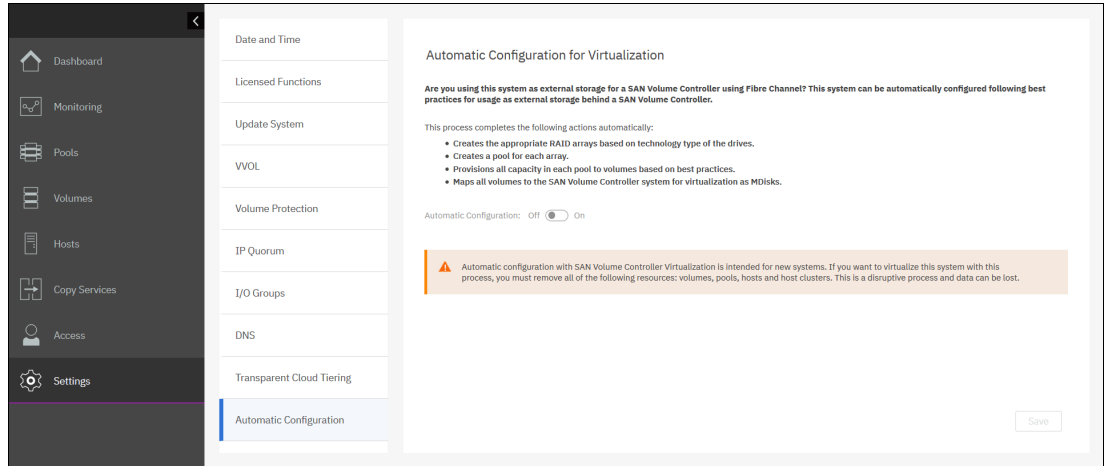


Figure 4-86 Automatic Configuration for Virtualization

## 4.10.6 Support menu

Use the Support pane to view and change Call Home settings, configure and manage connections, and upload support packages to the IBM Support Center.

The following options are available from the menu:

► Call Home

The Call Home feature transmits operational and event-related data to you and IBM through a Simple Mail Transfer Protocol (SMTP) server connection in the form of an event notification email. When configured, this function alerts IBM Support personnel about hardware failures and potentially serious configuration or environmental issues.

This view provides the following useful information about email notification and Call Home information (among others), as shown in Figure 4-87:

- IP of the email server (SMTP server) and port.
- Call Home email address.
- Email of one or more users set to receive one or more email notifications.
- Contact information of the person in the organization that is responsible for the system.

Call Home  
The support user receives call home events. Local users also receive event notifications.

Save Cancel

Transmission Settings  
 Send data using Call Home cloud service  
 Send data with email notifications

Call Home with cloud services  
Connection: ✔ Active  
Last Connection: Success at 10/25/2019 3:44:44 PM

Email Contact  
\* Contact Name \* Email Reply Address  
\* Telephone (Primary) \* Telephone (Alternate)  
\* Required

System Location  
\* Company Name \* Street Address  
\* City \* State or Province \* Postal Code  
\* Machine Location \* Country or Region  
United States  
\* Required

Figure 4-87 Call home settings

► Support Assistance

This option enables IBM Support personnel to access the system to complete troubleshooting and maintenance tasks. You can configure local Support Assistance, where IBM Support personnel visit your site to fix problems with the system, or Remote Support Assistance. Both local and Remote Support Assistance use secure connections to protect data exchange between the IBM Support Center and the system. More access controls can be added by the system administrator.

► Support Package

If Support Assistance is configured on your systems, you can automatically or manually upload new support packages to the IBM Support Center to help analyze and resolve errors on the system.

The menus are available by selecting **Settings** → **Support**, as shown in Figure 4-88 on page 185.

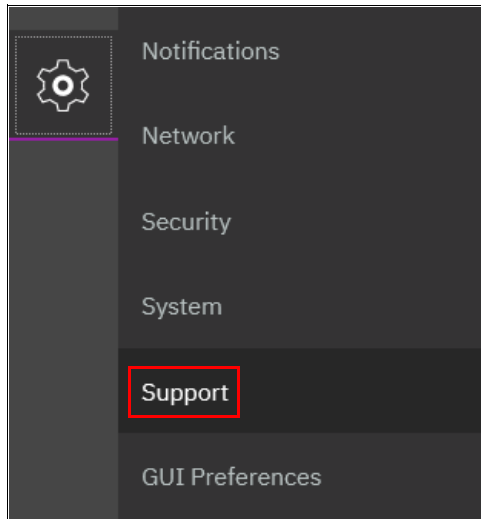


Figure 4-88 Support menu

For more information about how the Support menu helps with troubleshooting your system or how to back up your systems, see in Chapter 13, “Reliability, availability, and serviceability, and monitoring and troubleshooting” on page 753.

## 4.10.7 GUI Preferences menu

The **GUI Preferences** menu consists of two options:

- ▶ Login
- ▶ General

### Login message

IBM Spectrum Virtualize enables administrators to configure the welcome banner (login message). This message is a text message that appears in the GUI login window or at the CLI login prompt.

The content of the welcome message is helpful when you need to notify users about some important information about the system, such as security warnings or a location description. To define and enable the welcome message by using the GUI, edit the text area with the message content and click **Save** (see Figure 4-89).

Figure 4-89 Enabling the login message

The resulting login dialog box is shown in Figure 4-90. The system shows the welcome message in the GUI before login.

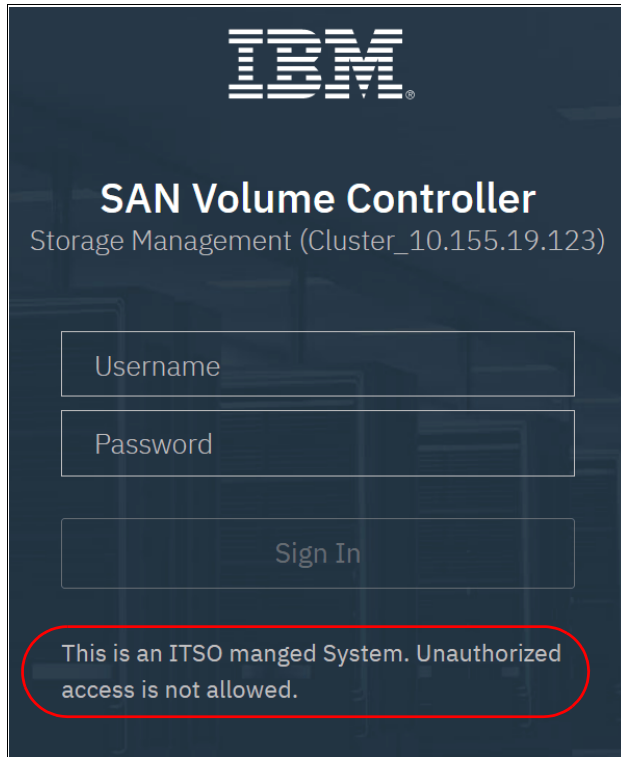


Figure 4-90 Welcome message in the GUI

Figure 4-91 shows the welcome message as it appears in the CLI.

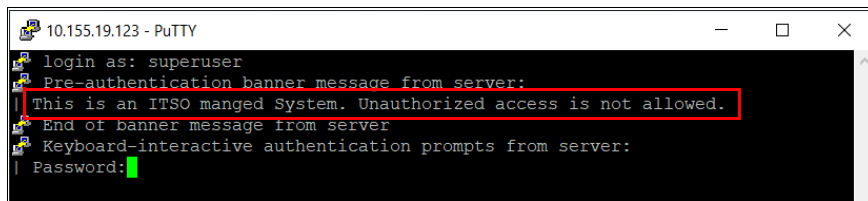


Figure 4-91 Welcome message in CLI

## General Settings

With the **General Settings** menu, you can refresh the GUI cache, set the low graphics mode option, and enable advanced pools settings.

To configure general GUI preferences, complete the following steps:

1. From the SAN Volume Controller Settings window, click **Settings** and select **GUI Preferences** (see Figure 4-92 on page 187).

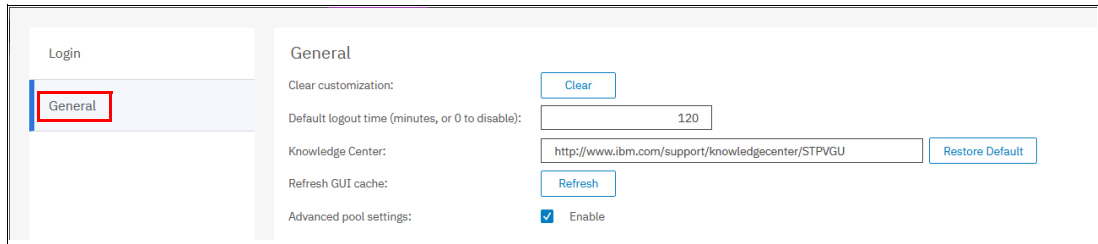


Figure 4-92 General GUI Preferences window

2. You can configure the following elements:

- Clear customizations

This option deletes all GUI preferences that are stored in the browser and restores the default preferences.

- Default logout time

Measured in minutes of inactivity in the established session.

- IBM Knowledge Center

You can change the URL of IBM Knowledge Center for IBM Spectrum Virtualize.

- Refresh GUI cache

This option causes the GUI to refresh all of its views and clears the GUI cache. The GUI looks up every object again. This option is useful if a value or object that is shown in the CLI is not being reflected in the GUI.

- Advanced pool settings

You can select the extent size during storage pool creation.

- Accessibility

Enables low graphics mode when the system is connected through a slower network.

## 4.11 Additional frequent tasks in the GUI

This section describes additional options and tasks that are available in the GUI of your SAN Volume Controller that are frequently used by administrators.

### 4.11.1 Renaming components

These sections provide guidance about how to rename your system and single nodes.

#### Renaming your SAN Volume Controller system

All objects in the SAN Volume Controller system have names that are user-defined or system-generated. Choose a meaningful name when you create an object. If you do not choose a name for the object, the system generates a name for you.

A well-chosen name serves both as a label for an object and as a tool for tracking and managing the object. Choosing a meaningful name is important if you decide to use configuration backup and restore.

When you choose a name for an object, apply the following naming rules:

- ▶ Names must begin with a letter.

**Important:** Do not start names by using an underscore (\_) character even though it is possible. Using an underscore as the first character of a name is a reserved naming convention that is used by the system configuration restore process.

- ▶ The first character cannot be numeric.
- ▶ The name can be a maximum of 63 characters, but there are exceptions. The name can be a maximum of 15 characters for RC relationships and groups. The `lsfabric` command displays long object names that are truncated to 15 characters for nodes and systems. (`lsrelationshipcandidate` or `lsrelationship` commands).
- ▶ Valid characters are uppercase letters (A - Z), lowercase letters (a - z), digits (0 - 9), the underscore (\_) character, a period (.), a hyphen (-), and a space.
- ▶ Names must not begin or end with a space.
- ▶ Object names must be unique within the object type. For example, you can have a volume that is called ABC and an MDisk called ABC, but you cannot have two volumes that are called ABC.
- ▶ The default object name is valid (an object prefix with an integer).
- ▶ Objects can be renamed to their current names.

To rename the system from the System pane, complete the following steps:

1. Click **Actions** in the upper-left corner of the SAN Volume Controller System pane, as shown in Figure 4-93.

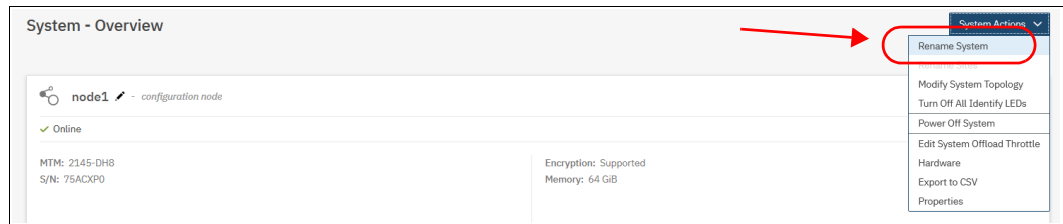


Figure 4-93 Actions on the System pane

2. The Rename System pane opens (see Figure 4-94). Specify a new name for the system and click **Rename**.

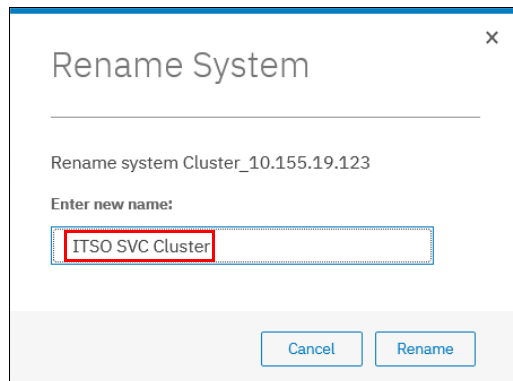


Figure 4-94 Renaming the system

**System name:** You can use the letters A - Z and a - z, the numbers 0 - 9, and the underscore (\_) character. The clustered system name can be 1 - 63 characters.

3. If you are certain that you want to change the system name click **Yes**.

**Warning:** When you rename your system, the iSCSI name automatically changes because it includes the system name by default. Therefore, this change needs more actions on iSCSI-attached hosts.

## Renaming a node

To rename a node, complete the following steps:

1. Go to the System View window and select one of the nodes. The Node Details pane for this node opens, as shown in Figure 4-95. Click **Rename**.

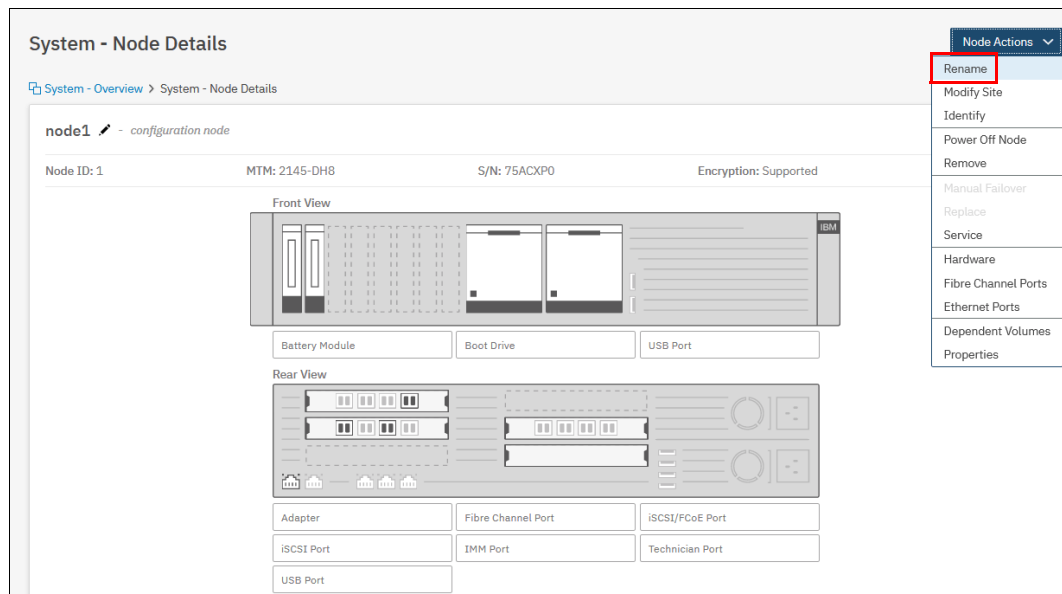


Figure 4-95 Renaming a node on the Node Details pane

2. Enter the new name of the node and click **Rename** (Figure 4-96).

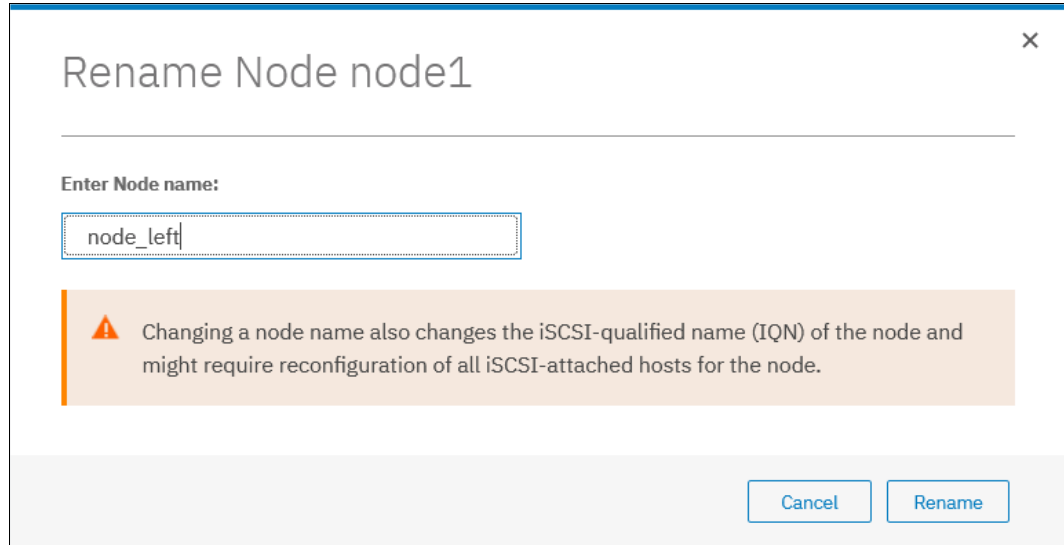


Figure 4-96 Entering the new node name

**Warning:** Changing the SAN Volume Controller node name causes an automatic IQN update and requires the reconfiguration of all iSCSI-attached hosts.

### Renaming sites

SAN Volume Controller supports configuring site settings that describe the location of the nodes and storage systems that are deployed in a stretched system configuration. This site information configuration is part of the configuration process for enhanced systems. The site information makes it possible for SAN Volume Controller to manage and reduce the amount of data that is transferred between the two sides of the system, which reduces the costs of maintaining the system.

**Note:** The renaming sites option is available only in a Stretched or HyperSwap topology.

Three site objects are automatically defined by SAN Volume Controller and numbered 1, 2, and 3. SAN Volume Controller creates the corresponding default names, `site1`, `site2`, and `site3` for each of the site objects. The `site1` and `site2` names are the two sites that make up the two halves of the enhanced system, and `site3` is the quorum disk. You can rename the sites to describe your data center locations.

To rename the sites, in the System pane, select **Actions** in the upper-left corner, and select **Rename Sites**, as shown in Figure 4-97.

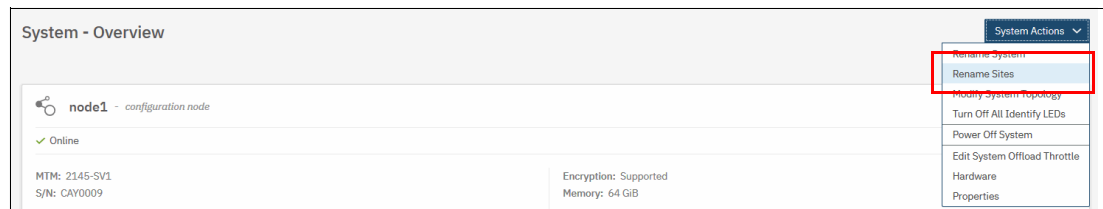


Figure 4-97 Rename Sites action



The Rename Sites window with the site information opens, as shown in Figure 4-98.

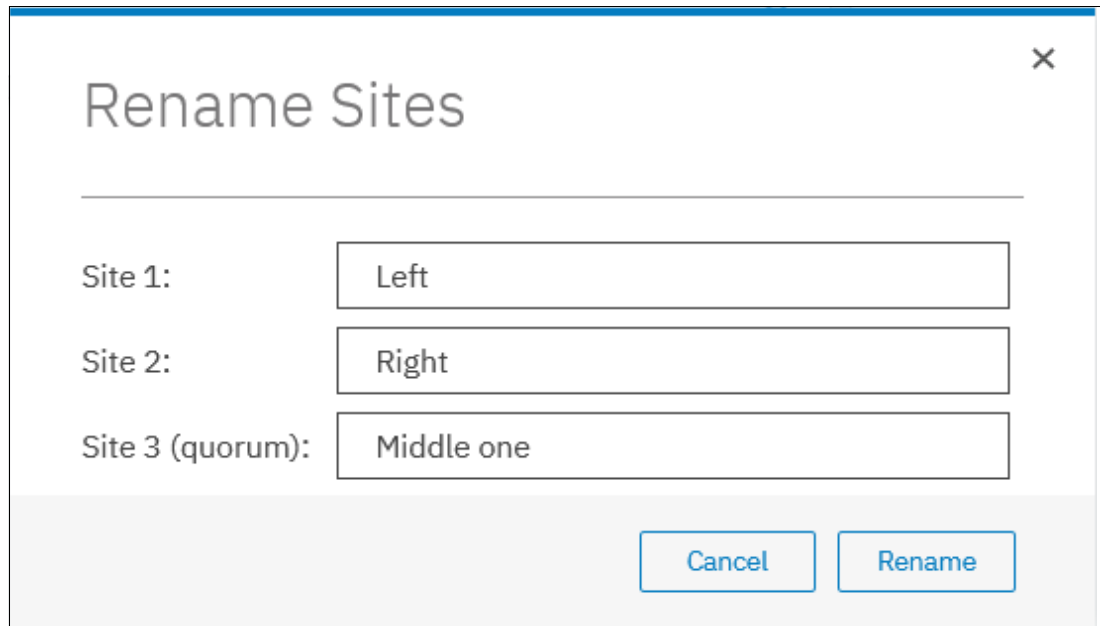


Figure 4-98 Rename Sites window

You can rename all or selected sites. Save your settings by clicking **Rename**.

### 4.11.2 Changing the system topology

You can create an enhanced resilient system configuration where each node on the system is physically on a different site. When used with mirroring technologies, such as volume mirroring or Copy Services, these configurations can be used to maintain access to data on the system in the event of power failures or site-wide outages.

There are two options for an enhanced resiliency configuration:

- ▶ A Stretched topology is ideal for a disaster recovery (DR) solution.
- ▶ HyperSwap fulfills high availability (HA) requirements.

If you prepared your infrastructure to support enhanced topology, you can proceed with the topology change by completing the following steps:

1. From the System pane, select **SystemActions** → **Modify System Topology**, as shown in Figure 4-99.

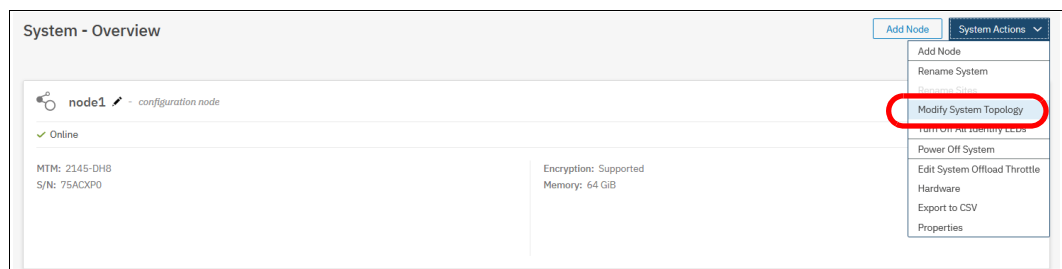


Figure 4-99 Modifying system topology

- The wizard opens and informs you about options to change the topology to either a Stretched Cluster or HyperSwap (see Figure 4-100).

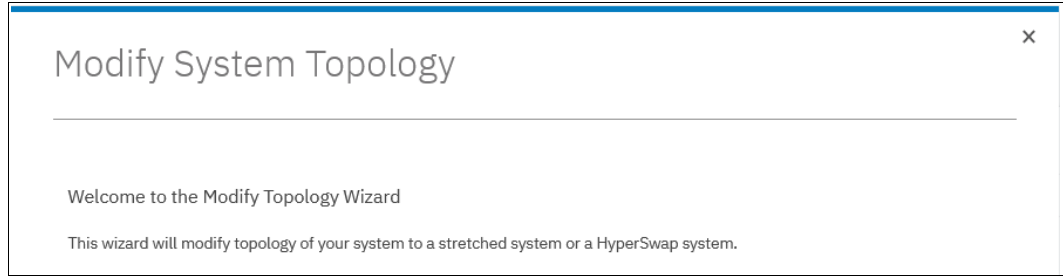


Figure 4-100 Topology wizard

- The system requires you to define three sites: Primary, Secondary, and Quorum sites. Assign reasonable names to the sites for easy identification, as shown in our example in Figure 4-101.

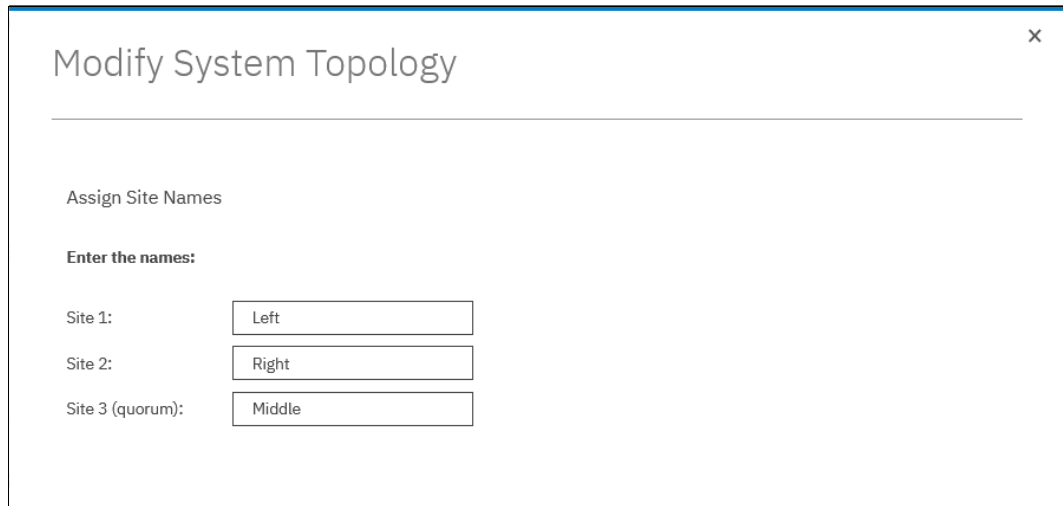


Figure 4-101 Assigning site names

4. Choose the topology. Although Stretched Cluster is optimal for DR solutions with asynchronous replication of primary volumes, HyperSwap is ideal for HA solutions with near-real-time replication. In our example, we decide to use a Stretched Cluster configuration (Figure 4-102).

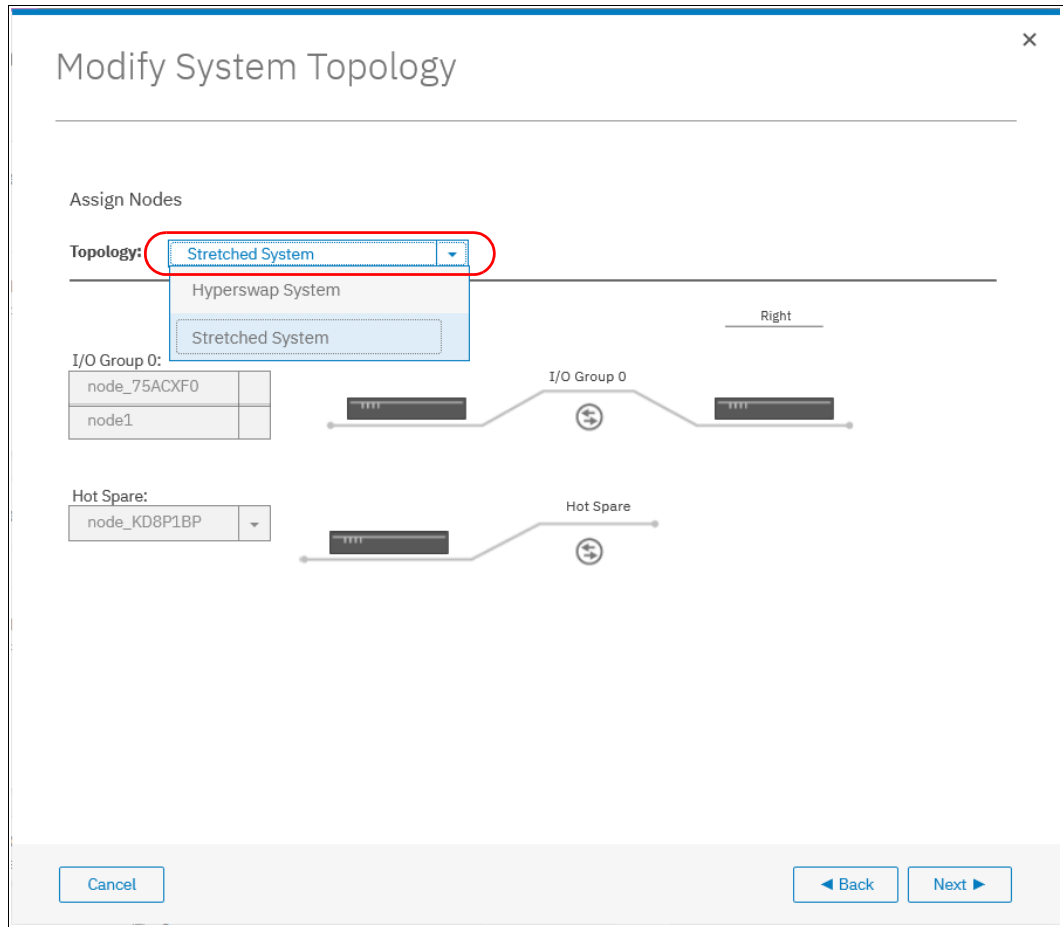


Figure 4-102 Changing the topology

- Assign hosts to primary sites. Right-click each host and modify sites for them one by one (see Figure 4-103). Also, assign offline hosts to primary sites because they might be down for maintenance or any other reason.

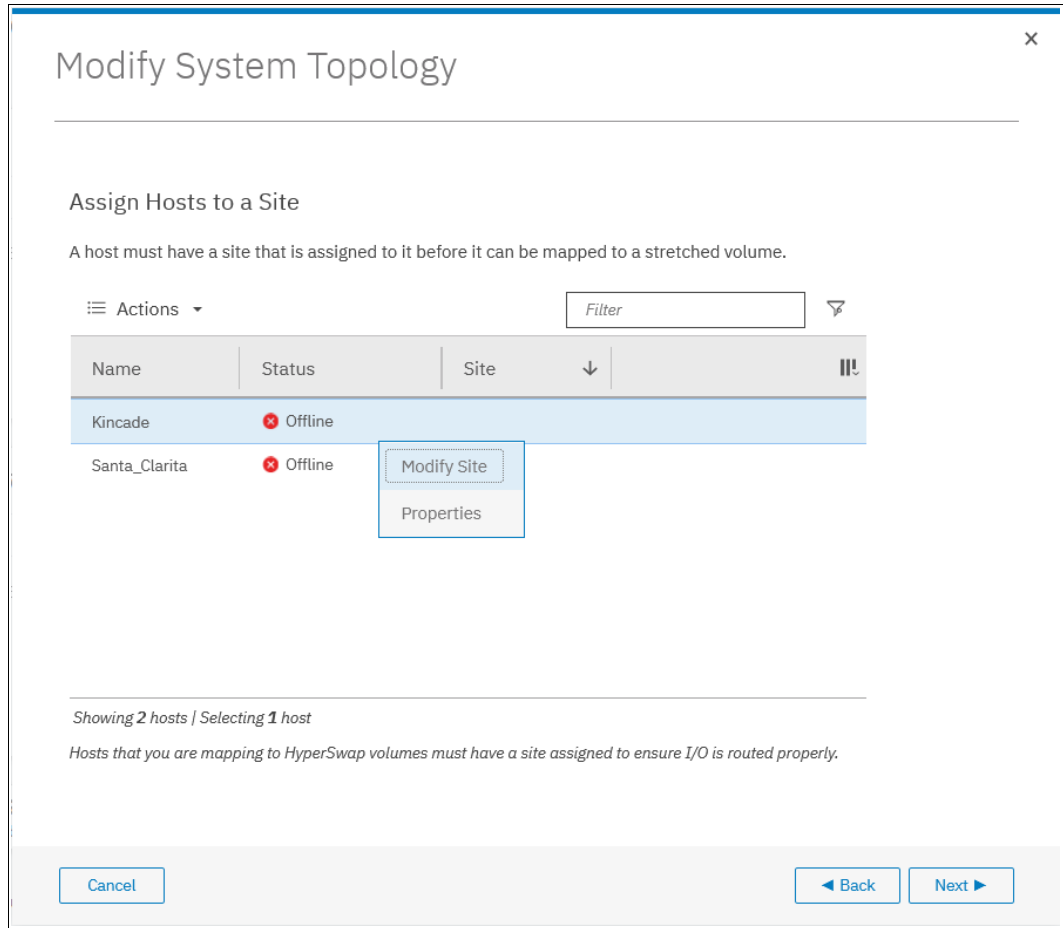


Figure 4-103 Assigning hosts to sites

The system view changes after you modify the Systems topology, as shown in Figure 4-104.

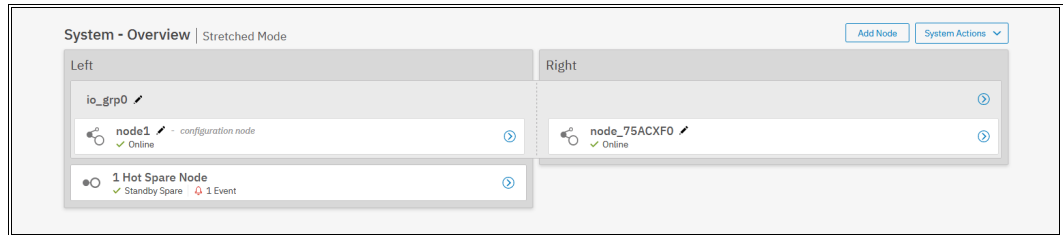


Figure 4-104 Modified system view

- Similarly, assign back-end storage to sites where the primary volumes will be provisioned (that is, where the hosts are primarily located), as shown in Figure 4-105. At least one storage device must be assigned to the site that is planned for Quorum volumes.

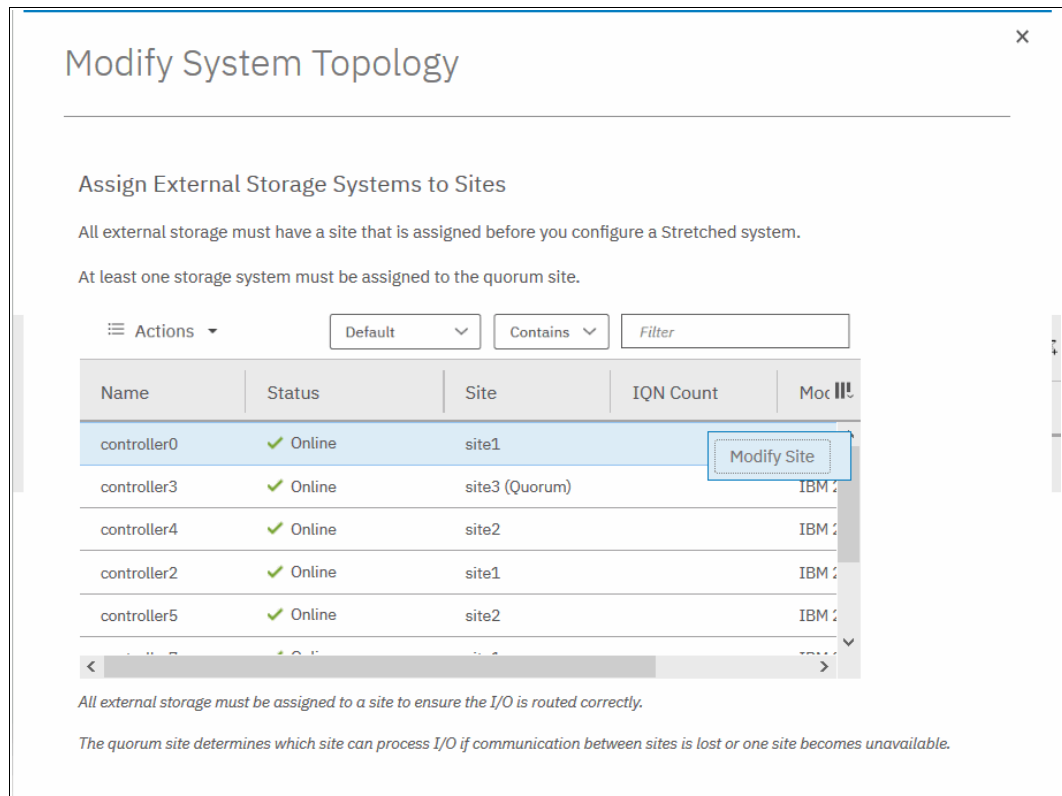


Figure 4-105 Assigning storage to sites

- After this process is complete, the Summary window opens (Figure 4-106). The system is ready to commit the changes. Click **Finish** to complete.

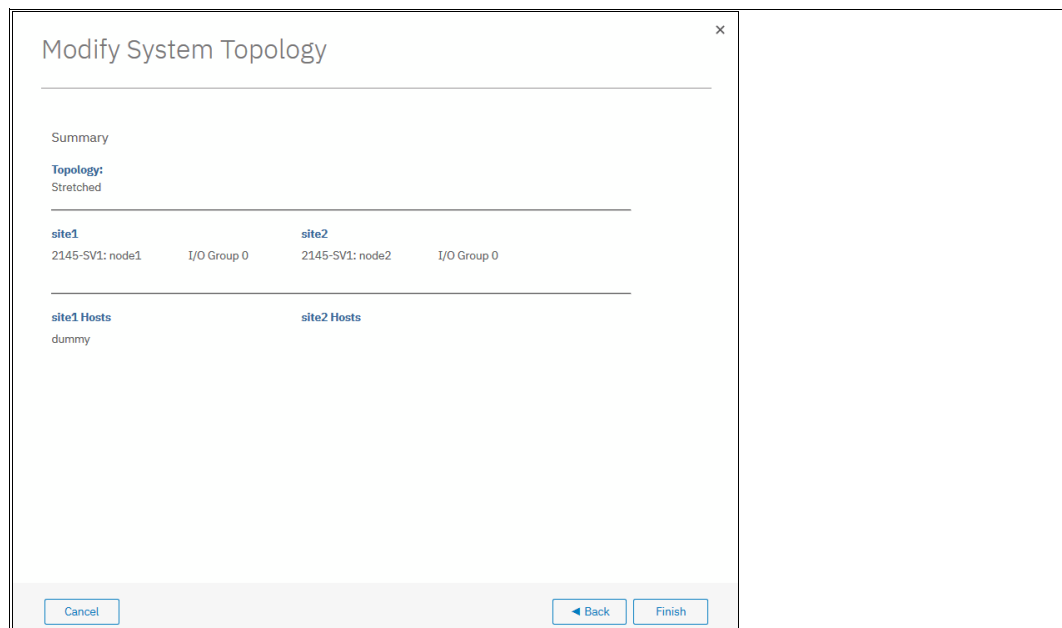


Figure 4-106 Summary of system topology changes

8. After the operation completes, select **System actions** → **Properties** to ensure that the topology now shows as Single Site, as shown in Figure 4-107.

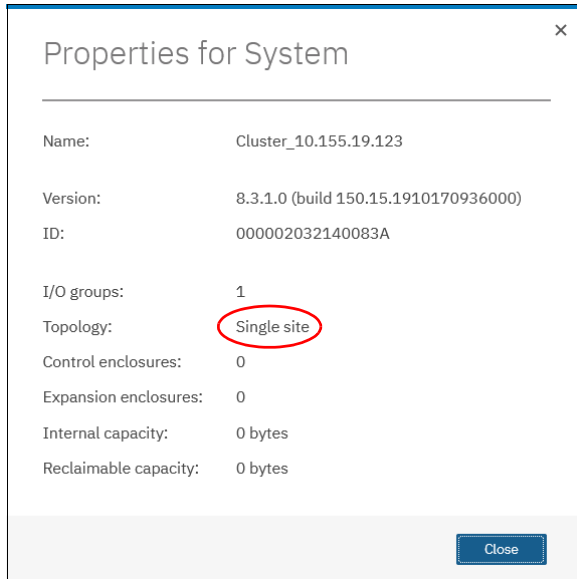


Figure 4-107 Checking system properties

As a validation step, verify that all hosts have the correctly mapped and active online volumes and that no error appears in the event log.

**Note:** If you want to change your Stretched Cluster configuration to a HyperSwap configuration or vice versa, you must first configure the topology as a single-site topology. Use the wizard, which provides only the option to modify the topology as a single site before you can select a HyperSwap or Stretched Cluster configuration.

For more information about resilient solutions for your SAN Volume Controller environment, see the following publications:

- ▶ *IBM Spectrum Virtualize and SAN Volume Controller Enhanced Stretched Cluster with VMware*, SG24-8211
- ▶ *IBM Storwize V7000, Spectrum Virtualize, HyperSwap, and VMware Implementation*, SG24-8317

### 4.11.3 Restarting the GUI service

The service that runs that GUI operates from the configuration node. Occasionally, you might need to restart this service if the GUI is not performing to your expectation (or you cannot connect). To do this task, complete the following steps:

1. Log in to the Service Assistant and identify the configuration node, as shown in Figure 4-108 on page 197.

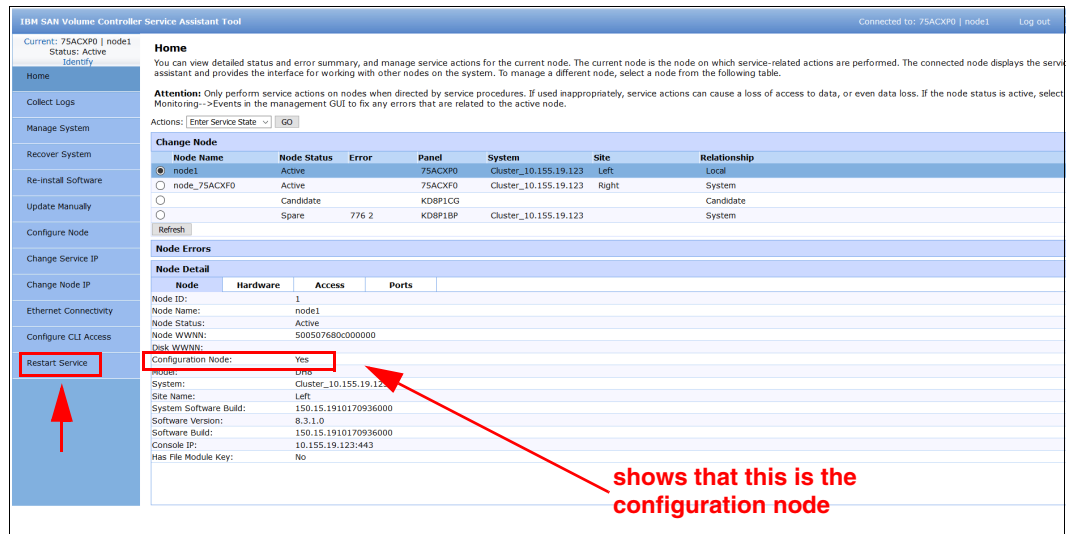


Figure 4-108 Identifying the config node on the service assistant

2. After the process completes, go to **Restart Service**, as shown in Figure 4-109.

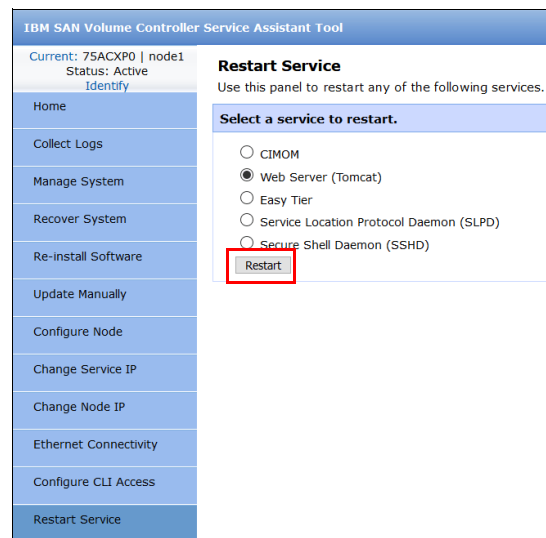


Figure 4-109 Restarting the Tomcat web server

3. Select **Web Server (Tomcat)**. Click **Restart**, and the web server that runs the GUI restarts. This task is a concurrent action, but the cluster GUI is unavailable while the server is restarting (the Service Assistant and CLI are not affected). After 5 minutes, check to see whether GUI access was restored.







## Storage pools

This chapter describes how IBM SAN Volume Controller manages physical storage resources. All storage resources that are under system control are managed by using *storage pools* or *managed disk (MDisk) groups (MDGs)*.

Storage pools aggregate internal and external capacity and provide the containers in which you can create volumes. Storage pools make it easier to dynamically allocate resources, maximize productivity, and reduce costs.

You can configure storage pools through the management GUI, either during initial configuration or later. Alternatively, you can configure the storage to your own requirements by using the command-line interface (CLI).

This chapter includes the following topics:

- ▶ 5.1, “Working with storage pools” on page 200
- ▶ 5.2, “Working with external controllers and MDisks” on page 217
- ▶ 5.3, “Working with internal drives and arrays” on page 231

## 5.1 Working with storage pools

Storage pools act as containers for MDisks, which provide storage capacity to the pool, and volumes that are provisioned from this capacity, which can be mapped to host systems. The system organizes storage in this fashion to ease storage management and make it more efficient.

MDisks can either be redundant array of independent disks (RAID) arrays that are created by using internal storage, such as drives and flash modules, or logical units (LUs) that are provided by external storage systems. A single storage pool can contain both types of MDisks, but a single MDisk can be part of only one storage pool. MDisks themselves are not visible to host systems.

Figure 5-1 provides an overview of how storage pools, MDisks, and volumes are related.

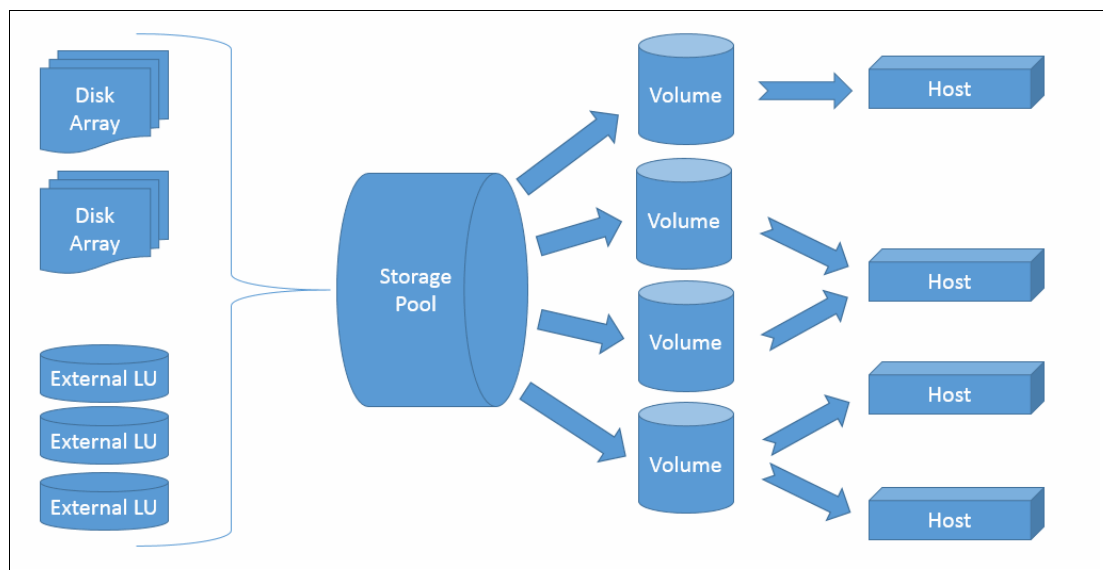


Figure 5-1 Relationship between MDisks, storage pools, and volumes

All MDisks in a pool are split into chunks of the same size, which are called *extents*. Volumes are created from the set of available extents in the pool. The extent size is a property of the storage pool and cannot be changed after the pool is created. The choice of extent size affects the total amount of storage that can be managed by the system.

It is possible to add MDisks to an existing pool to provide more usable capacity in the form of extents. The system automatically balances volume extents between the MDisks to provide the best performance to the volumes. It is also possible to remove extents from the pool by deleting an MDisk. The system automatically migrates extents that are in use by volumes to other MDisks in the same pool to make sure that the data on the extents is preserved.

A storage pool represents a failure domain. If one or more MDisks in a pool become inaccessible, all volumes (except for image mode volumes) in that pool are affected. Volumes in other pools are unaffected.

The system supports *standard pools* (parent pools and child pools) and *Data Reduction Pools* (DRPs).

Child pools are created from existing capacity that is assigned to a parent pool instead of being created directly from MDisks. When the child pool is created, the capacity for a child pool is reserved from the parent pool. This capacity is no longer reported as available capacity of the parent pool. In terms of volume creation and management, child pools are similar to parent pools.

DRPs use a set of techniques that can be used to reduce the amount of usable capacity that is required to store data, such as compression and deduplication. Data reduction can increase storage efficiency and performance, and reduce storage costs, especially for flash storage. DRPs automatically reclaim capacity that is no longer needed by host systems. This reclaimed capacity is given back to the pool as usable capacity and can be reused by other volumes. A DRP cannot be used as a parent pool to create a child pool.

For more information about DRP planning and implementation, see Chapter 9, “Advanced features for storage efficiency” on page 449 and *Introduction and Implementation of Data Reduction Pools and Deduplication*, SG24-8430.

In general, you manage storage as follows:

1. Create storage pools (standard or DRP), depending on your requirements and sizing.
2. Assign storage to these pools by using one or more of the following options:
  - Create array MDisks from internal drives or flash modules.
  - Add MDisks provisioned from external storage systems.
3. Create volumes in these pools and map them to hosts or host clusters.

You manage storage pools either in the Pools pane of the GUI or by using the CLI. To access the Pools pane, select **Pools** → **Pools**, as shown in Figure 5-2.

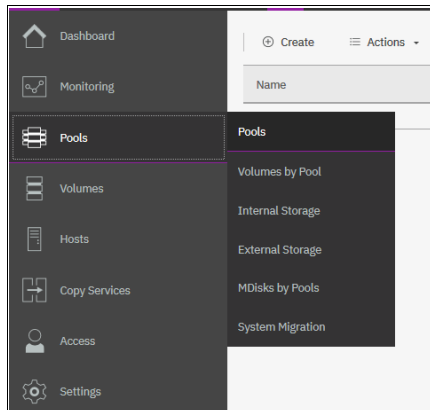


Figure 5-2 Accessing the Pools pane

The pane lists all storage pools and their major parameters. If a storage pool has child pools, they are also shown.

To see a list of configured storage pools by using the CLI, run the `lsmdiskgrp` command without any parameters, as shown in Example 5-1.

Example 5-1 The `lsmdiskgrp` output (some columns are not shown)

```

IBM_2145:ITS0-SV1:superuser>lsmdiskgrp
id name  status mdisk_count vdisk_count capacity extent_size free_capacity
0 Pool0 online 1 5 821.00GB 1024 771.00GB
1 Pool1 online 1 5 7.21TB 4096 6.14TB
  
```

## 5.1.1 Creating storage pools

To create a storage pool, complete the following steps:

1. Select **Pools** → **MDisks by Pools** and click **Create Pool**, or select **Pools** → **Pools** and click **Create**, as shown in Figure 5-3.

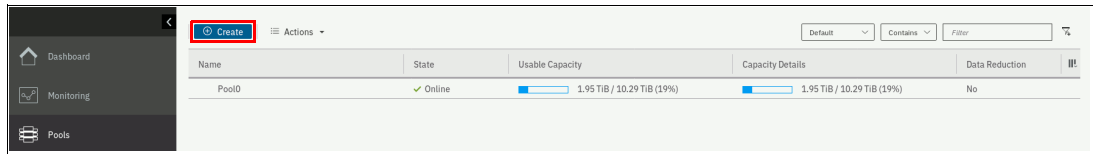


Figure 5-3 Option to create a storage pool in the Pools pane

Both alternatives open the dialog box that is shown in Figure 5-4.

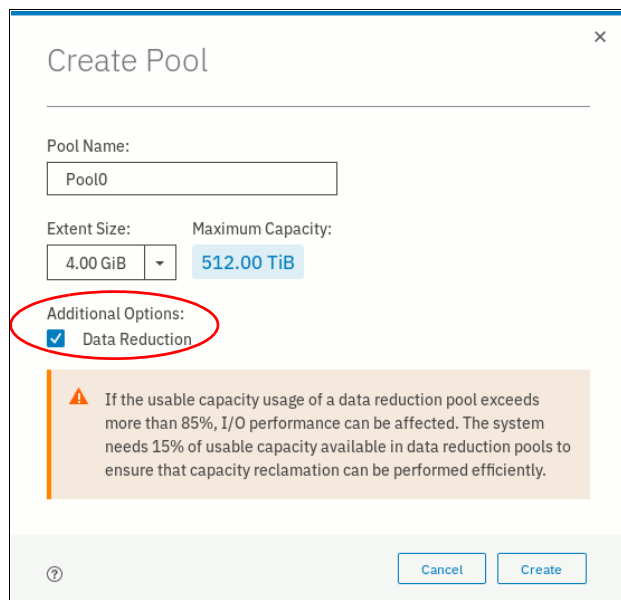


Figure 5-4 Create Pool dialog box

2. Select the **Data reduction** check box to create a DRP. Leaving it clear creates a standard storage pool.

**Note:** DRPs require careful planning and sizing. Limitations and performance characteristics of DRPs are different from standard pools.

A standard storage pool that is created by using the GUI has a default extent size of 1 GB. DRPs have a default extent size of 4 GB. The size of the extents is selected at creation time and cannot be changed later. The extent size controls the maximum total storage capacity that is manageable per system (across all pools). For DRPs, the extent size also controls the maximum capacity after reduction in the pool itself.

For more information about the differences between standard pools and DRPs and for extent size planning, see Chapter 2, “Planning” on page 53.

**Note:** Do not create DRPs with small extent sizes. For more information, see this [IBM Support alert](#).

When creating a standard pool, you cannot change the extent size by using the GUI by default. If you want to specify a different extent size, enable this option by selecting **Settings** → **GUI Preferences** → **General** and checking **Advanced pool settings**, as shown in Figure 5-5. Click **Save**.

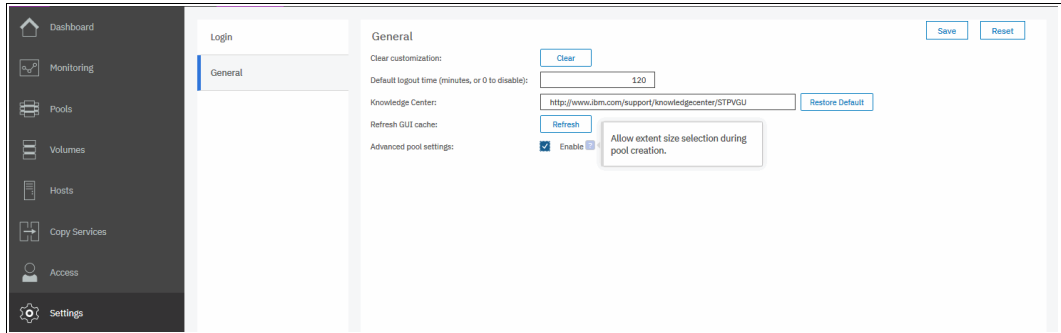


Figure 5-5 Advanced pool settings

When the **Advanced pool settings** option is enabled, you can also select an extent size for standard pools at creation time, as shown in Figure 5-6.

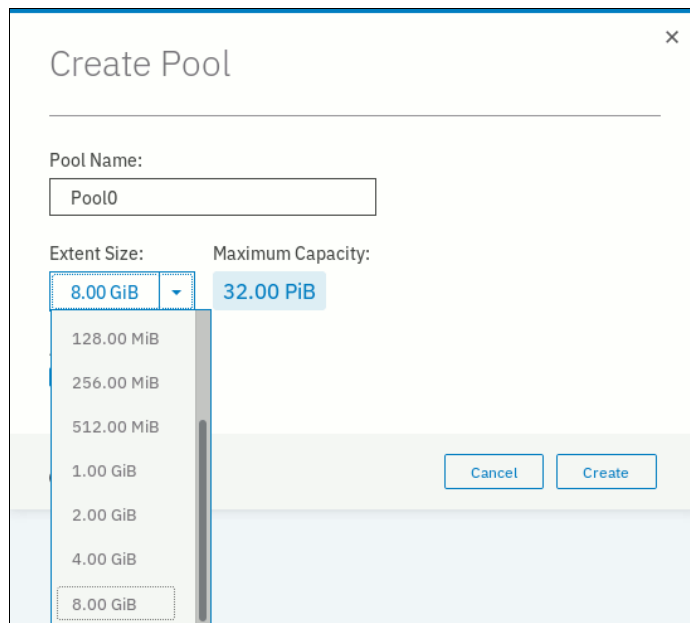


Figure 5-6 Creating a pool with Advanced settings selected

If an encryption license is installed and enabled, you can select whether the storage pool is encrypted, as shown in Figure 5-7 on page 204. The encryption setting of a storage pool is selected at creation time and cannot be changed later. By default, if encryption is enabled, encryption is selected. For more information about encryption and encrypted storage pools, see Chapter 12, “Encryption” on page 685.

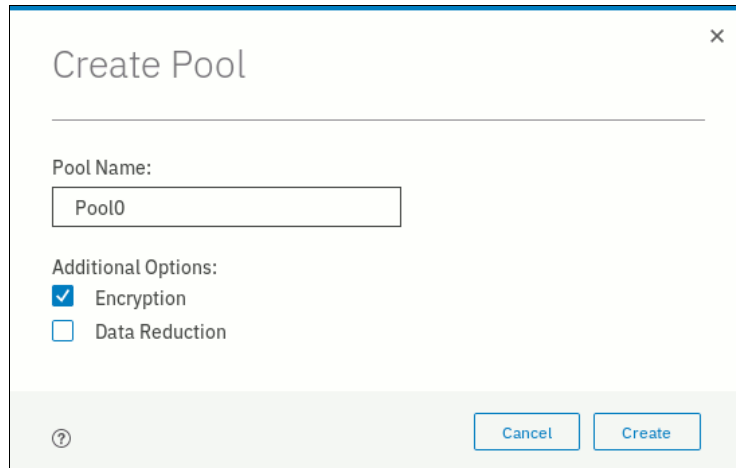


Figure 5-7 Creating a pool with encryption enabled

3. Enter the name for the pool and click **Create**.

**Naming rules:** When you choose a name for a pool, the following rules apply:

- ▶ Names must begin with a letter.
- ▶ The first character cannot be numerical.
- ▶ The name can be a maximum of 63 characters.
- ▶ Valid characters are uppercase letters (A - Z), lowercase letters (a - z), digits (0 - 9), underscore (\_), period (.), hyphen (-), and space.
- ▶ Names must not begin or end with a space.
- ▶ Object names must be unique within the object type. For example, you can have a volume that is named *ABC* and a storage pool that is called *ABC*, but not two storage pools that are both called *ABC*.
- ▶ The default object name is valid (object prefix with an integer).
- ▶ Objects can be renamed at a later stage.

The new pool is created and is included in the list of storage pools with zero bytes, as shown in Figure 5-8.

Name	State	Usable Capacity	Capacity Details	Data Reduction	!!!
Pool0	✓ Online	1.95 TiB / 10.29 TiB (19%)	1.95 TiB / 10.29 TiB (19%)	No	
Pool2	✓ Online	0 bytes	0 bytes	No	

Figure 5-8 Newly created empty pool

To perform this task by using the CLI, run the `mkmdiskgrp` command. The only required parameter is the extent size, which is specified by the `-ext` parameter and must have one of the following values: 16, 32, 64, 128, 256, 512, 1024, 2048, 4096, or 8192 (MB). To create a DRP, specify `-datareduction yes`. The minimum extent size of DRPs is 1024, and attempting to use a smaller extent size sets the extent size to 1024.

In Example 5-2, the command creates a DRP that is named “Pool2” with no MDisk in it.

### Example 5-2 The `mkmdiskgrp` command

```
IBM_2145:ITS0-SV1:superuser>mkmdiskgrp -name Pool2 -datareduction yes -ext 8192  
MDisk Group, id [2], successfully created
```

## 5.1.2 Managed disks in a storage pool

A storage pool is created as an empty container with no storage that is assigned to it. Storage is then added in the form of *MDisks*. An MDisk can be either a RAID array from internal storage (as an array of drives) or a LU from an external storage system. The same storage pool can include both internal and external MDisks.

Arrays are assigned to storage pools at creation time. Arrays cannot exist outside of a storage pool and they cannot be moved between storage pools. It is only possible to destroy an array by removing it from a pool and re-creating it within a new pool.

External MDisks can exist within a pool or outside of a pool. The MDisk object remains on a system if it is visible from external storage, but its access mode changes depending on whether it is assigned to a pool or not.

MDisks are managed by using the MDisks by Pools pane. To access the MDisks by Pools pane, select **Pools** → **MDisks by Pools**, as shown in Figure 5-9.

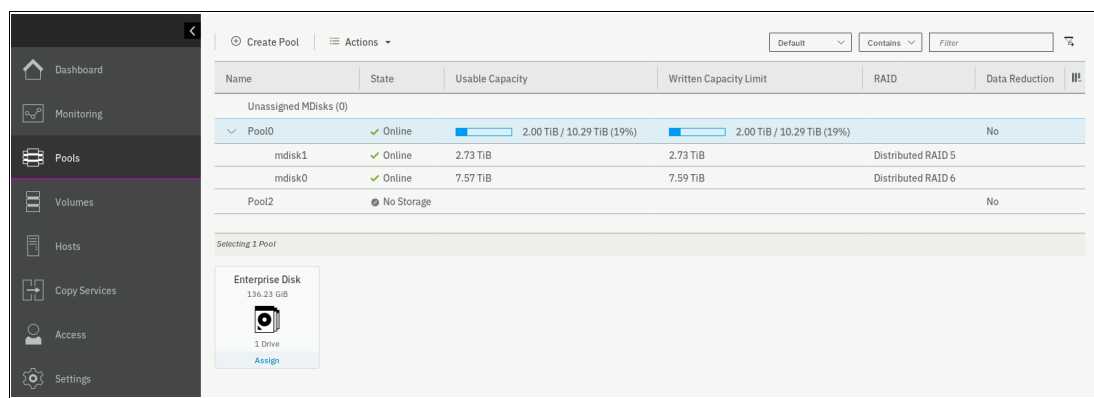


Figure 5-9 MDisks by Pools

The pane lists all the MDisks that are available in the system under the storage pool to which they belong. Unassigned MDisks are listed separately at the top. Both arrays and external MDisks are listed. For more information about operations with array MDisks, see 5.3, “Working with internal drives and arrays” on page 231. To implement a solution with external MDisks, see 5.2, “Working with external controllers and MDisks” on page 217.

To list all MDisks that are visible by the system by using the CLI, run the `lsmdisk` command without any parameters. If required, you can filter output to include only external or only array type MDisks.

## 5.1.3 Actions on storage pools

A number of actions can be performed on storage pools. To select an action, select the storage pool and click **Actions**, as shown in Figure 5-10. Alternatively, right-click the storage pool.

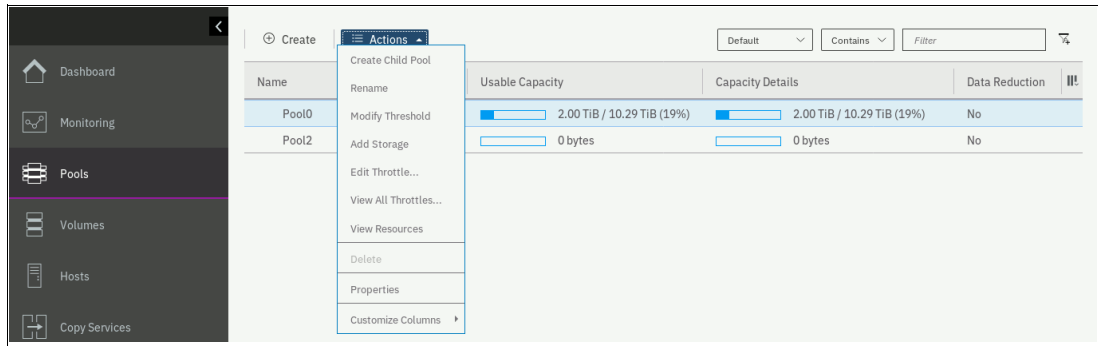


Figure 5-10 Pools actions menu

## Create Child Pool window

To create a child storage pool, click **Create Child Pool**. For more information about child storage pools and a detailed description of this wizard, see 5.1.4, “Child pools” on page 212. It is not possible to create a child pool from an empty pool or from a DRP.

## Rename window

To modify the name of a storage pool, click **Rename**. Enter the new name and click **Rename** in the dialog window.

To do this task by using the CLI, run the **chmdiskgrp** command. Example 5-3 shows how to rename Pool2 to StandardStoragePool. If successful, the command returns no output.

*Example 5-3 Using chmdiskgrp to rename a storage pool*

---

```
IBM_2145:ITS0-SV1:superuser>chmdiskgrp -name StandardStoragePool Pool2
IBM_2145:ITS0-SV1:superuser>
```

---

## Modify Threshold window

A warning event is generated when the amount of used capacity in the pool exceeds the warning threshold. When you use thin-provisioned volumes that auto-expand (automatically use available extents from the pool), monitor the capacity usage and get warnings before the pool runs out of free extents so that you can add storage before running out of space.

**Note:** The warning is generated only the first time that the threshold is exceeded by the used capacity in the storage pool.

To modify the threshold, select **Modify Threshold** and enter the new value. The default threshold is 80%. To disable warnings, set the threshold to 0%.

The threshold is visible in the pool properties and indicated by a red bar, as shown in Figure 5-11 on page 207.



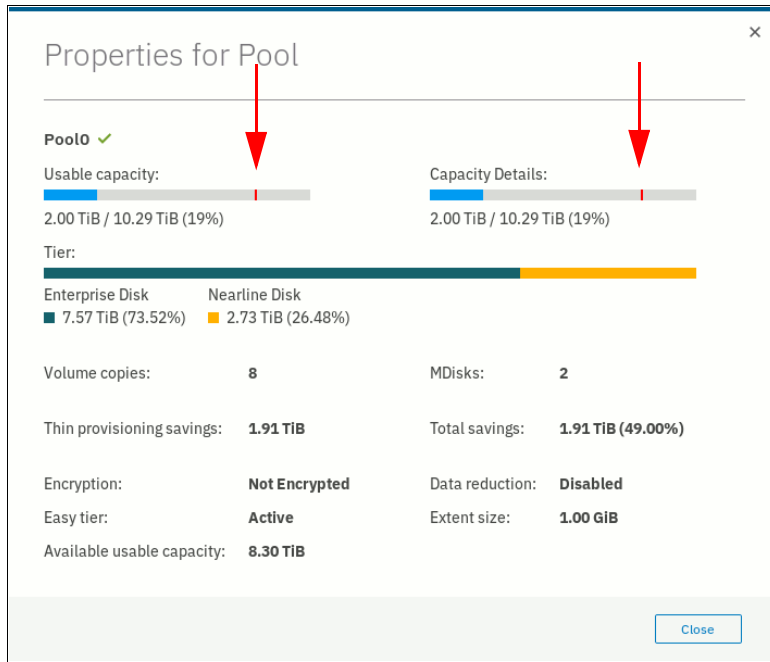


Figure 5-11 Pool properties with warning threshold

To do the task by using the CLI, run the **chmdiskgrp** command. You can specify the threshold by using a percentage. You can also set an exact value and specify a unit.

Example 5-4 shows the warning threshold set to 750 GB for Poo10.

*Example 5-4 Changing the warning threshold level by using the CLI*

---

```
IBM_2145:ITS0-SV1:superuser>chmdiskgrp -warning 750 -unit gb Poo10
IBM_2145:ITS0-SV1:superuser>
```

---

### Add Storage to Pool window

This action starts the configuration wizard, which assigns storage to the pool, as shown in Figure 5-12 on page 208.

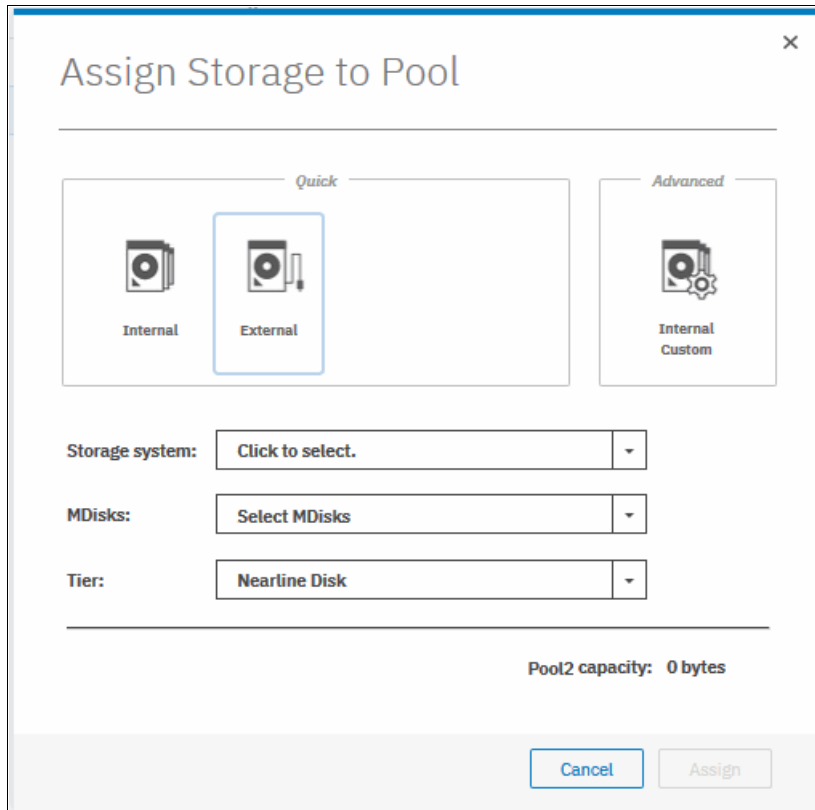


Figure 5-12 Add Storage to Pool wizard

If **Internal** or **Internal Custom** is chosen, the system guides you through array MDisk creation by using internal drives. If **External** is selected, the system guides you through the selection of external storage MDisks. If no external storage is attached, the **External** option is not shown.

## Quorum

With this option, you can reset the quorum configuration to its default, which is to return quorum assignments to a set of MDisks that is chosen by system.

## Edit Throttle for Pool window

Click this option to access the window where you set the pool's throttle configuration.

Throttles can be defined for storage pools to control I/O operations. If a throttle limit is defined, the system either processes the I/O for that object, or delays the processing of the I/O. Resources become free for more critical I/O operations.

You can use storage pool throttles to avoid overwhelming the back-end storage. Only parent pools support throttles because only parent pools contain MDisks from internal or external back-end storage. For volumes in child pools, the throttle of the parent pool is applied.

You can define a throttle for input/output operations per second (IOPS), bandwidth, or both, as shown in Figure 5-13 on page 209:

- ▶ **IOPS limit** indicates the limit of configured IOPS (for both reads and writes combined).
- ▶ **Bandwidth limit** indicates the bandwidth limit in megabytes per second (MBps). You can also specify the limit in gigabits per second (Gbps) or terabytes per second (TBps).

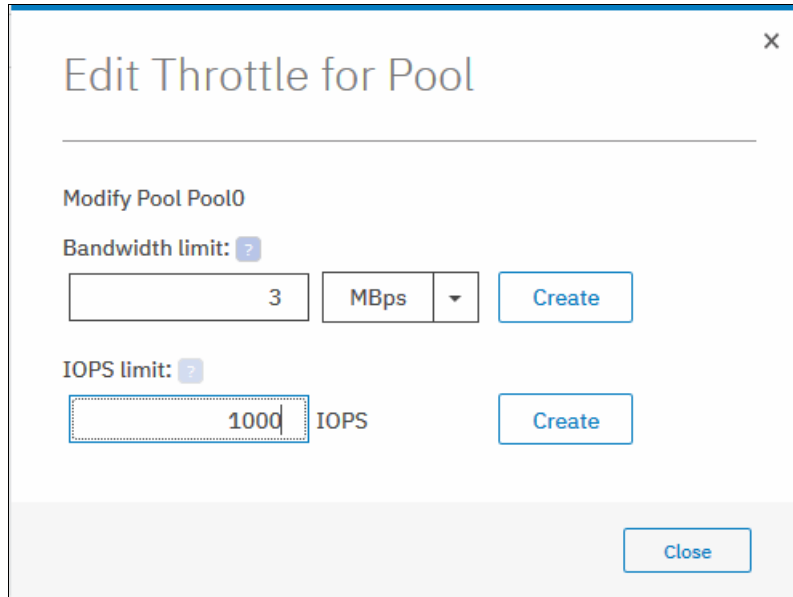


Figure 5-13 Edit throttle for Pool window

If more than one throttle applies to an I/O operation, the lowest and most stringent throttle is used. For example, if a throttle of 100 MBps is defined on a pool and a throttle of 200 MBps is defined on a volume of that pool, the I/O operations are limited to 100 MBps.

The throttle limit is a per node limit. For example, if a throttle limit is set for a volume at 100 IOPS, each node on the system that has access to the volume allows 100 IOPS for that volume. Any I/O operation that exceeds the throttle limit is queued at the receiving nodes. The multipath policies on the host determine how many nodes receive I/O operations and the effective throttle limit.

If a throttle exists for the storage pool, the dialog box that is shown in Figure 5-13 also shows the **Remove** button that is used to delete the throttle.

To set a storage pool throttle by using the CLI, run the **mkthrottle** command. Example 5-5 shows a storage pool throttle, named `iops_bw_limit`, that is set to 3 Mbps and 1000 IOPS on Pool0.

*Example 5-5 Setting a storage pool throttle by using the mkthrottle command*

---

```
IBM_2145:ITS0-SV1:superuser>mkthrottle -type mdiskgrp -iops 1000 -bandwidth 3
-name iops_bw_limit -mdiskgrp Pool0
Throttle, id [0], successfully created.
```

---

To remove a throttle by using the CLI, run the **rmthrottle** command. The command uses the throttle ID or throttle name as an argument, as shown in Example 5-6. The command returns no feedback if it runs successfully.

*Example 5-6 Removing a pool throttle by running the rmthrottle command*

---

```
IBM_2145:ITS0-SV1:superuser>rmthrottle iops_bw_limit
IBM_2145:ITS0-SV1:superuser>
```

---

## View All Throttles window

You can display the defined throttles by using the Pools pane. Right-click a pool and select **View all Throttles** to display the list of the pool's throttles. If you want to view the throttle of other elements (like **Volumes** or **Hosts**, for example), you can select **All Throttles** in the drop-down list, as shown in Figure 5-14.

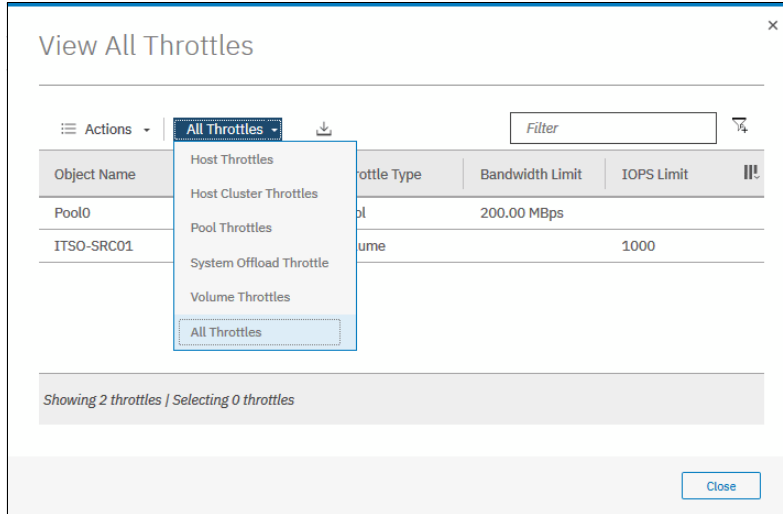


Figure 5-14 Viewing all throttles

To see a list of created throttles by using the CLI, run the `lsthrottle` command. When you run the command without arguments, it displays a list of all throttles on the system. To list only storage pool throttles, specify the `-filtervalue throttle_type=mdiskgrp` parameter.

## View Resources window

To browse a list of MDisks that are part of the storage pool, click **View Resources**, which opens the window that is shown in Figure 5-15.

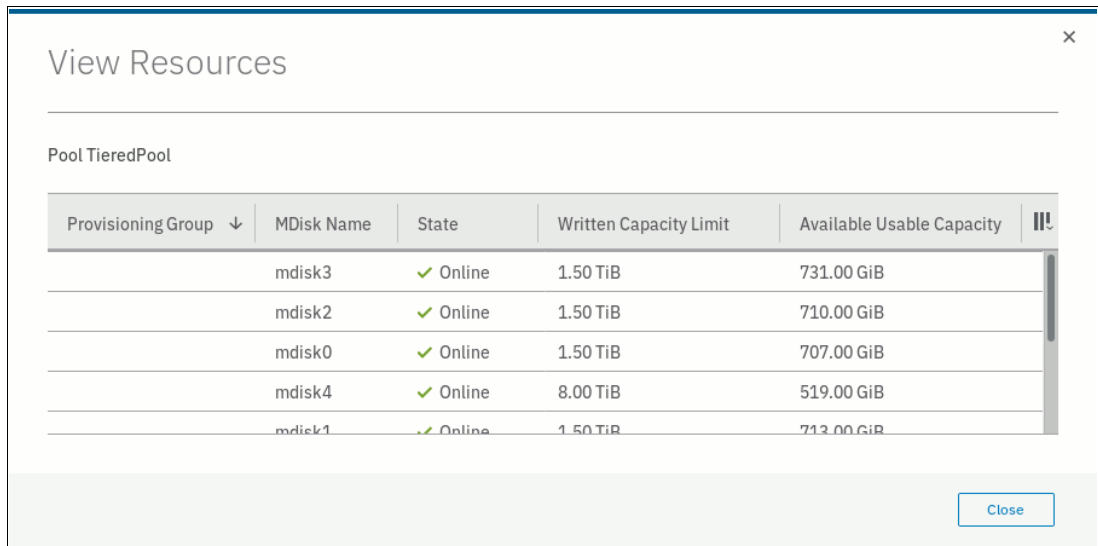


Figure 5-15 List of resources in the storage pool

To list storage pool resources by using the CLI, run the `lsmdisk` command. You can filter the output to display MDisk objects that belong only to a single MDisk group (storage pool), as shown in Example 5-7.

*Example 5-7 Using lsmdisk (some columns are not shown)*

```
IBM_2145:ITS0-SV1:superuser>lsmdisk -filtervalue mdisk_grp_name=Pool0
id name      status mode   mdisk_grp_id mdisk_grp_name  capacity
0 mdisk0  online managed 0           TieredPool      1.5TB
1 mdisk1  online managed 0           TieredPool      1.5TB
2 mdisk2  online managed 0           TieredPool      1.5TB
3 mdisk3  online managed 0           TieredPool      1.5TB
4 mdisk4  online managed 0           TieredPool      8.0TB
```

### Deleting a storage pool

A storage pool can be deleted by using the GUI only if no volumes are associated with it. Select **Delete** to delete the pool immediately without any additional confirmation.

If there are volumes in the pool, the **Delete** option is inactive and cannot be selected. Delete the volumes or migrate them to another storage pool before proceeding. For more information about volume migration and volume mirroring, see Chapter 6, “Volumes” on page 255.

After you delete a pool, the following actions occur:

- ▶ All the external MDisks in the pool return to a mode of *Unmanaged*.
- ▶ All the array mode MDisks in the pool are deleted and all member drives return to a status of *Candidate*.

To delete a storage pool by using the CLI, run the `rmmdiskgrp` command.

**Note:** Be careful when you run the `rmmdiskgrp` command with the `-force` parameter. Unlike the GUI, it does not prevent you from deleting a storage pool with volumes. This command deletes all volumes and host mappings on a storage pool, and they *cannot be recovered*.

### Properties for Pool window

Select **Properties** to display information about the storage pool. By hovering your cursor over the elements of the window and clicking **[?]**, you see a short description of each property, as shown in Figure 5-16 on page 212.

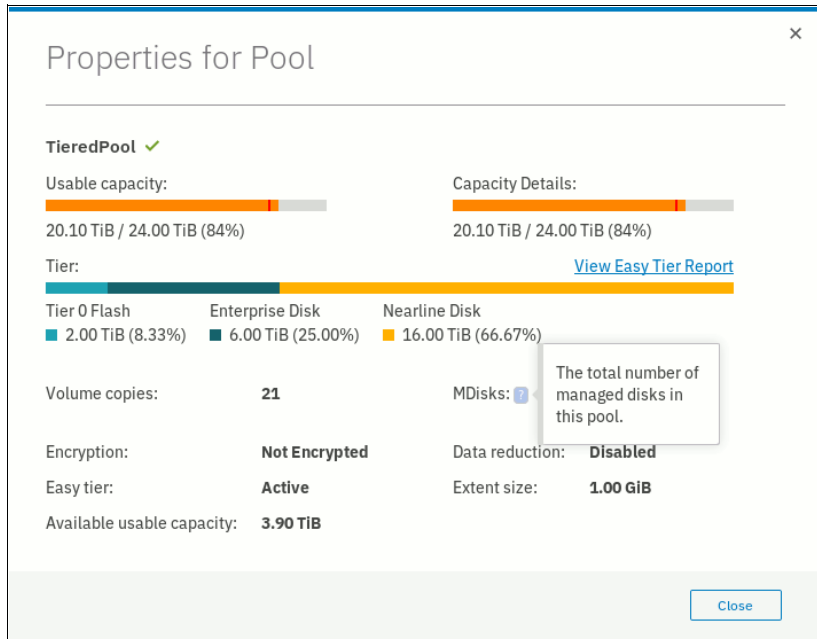


Figure 5-16 Pool properties and details

To display detailed information about the properties by using the CLI, run the `lsmdiskgrp` command with a storage pool name or ID as a parameter, as shown in Example 5-8.

Example 5-8 The `lsmdiskgrp` output (partially shown)

---

```

IBM_2145:ITS0-SV1:superuser>lsmdiskgrp TieredPool
id 0
name TieredPool
status online
mdisk_count 8
vdisk_count 21
capacity 24.00TB
extent_size 1024
free_capacity 3.90TB
<...>

```

---

## 5.1.4 Child pools

A *child pool* is a storage pool that is created within another storage pool. The storage pool in which the child storage pool is created is called the *parent storage pool*.

Unlike a parent pool, a child pool does not contain MDisks. Its capacity is provided exclusively by the parent pool in the form of extents. The capacity of a child pool is set at creation time, but can be modified later nondisruptively. The capacity must be a multiple of the parent pool extent size and must be smaller than the free capacity of the parent pool. Capacity that is assigned to a child pool is taken away from the capacity of the parent pool.

A child pool cannot be created from a DRP. Creating a child pool within another child pool is not possible either.

Child pools are useful when the capacity that is allocated to a specific set of volumes must be controlled. For example, child pools can be used with VMware vSphere Virtual Volumes

(VVOLs). Storage administrators can restrict access of VMware administrators to only a part of the storage pool and prevent volumes creation from affecting the rest of the parent storage pool.

Ownership groups can be used to restrict access to storage resources to a specific set of users, as described in Chapter 11, “Ownership groups” on page 673.

Child pools can also be useful when strict control over thin-provisioned volume expansion is needed. For example, you might create a child pool with no volumes in it to act as an emergency set of extents so that if the parent pool ever runs out of free extents, you can use the ones from the child pool.

On systems with encryption enabled, child pools can be created to migrate existing volumes in a non-encrypted pool to encrypted child pools. When you create a child pool after encryption is enabled, an encryption key is created for the child pool even when the parent pool is not encrypted. You can then use volume mirroring to migrate the volumes from the non-encrypted parent pool to the encrypted child pool. Child pools can also be used when a different encryption key is needed for different sets of volumes.

Child pools inherit most properties from their parent pools, and these properties cannot be changed. The inherited properties include:

- ▶ Extent size
- ▶ Easy Tier setting
- ▶ Encryption setting, but only if the parent pool is encrypted

## Creating a child storage pool

To create a child pool, complete the following steps:

1. Select **Pools** → **Pools**, right-click the parent pool that you want to create a child pool from, and select **Create Child Pool**, as shown in Figure 5-17.

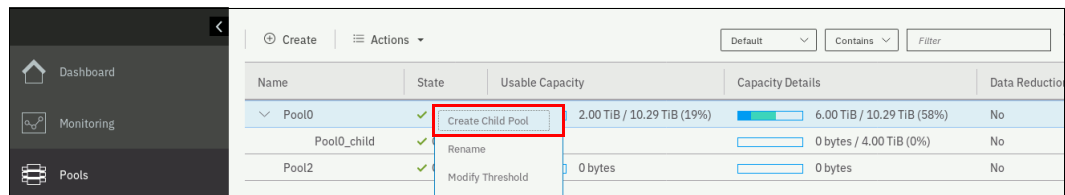


Figure 5-17 Creating a child pool

2. When the dialog box opens, enter the name and capacity of the child pool and click **Create**, as shown in Figure 5-18 on page 214.

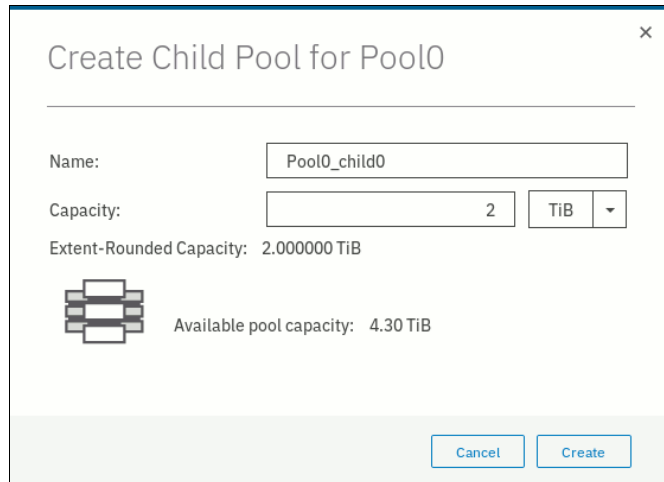


Figure 5-18 Creating a child pool

After the child pool is created, it is listed in the Pools pane under its parent pool. Toggle the sign to the left of the storage pool icon to either show or hide the child pools, as shown in Figure 5-19. The capacity that is assigned to the child pools is not usable in the parent pool, as shown by the gray area on the capacity details bar of the parent pool.

Name	State	Usable Capacity	Capacity Details	Data Reduction
Pool0	Online	2.00 TiB / 10.29 TiB (19%)	8.00 TiB / 10.29 TiB (78%)	No
Pool0_child	Online	0 bytes / 4.00 TiB (0%)		No
Pool0_child0	Online	0 bytes / 2.00 TiB (0%)		No
Pool2	Online	0 bytes	0 bytes	No

Figure 5-19 Listing parent and child pools

To create a child pool by using the CLI, run the `mkmdiskgrp` command. You must specify the parent pool for your new child pool and its size, as shown in Example 5-9. The size is in megabytes by default (unless the `-unit` parameter is used) and must be a multiple of the parent pool's extent size. In this case, it is  $100 * 1024 \text{ MB} = 100 \text{ GB}$ .

*Example 5-9 The `mkmdiskgrp` command to create child pools*

```
IBM_2145:ITS0-SV1:superuser>mkmdiskgrp -parentmdiskgrp Pool0 -size 102400 -name
Pool0_child0
MDisk Group, id [4], successfully created
```

### Actions for child storage pools

You can rename, resize, and delete a child pool. Also, it is possible to modify its warning threshold and assign it to an ownership group. To select an action, for example, **Resize**, complete the following steps:

1. Right-click the child storage pool, as shown in Figure 5-20 on page 215. Alternatively, select the storage pool and click **Actions**.



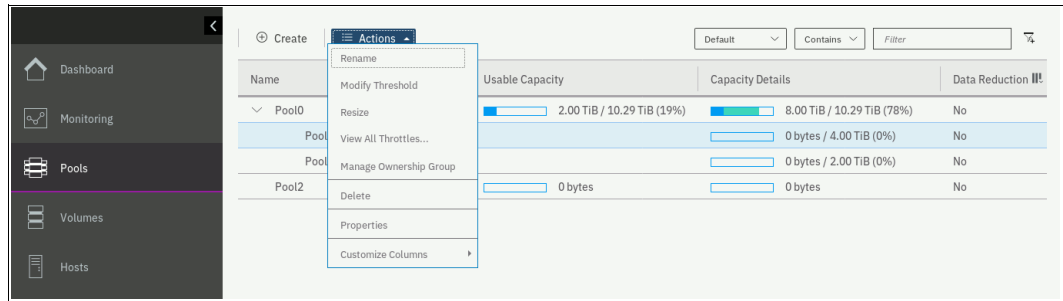


Figure 5-20 Actions for child storage pools

2. Select **Resize** to increase or decrease the capacity of the child storage pool, as shown in Figure 5-21. Enter the new pool capacity and click **Resize**.

**Note:** You cannot shrink a child pool below its real capacity. Thus, the new size of a child pool must be larger than the capacity that is used by its volumes.

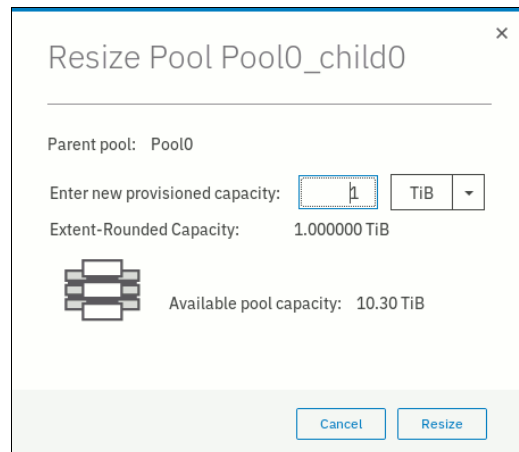


Figure 5-21 Resizing a child pool

When the child pool is shrunk, the system resets the warning threshold and issues a warning if the threshold is reached.

To rename and resize child pool by using the CLI, run the **chmdiskgrp** command. Example 5-10 renames the child pool `Pool0_child0` to `Pool0_child_new` and reduces its size to 44 GB. If successful, the command returns no feedback.

*Example 5-10 Running the chmdiskgrp command to rename a child pool*

---

```
IBM_2145:ITS0-SV1:superuser>chmdiskgrp -name Pool0_child_new -size 61440
Pool0_child0
IBM_2145:ITS0-SV1:superuser>
```

---

Deleting a child pool is a task that is like deleting a parent pool. As with a parent pool, the **Delete** action is disabled if the child pool contains volumes, as shown in Figure 5-22 on page 216.

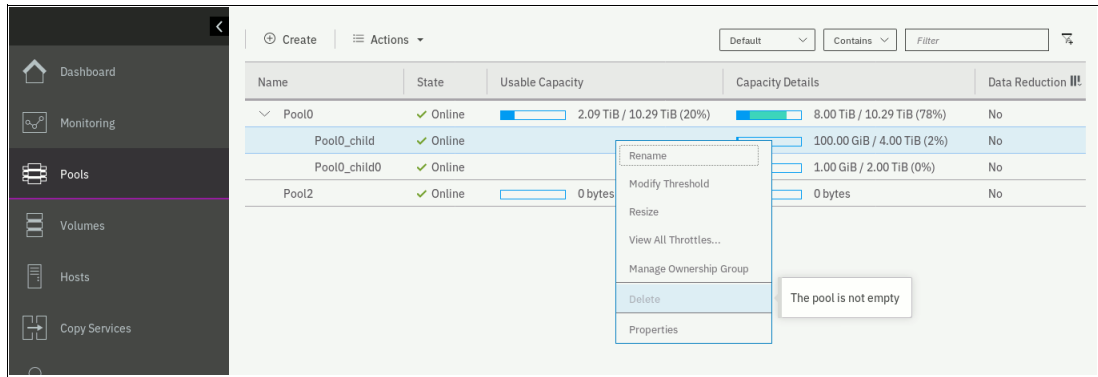


Figure 5-22 Deleting a child pool

After deleting a child pool, the extents that it occupied return to the parent pool as free capacity.

To delete a child pool by using the CLI, run the `rmm diskgrp` command.

To assign an existing ownership group to a child pool, click **Manage Ownership Group**, as shown in Figure 5-23. All volumes that are created in the child pool inherit the ownership group of the child pool. For more information, see Chapter 11, “Ownership groups” on page 673.

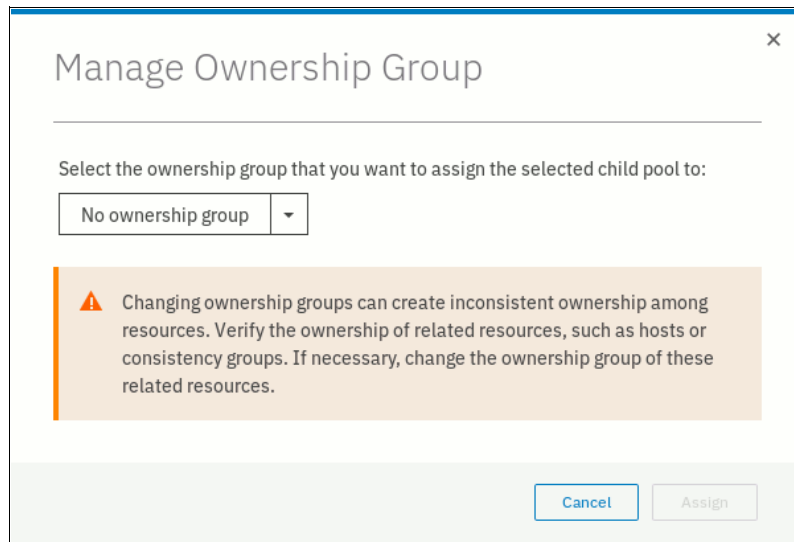


Figure 5-23 Managing the ownership group of a child pool

## Migrating volumes to and from child pools

To move a volume to another pool, you can use migration or volume mirroring in the same way that you use them for parent pools. For more information about volume migration and volume mirroring, see Chapter 6, “Volumes” on page 255.

The system supports migration of volumes between child pools within the same parent pool or migration of a volume between a child pool and its parent pool. Migrations between a source and target child pool with different parent pools are not supported. However, you can migrate the volume from the source child pool to its parent pool. Then, the volume can be migrated from the parent pool to the parent pool of the target child pool. Finally, the volume can be migrated from the target parent pool to the target child pool.

During a volume migration within a parent pool (between a child and its parent or between children with the same parent), there is no data movement, but there are extent reassignments.

Volume migration between a child storage pool and its parent storage pool can be performed by going to the **Volumes by Pool** page and clicking **Volumes**. Right-click a volume and select it to migrate it into a suitable pool.

In the example in Figure 5-24, the volume `vdisk0` was created in child pool `Pool0_child_new`. The child pools appear exactly like the parent pools in the Volumes by Pool pane.

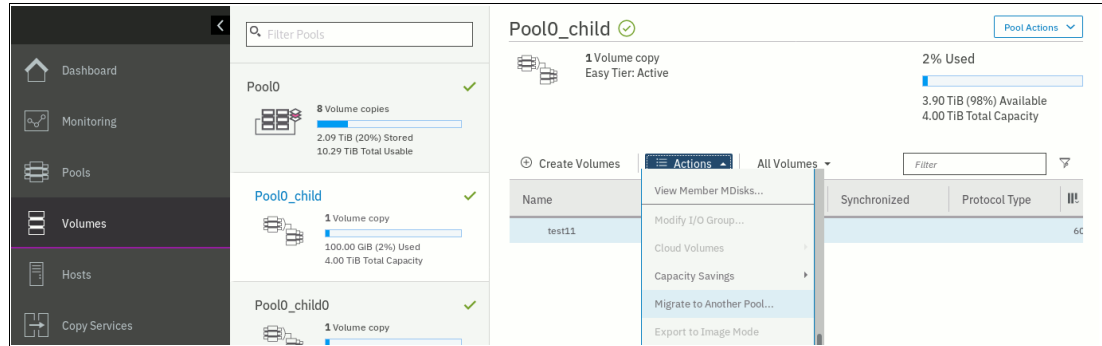


Figure 5-24 Actions menu in Volumes by Pool

For more information about CLI commands for migrating volumes to and from child pools, see Chapter 6, “Volumes” on page 255.

### 5.1.5 Encrypted storage pools

The system supports two types of encryption: hardware encryption and software encryption.

Hardware encryption is implemented at an array level, and software encryption is implemented at a storage pool level. For more information about encryption and encrypted storage pools, see Chapter 12, “Encryption” on page 685.

## 5.2 Working with external controllers and MDisks

*Controllers* are external storage systems that provide storage resources that are used as MDisks. The system supports external storage controllers that are attached through internet Small Computer Systems Interface (iSCSI) and through Fibre Channel (FC).

A key feature of the SAN Volume Controller is its ability to consolidate disk controllers from various vendors into storage pools. The storage administrator can manage and provision storage to applications from a single user interface and use a common set of advanced functions across all of the storage systems under the control of the SAN Volume Controller.

This concept is called *External Virtualization*, which makes your storage environment more flexible, cost-effective, and easy to manage. External Virtualization follows a differential licensing scheme that is based on Storage Capacity Units (SCU).

For more information about how to configure external storage systems, see Chapter 2, “Planning” on page 53.

## 5.2.1 External storage controllers

External storage controllers can be attached by using FC and iSCSI. The following sections describe how to attach external storage controllers to the system and how to manage them by using the GUI.

### System layers

A *system layer* affects how the system interacts with other external IBM Storwize or IBM FlashSystem family systems. SAN Volume Controller is always in the *replication* layer. A Storwize or FlashSystem family system is in either the *storage* layer (default) or the *replication* layer.

With these default settings, SAN Volume Controller can virtualize other Storwize or FlashSystem family systems. However, if the Storwize or FlashSystem family system was moved to the *replication layer*, it must be configured back to *storage layer* for the SAN Volume Controller to be able to use it as external storage.

The SAN Volume Controller system layer cannot be changed. The changes must be made on the external Storwize or FlashSystem family system instead.

**Note:** Before you change the system layer, the following conditions must be met:

- ▶ No host object can be configured with worldwide port names (WWPNs) from a Storwize or IBM FlashSystem family system.
- ▶ No system partnerships can be defined.
- ▶ No Storwize or IBM FlashSystem family system can be visible on the storage area network (SAN) fabric.

For more information about layers and how to change them, see [IBM Knowledge Center](#), and expand **Product overview** → **Technical overview** → **System layers**.

### Attachment by using Fibre Channel

A controller that is connected through FC is detected automatically by the system if the cabling, zoning, and system layer are configured correctly.

Any supported controller can be temporarily direct attached for migrating the data from that controller onto the system. For permanent virtualization, external storage controllers are connected through SAN switches. For more information about how to attach and zone back-end storage controllers to the system, see 2.6, “Fibre Channel SAN configuration planning” on page 58.

**Note:** At the time of writing, two back-end storage systems were supported as permanently directly attached: IBM DS8000 and IBM FlashSystem 900. For more information, see this [IBM Support web page](#).

If the external controller is not detected, ensure that the system is cabled and zoned into the same SAN as the external storage system. Check that layers are set correctly on both virtualizing and virtualized systems if they belong to the IBM Storwize or IBM FlashSystem family.

After the problem is corrected, rescan the FC network immediately by selecting **Pools** → **External Storage**, and then selecting **Actions** → **Discover Storage**, as shown in Figure 5-25 on page 219.

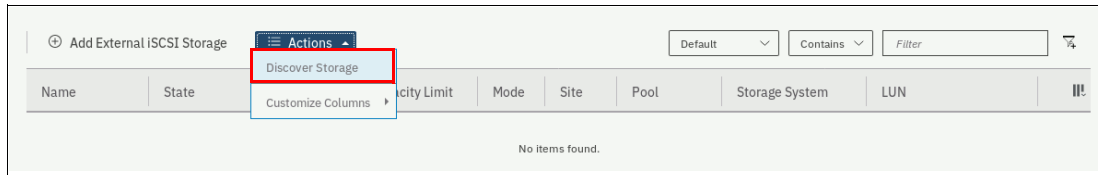


Figure 5-25 Discover Storage menu

This action runs the `detectmdisk` command. It returns no output. Although it might appear that the command completed, some extra time might be required for it to run. The command is asynchronous and returns a prompt while the command continues to run in the background.

## Attachment by using iSCSI

Unlike a Fibre Channel connection (FICON), you must manually configure iSCSI connections between the SAN Volume Controller and the external storage controller. Until you do this task, the controller is not listed in the External Storage pane. For more information about how to attach back-end storage controllers to the system, see Chapter 2, “Planning” on page 53.

To start virtualizing an iSCSI back-end controller, you must follow the documentation in [IBM Knowledge Center](#) to perform configuration steps that are specific to your back-end storage controller. You can see find the steps by selecting **Configuring** → **Configuring and servicing storage systems** → **External storage system configuration with iSCSI connections**.

For more information about configuring SAN Volume Controller to virtualize back-end storage controller with iSCSI, see *iSCSI Implementation and Best Practices on IBM Storwize Storage Systems*, SG24-8327.

## Managing external storage controllers

You can manage both FC and iSCSI storage controllers through the External Storage pane. To access the External Storage pane, select **Pools** → **External Storage**, as shown in Figure 5-26.

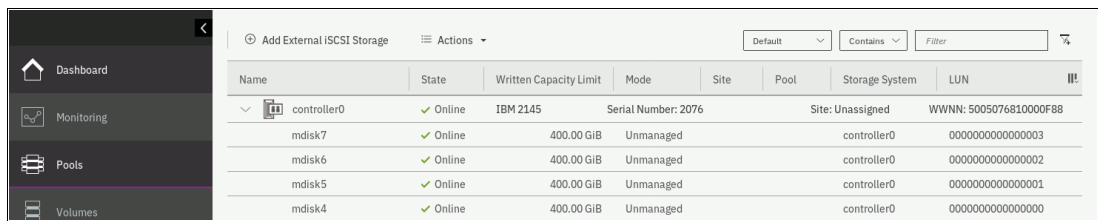


Figure 5-26 External Storage pane

**Note:** A controller that is connected through FC is detected automatically by the system. The cabling, the zoning, and system layers must be configured correctly. A controller that is connected through iSCSI must be added to the system manually.

Depending on the type of back-end system, it might be detected as one or more controller objects.

The External Storage pane lists the external controllers that are connected to the system and all the external MDisks that are detected by the system. The MDisks are organized by the external storage system that presents them. Toggle the sign to the left of the controller icon to show or hide the MDisks that are associated with the controller.

If you configured logical unit names on your external storage systems, it is not possible for the system to determine these names because they are local to the external storage system. However, you can use the LU unique identifiers (UIDs), or external storage system worldwide node names (WWNNs) and LU number to identify each device.

To list all visible external storage controllers with CLI, run the `lscontroller` command, as shown in Example 5-11.

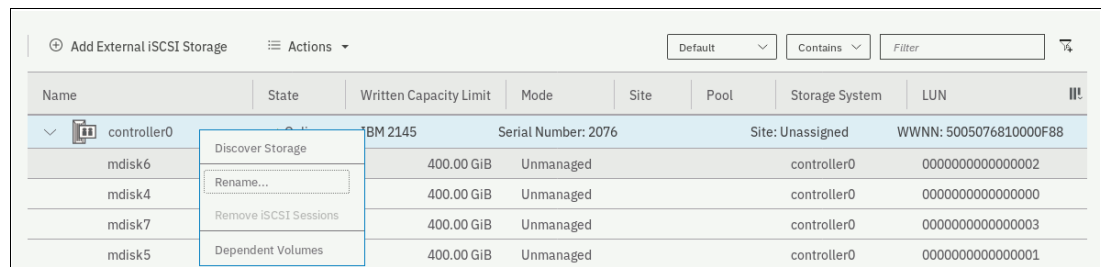
*Example 5-11 Listing controllers by using the CLI (some columns are not shown)*

```
IBM_2145:ITS0-SV1:superuser>lscontroller
id controller_name ctrl_s/n          vendor_id      product_id_low
0 controller0      2076          IBM            2145
1 controller1      2076          IBM            2145
2 controller2      2076          IBM            2145
<...>
```

## 5.2.2 Actions on external storage controllers

You can perform many actions on external storage controllers. Some actions are available for external iSCSI controllers only.

To select any action, select **Pools** → **External Storage** and right-click the controller, as shown in Figure 5-27. Alternatively, select the controller and click **Actions**.



*Figure 5-27 Actions for external storage*

### Discover Storage

When you create or remove LUs on an external storage system, the change might not be detected immediately. In this case, click **Discover Storage** so that the system can rescan the FC or iSCSI network. In general, the system automatically detects disks when they appear on the network. However, some FC controllers do not send the required SCSI primitives that are necessary to automatically discover the new disks.

The rescan process discovers any new MDisks that were added to the system and rebalances MDisk access across the available ports. It also detects any loss of availability of the controller ports.

This action runs the `detectmdisk` command.

## Rename

To modify the name of an external controller to simplify administration tasks, click **Rename**. The naming rules are the same as for storage pools, and they can be found in 5.1.1, “Creating storage pools” on page 202.

To rename a storage controller by using the CLI, run the **chcontroller** command. Example 5-12 shows a use case.

## Removing iSCSI sessions

This action is available only for external controllers that are attached with iSCSI. To remove the iSCSI session that is established between the source and target port, right-click the session and select **Remove**.

For more information about the CLI commands and detailed instructions, see *iSCSI Implementation and Best Practices on IBM Storwize Storage Systems*, SG24-8327.

## Modifying a site

This action is available only for systems that are configured as an Enhanced Stretched Cluster (ESC) or HyperSwap topology. To change the site with which the external controller is associated, select **Modify Site**, as shown in Figure 5-28.

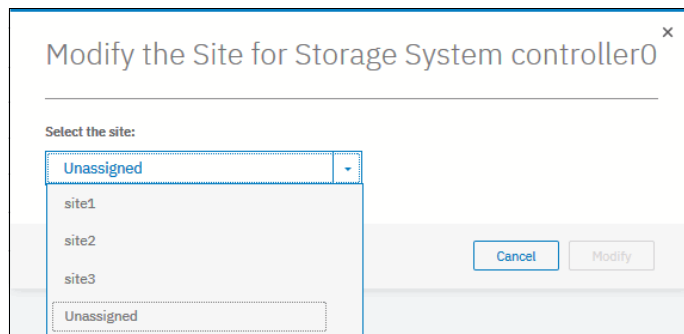


Figure 5-28 Modifying the site of an external controller

To change the controller site assignment by using the CLI, run the **chcontroller** command. Example 5-12 shows that controller0 was renamed and reassigned to a different site.

### Example 5-12 Changing a controller's name and site

---

```
IBM_2145:ITS0-SV1:superuser>chcontroller -name site3_controller -site site3
controller0
IBM_2145:ITS0-SV1:superuser>
```

---

## 5.2.3 Working with external MDisks

After an external back-end storage controller is configured, attached to the system, and detected as a controller, you can work with LUs that are provisioned from it. Each LU is represented by an MDisk object.

External MDisks can have one of the following modes:

► *Unmanaged*

External MDisks are initially discovered by the system as unmanaged MDisks. An unmanaged MDisk is not a member of any storage pool. It is not associated with any volumes, and has no metadata that is stored on it. The system does not write to an MDisk that is in unmanaged mode except when it attempts to change the mode of the MDisk to one of the other modes. Removing an external MDisk from a pool returns it to unmanaged mode.

► *Managed*

When unmanaged MDisks are added to storage pools, they become managed. Managed mode MDisks are always members of a storage pool, and their extents contribute to the storage pool. This mode is the most common and normal mode for an MDisk.

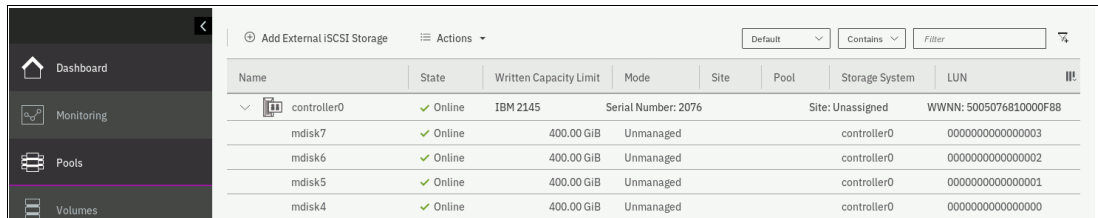
► *Image*

Image mode provides a direct block-for-block conversion from the MDisk to a volume. This mode is provided to satisfy the following major usage scenarios:

- Presenting existing data on an MDisk through the system to an attached host
- Importing existing data on an MDisk into the system
- Exporting data on a volume by performing a migration to an image mode MDisk

### Listing external MDisks

You can manage external MDisks by using the External Storage pane, which is accessed by selecting **Pools** → **External Storage**, as shown in Figure 5-29.



Name	State	Written Capacity Limit	Mode	Site	Pool	Storage System	LUN
controller0	Online	IBM Z145	Serial Number: 2076	Site: Unassigned		WWNN: 5005076810000F88	
mdisk7	Online	400.00 GiB	Unmanaged		controller0	0000000000000003	
mdisk6	Online	400.00 GiB	Unmanaged		controller0	0000000000000002	
mdisk5	Online	400.00 GiB	Unmanaged		controller0	0000000000000001	
mdisk4	Online	400.00 GiB	Unmanaged		controller0	0000000000000000	

Figure 5-29 External Storage pane

To list all MDisks that are visible by the system by using the CLI, run the `lsmdisk` command without any parameters. If required, you can filter output to include only external or only array type MDisks.

### Assigning MDisks to pools

You can add *Unmanaged* MDisks to an existing pool or create a pool to include them. If no storage pool exists yet, follow the procedure that is outlined in 5.1.1, “Creating storage pools” on page 202.



Figure 5-30 shows how to add selected MDisk to an existing storage pool. Click **Assign** under the **Actions** menu or right-click the MDisk and select **Assign**.

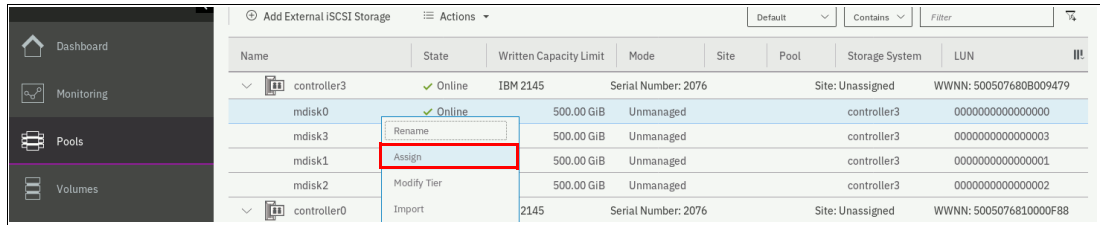


Figure 5-30 Assigning an unmanaged MDisk

After you click **Assign**, a dialog box opens, as shown in Figure 5-31. Select the target pool, MDisk storage tier, and external encryption setting.

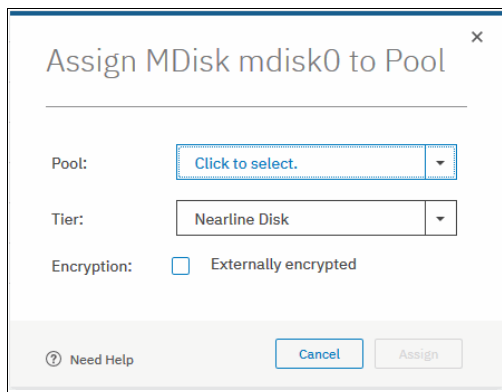


Figure 5-31 Assign MDisk dialog box

When you add MDisks to pools, you must assign them to the correct storage tiers. It is important to set the tiers correctly if you plan to use the Easy Tier feature. Using an incorrect tier can mean that the Easy Tier algorithm might make wrong decisions and thus affect system performance. For more information about storage tiers, see Chapter 9, “Advanced features for storage efficiency” on page 449.

The storage tier setting can also be changed after the MDisk is assigned to the pool.

Select the **Externally encrypted** check box if your back-end storage performs data encryption. For more information about SAN Volume Controller encryption, see Chapter 5, “Storage pools” on page 199.

After the task completes, click **Close**.

**Note:** If the external storage LUs that are presented to the SAN Volume Controller contain data that must be retained, do not use the **Assign** option to add the MDisks to a pool. This option destroys the data on the LU. Instead, use the **Import** option to create an image mode MDisk. For more information, see Chapter 8, “Storage migration” on page 429.

To see the external MDisks that are assigned to a pool within the system, select **Pools** → **MDisks by Pools**.

When a new MDisk is added to a pool that already contains MDisks and volumes, the Easy Tier feature automatically balances volume extents between the MDisks in the pool as a background process. The goal of this process is to distribute extents in a way that provides the best performance to the volumes. It does *not* attempt to balance the amount of data evenly between all MDisks.

The data migration decisions that Easy Tier makes between tiers of storage (inter-tier) or within a single tier (intra-tier) are based on the I/O activity that is measured. Therefore, when you add an MDisk to a pool, extent migrations are not necessarily performed immediately. No migration of extents occurs until there is sufficient I/O activity to trigger it.

If Easy Tier is turned off, no extent migration is performed. Only newly allocated extents are written to a new MDisk.

For more information about the Easy Tier feature, see Chapter 9, “Advanced features for storage efficiency” on page 449.

To assign an external MDisk to a storage pool by using the CLI, run the `addmdisk` command. You must specify the MDisk name or ID, MDisk tier, and target storage pool, as shown in Example 5-13. The command returns no feedback.

*Example 5-13 The addmdisk command*

```
IBM_2145:ITS0-SV1:superuser>addmdisk -mdisk mdisk3 -tier enterprise Pool0
IBM_2145:ITS0-SV1:superuser>
```

## 5.2.4 Actions for external MDisks

External MDisks support specific actions that are not supported on RAID arrays that are made from internal storage. Some actions are supported only on unmanaged external MDisks, and some are supported only on managed external MDisks.

To choose an action, select **Pools** → **External Storage** or **Pools** → **MDisks by Pools**, select the external MDisk, and click **Actions**, as shown in Figure 5-32. Alternatively, right-click the external MDisk.

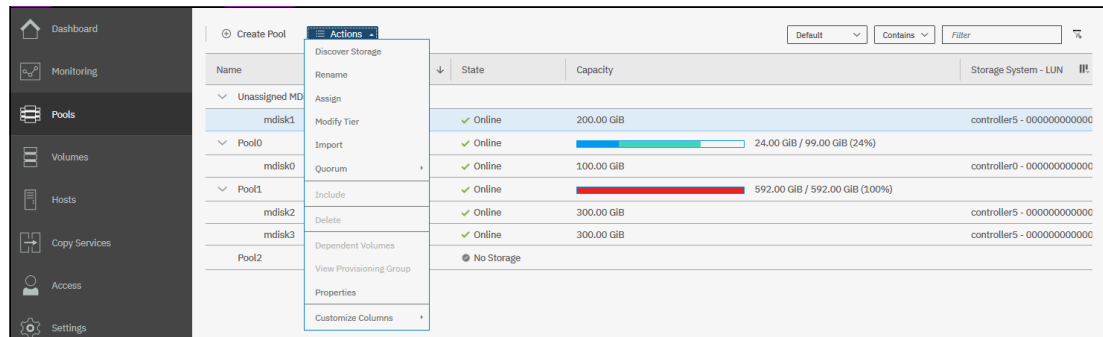


Figure 5-32 Actions for MDisks

### Discover Storage

This option is available even if no MDisks are selected. By running it, you cause the system to rescan the iSCSI and FC network for these purposes:

- ▶ Find any new MDisks that might be added.
- ▶ Rebalance MDisk access across all available controller device ports.

This action runs the `detectmdisk` command.

## Assign

This action is available only for unmanaged MDisks. Select **Assign** to open the dialog box that is explained in “Assigning MDisks to pools” on page 222.

## Modify Tier

To modify the tier to which the external MDisk is assigned, select **Modify Tier**, as shown in Figure 5-33. This setting is adjustable because the system cannot always detect the tiers that are associated with external storage automatically, unlike with internal arrays.

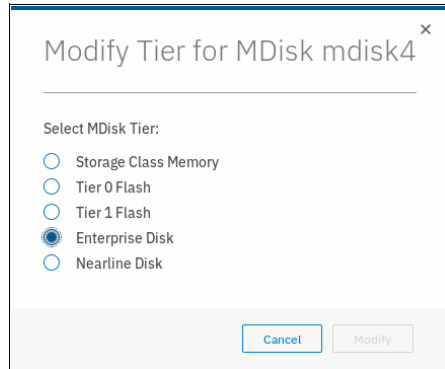


Figure 5-33 Modifying an external MDisk tier

For more information about storage tiers and their importance, see 9.1, “IBM Easy Tier” on page 450.

To change the external MDisk storage tier, run the `chmdisk` command. Example 5-14 shows setting the new tier to `mdisk2`. No feedback is returned.

*Example 5-14 Changing the tier setting by using the CLI*

---

```
IBM_2145:ITS0-SV1:superuser>chmdisk -tier tier1_flash mdisk2
IBM_2145:ITS0-SV1:superuser>
```

---

## Modify Encryption option

To modify the encryption setting for the MDisk, select **Modify Encryption**. This option is available only when encryption is enabled.

If the external MDisk is already encrypted by the external storage system, change the encryption state of the MDisk to **Externally encrypted**. This setting stops the system from encrypting the MDisk again if the MDisk is part of an encrypted storage pool.

For more information about encryption, encrypted storage pools, and self-encrypting MDisks, see Chapter 12, “Encryption” on page 685.

To perform this task by using the CLI, run the `chmdisk` command, as shown in Example 5-15.

*Example 5-15 Using chmdisk to modify encryption*

---

```
IBM_2145:ITS0-SV1:superuser>chmdisk -encrypt yes mdisk5
IBM_2145:ITS0-SV1:superuser>
```

---

## Import

This action is available only for unmanaged MDisks. Importing an unmanaged MDisk enables you to preserve the existing data on the MDisk. You can migrate the data to a new volume or keep the data on the external system.

MDisks are imported for storage migration. The SAN Volume Controller provides a migration wizard to help with this process, which is described in Chapter 8, “Storage migration” on page 429.

**Note:** The wizard is the preferred method to migrate data from legacy storage to the SAN Volume Controller. When an MDisk is imported, the data on the original LU is not modified. The system acts as a pass-through, and the extents of the imported MDisk do not contribute to storage pools.

To choose one of the following migration methods, select **Import**:

► **Import to temporary pool as image mode volume**

This method does not migrate data from the source MDisk. It creates an *image mode volume* that has a direct block-for-block conversion of the MDisk. The existing data is preserved on the external storage controller, but it is also accessible from the SAN Volume Controller system.

In this method, the image mode volume is created in a temporary migration pool and is presented through the SAN Volume Controller. Choose the extent size of the temporary pool and click **Import**, as shown in Figure 5-34.

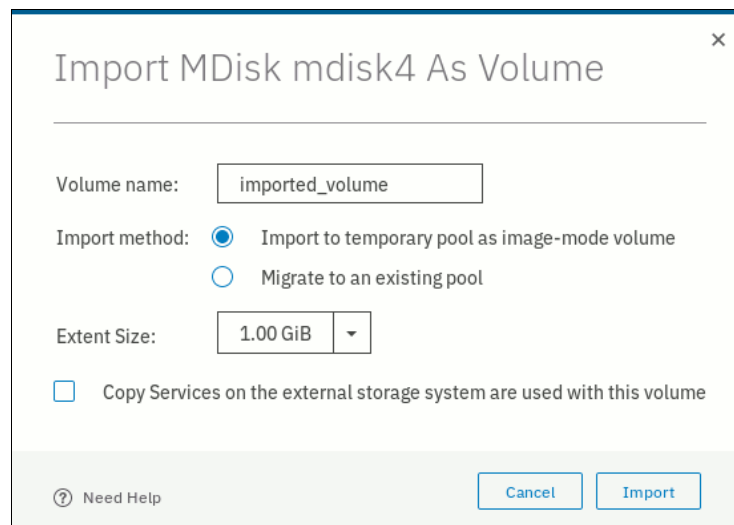


Figure 5-34 Importing an unmanaged MDisk

The MDisk is imported and listed as an image mode MDisk in the temporary migration pool, as shown in Figure 5-35.

Name	State	Usable Capacity	Written Capacity Limit	Storage System - LUN	Data Reduction
Unassigned MDisks (0)					
MigrationPool_1024	Online	<div style="width: 100%; height: 10px; background-color: red;"></div>	400.00 GiB / 400.00 GiB (100%)		No
mdisk4	Online	400.00 GiB		controller0 - 0000000000000000	

Figure 5-35 Image mode imported MDisk



The data migration begins automatically after the MDisk is imported successfully as an image mode volume. You can check the migration progress by clicking the task under **Running Tasks**, as shown in Figure 5-38.

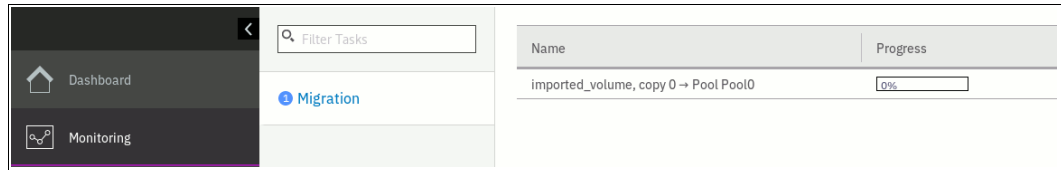


Figure 5-38 MDisk migration in the Running Tasks pane

After the migration completes, the volume is available in the chosen destination pool. This volume is no longer an image mode volume. It is now virtualized by the system.

All data is migrated off the source MDisk, and the MDisk switched its mode, as shown in Figure 5-39.

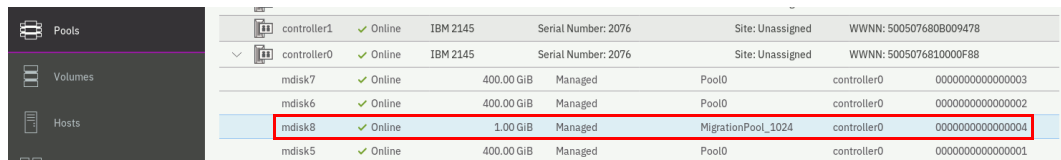


Figure 5-39 Imported MDisks appear as “Managed”

The MDisk can be removed from the migration pool. It returns to the list of external MDisks as Unmanaged. The MDisk can now be used as a regular managed MDisk in a storage pool, or it can be decommissioned.

Alternatively, importing and migrating external MDisks to another pool can be done by selecting **Pools** → **System Migration** to start the system migration wizard. For more information, see Chapter 8, “Storage migration” on page 429.

## Quorum

This menu option enables you to introduce a new set of quorum disks. When three online managed MDisks are selected, the **Quorum** → **Modify Quorum Disks** menu becomes available, as shown on Figure 5-40. For HyperSwap and ESC configurations, MDisks must belong to three different sites.

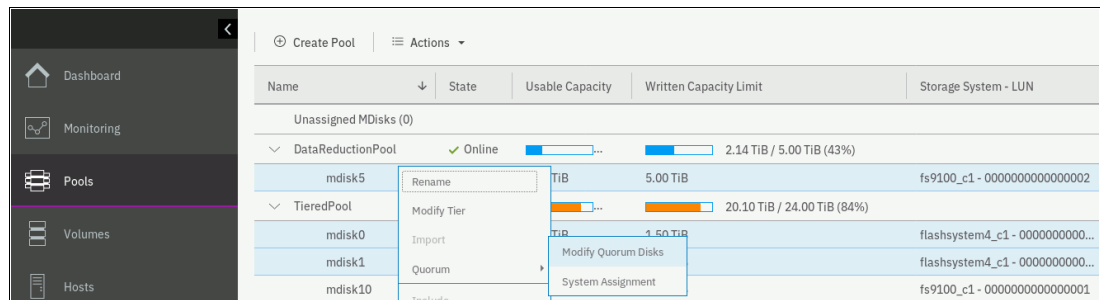


Figure 5-40 Selecting new quorum disks

To list and change the quorum configuration, run the `lquorum` and `chquorum` commands.

## Include

The system can exclude an MDisk from its storage pool if it has multiple I/O failures or has persistent connection errors. Exclusion ensures that there is no excessive error recovery that might impact other parts of the system. If an MDisk is automatically excluded, run the DMP to resolve any connection and I/O failure errors.

If no error event is associated with the MDisk in the log and the external problem is corrected, click **Include** to add the excluded MDisk back to the storage pool.

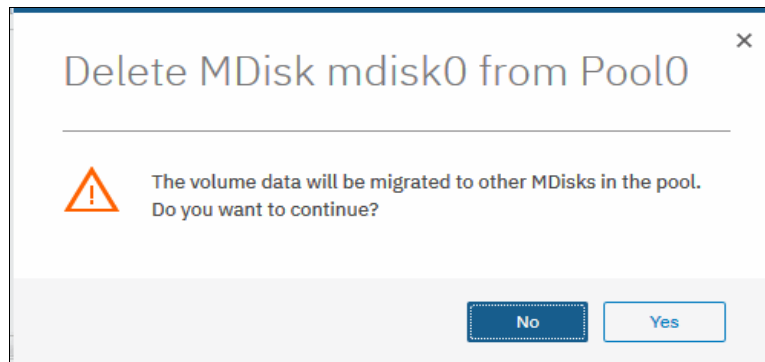
The `includemdisk` command performs the same task. The command needs the MDisk name or ID to be provided as a parameter, as shown in Example 5-16.

*Example 5-16 Including a degraded MDisk by using the CLI*

```
IBM_2145:ITS0-SV1:superuser>includemdisk mdisk3
IBM_2145:ITS0-SV1:superuser>
```

## Remove

In some cases, you might want to remove external MDisks from their storage pool. To remove the MDisk from the storage pool, click **Remove**. After the MDisk is removed, it goes back to the Unmanaged state. If there are no volumes in the storage pool to which this MDisk is allocated, the command runs immediately without more confirmation. If there are volumes in the pool, you are prompted to confirm the action, as shown in Figure 5-41.



*Figure 5-41 Removing an MDisk from a pool*

Confirming the action starts a migration of volumes to extents on that MDisk to other MDisks in the pool. During this background process, the MDisk remains a part of the storage pool. Only when the migration completes is the MDisk removed from the storage pool and returns to Unmanaged mode.

Ensure that you have enough available capacity remaining in the storage pool to allocate the data that is being migrated from the removed MDisk, or this command fails.

**Important:** The MDisk that you are removing must remain accessible to the system while all data is copied to other MDisks in the same storage pool. If the MDisk is unmapped before the migration finishes, all volumes in the storage pool go offline and remain in this state until the removed MDisk is connected again.

To remove an MDisk from a storage pool by using the CLI, run the `rmmdisk` command. You must use the `-force` parameter if you must migrate volume extents to other MDisks in a storage pool.

The command fails if you do not have enough available capacity remaining in the storage pool to allocate the data that you are migrating from the removed array.

### Dependent Volumes

A volume depends on an MDisk if the MDisk becoming unavailable results in a loss of access or a loss of data for that volume. Use this option before you do maintenance operations to confirm which volumes (if any) are affected. Selecting an MDisk and clicking **Dependent Volumes** lists the volumes that depend on that MDisk. An example is shown in Figure 5-42.

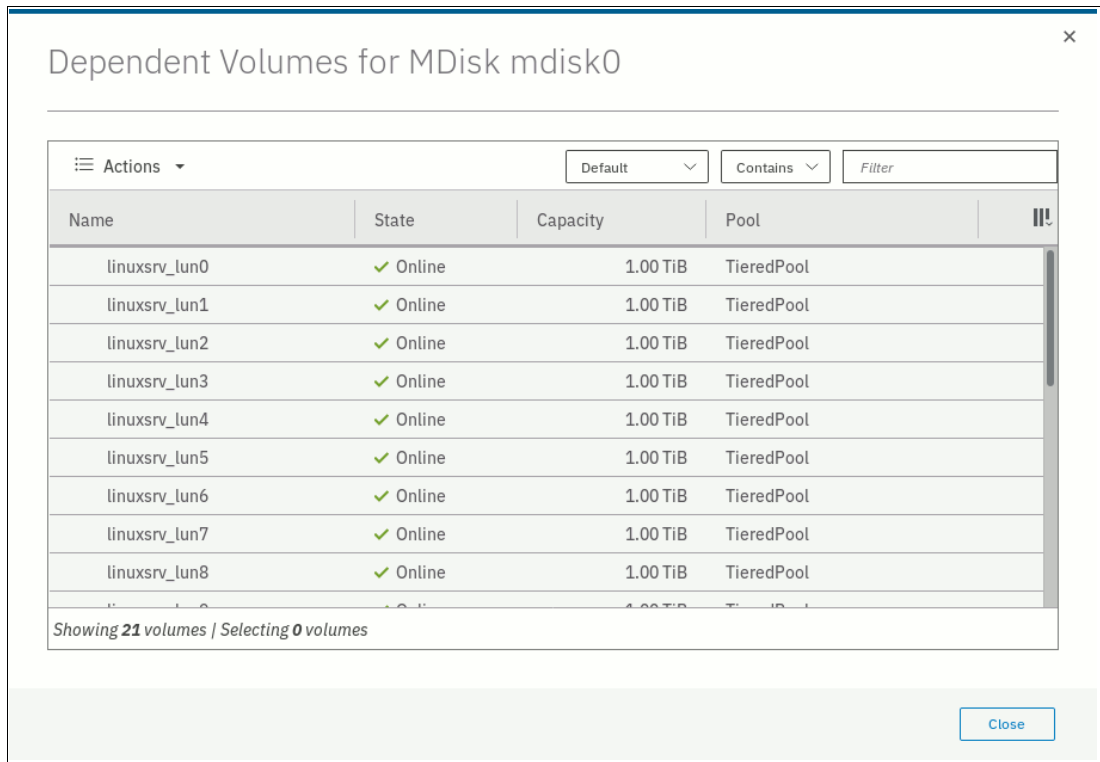


Figure 5-42 Dependent volumes for an MDisk

You can get the same information by running the `lsdependentvdisks` command (see Example 5-17).

*Example 5-17 Listing VDisks that depend on an MDisk by using the CLI*

```
IBM_2145:ITS0-SV1:superuser>lsdependentvdisks -mdisk mdisk0
vdisk_id vdisk_name
0        linuxsrv_lun0
1        linuxsrv_lun1
2        linuxsrv_lun2
3        linuxsrv_lun3
4        linuxsrv_lun4
<...>
```



## View Provisioning Groups

Provisioning groups are used for capacity reporting and monitoring of overprovisioned external storage controllers. Each overprovisioned MDisk is part of a provisioning group that defines the physical storage resources that are available to a set of MDisks. Storage controllers report the usable capacity of an overprovisioned MDisk based on its provisioning group. If multiple MDisks are part of the same provisioning group, then these MDisks share the physical storage resources and report the same usable capacity. However, this usable capacity is not available to each MDisk individually because it is shared among all these MDisks.

To know the usable capacity that is available to the system or to a pool when overprovisioned storage is used, you must account for the usable capacity of each provisioning group. To show a summary of overprovisioned external storage, including controllers, MDisks, and provisioning groups, click **View Provisioning Groups**, as shown in Figure 5-43.

MDisk Name	State	Written Capacity Limit	Pool	Storage System - LUN
mdisk5	Online	500.00 GiB	Pool1	flashsystem4_c1 - 000000000...
mdisk6	Online	500.00 GiB	Pool1	flashsystem4_c1 - 000000000...

Figure 5-43 View Provisioning Groups

## 5.3 Working with internal drives and arrays

An *array* is a type of MDisk that is made up of disk drives (or flash drive modules). These drives are members of the array. A RAID is a method of configuring member drives to create high availability (HA) and high-performance groupings of drives. The system supports nondistributed (traditional) and distributed array configurations.

### 5.3.1 Working with drives

This section describes how to manage internal storage disk drives and configure them to be used in arrays.

## Listing disk drives

The system provides an Internal Storage pane for managing all internal drives. To access the Internal Storage pane, click **Pools** → **Internal Storage**, as shown in Figure 5-44.

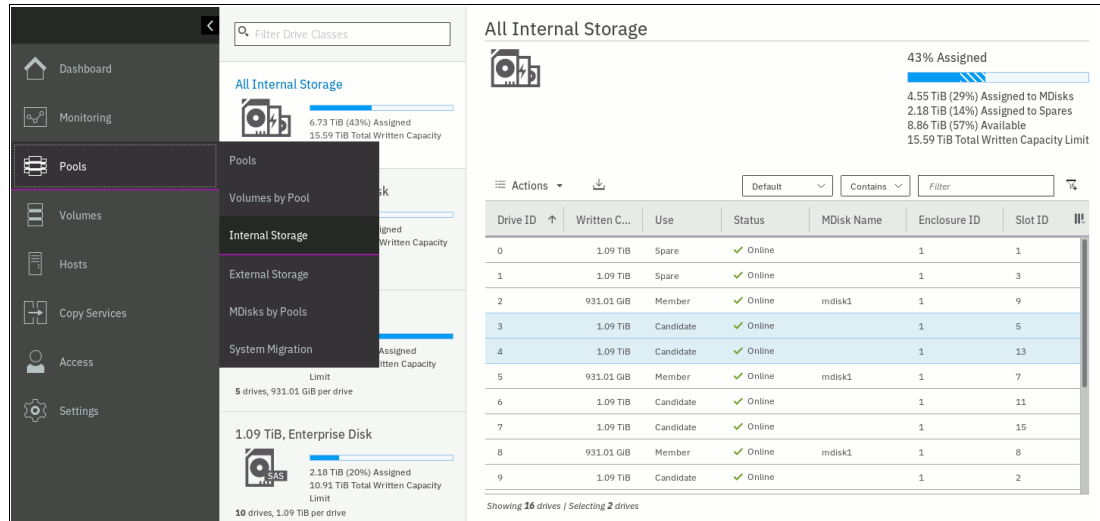


Figure 5-44 Internal storage pane

The pane gives an overview of the internal drives in the system. Select **All Internal Storage** in the drive class filter to display all drives that are managed in the system, including all I/O groups and expansion enclosures.

Alternatively, you can filter the drives by their type or class. For example, you can choose to show only Enterprise disks, Nearline (NL) disks, or Flash. Select the class on the left side of the pane to filter the list and display only the drives of the selected class.

You can find information about the capacity allocation of each drive class in the upper right corner, as shown in Figure 5-44:

- ▶ **Assigned to MDisks** shows the storage capacity of the selected drive class that is assigned to MDisks.
- ▶ **Assigned to Spares** shows the storage capacity of the selected drive class that is used for spare drives.
- ▶ **Available** shows the storage capacity of the selected drive class that is not yet assigned to either MDisks or Spares.
- ▶ **Total Written Capacity Limit** shows the total amount of storage capacity of the drives in the selected class.

If **All Internal Storage** is selected under the drive class filter, the values that are shown refer to the entire internal storage.

The percentage bar indicates how much of the total written capacity limit is assigned to MDisks and Spares. MDisk capacity is represented by the solid portion and spare capacity by the shaded portion of the bar.

To list all internal drives available in the system, run the `lsdrive` CLI command. If needed, you can filter output to list only drives that belong to particular enclosure, only that have specific capacity, or by another attributes (see Example 5-18).

*Example 5-18 lsdrive output (some lines and columns are not shown)*

```
IBM_2145:ITS0-SV1:superuser>lsdrive
id status error_sequence_number use      tech_type      capacity mdisk_id
0  online                spare      tier_enterprise 1.1TB
1  online                spare      tier_enterprise 1.1TB
2  online                member     tier_nearline  931.0GB  16
3  online                candidate  tier_enterprise 1.1TB
4  online                candidate  tier_enterprise 1.1TB
5  online                member     tier_nearline  931.0GB  16
<...>
```

The drive list shows Status of each drive. A drive can be Online, which means that drive is fully accessible by both nodes in the I/O group. A Degraded drive is accessible only by one of the two nodes. An Offline status indicates that the drive is not accessible by any of the nodes; for example, because it was physically removed from the enclosure or because it is unresponsive or failing.

The drive Use attribute describes the role that it plays in the system. The values and meanings are:

- ▶ **Unused:** The system can access the drive, but was not told to take ownership of it. Most actions on the drive are not permitted. This state is a safe state for newly added hardware.
- ▶ **Candidate:** The drive is owned by the system, and is not currently part of the RAID configuration. It is available to be used in an array MDisk.
- ▶ **Spare:** The drive is a hot spare protecting nondistributed (traditional) RAID arrays. If any member of such an array fails, a spare drive is taken and becomes a Member for rebuilding the array.
- ▶ **Member:** The drive is part of a RAID array.
- ▶ **Failed:** The drive is owned by the system and is diagnosed as faulty. It is waiting for a service action.

The drive use can transition between different values, but not all transitions are valid, as shown in Figure 5-45.

		To				
		unused	candidate	failed	member	spare
From	unused	yes	yes	no	no	no
	candidate	yes	yes	yes	no	yes
	failed	yes	yes	yes	no	no
	member	no	no	yes	no	no
	spare	no	yes	yes	no	yes

*Figure 5-45 Drive use transitions*

The system automatically sets the drive use to Member when creating a RAID array. Changing the drive use from Member to Failed is allowed only if the array is not dependent on the drive and more confirmation is required when taking a drive offline when no spare is available. Transitioning a Candidate drive to Failed is possible only by using the CLI.

**Note:** To start configuring arrays in a new system, all Unused drives must be configured as Candidates. Initial setup or Assign storage GUI wizards do that automatically.

Several actions can be performed on internal drives. To perform any action, select one or more drives and right-click the selection, as shown in Figure 5-46. Alternatively, select the drives and click **Actions**.

The screenshot shows a storage management interface. On the left, there are three drive details panels:
 

- 136.23 GiB, Enterprise Disk: 0 bytes (0%) Assigned, 681.16 GiB Total Written Capacity Limit, 5 drives, 136.23 GiB per drive.
- 931.01 GiB, Nearline Disk: 3.64 TiB (80%) Assigned, 4.55 TiB Total Written Capacity Limit, 5 drives, 931.01 GiB per drive.
- 1.09 TiB, Enterprise Disk.

 On the right, there is a table of drives with columns: Drive ID, Written C..., Use, Status, MDisk Name, Enclosure ID, and Slot ID. Drive 1 is selected, and a context menu is open over it, showing options: Fix Error, Take Offline, Mark as..., Identify, Upgrade, Dependent Volumes, and Properties.

Drive ID	Written C...	Use	Status	MDisk Name	Enclosure ID	Slot ID
0	1.09 TiB	Member	✓ Online	mdisk0	1	1
1	1.09 TiB	Candidate	✓ Online		1	3
2	931.01 GiB	Member	✓ Online		1	9
3	1.09 TiB	Candidate	✓ Online		1	5
4	1.09 TiB	Member	✓ Online		1	13
5	931.01 GiB	Member	✓ Online		1	7
6	1.09 TiB	Member	✓ Online		1	11
7	1.09 TiB	Member	✓ Online		1	15
8	931.01 GiB	Member	✓ Online		1	8

Figure 5-46 Actions on internal storage

The actions that are available in the drop-down menu depend on the status and use of the drives that are selected. Some actions can be performed only on drives in a specific state, and some are possible only when a single drive is selected.

### Action: Fix error

This action is available only if the selected drive has an error event associated with it. Select **Fix Error** to start the DMP for the selected drive. For more information about DMPs, see Chapter 13, “Reliability, availability, and serviceability, and monitoring and troubleshooting” on page 753.

### Action: Take offline

If a problem is identified with a specific drive, you can select **Take Offline** to take the drive offline. You must confirm the action, as shown in Figure 5-47.

The screenshot shows a dialog box titled "Take Drive Offline (drive 3)". It contains a warning icon and the following text:
 

- This action takes the drive offline and is to be used only when there are problems on the drive.
- The system prevents the drive from being taken offline if there will be resulting data loss.
- Only take a drive offline if a spare drive is available.
- Take the drive offline even if redundancy is lost on the array.

 At the bottom, there are "Cancel" and "OK" buttons.

Figure 5-47 Taking a drive offline

The system prevents you from taking the drive offline if one of the following conditions is true:

- ▶ Taking the drive offline results in a loss of access to data.
- ▶ The first option was selected and no suitable spares are available to start a rebuild.

If a spare is available and the drive is taken offline, the associated MDisk remains Online and the RAID array starts a rebuild by using a suitable spare. If no spare is available and the drive is taken offline by using the second option, the status of the associated MDisk becomes Degraded. The status of the storage pool to which the MDisk belongs becomes Degraded as well.

A drive that is taken offline is considered Failed, as shown in Figure 5-48.

Drive ID	Written C...	Use	Status	MDisk Name	Enclosure ID	Slot ID
0	1.09 TiB	Member	Online	mdisk0	1	1
1	1.09 TiB	Member	Online	mdisk0	1	3
2	931.01 GiB	Member	Online	mdisk1	1	9
3	1.09 TiB	Failed	Offline		1	5

Figure 5-48 An offline drive is marked as failed

To take a drive offline with the CLI, run the `chdrive` command (see Example 5-19). This command returns no feedback. Use the `-allowdegraded` parameter to set a member drive offline, even if no suitable spare is available.

*Example 5-19 Setting drive offline with CLI*

```
IBM_2145:ITS0-SV1:superuser>chdrive -use failed 3
IBM_2145:ITS0-SV1:superuser>
```

The system prevents you from taking a drive offline if the RAID array depends on that drive and doing so results in a loss of access to data, as shown in Figure 5-49.

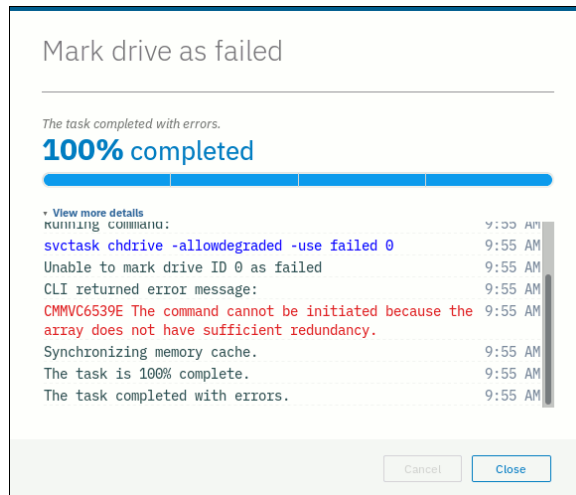


Figure 5-49 Taking a drive offline fails if it might result in a loss of access to data

## Action: Mark as

Select **Mark as** to change the use that is assigned to the drive, as shown in Figure 5-50. The list of available options depends on the current drive use and state. Refer to the allowed state transitions that are shown in Figure 5-45 on page 233.

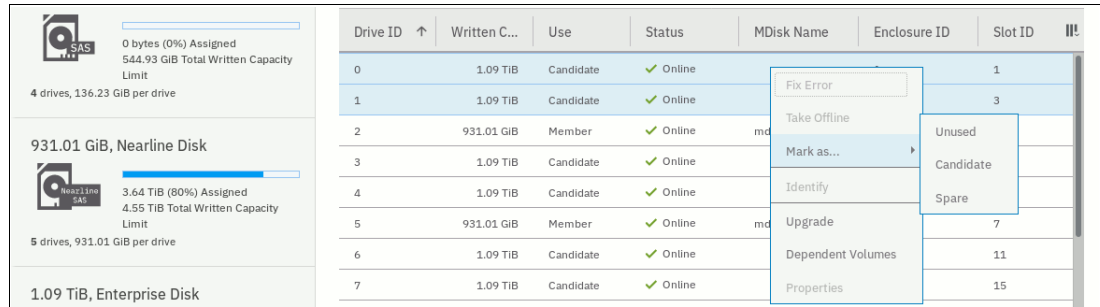


Figure 5-50 A drive can be marked as Unused, Candidate, or Spare

To change the drive role with the CLI, run the **chdrive** command (see Example 5-20). It shows the drive that was set offline with a previous command is set to spare. Notice that it cannot go from Failed to Spare use in one step. It needs to be assigned to a Candidate role before.

### Example 5-20 Changing drive role with CLI

```
IBM_2145:ITS0-SV1:superuser>lsdrive -filtervalue status=offline
id status error_sequence_number use tech_type capacity mdisk_id
3 offline failed tier_enterprise 558.4GB
IBM_2145:ITS0-SV1:superuser>chdrive -use spare 3
CMMVC6537E The command cannot be initiated because the drive that you have
specified has a Use property that is not supported for the task.
IBM_2145:ITS0-SV1:superuser>chdrive -use candidate 3
IBM_2145:ITS0-SV1:superuser>chdrive -use spare 3
IBM_2145:ITS0-SV1:superuser>
```

## Action: Identify

Select **Identify** to turn on the light-emitting diode (LED) light of the enclosure slot of the selected drive. This allows you to easily locate a drive that must be replaced or that you want to troubleshoot. A dialog box opens, which confirms that the LED was turned, on as shown in Figure 5-51.

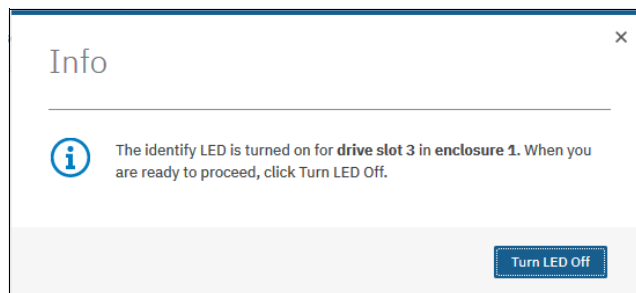


Figure 5-51 Identifying an internal drive

This action makes the amber LED that is associated with the drive that had this action performed flash continuously.

Click **Turn LED Off** when you are finished. The LED is returned to its initial state.

On the CLI, run the **chenclosureslot** command to turn on the LED. See Example 5-21 for commands to locate the enclosure and slot for drive 21 and then to turn the identification LED of slot 3 in enclosure 1 on and off again.

*Example 5-21 Changing slot LED to identification mode with CLI*

---

```
IBM_2145:ITS0-SV1:superuser>lsdrive 21
id 21
<...>
enclosure_id 1
slot_id 3
<...>
IBM_2145:ITS0-SV1:superuser>chenclosureslot -identify yes -slot 3 1
IBM_2145:ITS0-SV1:superuser>lsenclosureslot -slot 3 1
enclosure_id 1
slot_id 3
fault_LED slow_flashing
powered yes
drive_present yes
drive_id 21
IBM_2145:ITS0-SV1:superuser>chenclosureslot -identify no -slot 3 1
```

---

### Action: Upgrade

Selecting **Upgrade** allows the user to update the drive firmware, as shown in Figure 5-52. You can choose to update an individual drive, selected drives, or all the drives in the system.

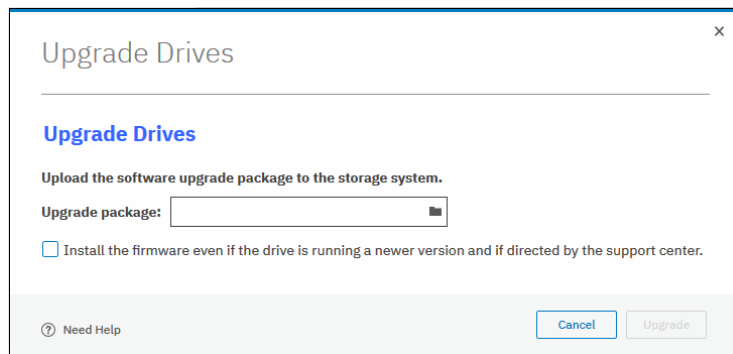


Figure 5-52 Upgrading a drive or a set of drives

For more information about updating drive firmware, see Chapter 13, “Reliability, availability, and serviceability, and monitoring and troubleshooting” on page 753.

### Action: Dependent volumes

Select **Dependent Volumes** to list the volumes that depend on the selected drives. A volume depends on a drive or a set of drives when removal or failure of that drive or set of drives results in a loss of access or a loss of data for that volume. Use this option before you perform maintenance operations to confirm which volumes (if any) are affected.

Figure 5-53 shows the list of volumes that depend on a set of three drives that belong to the same MDisk. All listed volumes go offline if *all* selected drives go offline at the same time. This does *not* mean that volumes go offline if a single drive or two of the three drives were to go offline.

Name	State	Capacity	Pool
vdisk0	✓ Online	100.00 GiB	mdiskgrp1
vdisk1	✓ Online	10.00 GiB	mdiskgrp1
vdisk2	✓ Online	10.00 GiB	mdiskgrp1
vdisk3	✓ Online	10.00 GiB	mdiskgrp1
vdisk4	✓ Online	10.00 GiB	mdiskgrp1
vdisk5	✓ Online	10.00 GiB	mdiskgrp1
vdisk6	✓ Online	10.00 GiB	mdiskgrp1

Showing 7 volumes | Selecting 0 volumes

Figure 5-53 List of volumes dependent on disks 7, 8, 9

Whether there are dependent volumes depends on the redundancy of the RAID array at a moment. It is based on the RAID level, state of the array, and state of the other member drives in the array. For example, it takes three or more drives going offline at the same time in healthy RAID 6 array to have dependent volumes.

**Note:** A lack of dependent volumes does not imply that there are no volumes using the drive. Volume dependency shows the list of volumes that become unavailable if the drive or the set of selected drives becomes unavailable.

You can get the same information by running the CLI command `lsdependentvdisks`. Use the parameter `-drive` with a list of drive IDs that you are checking, separated with a colon (:), as shown in Example 5-22.

*Example 5-22 Listing volumes dependent on drives with the CLI*

```
IBM_2145:ITS0-SV1:superuser>lsdependentvdisks -drive 7:8:9
vdisk_id vdisk_name
0        vdisk0
1        vdisk1
2        vdisk2
3        vdisk3
4        vdisk4
5        vdisk5
6        vdisk6
```



## Action: Properties

Select **Properties** to view more information about the drive, as shown in Figure 5-54.

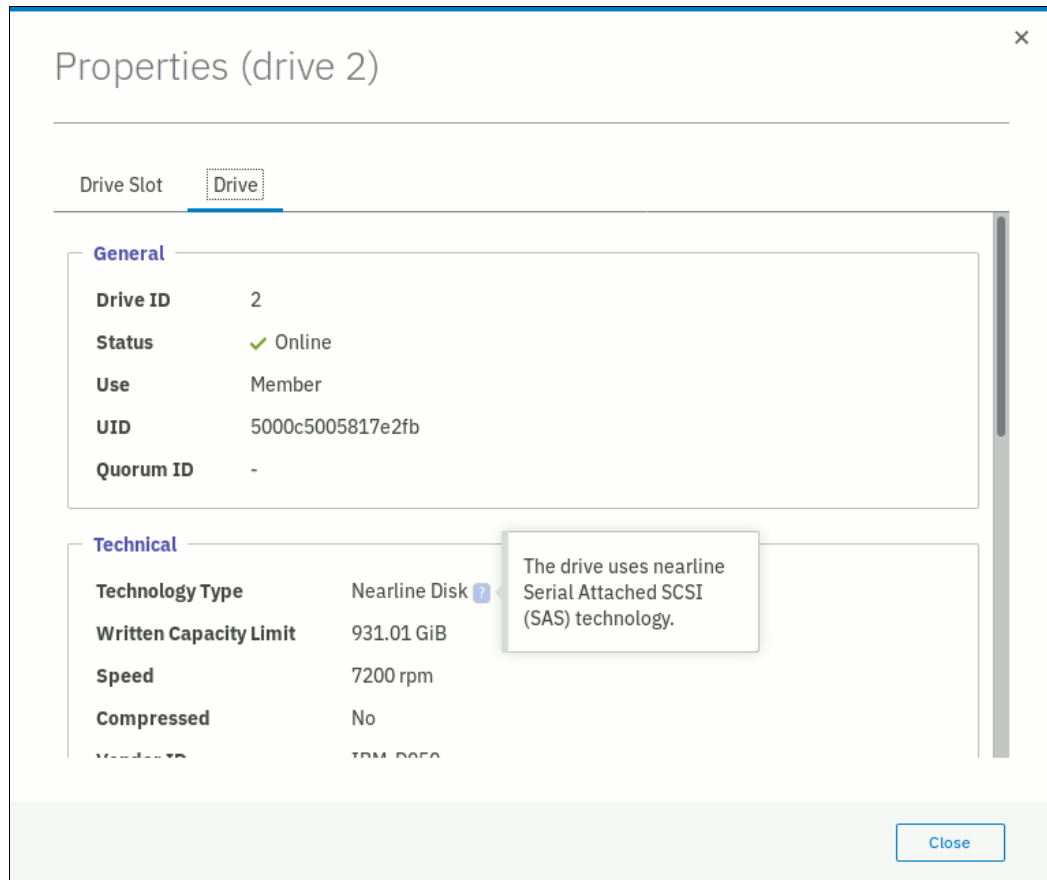


Figure 5-54 Drive properties

You can find a short description of each drive property by hovering on it and then clicking [?]. You can also display drive slot details by changing to the **Drive Slot** tab.

To get all available information about the particular drive, run the CLI command `l drive` with drive ID as the parameter. To get slot information, run the `l enclosureslot` command.

## 5.3.2 RAID and distributed RAID

To use internal SAN Volume Controller disks in storage pools, they must be joined into RAID arrays to form array mode MDisks.

RAID provides the following key design goals:

- ▶ Increased data reliability
- ▶ Increased input/output (I/O) performance

### Introduction to RAID technology

RAID technology can provide better performance for data access, HA for the data, or a combination of both. RAID levels define a tradeoff between HA, performance, and cost.

When multiple physical disks are set up to use the RAID technology, they are in a *RAID array*. The SAN Volume Controller provides multiple, traditional RAID (TRAIID) levels:

- ▶ RAID 0
- ▶ RAID 1
- ▶ RAID 5
- ▶ RAID 6
- ▶ RAID 10

**Note:** RAID 0 does not provide any redundancy. A single drive failure in a RAID 0 array causes data loss.

In a TRAIID approach, whether it is RAID 10, RAID 5, or RAID 6, data is spread among up to 16 drives in an array. There are separate spare drives that do not belong to an array and might protect multiple arrays. When one of the drives within the array fails, the system rebuilds the array by using a spare drive:

- ▶ For RAID 10, all data is read from the mirrored copy.
- ▶ For RAID 5 or RAID 6, data is calculated from remaining data stripes and parity.

This data is then written to a spare drive. The spare becomes a member of the array when the rebuild starts. After the rebuild is complete and the failed drive is replaced, a member exchange is performed to add the replacement drive to the array and to restore the spare to its original state so it can act as a hot spare again for another drive failure in the future.

During a rebuild of a TRAIID array, writes are submitted to a single spare drive that can become a bottleneck and might affect I/O performance. With increasing drive capacity, the rebuild time increases significantly. The probability of a second failure during the rebuild process also becomes more likely. Outside of any rebuild activity, the spare drives are idle and do not process I/O requests for the system.

Distributed Redundant Array of Independent Disks (DRAID) addresses these shortcomings.

## Distributed RAID

In DRAID, there are no dedicated spare drives being idle most of the time. All of the 4 - 128 drives in the array process I/O requests at all times instead, which improves the overall I/O performance. Spare capacity is spread across all member drives to form one or more *rebuild areas*. During a rebuild, write workload is distributed across all drives, which removes the single drive bottleneck of traditional arrays.

By using this approach, DRAID reduces the rebuild time, the impact on I/O performance during rebuild and the probability of a second failure during rebuild. As with TRAIID, a DRAID 6 array can tolerate two drive failures and survive. If another drive fails in the same array before the array is rebuilt, the MDisk and the storage pool go offline. That is, DRAID has the same redundancy characteristics as TRAIID.

A rebuild after a drive failure reconstructs the data on the failed drive and distributes it across all drives in the array by using a rebuild area. After the failed drive is replaced, a copyback process copies the data to the replacement drive and to free up the rebuild area so that it can be used for another drive failure in the future.

The following DRAID types are available:

- ▶ DRAID 5
- ▶ DRAID 6

Figure 5-55 shows an example of a DRAID 6 with 10 disks. The capacity on the drives is divided into many packs. The reserved spare capacity (marked in yellow) is equivalent to two spare drives, but the capacity is distributed across all of the drives (depending on the pack number) to form two rebuild areas. The data is striped similar to a TRAIID array, but the number of drives in the array can be larger than the stripe width.

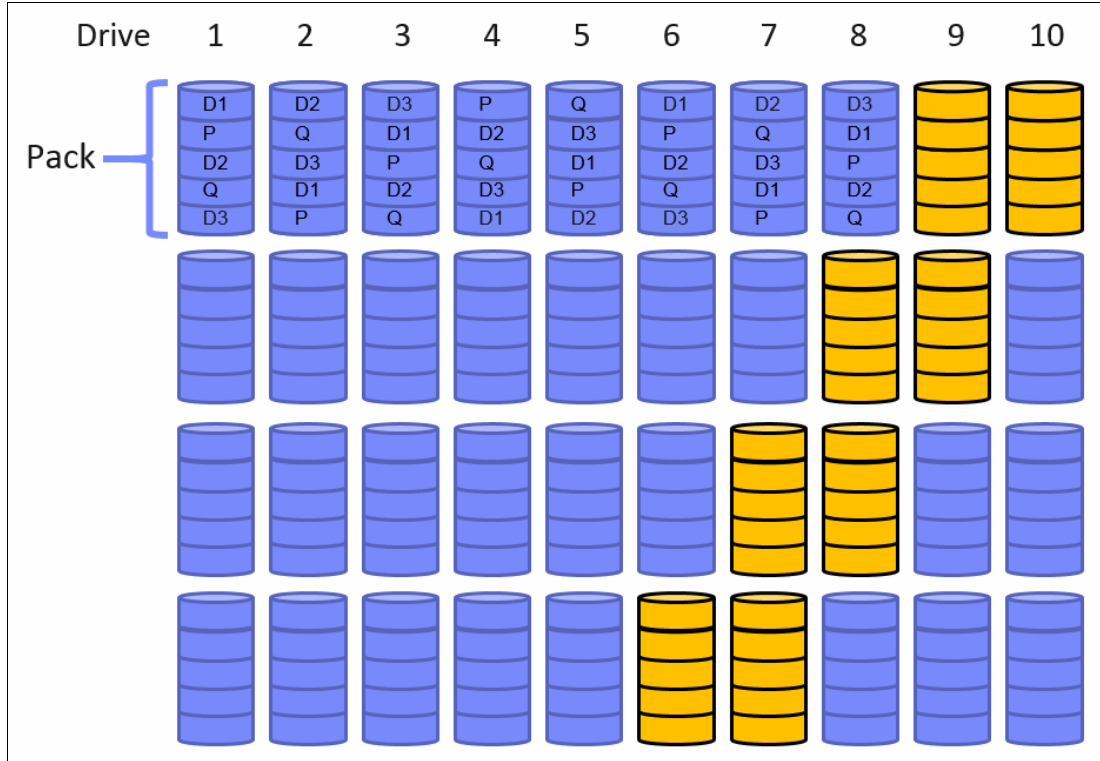


Figure 5-55 Distributed RAID 6 (for simplification, not all packs are shown)

Figure 5-56 on page 242 shows what happens after a single drive failure in this DRAID 6 array. Drive 3 failed and the array is using half of the spare capacity in each pack (marked in green) to rebuild the data of the failed drive. All drives are involved in the rebuild process, which significantly reduces the rebuild time. One of the two distributed rebuild areas is in use, but the second rebuild area can be used to rebuild the array once more after another failure.

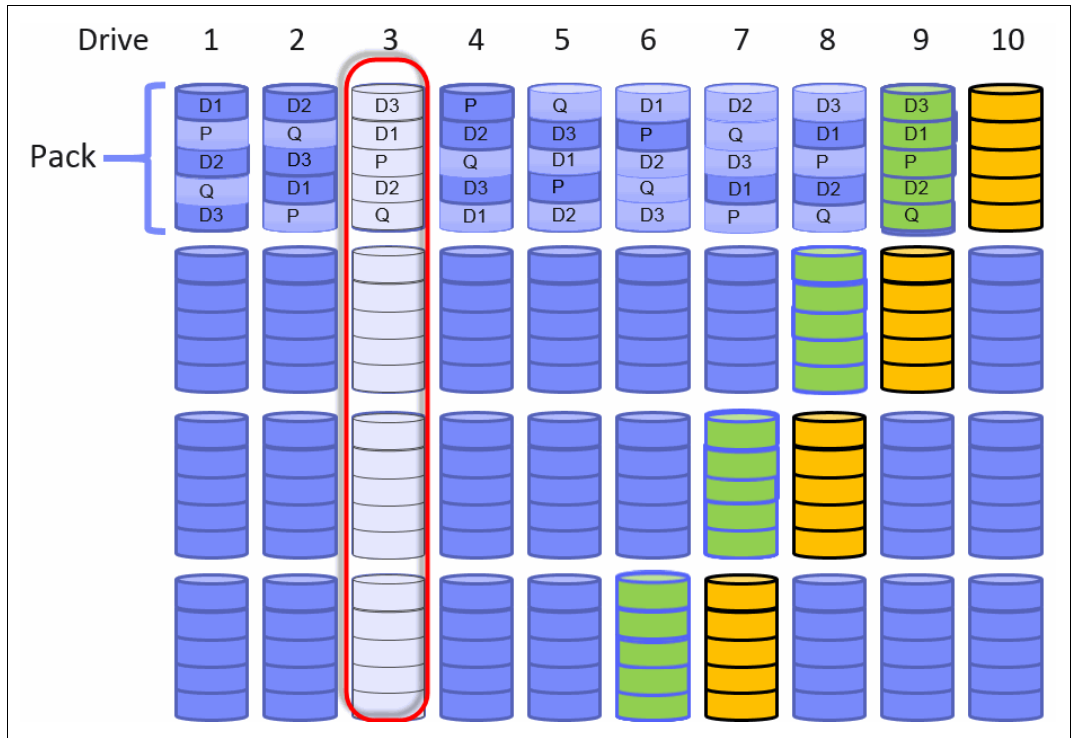


Figure 5-56 Single drive failure with DRAID 6 (for simplification, not all packs are shown)

After the rebuild completes, the array can sustain two more drive failures, even before drive 3 is replaced. If no rebuild area is available to perform a rebuild after another drive failure, the array becomes Degraded until a rebuild area is available again and the rebuild can start.

After drive 3 is replaced, a copyback process copies the data from the occupied rebuild area to the replacement drive to empty the rebuild area and make sure it can be used again for a new rebuild.

DRAID addresses the main disadvantages of TRAITD while providing the same redundancy characteristics:

- ▶ In case of a drive failure, data is read from many drives and written to many drives. This process minimizes the impact on performance during the rebuild process. Also, it significantly reduces rebuild time. Depending on the distributed array configuration and drive sizes, the rebuild process can be up to 10 times faster.
- ▶ Spare space is distributed throughout the array, which means more drives are processing I/O, and no dedicated spare drives are idling.

The DRAID implementation has the following other advantages:

- ▶ Arrays can be much larger than before and can span many more drives, which improves the performance of the array. The maximum number of drives a DRAID can contain is 128.
- ▶ Existing distributed arrays can be expanded by adding one or more drives. Traditional arrays cannot be expanded.
- ▶ Distributed arrays use all node CPU cores to improve performance, especially in configurations with a small number of arrays.

The following minimum number of drives are needed to build a Distributed Array:

- ▶ Six drives for a DRAID 6 array
- ▶ Four drives for a DRAID 5 array

### 5.3.3 Creating arrays

Only RAID arrays (array mode MDisks) can be added to a storage pool. It is not possible to add a Just A Bunch Of Disks (JBOD) or a single drive. It is also not possible to create a RAID array without assigning it to a storage pool.

**Note:** It is recommended to use DRAID 6 whenever possible. DRAID technology dramatically reduces rebuild times, decreases the exposure volumes have to the extra load of recovering redundancy, and improves performance.

To create a RAID array from internal storage, right-click the storage pool that you want to add it to and select **Add Storage**, as shown in Figure 5-57.

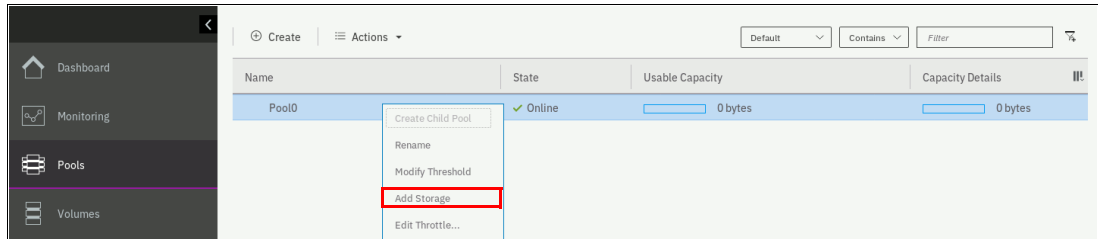


Figure 5-57 Adding storage to a pool

This action starts the configuration wizard that is shown in Figure 5-58. If any of the drives are found with an Unused role, it is suggested to reconfigure them as Candidates to be included into the configuration.

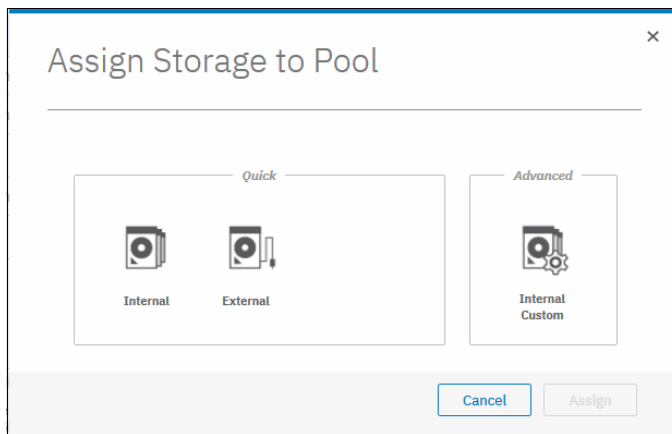


Figure 5-58 Assigning storage to a pool

If **Internal** or **Internal Custom** is chosen, the system guides you through the array MDisk creation process. If **External** is selected, systems guide you through the selection of external storage.

**Add Storage** dialog provides the following two options to configure internal arrays:

- ▶ **Quick** → **Internal**, which means the system automatically chooses the recommended RAID configuration for drives
- ▶ **Advanced** → **Internal Custom**, which provides more flexibility and lets you choose RAID parameters manually

### Quick Internal configuration

When selecting **Internal**, the system automatically recommends the RAID type (traditional or distributed) and level (such as RAID 6), number of drives, stripe width, and the number of rebuild areas (or spares for TRAITD) for each drive class. The number of drives and the stripe width can be adjusted before the array is created. Depending on the adjustments that are made, the system might select a different RAID type and level. A summary view can be expanded to preview the details of the arrays that are going to be created.

**Note:** It is not possible to change the RAID level or stripe width of an existing array. You also cannot change the drive count of a traditional array. If you need to change these properties, you must delete the array MDisk and re-create it with the required settings.

In the example that is shown in Figure 5-59, the dialog box suggests creating one DRAID 6 array with all of the 10 10 K enterprise drives by using a stripe width of 9 and to create a DRAID 5 array with a stripe width of 4 by using the five NL drives. An insufficient number of 15 K enterprise drives are available to create an array.

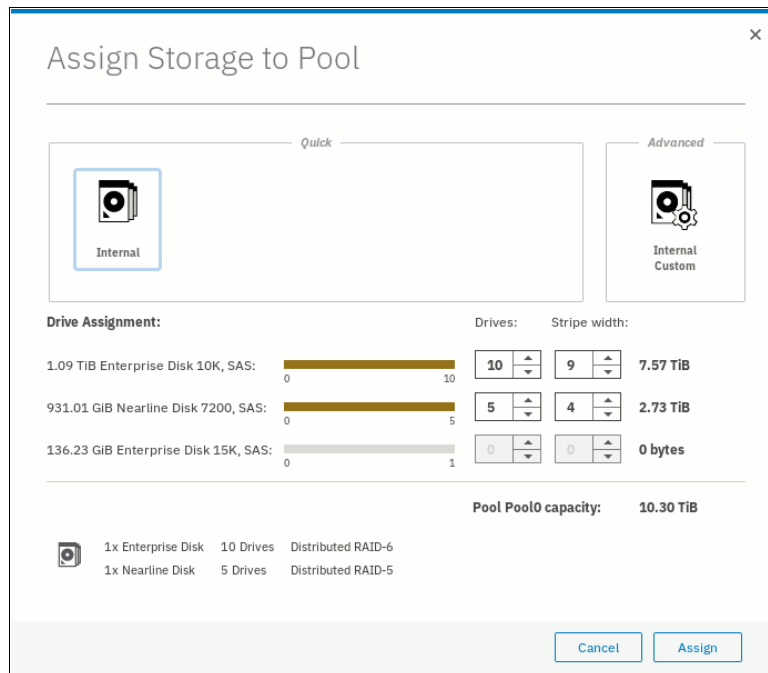


Figure 5-59 Assign Storage to Pool dialog

If you select two drives only, the system automatically updates the recommendation to a RAID 1 array. For more control of the array creation steps, you can select the **Internal Custom** option.

By default, if sufficient candidate drives are present, the system might recommend traditional arrays. However, switch to DRAID when possible, with the **Internal Custom** option.

If the system has multiple drive classes (such as Flash and Enterprise disks), the following default option is available:

1. Create multiple arrays of different tiers.
2. Assign them to the pool to take advantage of the Easy Tier functionality.

However, this configuration can be adjusted by setting the number of drives of different classes to zero. For more information about Easy Tier, see Chapter 9, “Advanced features for storage efficiency” on page 449.

If you are adding storage to a pool with storage already assigned, the existing storage configuration is considered for the recommendation. The system aims to achieve a balanced configuration, so some properties are inherited from existing arrays in the pool for a specific drive class. It is not possible to add RAID arrays significantly different from existing arrays in a pool when the GUI is used.

For example, if the pool has an existing DRAID 6 array made of 16 drives, you cannot add a two drive RAID 1 array to the same pool. Otherwise, an imbalanced storage pool is created.

You can still add any array of any configuration to an existing pool by using the CLI.

When you are satisfied with the configuration presented, click **Assign**. The RAID arrays are then created, added as array mode MDisks to the pool, and initialized in the background. You can monitor the progress of the initialization by selecting the corresponding task under **Running Tasks** in the upper-right corner of the GUI, as shown in Figure 5-60. The array is available for I/O during this process and there is no need to wait for it to complete.

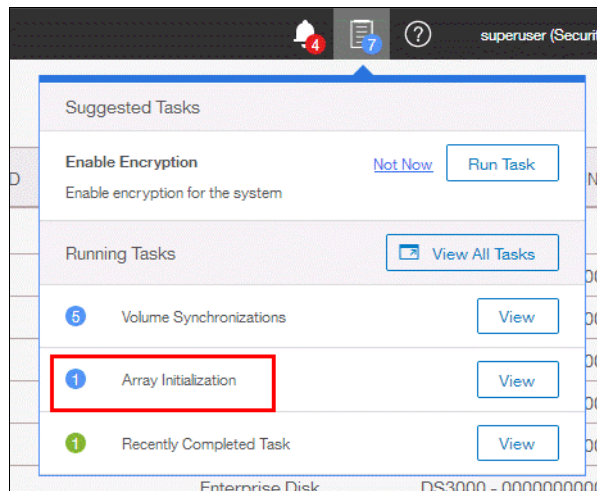


Figure 5-60 Array Initialization task

Click **View** in the Running tasks list to see the initialization progress and the time remaining, as shown in Figure 5-61 on page 246. The time it takes to initialize an array depends on the type of drives of which it consists. For example, an array of Flash drives is much quicker to initialize than NL-serial-attached SCSI (SAS) drives.

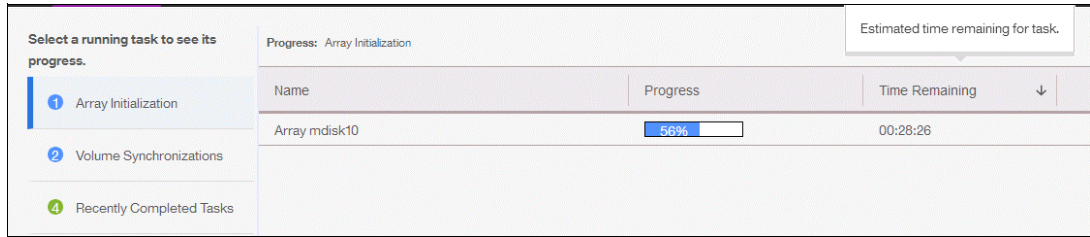


Figure 5-61 Array initialization task progress information

## Advanced Custom configuration

Select **Internal Custom** to customize the configuration of MDisks made out of internal drives. The following values can be customized:

- ▶ RAID type and level
- ▶ Number of spares (or spare areas)
- ▶ Array width
- ▶ Stripe width
- ▶ Number of drives of each class

Figure 5-62 shows an example with 10 drives ready to be configured as DRAID 6, with two rebuild areas distributed over all drives as spare capacity.

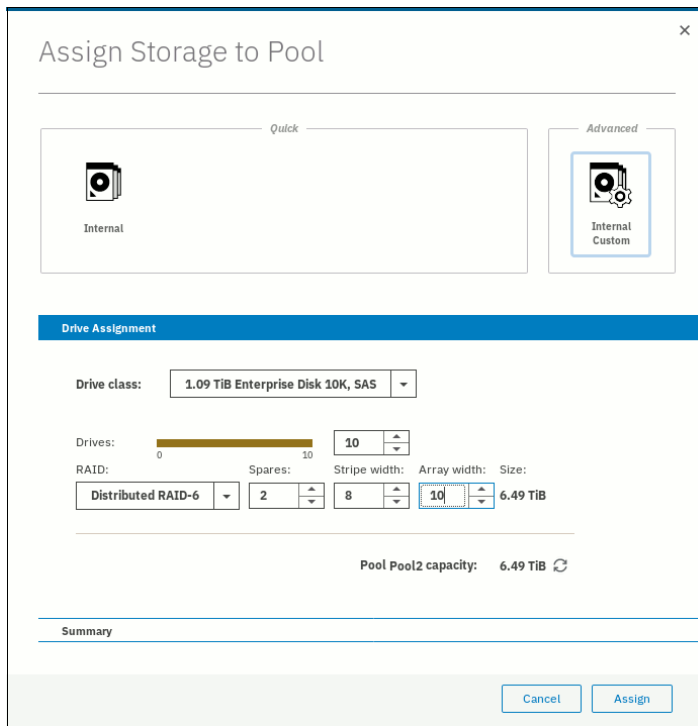


Figure 5-62 Adding internal storage to a pool using the Advanced option

To return to the default settings, click **Refresh** next to the pool capacity. To create and assign the arrays, click **Assign**.

Unlike automatic assignment, custom internal configuration does not create multiple arrays. Each array must be created separately.



As with automatic assignment, the system does not allow you to add significantly different arrays to an existing and populated pool because it aims to achieve balance across MDisks inside one pool. You can still add any array of any configuration to an existing pool by using the CLI.

**Note:** Spare drives are not assigned when TRAIID arrays are created by using the Internal Custom configuration wizard. You must set them up manually.

## Configuring arrays with the CLI

When working with the CLI, run the `mkarray` command to create TRAIID and the `mkdistributedarray` command to create DRAID. For this process, it is required to retrieve a list of drives that are ready to become array members. For more information about how to list all available drives, and read and change their use modes, see 5.3.1, “Working with drives” on page 231.

To get the recommended array configuration on the CLI, run the `lsdriveclass` command to list the available drive classes and the `lsarrayrecommendation` commands as shown in Example 5-23. The recommendations are listed in the order of preference.

*Example 5-23 Listing array recommendations using the CLI*

---

```
IBM_2145:ITS0-SV1:superuser>lsdriveclass
id RPM   capacity tech_type      block_size candidate_count
0  10000 1.1TB   tier_enterprise 512         10
1  7200  931.0GB tier_nearline  512         5
2  15000 136.2GB tier_enterprise 512         1
IBM_Storwize:ITS0V7K:superuser>lsarrayrecommendation -driveclass 0 -drivecount 10
Pool2
raid_level distributed stripe_width rebuild_areas drive_count array_count capacity
raid6      yes       9           1           10          1           7.6TB
raid6      no        10          0           10          1           8.7TB
raid5      yes       9           1           10          1           8.7TB
raid5      no        9           0           9           1           8.7TB
raid10     no        8           0           8           1           4.4TB
raid1      no        2           0           2           5           5.5TB
```

---

To create the recommended DRAID 6 array, specify the RAID level, drive class, number of drives, stripe width, number of rebuild areas, and the storage pool. The system automatically chooses drives for the array from the available drives in the class. As shown in Example 5-24, a DRAID 6 array is created out of 10 drives of class 0 by using a stripe width of 9 and a single rebuild area and adds it to Pool2.

*Example 5-24 Creating DRAID with `mkdistributedarray`*

---

```
IBM_2145:ITS0-SV1:superuser>mkdistributedarray -level raid6 -driveclass 0
-drivecount 10 -stripewidth 9 -rebuildareas 1 Pool2
MDisk, id [0], successfully created
```

---

There are default values for the stripe width and the number of rebuild areas, depending on RAID level and the drive count. In this example, it was required to specify the stripe width because for DRAID 6 it is 12 by default. The drive count value must equal or be greater than the sum of the stripe width and the number of rebuild areas.

To create a TRAIID 10 MDisk instead, you must specify a list of drives that you want to add as members, its RAID level, and the storage pool name or ID to which you want to add this array.

As shown in Example 5-25, a RAID-10 array is created and Pool2 is added to it. A spare drive also is designated.

*Example 5-25 Creating TRAIID with mkarray*

```
IBM_2145:ITS0-SV1:superuser>mkarray -level raid10 -drive 0:1:2:3:4:5:6:7 Pool2
MDisk, id [0], successfully created
IBM_Storwize:ITS0V7K:superuser>chdrive -use spare 8
```

**Note:** Do not forget to designate some of the drives as spares when creating traditional arrays. Spare drives are required to perform a rebuild immediately after a drive failure.

The storage pool must exist (see 5.1.1, “Creating storage pools” on page 202). To check array initialization progress with the CLI, run the `lsarrayinfo progress` command.

### 5.3.4 Actions on arrays

MDisks that are created from internal storage support specific actions that are not supported on external MDisks. Some actions that are supported on TRAIID arrays are not supported on DRAID arrays and vice versa.

To choose an action, open **Pools** → **MDisks by Pools**, select the array (MDisk), and click **Actions**. Alternatively, right-click the array, as shown in Figure 5-63.

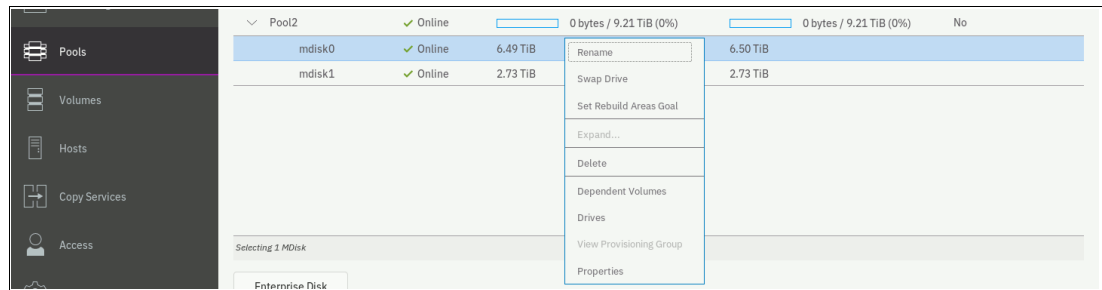


Figure 5-63 Actions on arrays

#### Rename

Select this option to change the name of an MDisk.

The CLI command for this operation is `charray` (see Example 5-26). No feedback is returned.

*Example 5-26 Renaming array MDisk with charray*

```
IBM_2145:ITS0-SV1:superuser>charray -name Distributed_array mdisk1
IBM_2145:ITS0-SV1:superuser>
```

#### Swap drive

Select **Swap Drive** to replace a drive in the array with another drive. The other drive must have use of Candidate or Spare. Use this action to perform proactive drive replacement to replace a drive that is not failed but is expected to fail soon, for example, as indicated by an error message in the event log.

Figure 5-64 shows the dialog box that opens. Select the member drive to be replaced and the replacement drive, and click **Swap**.

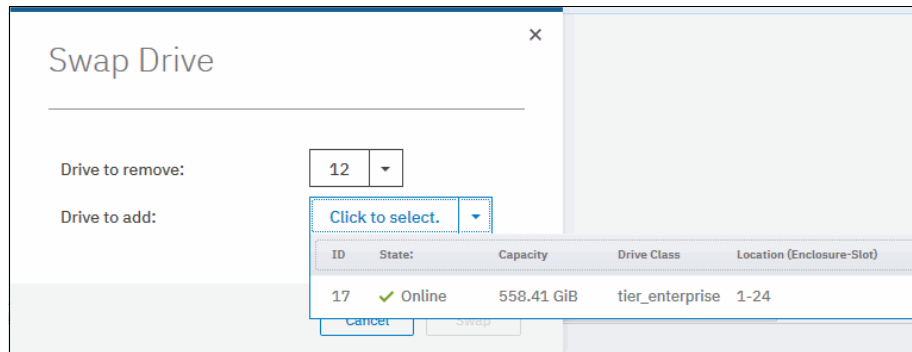


Figure 5-64 Swapping array member with another candidate or spare drive

The exchange of the drives starts running in the background. The volumes on the affected MDisk remain accessible during the process.

Swapping a drive in a traditional array performs a concurrent member exchange that does not reduce the redundancy of the array. The data of the old member is copied to the new member and after completion the old member is removed from the array.

In a distributed array, the system immediately removes the old member from the array and performs a rebuild. After the rebuild is complete, a copyback is started to copy the data to the new member drive. This process is nondisruptive, but reduces the redundancy of the array during the rebuild process.

The **charraymember** CLI command is run to perform this task. Example 5-27 shows the replacement of array member ID 7 that was assigned to drive ID 12, with drive ID 17. Notice that the **-immediate** parameter is required for distributed arrays to acknowledge that a rebuild starts.

Example 5-27 Replacing array member with CLI (some columns are not shown)

```
IBM_2145:ITS0-SV1:superuser>lsarraymember 16
mdisk_id mdisk_name      member_id drive_id new_drive_id spare_protection
16      Distributed_array 6         18         1
16      Distributed_array 7         12         1
16      Distributed_array 8         15         1
<...>
IBM_2145:ITS0-SV1:superuser>lsdrive
id status error_sequence_number use      tech_type      capacity
16 online                member    tier_enterprise 558.4GB 16
17 online                spare     tier_enterprise 558.4GB
18 online                member    tier_enterprise 558.4GB 16
<...>
IBM_2145:ITS0-SV1:superuser>charraymember -immediate -member 7 -newdrive 17
Distributed_array
IBM_2145:ITS0-SV1:superuser>
```

### Set Spare Goal or Set Rebuild Areas Goal option

Select this option to set the number of spare drives (on TRAITD) or rebuild areas (on DRAID) that are expected to protect the array from drive failures.

If the number of rebuild areas that is available does not meet the configured goal, an error is logged in the event log, as shown in Figure 5-65. This error can be fixed by replacing failed drives in the DRAID array.

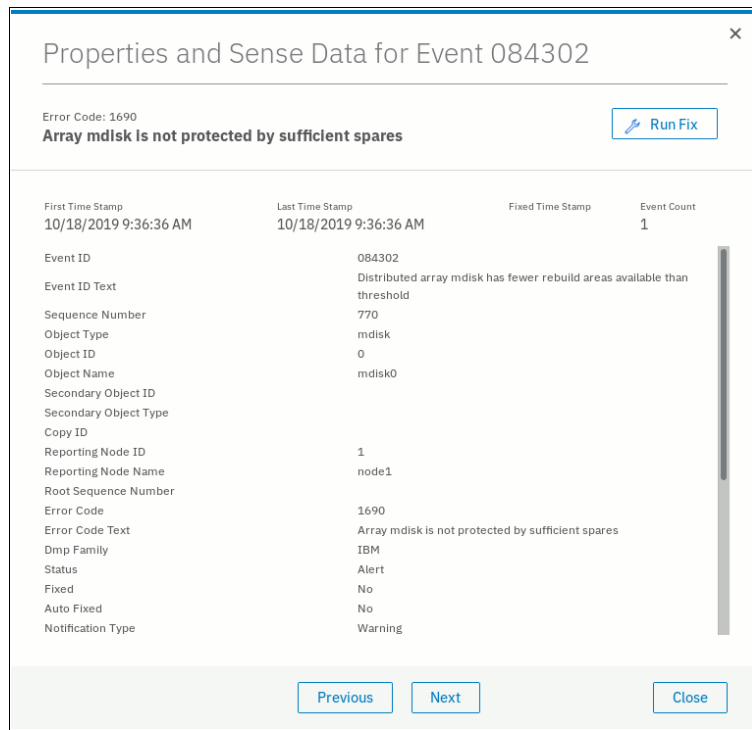


Figure 5-65 An error 1690 is logged if insufficient rebuild areas are available

**Note:** This option does not change the number of rebuild areas or spares that is available to the array. It specifies only at which point a warning event is generated. Setting the goal to 0 does not prevent the array from rebuilding.

In the CLI, this task is performed by using the **charray** command (see Example 5-28).

*Example 5-28 Adjusting array goals with **charray** (some columns are not shown)*

```
IBM_2145:ITS0-SV1:superuser>lsarray
mdisk_id mdisk_name      status mdisk_grp_id mdisk_grp_name distributed
0        mdisk0            online 0             mdiskgrp0 no
16       Distributed_array online 1             mdiskgrp1 yes
IBM_2145:ITS0-SV1:superuser>charray -sparegoal 2 mdisk0
IBM_2145:ITS0-SV1:superuser>charray -rebuildareasgoal 0 Distributed_array
```

## Expand

Select **Expand** to expand the array by adding drives to it to increase the available capacity of the array or to create rebuild areas. Only distributed arrays can be expanded, and the option is not available for traditional arrays.

Candidate drives of a drive class that is compatible to the drive class of the array must be available in the system; otherwise, an error message is shown and the array cannot be expanded. A drive class is compatible to another if its characteristics, such as capacity and performance, are an exact match or are superior. In most cases, drives of the same class should be used to expand an array.

The dialog box that is shown in Figure 5-66 is displayed to give the user an overview of the current size of the array, the number of available candidate drives in the selected drive class, and the new array capacity after the expansion. The drive class, and the number of drives to add can be modified as required, and the projected new array capacity are updated as needed. To add rebuild areas to the array, click **Advanced Settings** and modify the number of extra spares.

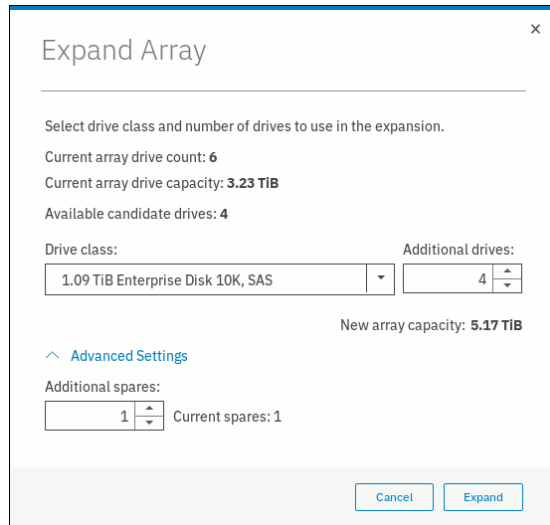


Figure 5-66 Expand a distributed array

Clicking **Expand** starts a background process that adds the selected number of drives to the array. As part of the expansion, the system automatically migrates data for optimal performance for the new expanded configuration.

You can monitor the progress of the expansion by clicking the **Running Tasks** icon in the upper-right corner of the GUI, or by selecting **Monitoring** → **Background tasks**, as shown in Figure 5-67.



Figure 5-67 Array expansion progress

On the CLI, this task is performed by running the `expandarray` command. To get a list of compatible drive classes, the `lscompatibledriveclasses` command can be run (see Example 5-29).

*Example 5-29 Expand an array using the CLI*

```
IBM_2145:ITS0-SV1:superuser>lsarray 0
<..>
capacity 3.2TB
<..>
drive_class_id 0
```

```

drive_count 6
<..>
rebuild_areas_total 1
IBM_2145:ITS0-SV1:superuser>lscompatibledriveclasses 0
id
0
IBM_2145:ITS0-SV1:superuser>expandarray -driveclass 0 -totaldrivecount 10
-totalrebuildareas 2 0
IBM_2145:ITS0-SV1:superuser>lsarrayexpansionprogress
progress estimated_completion_time target_capacity additional_capacity_remaining
29          191018233758             5.17TB          1.38TB

```

**Note:** The `expandarray` command on the CLI expects the total drive count *after* the expansion as a parameter, including the number of new drives and the number of drives in the array before the expansion. The same is true for the number of rebuild areas.

## Delete

Select **Delete** to remove the array from the storage pool and delete it. An array MDisk does not exist outside of a storage pool. Therefore, an array cannot be removed from the pool without being deleted. All drives that belong to the deleted array return into Candidate.

If no volumes use extents from this array, the command runs immediately without more confirmation. If volumes use extents from this array, you are prompted to confirm the action, as shown in Figure 5-68.

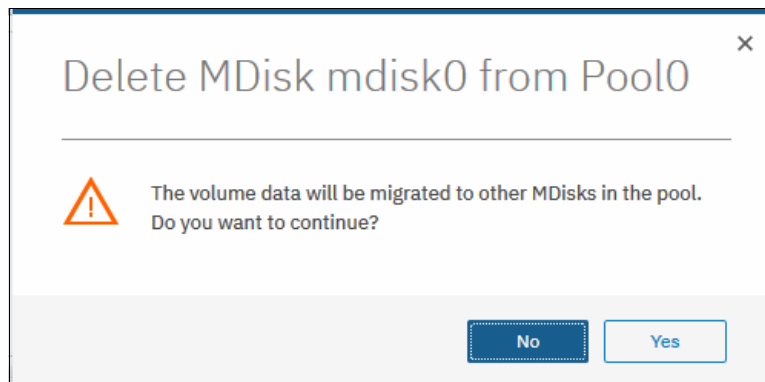


Figure 5-68 Deleting an array from a non-empty storage pool

Confirming the deletion starts a background process that migrates used extents on the MDisk to other MDisks in the same storage pool. After completion of that process, the array is removed from the storage pool and deleted.

**Note:** The command fails if not enough available capacity remains in the storage pool to allocate the capacity that is being migrated from the removed array.

To delete the array with the CLI, run the `rmarray` command. The `-force` parameter is required if volume extents must be migrated to other MDisks in a storage pool.

Use the **Running Tasks** section in the GUI or run the `lsmigrate` command on the CLI to monitor the progress of the migration. The MDisk continues to exist until the migration completes.

## Dependent volumes

A volume is dependent on an MDisk if the MDisk that is becoming unavailable results in a loss of access or a loss of data for that volume. Use this option before you conduct maintenance operations to confirm which volumes (if any) are affected.

If an MDisk in a storage pool goes offline, the entire storage pool goes offline. This means all volumes in a storage pool usually depend on each MDisk in the same pool, even if the MDisk does not have extents for each of the volumes. Clicking the **Dependent Volumes Action** menu of an MDisk lists the volumes that depend on that MDisk, as shown in Figure 5-69.

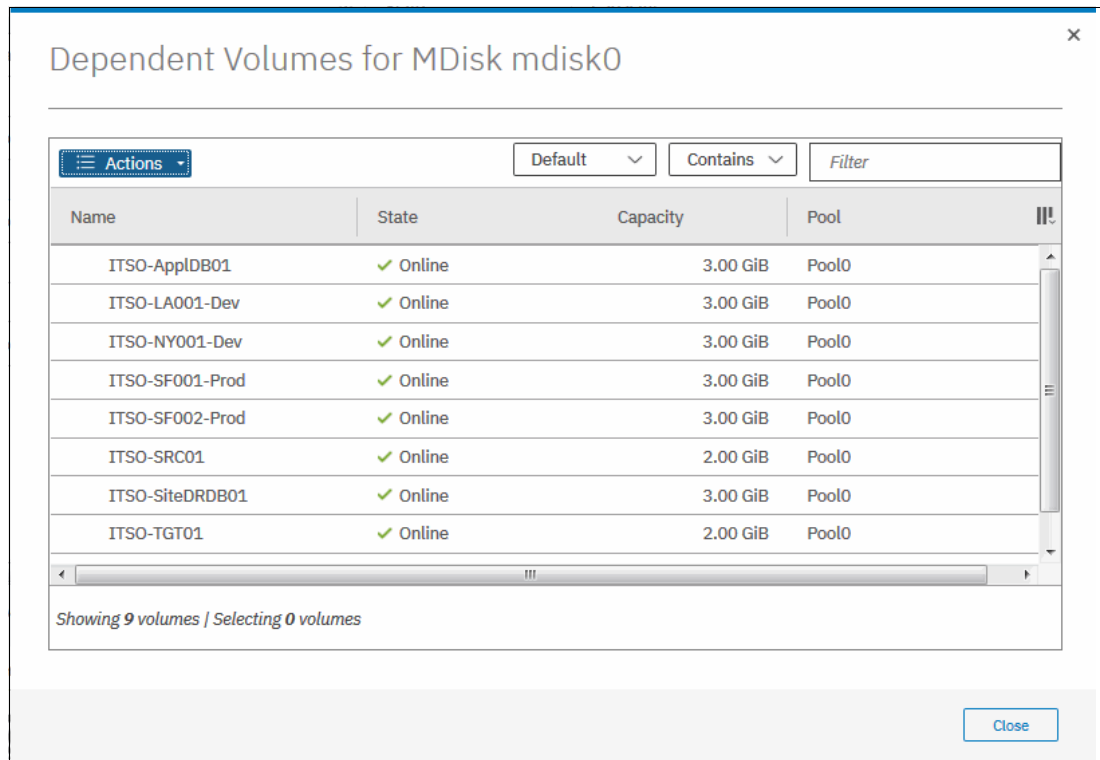


Figure 5-69 Dependent volumes for MDisk mdisk0

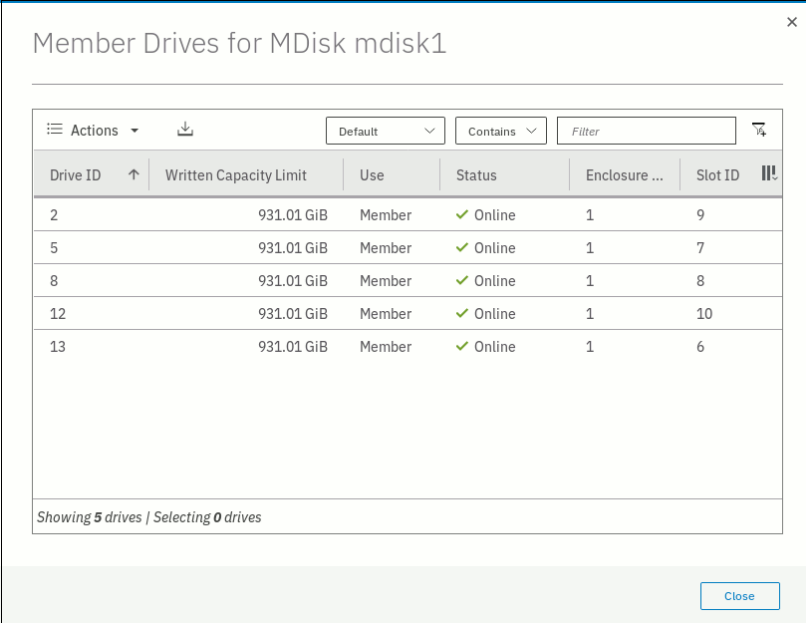
You can get the same information by running the CLI command `lsdependentvdisks` (see Example 5-30).

*Example 5-30 Listing VDisks that depend on MDisk with CLI*

```
IBM_2145:ITSO-SV1:superuser>lsdependentvdisks -mdisk mdisk0
vdisk_id vdisk_name
0        ITSO-SRC01
1        ITSO-TGT01
2        ITSO-App1DB01
<...>
```

## Drives

Select **Drives** to see information about the member drives that are included in the array, as shown in Figure 5-70.



The screenshot shows a window titled "Member Drives for MDisk mdisk1". It contains a table with the following columns: Drive ID, Written Capacity Limit, Use, Status, Enclosure ..., and Slot ID. The table lists five drives, all with a capacity of 931.01 GiB and a status of "Online".

Drive ID	Written Capacity Limit	Use	Status	Enclosure ...	Slot ID
2	931.01 GiB	Member	✓ Online	1	9
5	931.01 GiB	Member	✓ Online	1	7
8	931.01 GiB	Member	✓ Online	1	8
12	931.01 GiB	Member	✓ Online	1	10
13	931.01 GiB	Member	✓ Online	1	6

Showing 5 drives | Selecting 0 drives

Close

Figure 5-70 List of drives in an array

You can run the CLI command **lsarraymember** to get the same information with the CLI. Provide an array name or ID as the parameter to filter output by the array. If run without arguments, the command lists all members of all configured arrays.

## Properties

This section shows all available array MDisk parameters: its state, capacity, RAID level, and others.

Run the CLI command **lsarray** to get a list of all configured arrays and **lsarray** with array name or ID as the parameter to get extended information about the selected array, as shown in Example 5-31.

### Example 5-31 lsarray output (truncated)

```
IBM_2145:ITS0-SV1:superuser>lsarray
mdisk_id mdisk_name      status mdisk_grp_id mdisk_grp_name capacity
0         mdisk0             online 0             mdiskgrp0     1.3TB
16        Distributed_array online 1             mdiskgrp1     2.2TB
IBM_2145:ITS0-SV1:superuser>lsarray 16
mdisk_id 16
mdisk_name Distributed_array
status online
mode array
mdisk_grp_id 1
mdisk_grp_name mdiskgrp1
capacity 2.2TB
<...>
```





# Volumes

In IBM Spectrum Virtualize, a *volume* is an amount of storage space that is provisioned out of a storage pool and presented to a host as a Small Computer System Interface (SCSI) logical unit (LU); that is, a logical disk.

This chapter describes how to create and provision volumes on IBM Spectrum Virtualize systems. The first part of this chapter provides a brief overview of IBM Spectrum Virtualize volumes, the classes of volumes available, and available volume customization options.

The second part of this chapter describes how to create, modify, and map volumes by using the GUI.

The third part of this chapter provides an introduction to volume manipulation from the command-line interface (CLI).

This chapter includes the following topics:

- ▶ 6.1, “Introduction to volumes” on page 256
- ▶ 6.2, “Volume characteristics” on page 256
- ▶ 6.3, “Virtual volumes” on page 273
- ▶ 6.4, “Volumes in multi-site topologies” on page 274
- ▶ 6.5, “Operations on volumes” on page 276
- ▶ 6.6, “Volume operations by using the CLI” on page 322

## 6.1 Introduction to volumes

A volume is a logical disk that the system presents to attached hosts. For the attached host, a volume is a SCSI target that uses a 512-byte block size.

For an IBM Spectrum Virtualize system cluster, the volume that is presented to a host is internally represented as a virtual disk (VDisk). A VDisk is an area of usable storage that was allocated out of a pool of storage that is managed by IBM Spectrum Virtualize cluster. The term *virtual* is used because the volume that is presented does not necessarily exist on a single physical entity.

**Note:** Volumes are made out of extents that are allocated from a storage pool. Storage pools group managed disks (MDisks), which are Redundant Array of Independent Disks (RAID) arrays from internal storage, or LUs that are presented to and virtualized by IBM Spectrum Virtualize system. Each MDisk is divided into sequentially numbered extents (0 based indexing). The extent size is a property of a storage pool, and is used for all MDisks comprising the storage pool.

MDisks are internal objects that are used for storage management. They are not directly visible to or used by host systems.

A volume is presented to hosts by an I/O group, and within that group has a preferred node; that is, a node that by default serves I/O requests to that volume. When host requests an I/O operation to a volume, the multipath driver on the host identifies the preferred node for the volume and by default uses only paths to this node to for I/O requests.

## 6.2 Volume characteristics

The following parameters characterize each volume. They should be set correctly to match the requirements of the storage consumer (an application running on a host):

- ▶ Size
- ▶ Performance (input/output operations per second [IOPS], response time, and bandwidth)
- ▶ Resiliency
- ▶ Storage efficiency
- ▶ Security (data-at-rest encryption)
- ▶ Extent allocation policy
- ▶ Management mode

Also, volumes can be configured as VMware vSphere Virtual Volumes (VVOLs).

VVOLs change the approach to VMware virtual machines (VMs) disk configuration. Instead, files on a VMware VM file system (VMFS) distributed file system that is created on a single large volume, VVOLs introduce one-to-one mapping between VM disks and storage volumes. VVOLs can be managed by VMware infrastructure, which allows storage system administrators to delegate VM disk management to VMware infrastructure specialists.

To provide storage consumers adequate service, all the parameters need to be correctly set. Importantly, the various parameters can be interdependent or setting one of them might affect other aspects of the volume. The volume parameters and their interdependencies are covered in the following topics.

## 6.2.1 Volume type

The *type* attribute of a volume defines the method of allocation of extents that make up the volume copy:

- ▶ A *striped* volume contains a volume copy that has extents that are allocated from multiple MDisks from the storage pool. By default, extents are allocated by using round-robin algorithm from all MDisks in the storage pool. However, it is possible to supply a list of MDisks to use for volume creation.

**Attention:** By default, striped volume copies are striped across all MDisks in the storage pool (see Figure 6-1). If some of the MDisks are smaller than others, the extents on the smaller MDisks are used up before the larger MDisks run out of extents. Manually specifying the stripe set in this case might result in the volume copy not being created.

If you are unsure if sufficient free space is available to create a striped volume copy, select one of the following options:

- ▶ Check the free space on each MDisk in the storage pool by running the **lsfreextents** command.
- ▶ Allow the system to automatically create the volume copy by not supplying a specific stripe set.

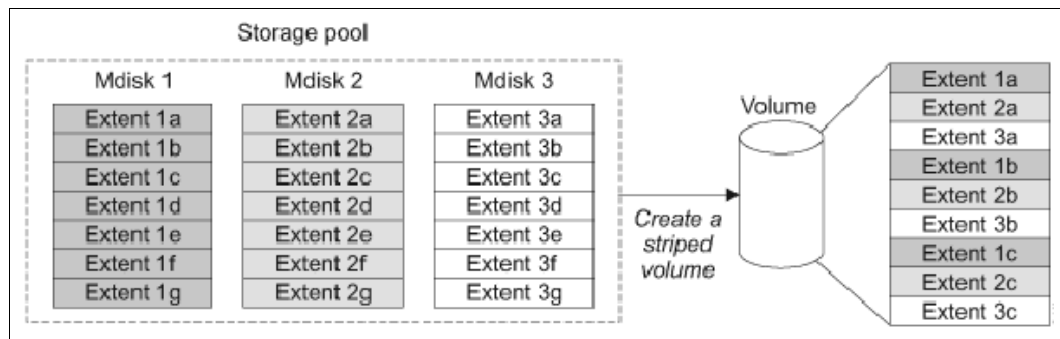


Figure 6-1 Striped extent allocation

- ▶ A *sequential* volume contains a volume copy that has extents that are allocated sequentially on one MDisk.
- ▶ *Image-mode* volume is a special type of volume that has a direct one-to-one mapping to one (image mode) MDisk.

For striped volumes, the extents are allocated from the set of MDisks (by default, all MDisks in the storage pool). Consider the following points:

- ▶ An MDisk is picked by using a pseudo-random algorithm and an extent is allocated from this MDisk. This approach minimizes the probability of triggering the *striping effect*, which might lead to poor performance for workloads that generate many metadata I/O, or that create multiple sequential streams.
- ▶ All subsequent extents (if required) are allocated from the MDisk set by using round-robin algorithm.
- ▶ If an MDisk has no free extents when its turn arrives, the algorithm moves to the next MDisk in the set that has a free extent.

**Note:** The *striping effect* occurs when multiple logical volumes that are defined on a set of physical storage devices (MDisks) store their metadata or file system transaction log on the same physical device (MDisk).

Because of the way the file systems work, many I/O to file system metadata disk regions exist. For example, for a journaling file system, a write to a file might require two or more writes to the file system journal; at minimum, one to make a note of the intended file system update, and one marking successful completion of the file write.

If multiple volumes (each with its own file system) are defined on the same set of MDisks, and all (or most) of them store their metadata on the same MDisk, a disproportionately large I/O load is generated on this MDisk, which can result in suboptimal performance of the storage system. Pseudo-randomly allocating first MDisk for new volume extent allocation minimizes probability, that multiple file systems that are created on these volumes places their metadata regions on the same physical MDisk.

## 6.2.2 Managed mode and image mode

Volumes are configured within IBM Spectrum Virtualize by allocating a set of extents off one or more managed mode MDisks in the storage pool. Extents are the smallest allocation unit at the time of volume creation, so each MDisk extent maps to exactly one volume extent.

**Note:** An MDisk extent maps to exactly one volume extent. For volumes with two copies, one volume extent maps to two MDisk extents (one for each volume copy).

Figure 6-2 shows this mapping. It also shows a volume that consists of several extents that are shown as V0 - V7. Each of these extents is mapped to an extent on one of the MDisks: A, B, or C. The mapping table stores the details of this indirection.

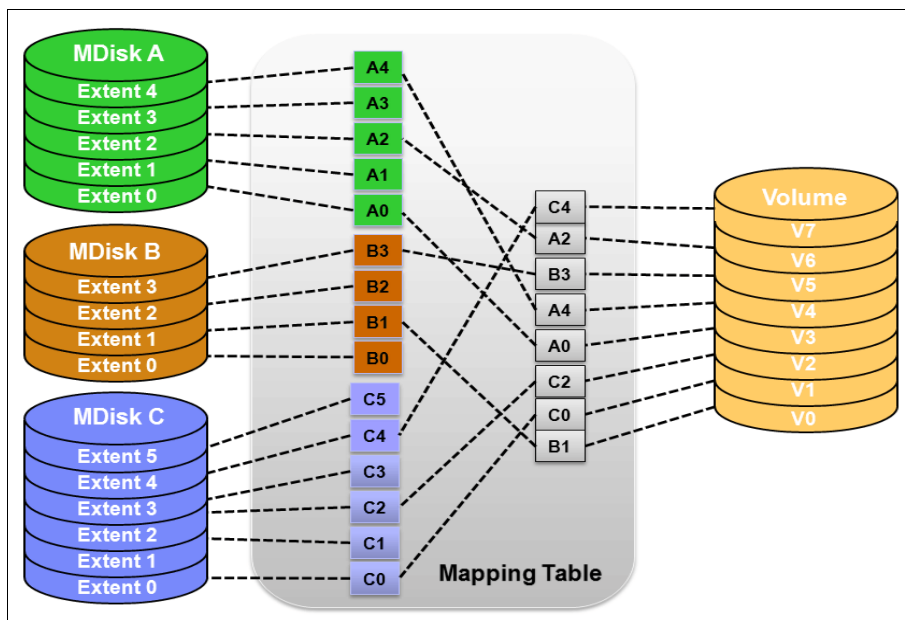


Figure 6-2 Simple view of block virtualization

Several of the MDisk extents are unused; that is, no volume extent maps to them. These unused extents are available for use in creating volumes, migration, expansion, and so on.

The default and most common type of volumes in IBM Spectrum Virtualize are managed mode volumes. Managed mode volumes are allocated from a set of MDisk that belong to a storage pool and can be subject of the full set of virtualization functions.

In particular, they offer full flexibility in mapping between logical volume representation (logical blocks) and physical storage that is used to store these blocks. This requires that physical storage (MDisks) are fully managed by IBM Spectrum Virtualize, which means that the LUs that are presented to the IBM Spectrum Virtualize by the back-end storage systems does not contain any data when it is added to the storage pool.

Image mode volumes enable IBM Spectrum Virtualize to work with LUs that were previously directly mapped to hosts. This is often required when IBM Spectrum Virtualize is introduced into a storage environment and image mode volumes are used to enable seamless migration of data and smooth transition to virtualized storage.

Image mode creates one-to-one mapping of logical block addresses (LBAs) between a volume and a single MDisk (LU presented by the virtualized storage). Image mode volumes have the minimum size of one block (512 bytes) and always occupy at least one extent. An image mode MDisk cannot be used as a quorum disk and no IBM Spectrum Virtualize system metadata extents are allocated from it; however, all of the IBM Spectrum Virtualize copy services functions can be applied to image mode disks. The difference between a managed mode volume (with striped extent allocation) and an image mode volume is shown in Figure 6-3.

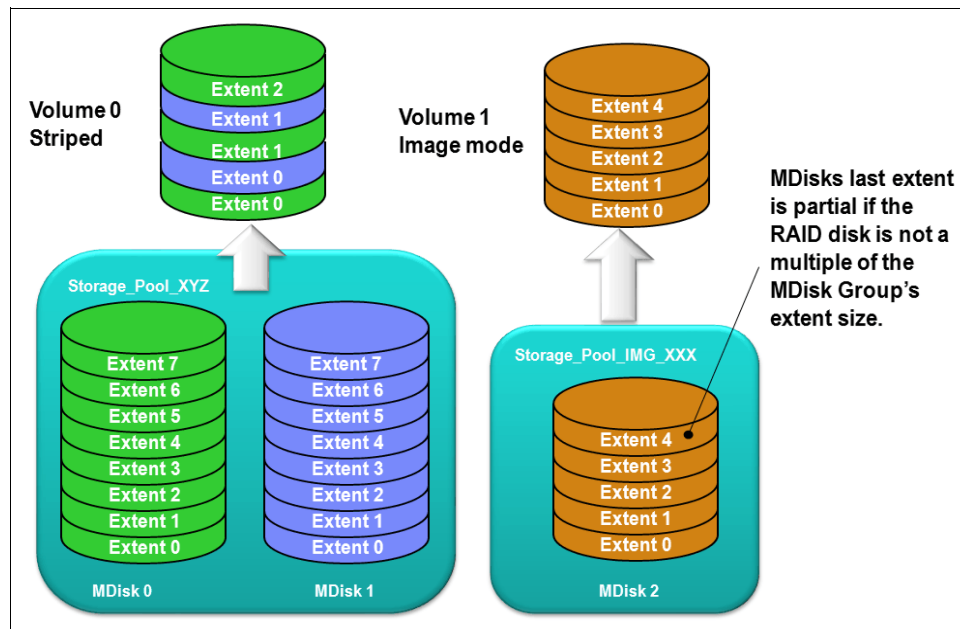


Figure 6-3 Image mode volume versus striped volume

An image mode volume is mapped to only one image mode MDisk and is mapped to the entirety of this MDisk. Therefore, the image mode volume capacity is equal to the size of the corresponding image mode MDisk. If the size of the (image mode) MDisk is not a multiple of the MDisk group (MDG) extent size, the last extent is marked as partial (not filled).

When you create an image mode volume, you map it to an MDisk that must be in unmanaged mode and must not be a member of a storage pool. As the image mode volume is configured, the MDisk is made a member of the specified storage pool. It is a good practice to use a dedicated pool for image mode MDisks, with a name indicating its role, such as Storage Pool\_IMG\_XXX.

IBM Spectrum Virtualize also supports the reverse process, in which a managed mode volume can be migrated to an image mode volume. The extent size that is chosen for this specific storage pool must be the same as the extent size of the storage pool into which you plan to migrate the data off the image mode volumes. If a volume is migrated to another MDisk, it is represented as being in managed mode during the migration. Its mode changes to “image” only after the process completes.

Volumes have the following characteristics or attributes:

- ▶ Volumes can be configured as managed or in image mode.
- ▶ Volumes can have extents that are allocated in striped or sequential mode.
- ▶ Volumes can be created as standard-provisioned or thin-provisioned. A conversion from a standard-provisioned to a thin-provisioned volume and vice versa can be done at run time.
- ▶ Volumes can be configured to have single data copy, or be mirrored to make them resistant to disk subsystem failures or to improve the read performance.
- ▶ Volumes can be compressed.
- ▶ Volumes can be configured as VVols, which allows VMware vCenter to manage storage objects. The storage system administrator can create these objects and assign ownership to VMware administrators to allow them to manage these objects.

Volumes have two major modes: Managed mode and image mode. For managed mode volumes, the sequential policy and the striped policy define how the extents of the volume are allocated from the storage pool.

- ▶ Volumes can be replicated synchronously or asynchronously. An IBM Spectrum Virtualize system can maintain active volume mirrors with a maximum of three other IBM Spectrum Virtualize systems, but not from the same source volume.

These volume characteristics are described next.

### 6.2.3 Size

Each volume has the following associated values that describe its size:

- ▶ The real (physical) capacity is the size of storage space that is allocated to the volume from the storage pool. It determines how many MDisk extents are allocated to form the volume.
- ▶ The virtual capacity is reported to the host and to IBM Spectrum Virtualize components (for example, FlashCopy, cache, and remote copy).

In a *standard-provisioned* volume, these two values are the same. In a *thin-provisioned* volume, the real capacity can be as little as a few percent of virtual capacity. The *real capacity* is used to store the user data and in the case of thin-provisioned volumes, metadata of the volume. The real capacity can be specified as an absolute value, or as a percentage of the virtual capacity. Volume size can be specified in units down to 512-byte blocks (see Figure 6-4 on page 261).

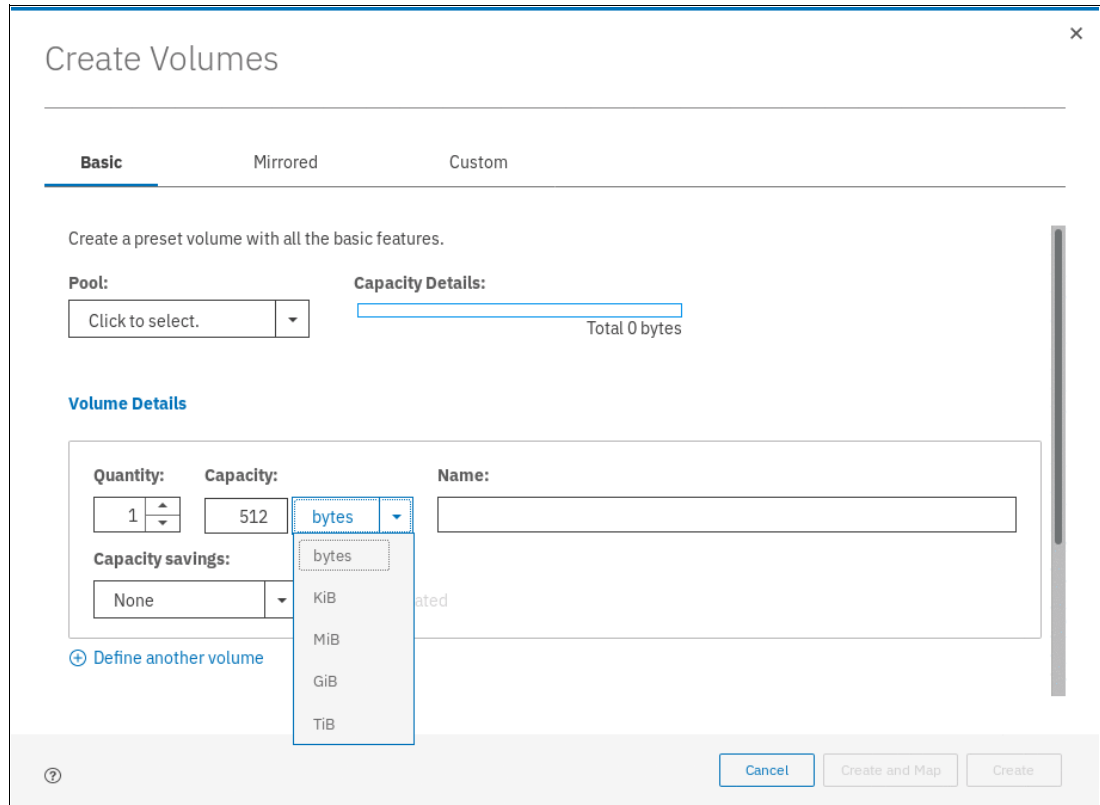


Figure 6-4 Smallest possible volume size

However, a volume is composed of storage pool extents. Therefore, it is not possible to allocate less than one extent to create a volume. Effectively, the internal unit of volume size is the extent size of the pools in which the volume is created. For example, a basic volume of a size 512 bytes created in a pool with the default extend size (1024 MiB) uses 1024 MiB of the pool space because an entire extent must be allocated to provide the space for the volume.

In practice, this rounding up of volume size to the whole number of extents has little effect on storage use efficiency, unless the storage system serves large number of small volumes. For more information about storage pools and extents, see Chapter 5, “Storage pools” on page 199.

## 6.2.4 Performance

The basic metrics of volume performance are the number of IOPS the volume can provide, typical time to service an IO request, and bandwidth of data served to a host.

Generally, volume performance is defined by the pool or pools that are used to create the volume. The pool determines the media bus (Non-Volatile Memory Express [NVMe] or Serial Advanced Technology Attachment [SATA]), media type (FlashCore Module [FCM] modules, SSDs, or hard disk drives [HDDs]), RAID level and number of drives per RAID array, and the possibility for IBM Easy Tier function to optimize performance of a volume. Also, storage efficiency, security, and allocation policy configuration settings of a module might change these characteristics.

Other aspect for performance of volumes that are served by IBM SAN Volume Controller is the back-end storage that is used to serve the volume. Back-end storage cache size controller performance and the load level on the storage also are factors that affect the performance of a volume.

## 6.2.5 Volume copies

A volume can have one or two physical copies. Although this is not required, the two copies of a mirrored volume typically are allocated from different storage pools that are backed by different physical hardware to increase volume resiliency. Each copy of the volume has the same virtual capacity, but the two copies can have different characteristics, including different real capacity.

However, each volume copy is not a separate object and can be manipulated only in the context of the volume. A mirrored volume behaves in the same way as any other volume. For example, all of its copies are expanded or shrunk when the volume is resized, it can participate in FlashCopy and RC relationships, is serviced by an I/O group, and includes a preferred node.

Volume copies are identified in the GUI by a copy ID, which can have value 0 or 1. Copies of the volume can be split, which provides a point-in-time (PiT) copy of a volume. An overview of volume mirroring is shown in Figure 6-5.

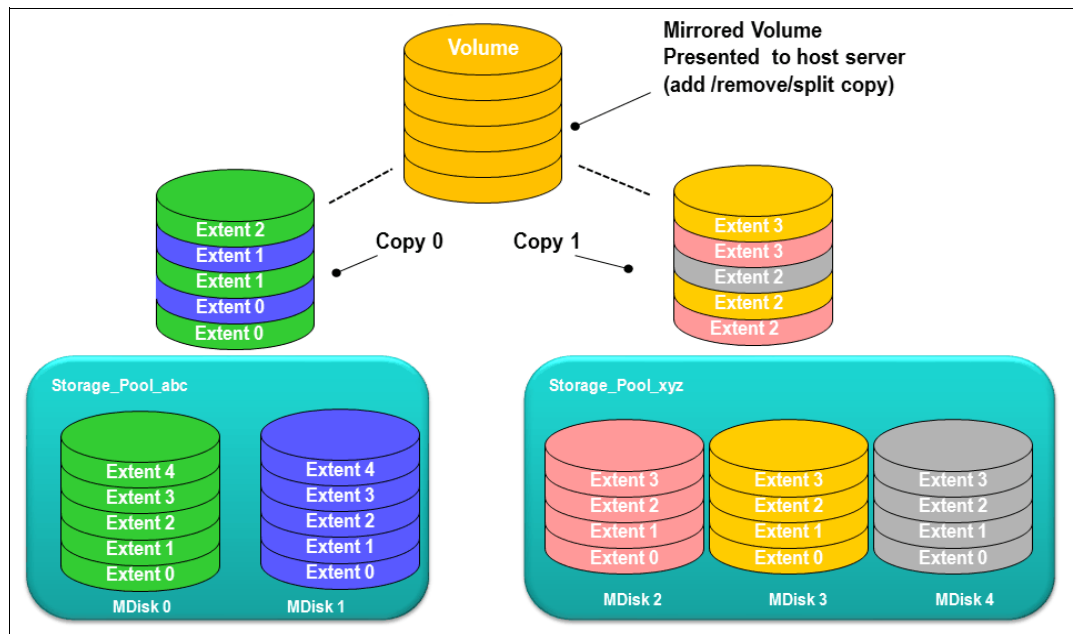


Figure 6-5 Volume mirroring overview

A copy can be added to a volume with a single copy or removed from a volume with two copies. Internal safety mechanisms prevent accidental removal of the only remaining copy of a volume.

A newly created, unformatted volume with two copies initially has the two copies in an out-of-synchronization state. The primary copy is defined as “fresh” and the secondary copy is defined as “stale”, and the volume is immediately available for use.



The synchronization process updates the secondary copy until it is fully synchronized; that is, data that is stored on the secondary copy matches data on the primary copy. This update is done at the *synchronization rate* that is defined when the volume is created, and can be modified after volume creation. The synchronization status for mirrored volumes is recorded on the storage system quorum disk.

If a mirrored volume is created by using the **format** parameter, both copies are formatted in parallel. The volume comes online when both operations are complete with the copies in sync.

If it is known that MDisk space (which is used for creating volume copies) is formatted or if the user does not require read stability, a `no_synchronization` option can be selected that declares the copies as synchronized even when they are not.

Creating more volume copy is beneficial in multiple scenarios, as shown in the following examples:

- ▶ Improving volume resilience by protecting them from a single storage system failure. This requires that each volume copy be configured on a different storage system.
- ▶ Providing concurrent maintenance of a storage system that does not natively support concurrent maintenance (for volumes on external virtualized storage).
- ▶ Providing an alternative method of data migration with improved availability characteristics. While a volume is migrated by using the data migration feature, it is vulnerable to failures on the source and target storage pool. Volume mirroring provides an alternative migration method that is not affected by the destination volume pool availability. For more information about this volume migration method, see “Volume migration by adding a volume copy” on page 319.

**Note:** For migrating volumes to a Data Reduction Pool (DRP), volume mirroring is the only migration method because DRPs do not support `migrate` commands.

- ▶ Converting between standard-provisioned volumes and thin-provisioned volumes (in either direction).

If one of the mirrored volume copies is temporarily unavailable (for example, because the storage system that provides the pool is unavailable), the volume remains accessible to servers. The storage system remembers which areas of the volume were modified after loss of access to a volume copy and resynchronizes only these areas when both copies are available.

**Note:** Volume mirroring is not a disaster recovery (DR) solution because both copies are accessed by the same node pair and addressable by only a single cluster. However, if correctly planned, it can improve availability.

The storage system tracks the synchronization status of volume copies by dividing the volume into 256 kibibyte (KiB) grains and maintaining a bitmap of stale grains (on the quorum disk), mapping 1 bit to one grain of the volume space. If the mirrored volume needs resynchronization, the system copies to the out-of-sync volume copy only these grains that were written to (changed) since the synchronization was lost. This approach is known as an *incremental synchronization* and it minimizes the time that is required to synchronize the volume copies.

**Important:** Mirrored volumes can be taken offline if no quorum disk is available. This behavior occurs because the synchronization status of mirrored volumes is recorded on the quorum disk.

A volume with more than one copy can be checked to see whether all of the copies are identical or consistent. If a medium error is encountered while it is reading from one copy, it is repaired by using data from the other copy. This consistency check is performed asynchronously with host I/O.

As mirrored volumes use bitmap space at a rate of 1 bit per 256 KiB grain, 1 MiB of bitmap space supports up to 2 TiB of mirrored volumes. The default size of bitmap space is 20 MiB, which allows configuration of up to 40 TiB of mirrored volumes. If all 512 MiB of variable bitmap space is allocated to mirrored volumes, 1 PiB of mirrored volumes can be supported.

Table 6-1 lists the bitmap space configuration options.

Table 6-1 *Bitmap space default configuration*

Copy service	Minimum allocated bitmap space	Default allocated bitmap space	Maximum allocated bitmap space	Minimum <sup>a</sup> capacity when using the default values
RC <sup>b</sup>	0	20 MiB	512 MiB	40 TiB of remote mirroring volume capacity
FlashCopy <sup>c</sup>	0	20 MiB	2 GiB	<ul style="list-style-type: none"> <li>▶ 10 TiB of FlashCopy source volume capacity</li> <li>▶ 5 TiB of incremental FlashCopy source volume capacity</li> </ul>
Volume mirroring	0	20 MiB	512 MiB	40 TiB of mirrored volumes
RAID	0	40 MiB	512 MiB	<ul style="list-style-type: none"> <li>▶ 80 TiB array capacity using RAID 0, 1, or 10</li> <li>▶ 80 TiB array capacity in three-disk RAID 5 array</li> <li>▶ Slightly less than 120 TiB array capacity in five-disk RAID 6 array</li> </ul>

a. The amount of available capacity might increase based on settings, such as grain size and strip size. RAID is subject to a 15% margin of error.

b. RC includes Metro Mirror (MM), Global Mirror (GM), and active-active relationships.

c. FlashCopy includes the FlashCopy function, Global Mirror with Change Volumes (GMCV), and active-active relationships.

The sum of all bitmap memory allocation for all functions except FlashCopy must not exceed 552 MiB.

## 6.2.6 I/O operations data flow

Although a mirrored volume looks to its users the same as a volume with a single copy, some differences exist in how I/O operations are performed internally for volumes with single or two copies.

## Read I/O operations data flow

If the volume is mirrored (that is, two copies of the volume exist), one copy is known as the *primary copy*. If the primary is available and synchronized, reads from the volume are directed to that copy. The user selects which copy is primary when the volume is created, but can change this setting at any time. In the management GUI, an asterisk indicates the primary copy of the mirrored volume. Placing the primary copy on a high-performance controller maximizes the read performance of the volume.

For non-mirrored volumes, only one volume copy exists; therefore, no choice exists for read source, and all reads are directed to the single volume copy.

## Write I/O operations data flow

For write I/O operations to a mirrored volume, the host sends the I/O request to the preferred node, which is responsible for destaging the data from cache. Figure 6-6 shows the data flow for this scenario.

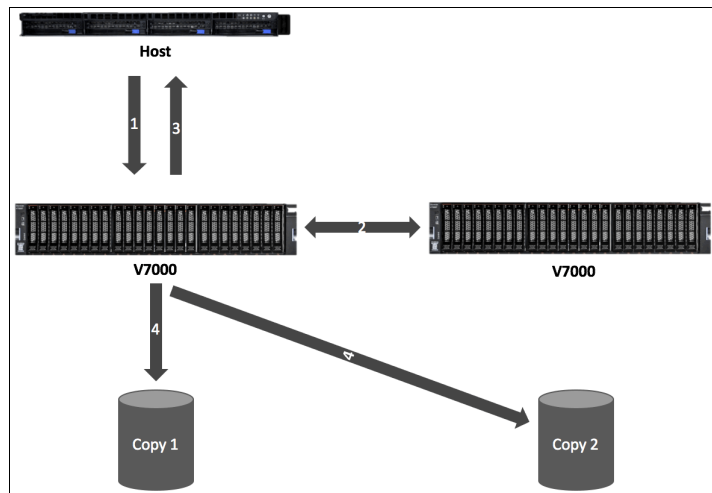


Figure 6-6 Data flow for write I/O processing in a mirrored volume

As shown in Figure 6-7 on page 266, the writes are sent by the host to the preferred node for the volume (1). Then, the data is mirrored to the cache of the partner node in the I/O group (2), and acknowledgment of the write operation is sent to the host (3). The preferred node then destages the written data to all volume copies (4). The example that is shown in Figure 6-7 on page 266 shows a case with a destage to a mirrored volume; that is, one with two physical data copies.

With V7.3, the cache architecture changed from an upper-cache design to a two-layer cache design. With this change, the data is only written once, and is then directly destaged from the controller to the disk system.

Figure 6-7 shows the data flow in a stretched environment.

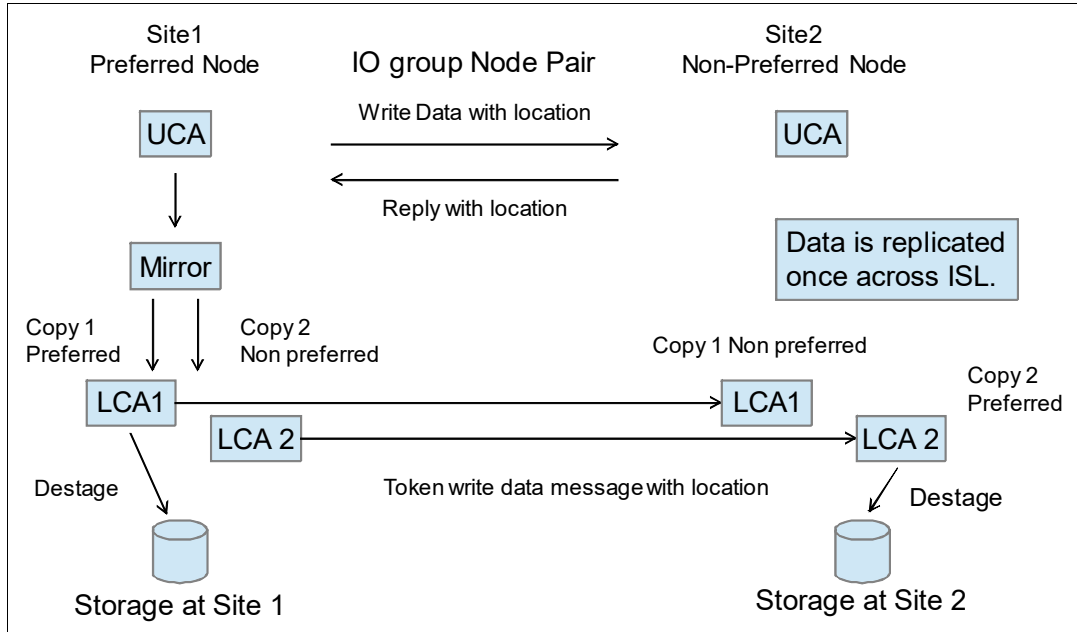


Figure 6-7 Design of an Enhanced Stretched Cluster

## 6.2.7 Storage efficiency

IBM Spectrum Virtualize storage systems allow you to configure DRPs that provide several technologies that increase efficiency of physical storage use:

- ▶ Thin provisioning
- ▶ Deduplication (block-level and pattern-matching)
- ▶ Compression
- ▶ SCSI UNMAP support

**Note:** Storage efficiency options can require more license and hardware components, depending on the model and configuration of your storage system.

Implementation of DRPs requires careful planning and sizing. For more information, see Chapter 2, “Planning” on page 53, and Chapter 9, “Advanced features for storage efficiency” on page 449.

DRPs use multithreading and hardware acceleration (where available) to provide storage efficiency functions on IBM Spectrum Virtualize storage systems. The use of storage efficiency options increases the number of IO operations the storage system must realize compared to access to a basic volume. This is because space-efficient volumes require the storage system to both to write the data that is sent by the host and the metadata that is required to maintain a space-efficient volume.

**Note:** FCM modules include compression hardware; therefore, they provide data set size reduction with no performance penalty.

For more information about storage efficiency functions of IBM Spectrum Virtualize, see Chapter 5, “Storage pools” on page 199 *Introduction and Implementation of Data Reduction Pools and Deduplication*, SG24-8430.

It is possible to benefit from compression and data-at-rest encryption because encryption is done after compression. However, the size of data that is encrypted at the host level is unlikely to be reduced by compression or deduplication at the storage system.

## Standard and thin-provisioned volumes

A standard-provisioned volume directly maps logical blocks on the virtual volume to physical blocks on storage media. Therefore, its virtual and physical capacity are identical.

A thin-provisioned volume has virtual capacity larger than physical capacity. Thin-provisioning is the base technology for all space-efficient volumes. When a thin-provisioned volume is created, a small amount of the real capacity is used for initial metadata. This metadata holds a mapping of an LBA in the volume to a *grain* on a physically allocated extent.

When a write request comes from a host, the block address for which the write is requested is checked against the mapping table. If a previous write to a block on the same grain exists as the incoming request, then physical storage was allocated for this LBA and can be used to service the request; otherwise, new physical grain is allocated to store the data, and the mapping table is updated to record that allocation.

**Note:** If you use of thin-provisioned volumes, it is *strongly* recommended to closely monitor available space in the pool that contains these volumes. If a thin-provisioned volume does not have enough real capacity for a write operation, the volume is taken offline and an error is logged. There is limited ability to recover with unmap. Also, consider creating a fully allocated sacrificial emergency space volume.

The grain size is defined when the volume is created and cannot be changed later. The grain size can be 32 KiB, 64 KiB, 128 KiB, or 256 KiB. The default grain size is 256 KiB, which is the preferred option; however, the following factors must be considered when deciding on the grain size:

- ▶ Smaller grain size helps to save space. If a 16 KiB write I/O requires a new physical grain to be allocated, the used space is 50% of a 32 KiB grain, but just over 6% of 256 KiB grain. If no subsequent writes to other blocks of the grain occur, the volume provisioning is less efficient for volumes with larger grain.
- ▶ Smaller grain size requires more metadata I/O to be performed, which increases the load on the physical back-end storage systems.
- ▶ When a thin-provisioned volume is a FlashCopy source or target volume, specify the same grain size for FlashCopy and the thin-provisioned volume configuration. Use 256 KiB grain to maximize performance.
- ▶ Grain size affects maximum size of the thin-provisioned volume. For 32 KiB size, the volume size cannot exceed 260 TiB.

Figure 6-8 on page 268 shows the thin-provisioning concept.

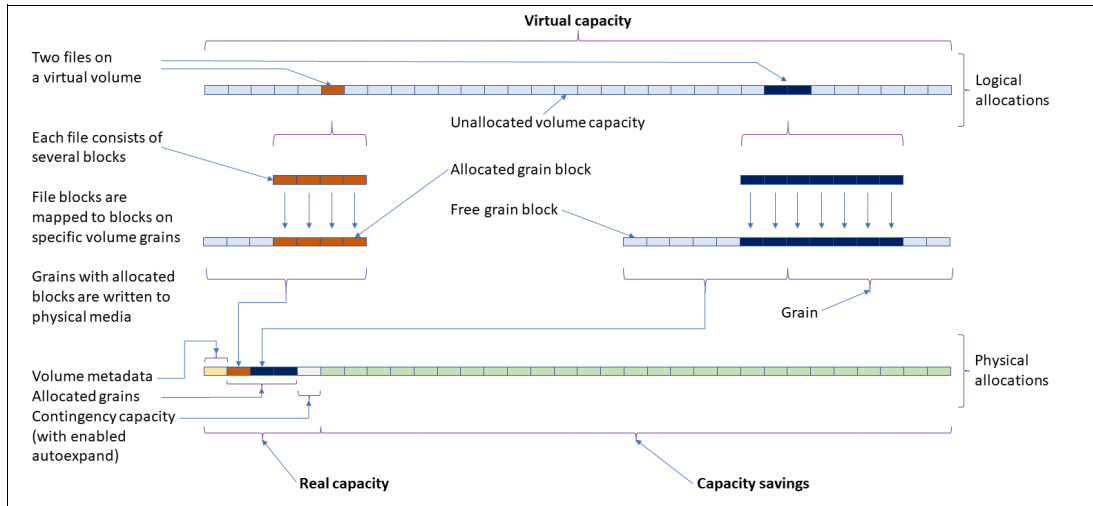


Figure 6-8 Conceptual diagram of thin-provisioned volume

Thin-provisioned volumes use metadata to enable capacity savings, and each grain of user data requires metadata to be stored. Therefore, the I/O rates that are obtained from thin-provisioned volumes are lower than the I/O rates that are obtained from standard-provisioned volumes.

The metadata storage that is used is never greater than 0.1% of the user data. The resource usage is independent of the virtual capacity of the volume.

**Thin-provisioned volume format:** Thin-provisioned volumes do not need formatting. A read I/O, which requests data from not allocated data space, returns zeros. When a write I/O causes space to be allocated, the grain is “zeroed” before use. Also, when a full-grain write consists of “all zeros”, no space is physically allocated on disk.

The real capacity of a thin-provisioned volume can be changed if the volume is not in image mode. Thin-provisioned volumes use the grains of real capacity that is provided in ascending order as new data is written to the volume. If the user initially assigns too much real capacity to the volume, the real capacity can be reduced to free storage for other uses.

A thin-provisioned volume can be configured to *autoexpand*. This feature causes the IBM Spectrum Virtualize to automatically add a fixed amount of extra real capacity to the thin-provisioned volume as required. Autoexpand does not cause the real capacity to grow much beyond the virtual capacity. Instead, it attempts to maintain a fixed amount of unused real capacity for the volume, which is known as the *contingency capacity*.

The contingency capacity is initially set to the real capacity that is assigned when the volume is created. If the user modifies the real capacity, the contingency capacity is reset to be the difference between the used capacity and real capacity.

A volume that is created without the autoexpand feature, and therefore has a zero contingency capacity, goes offline when the real capacity is used and it must expand.

To facilitate management of the auto expansion of thin-provisioned volumes, a capacity warning should be set for the storage pools from which they are allocated. When the used capacity of the pool exceeds the warning capacity, a warning event is logged. For example, if a warning of 80% is specified, an event is logged when 20% of the pool capacity remains free.

A thin-provisioned volume can be converted nondisruptively to a standard-provisioned volume (or vice versa) by using the volume mirroring function. You can create a thin-provisioned copy to a standard-provisioned primary volume and then remove the standard-provisioned copy from the volume after they are synchronized.

The standard-provisioned-to-thin-provisioned migration procedure uses a zero-detection algorithm so that grains that contain all zeros do not cause any real capacity to be used.

Thin-provisioned volumes can be used as volumes that are assigned to the host by FlashCopy to implement thin-provisioned FlashCopy targets. When creating a mirrored volume, a thin-provisioned volume can be created as a second volume copy, whether the primary copy is a fully or thin-provisioned volume.

## Deduplicated volumes

Deduplication is a specialized data set reduction technique. However, in contrast to the standard file-compression tools that work on single files or sets of files, deduplication is a technique that is applied on a larger scale, such as a file system or volume. In IBM Spectrum Virtualize, deduplication can be enabled for thin provisioned and compressed volumes that are created in DRPs.

IBM Spectrum Virtualize uses two techniques to detect duplicate data:

- ▶ Pattern matching
- ▶ Data signature (hash)

Deduplication works by identifying repeating chunks in the data that is written to the storage system. Pattern matching looks for a known data pattern (for example “all ones”), while data signature-based algorithm calculates a signature for each data chunk (by using a hash function), and checks if the calculated signature is present in the deduplication database.

If a known pattern or a signature match is found, the data chunk is replaced by a reference to a stored chunk, which reduces storage space that is required for storing the data. Conversely, if no match is found, the data chunk is stored without modification and its signature is added to the deduplication database.

To maximize the space that is available for the deduplication database, the system distributes it between all nodes in the I/O groups that contain deduplicated volumes. Each node holds a distinct portion of the records that are stored in the database. If nodes are removed or added to the system, the database is redistributed between the anodes to ensure optimal use of available resources.

Depending on the data type that is stored on the volume, the capacity savings can be significant. Examples of use cases that typically benefit from deduplication are virtual environments with multiple VMs running the same operating system and backup servers. In both cases, it is expected that multiple copies of identical files exist, such as components of the standard operating system or applications that are used in the organization. However, data that is encrypted at the file system level might not benefit from deduplication because encryption of the same file with different keys provides different output, which makes deduplication impossible.

Although deduplication (as with other features of IBM Spectrum Virtualize) is not apparent to users and applications, it must be planned for and understood before implementation because it might reduce the redundancy of a solution. For example, if an application stores two copies of a file to reduce chances of data corruption by way of a random event, the copies are deduplicated and thus the intended redundancy is removed from the system if these copies are on the same volume.

When planning the use of deduplicated volumes, be aware of update and performance considerations and the following software and hardware requirements:

- ▶ Code level V8.1.2 or higher is needed for DRPs.
- ▶ Code level V8.1.3 or higher is needed for deduplication.
- ▶ Nodes must have at least 32 GB memory to support deduplication. Nodes that have more than 64 GB memory can use a bigger deduplication fingerprint database, which might lead to better deduplication.
- ▶ Avoid the use of GMCV to or from a deduplicated volume.
- ▶ You must run supported hardware. For more information about the valid hardware and features combinations, see [IBM Knowledge Center](#) and select **Planning** → **Storage configuration planning**.

## Compressed volumes

Volumes that are created in a DRP can be compressed. This causes data that is written to the volume to be compressed before committing them to back-end storage, which reduces the physical capacity that is required to store the data. Because enabling compression does not incur more metadata handling penalties, in most cases it is recommended to enable compression on thin-provisioned volumes.

**Note:** Consider the following points:

- ▶ You can use the management GUI or the command-line interface to run the built-in compression estimation tool. This tool can be used to determine the capacity savings that are possible for data on the system by using compression.
- ▶ Data compression for volumes that are backed by flash-based storage results in the extra benefit of a reduction of write amplification, which has a beneficial effect on media longevity.

## Capacity reclamation

File deletion in modern file systems is realized by updating file system metadata and marking the physical storage space that is used by the removed file as unused. The data of the removed file is not overwritten. This improves file system performance by reducing the number of I/O operations on physical storage that is required to perform file deletion.

However, this approach affects the management of real capacity of volumes with enabled capacity savings. File system deletion frees space at the file system level, but physical data blocks that are allocated by the storage for the file still take up the real capacity of a volume.

To address this issue, file systems added support for SCSI **UNMAP** command, which can be issued after file deletion, and informs the storage system that physical blocks that are used by the removed file should be marked as no longer in use and can be freed up. Modern operating systems issue SCSI **UNMAP** commands only to storage that advertises support for this feature.

V8.1.0 and later releases support the SCSI **UNMAP** command on Spectrum Virtualize systems. This enables hosts to notify the storage controller of capacity that is no longer required and can be reused or de-allocated, which might improve capacity savings.

**Note:** For volumes that are outside DRPs, the complete stack from the operating system down to back-end storage controller must support **unmap** to enable the capacity reclamation. SCSI **unmap** is passed only to specific back-end storage controllers.



Consider the following points:

- ▶ V8.1.2 can also reclaim capacity in DRPs when a host issues SCSI **unmap** commands.
- ▶ By default, V8.2.1 does not advertise to hosts support for SCSI **unmap**.
- ▶ In V8.3.1 support for the host SCSI **unmap** command is enabled by default.

Before enabling SCSI unmap, see this IBM Support [web page](#).

Analyze your storage stack to optimally balance advantages and costs of data reclamation.

### Data reduction at two levels

It is possible to design and implement a solution where data reduction technologies are applied at both the SAN Volume Controller and the back-end storage. Such solutions allow small extra savings from compressing metadata

However, such cases require meticulous planning. Incorrect implementation may lead for example to out-of-space event being triggered by DRP garbage collection process or Easy Tier hot data migrations.

**Note:** You can find information required to correctly plan advanced data reduction architectures in:

- ▶ Planning for deduplicated volumes chapter at **Planning** → **Storage configuration planning** of the SAN Volume Controller Knowledge Center.
- ▶ *Introduction and Implementation of Data Reduction Pools and Deduplication*, SG24-8430
- ▶ [IBM Spectrum Virtualize Data Reduction Best Practices whitepaper](#)

## 6.2.8 Encryption

IBM Spectrum Virtualize systems can be configured to enable data-at-rest encryption. This functionality is realized in hardware (self-encrypting drives or in serial-attached SCSI (SAS) controller for drives that do not support self-encryption and are connected by way of the SAS bus) or in software (external virtualized storage).

See Chapter 12, “Encryption” on page 685 for more information about creating and managing encrypted volumes.

## 6.2.9 Cache mode

Another volume parameter is its cache characteristics. Under normal conditions, a volume’s read and write data is held in the cache of its preferred node, with a mirrored copy of write data held in the partner node of the same I/O group. However, it is possible to create a volume with different cache characteristics, if this is required.

Cache setting of a volume can have the following values:

- ▶ *readwrite*: All read and write I/O operations that are performed by the volume are stored in cache. This is the default cache mode for all volumes.
- ▶ *readonly*: Only read I/O operations that are performed on the volume are stored in cache. Writes to the volume are not cached.
- ▶ *disabled*: No I/O operations on the volume are stored in cache. I/Os are passed directly to the back-end storage controller rather than being held in the node’s cache.

Having cache-disabled volumes makes it possible to use the native copy services in the underlying RAID array controller for MDisks (logical unit numbers (LUNs)) that are used as IBM Spectrum Virtualize image mode volumes. However, use of IBM Spectrum Virtualize Copy Services rather than the underlying disk controller copy services gives better results.

**Note:** Disabling volume cache is a prerequisite for the use of native copy services on image mode volumes that are defined on storage systems that are virtualized by IBM Spectrum Virtualize. Contact IBM Support before turning off the cache for volumes in your production environment to avoid performance degradation.

## 6.2.10 I/O throttling

You can set a limit on the number of I/O operations that are realized by a volume. This limitation is called *I/O throttling* or *governing*.

The limit can be set in terms of number of IOPS or bandwidth (MBps, GBps, or TBps). By default, I/O throttling is disabled, but each volume can have up to two throttles defined: one for bandwidth and one for IOPS.

When deciding between the use of IOPS or bandwidth as the I/O governing throttle, consider the disk access profile of the application that is the primary volume user. Database applications generally issue large amounts of I/O operations, but transfer a relatively small amount of data. In this case, setting an I/O governing throttle that is based on bandwidth might not achieve much. A throttle based on IOPS would be better suited to this use case.

Conversely, a video streaming application generally issues a small amount of I/O, but transfers large amounts of data. Therefore, it is better to use a bandwidth throttle for the volume in this case.

An I/O governing rate of 0 does not mean that zero IOPS or bandwidth can be achieved for this volume; rather, it means that no throttle is set for this volume.

**Note:** Consider the following points:

- ▶ I/O governing does not affect FlashCopy and data migration I/O rates.
- ▶ I/O governing on MM or GM secondary volumes does not affect the rate of data copy from the primary volume.

For more information about how to configure I/O throttle on a volume, see 6.5.4, “I/O throttling” on page 295.

## 6.2.11 Volume protection

Volume protection prevents volumes or host mappings from being deleted if the system detects recent I/O activity. This global setting is enabled by default on new systems. You can set this value to apply to all volumes that are configured on your system, or control whether the system-level volume protection is enabled or disabled on specific pools.

There are two levels at which the volume protection needs to be enabled to be effective: system level and pool level. Both must be enabled for protection to be active on a pool. The pool-level protection depends on the system-level setting to ensure that protection is applied consistently for volumes within that pool. If system-level protection is enabled, but pool-level protection is not enabled, any volumes in the pool can be deleted.

When you enable volume protection at the system level, you specify a time period in minutes that the volume must be idle before it can be deleted. If volume protection is enabled and the time period is not expired, the volume deletion fails, even if the **-force** parameter is used. To prevent volumes that are configured in a pool from inadvertent deletion, enable volume protection at the pool level.

The following CLI commands and the corresponding GUI activities are affected by volume protection setting:

- ▶ `rmvdisk`
- ▶ `rmvdiskcopy`
- ▶ `rmvvolume`
- ▶ `rmvdiskhostmap`
- ▶ `rmvolumehostclustermap`
- ▶ `rmmdiskgrp`
- ▶ `rmhostiogr`
- ▶ `rmhost`
- ▶ `rmhostcluster`
- ▶ `rmhostport`
- ▶ `mkrcrelationship`

Volume protection can be set both from GUI (new in V.8.3.1, see 6.5.5, “Volume protection” on page 300), and CLI (see 6.6.9, “Volume protection” on page 336).

## 6.2.12 Secure data deletion

The system provides methods to securely erase data from a drive or from a boot drive when a control enclosure is decommissioned.

Secure data deletion effectively erases or overwrites all traces of data from a data storage device. The original data on that device becomes inaccessible and cannot be reconstructed. You can securely delete data on individual drives and on a boot drive of a control enclosure. The methods and commands that are used to securely delete data enable the system to be used in compliance with European Regulation EU2019/424.

For more information about this procedure, see [IBM Knowledge Center](#).

## 6.3 Virtual volumes

IBM Spectrum Virtualize V7.6 introduced support for *virtual volumes*. These volumes enable support for VVOLs, which allow VMware vCenter to manage system objects, such as volumes and pools. The IBM Spectrum Virtualize system administrators can create volume objects of this class, and assign ownership to VMware administrators to simplify management.

For more information about configuring VVol with IBM Spectrum Virtualize, see *Configuring VMware Virtual Volumes for Systems Powered by IBM Spectrum Virtualize*, SG24-8328.

## 6.4 Volumes in multi-site topologies

IBM Spectrum Virtualize system can be set up in a multi-site configuration, which makes the system aware of which system components (I/O groups, nodes, and back-end storage) are at which site. For defining the storage topology, a site is defined as an independent failure domain, which means that if one site fails, the other site can continue to operate without disruption.

The sites can be in the same data center room or across rooms in the data center, in buildings on the same campus, or buildings in different cities, depending on the type and scale of a failure the solution must survive.

The following topologies are available:

- ▶ *Standard* topology, which is intended for single-site configurations that does not allow site definition and assumes all components of the solution to be at a single site. Global Mirror (GM) or Metro Mirror (MM) can be used to maintain a copy of a volume on a different system at a remote site, which can be used for DR.
- ▶ *HyperSwap* topology, which is a three-site highly available (HA) configuration in which each I/O group is at a different site. A volume can be active on two I/O groups so that if one site is not available, it can immediately be accessed by way of the other site.
- ▶ *Stretched* topology (Enhanced Stretched Cluster [ESC]). When set-up in the stretched topology, each node of an I/O group of the storage system is at a different site and volumes have a copy at each site. Access to a volume can continue when one site is not available but with reduced performance. This topology is also known as an *enhanced stretched system*.

The reason for the use of an enhanced stretched system rather than the HyperSwap topology can be the use of GM or MM to a third site, or because the system was configured as an enhanced stretched system before HyperSwap being released.

The stretched topology uses fewer system resources compared to HyperSwap, which allows a greater number of HA volumes to be configured. However, during a disaster that makes one site unavailable, the system cache on the nodes of the surviving site is disabled. The HyperSwap topology uses extra system resources to support a full independent cache on each site, which allows full performance (even if one site is lost). In some environments, a HyperSwap topology provides better performance than a stretched topology.

If the objective of your solution design is HA, it is better to use an IBM HyperSwap topology. However, if the objectives include DR, complex copy services, or highest scalability, review the Planning for high availability section of the SAN Volume Controller Knowledge Center (**IBM SAN Volume Controller → Planning → Planning for high availability**) before choosing the topology.

**Note:** Multi-site topologies of IBM Spectrum Virtualize use two sites as component locations (nodes and back-end storage), and a third site as a location for a tie-breaker component that is used to resolve split-brain scenarios where the storage system components lose communication with each other.

For more information about ESC and HyperSwap, see the white paper *IBM Spectrum Virtualize HyperSwap configuration*, [WP102538](#).

The Create Volumes menu provides the following options, depending on the configured system topology:

- ▶ With standard topology, the available options are: Basic, Mirrored, and Custom.
- ▶ With HyperSwap topology, the options are: Basic, HyperSwap, and Custom.
- ▶ With stretched topology, the options are: Basic, Stretched, and Custom.

The multi-site topologies provide HA volumes accessible through two sites at up to 300 km (186.4 miles) apart. A fully independent copy of the data is maintained at each site.

**Note:** The determining factor for HyperSwap configuration validity is the time that it takes to send the data between the sites. Therefore, the distance should be estimated and allow for the fact that the distance between the sites that is measured along the data path is longer than the geographic distance. Also, each device on the data path that adds latency increases the effective distance between the sites.

## 6.4.1 HyperSwap topology

In the HyperSwap topology both nodes of an I/O group are in the same site. Therefore, to get a volume resiliently stored on both sites, the storage system must be composed of at least two I/O groups.

When data is written by hosts at either site, both copies are synchronously updated before the write operation completion is reported to the host. The HyperSwap function automatically optimizes itself to minimize data transmitted between sites and to minimize host read and write latency.

If the nodes or storage at either site go offline, the HyperSwap function automatically fails over access to the other copy. The HyperSwap function also automatically resynchronizes the two copies when possible.

The HyperSwap function is built on the foundation of two earlier technologies: Non-Disruptive Volume Move (NDVM) function that was introduced in IBM Spectrum Virtualize V6.4, and the RC features that include MM, GM, and GMCV.

HyperSwap volume configuration is possible only after the IBM Spectrum Virtualize system is configured in HyperSwap topology. After this topology change, the GUI presents an option to create a HyperSwap volume and creates them by running the `mkvolume` command instead of the usual `mkvdisk` command. The GUI continues to run the `mkvdisk` command when all other classes of volumes are created.

**Note:** It is still possible to create HyperSwap volumes in V7.5, as described in [this IBM Support white paper](#).

For more information, see *IBM Storwize V7000, Spectrum Virtualize, HyperSwap, and VMware Implementation*, SG24-8317.

From the perspective of a host or a storage administrator, a HyperSwap volume is a single entity; however, it is realized by using four volumes and a set of FlashCopy maps and an RC relationship (see Figure 6-9 on page 276).

## What does a HyperSwap volume look like?

- Composite object, comprised of:
  - 4 VDisks
  - 1 “active-active” relationship
  - 4 FlashCopy maps (for Change Volumes)
  - Additional access I/O group
- Remote Copy relationship and FlashCopy maps are automatically controlled
- Can be thick/thin/compressed
- Can be encrypted

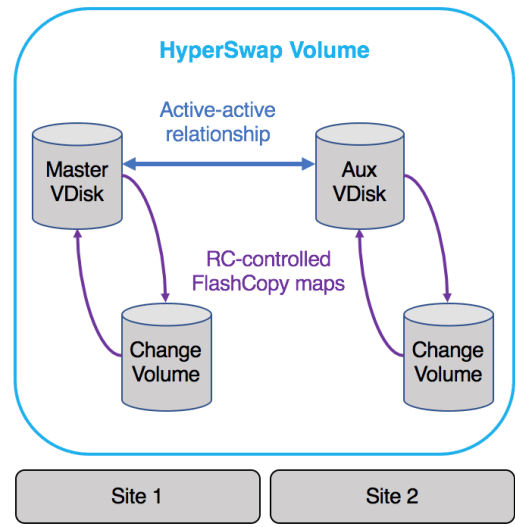


Figure 6-9 What makes up a HyperSwap Volume

The GUI simplifies the HyperSwap volume creation process by asking about required volume parameters only, and automatically configuring all the underlying volumes, FlashCopy maps, and volume replications relationships.

An example of a HyperSwap volume configuration is shown in Figure 6-29 on page 292.

## 6.4.2 Stretched topology

In the stretched topology, each I/O group in the system has one node on one site, and one node on the other site. The topology works with any number of I/O groups of 1 - 4.

The ESC function should be used if:

- ▶ You need synchronous or asynchronous replication by using MM or GM to a third site. It is possible to extend the configuration to four sites if the remote system is also configured with a stretched system topology.
- ▶ You need the absolute maximum number of HA volumes in the system.

For more information about creating volumes in stretched topology, see “Stretched volumes” on page 293

## 6.5 Operations on volumes

This section describes how to perform operations on volumes using GUI. The following operations can be performed on a volume:

- ▶ Volumes can be created and deleted.
- ▶ Volumes can have their characteristics modified, including:
  - Size (expanding or shrinking)
  - Number of copies (adding or removing a copy)
  - I/O throttling

- Protection
- ▶ Volumes can be migrated at run time to another MDisk or storage pool.
- ▶ PiT volume can be created by using FlashCopy. Multiple snapshots and quick restore from snapshots (reverse FlashCopy) are supported.
- ▶ Volumes can be mapped to (and unmapped from) hosts.

**Note:** With V7.4 and later, it is possible to prevent accidental deletion of volumes if they recently performed any I/O operations. This feature is called *Volume protection*, and it prevents active volumes, or host mappings, from being deleted inadvertently. This process is done by using a global system setting. For more information, see 6.6.9, “Volume protection” on page 336, and [IBM Knowledge Center](#).

## 6.5.1 Creating volumes

This section focuses on the use of the Create Volumes menu to create Basic and Mirrored volumes in a system with standard topology. Volume creation is available on the following volume classes:

- ▶ Basic
- ▶ Mirrored
- ▶ Custom

To create a volume, complete the following steps:

1. Click the **Volumes** menu and click the **Volumes** option of the IBM Spectrum Virtualize GUI, as shown in Figure 6-10.

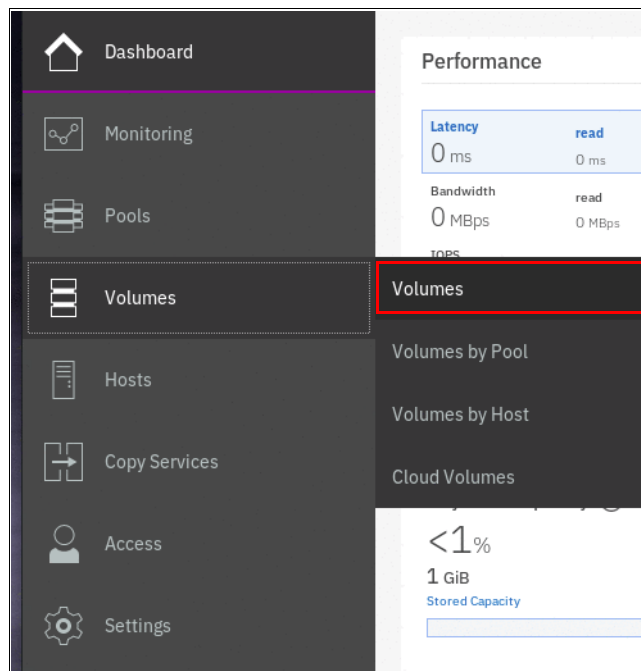


Figure 6-10 Volumes menu

A list of volumes, their state, capacity, and associated storage pools, is displayed.

2. To create a volume, click **Create Volumes**, as shown in Figure 6-11.

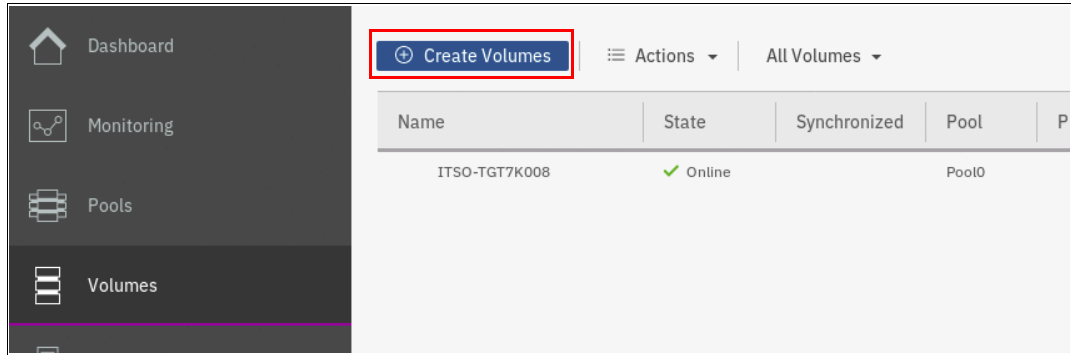


Figure 6-11 Create Volumes button

The Create Volumes tab opens the Create Volumes window, which displays the available creation methods.

**Note:** The volume classes that are displayed in the Create Volumes window depend on the topology of the system.

The Create Volumes window for standard topology is shown in Figure 6-13 on page 280.

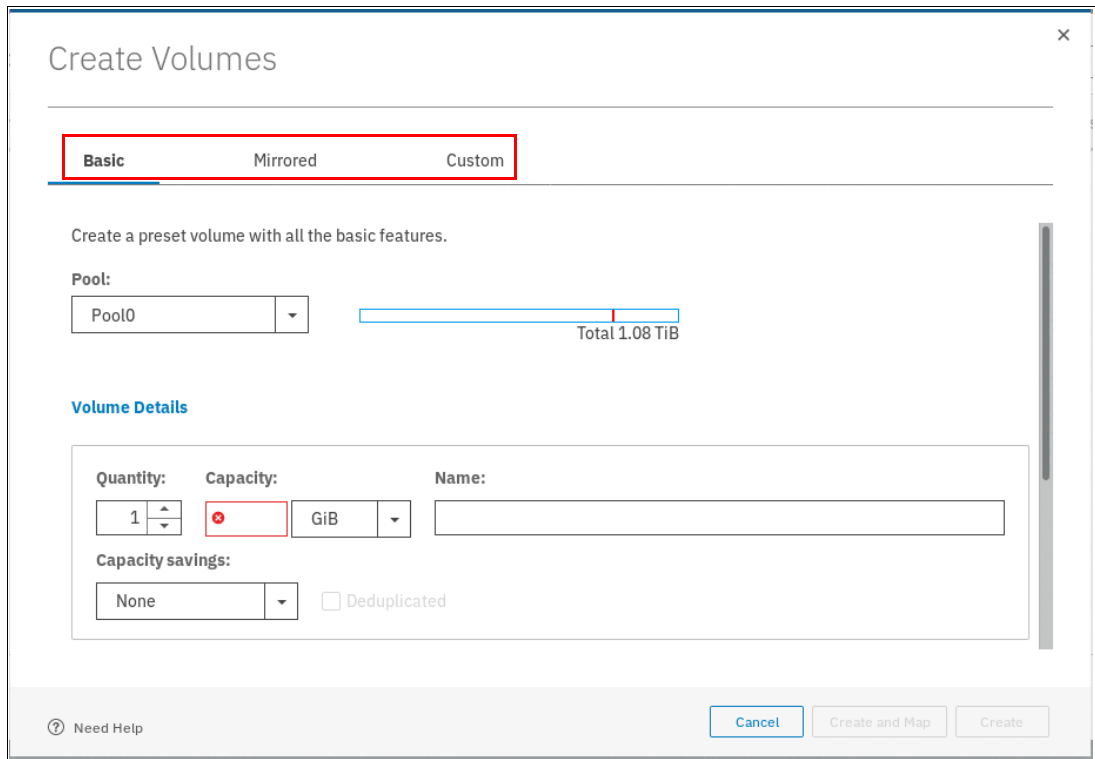


Figure 6-12 Basic, Mirrored, and Custom Volume Creation options



**Note:** Consider the following points:

- ▶ A *Basic volume* is a volume that has only one physical copy, uses storage that is allocated from a single pool on one site, and uses read/write cache mode.
- ▶ A *Mirrored volume* is a volume with two physical copies, where each volume copy can belong to a different storage pool.
- ▶ A *Custom volume* (in the context of this menu) is a Basic or Mirrored volume with values of some of its parameters that are changed from the defaults.
- ▶ The Create Volumes window also provides (by way of the Capacity Savings parameter) the ability to change the default provisioning of a Basic or Mirrored Volume to Thin-provisioned or Compressed.

For more information, see “Capacity savings option” on page 284.

## Creating basic volumes

A basic volume is a volume that has only one physical copy. Basic volumes are supported in any system topology and are common to all configurations. Basic volumes can be of any type of virtualization: striped, sequential, or image. They can also use any type of capacity savings: thin-provisioning, compressed, or none. Deduplication can be configured with thin-provisioned and compressed volumes in DRPs for added capacity savings.

To create a basic volume, click **Basic** as shown in Figure 6-13 on page 280. This action opens Basic volume menu in which you can define the following parameters:

- ▶ Pool: The Pool in which the volume is created (drop-down menu).
- ▶ Quantity: Number of volumes to be created (numeric up or down).
- ▶ Capacity: Size of the volume in specified units (drop-down menu).
- ▶ Capacity Savings (drop-down menu):
  - None
  - Thin-provisioned
  - Compressed
- ▶ Name: Name of the volume (cannot start with a number).
- ▶ I/O group.

The Basic Volume creation window is shown in Figure 6-13.

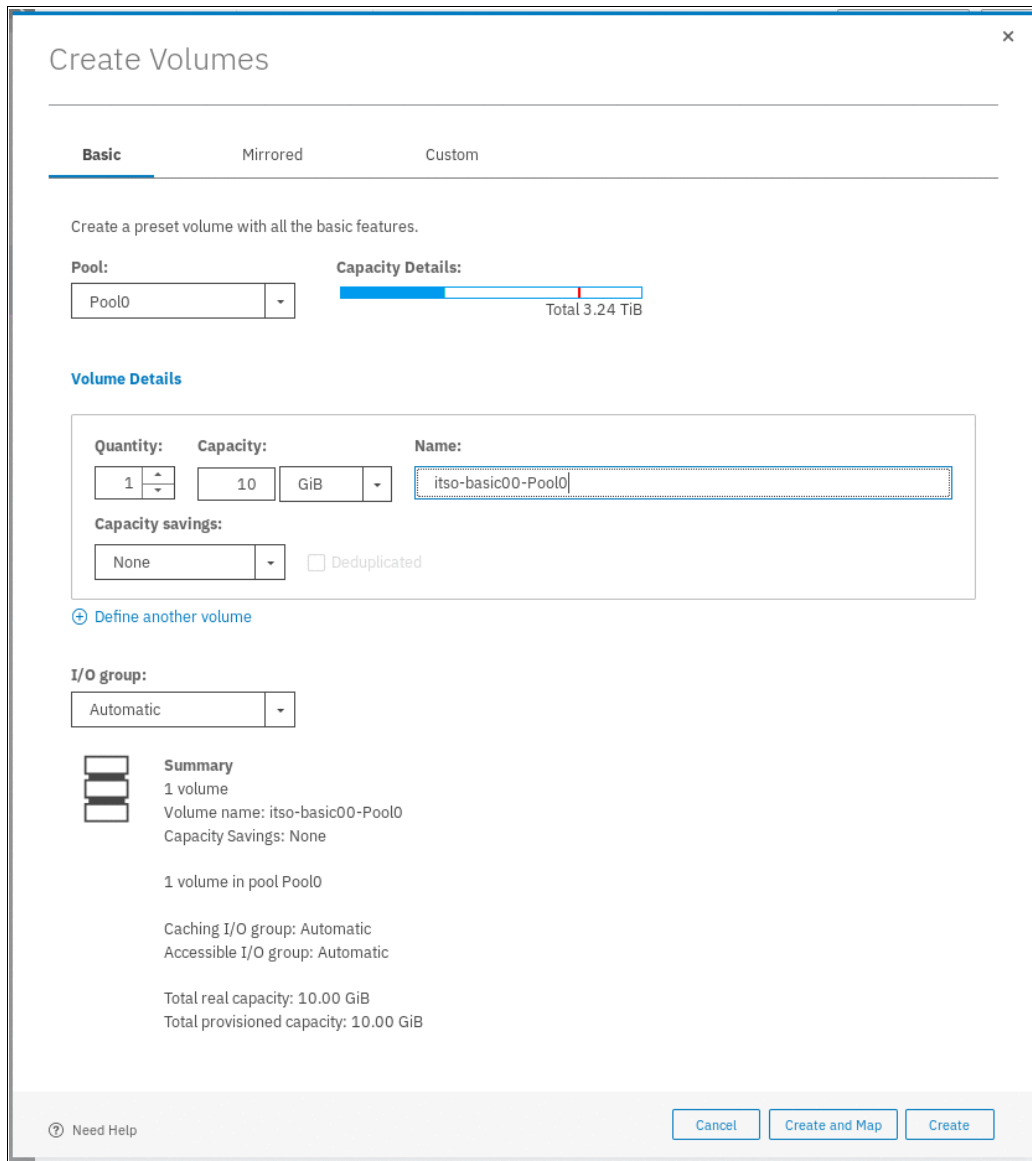


Figure 6-13 Create Volumes window

Define and consistently use a suitable volume naming convention to facilitate easy identification. For example, a volume name can contain the name of the pool or some tag that identifies the underlying storage subsystem, the host or cluster name that the volume is mapped to, and the content of this volume, such as the name of the applications that use the volume.

When all of the characteristics of the basic volume are defined, it can be created by selecting one of the following options:

- ▶ **Create**
- ▶ **Create and Map**

**Note:** The Plus sign (+) icon can be used to create more volumes in the same instance of the volume creation wizard.

In the example, the Create option was selected. The volume-to-host mapping can be performed later, as described in section 6.5.8, “Mapping a volume to a host” on page 312.

When the operation completes, the volume is seen in the Volumes pane in state `Online` (formatting), as shown in Figure 6-14.

Name	↓	State	Synchronized	Pool
itso-basic00-Pool0		✓ Online (formatting)		Pool0
Vdisk-compr-dedup3		✓ Online		Pool1
Vdisk-compr-dedup2		✓ Online		Pool1

Figure 6-14 Basic volume formatting

By default GUI does not display any information about the commands it runs to complete a task. However, while a command runs, you can click **View more details** to see the underlying CLI commands that are run to create the volume and a report of completion of the operation, as shown in Figure 6-15.

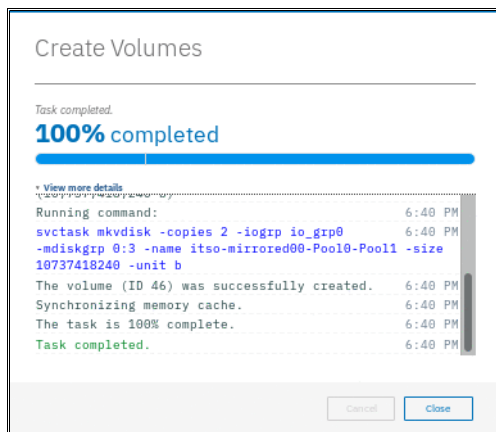


Figure 6-15 View more details of volume creation

**Note:** Consider the following points:

- ▶ standard-provisioned volumes are automatically formatted through the quick initialization process after the volume is created. This process makes standard-provisioned volumes available for use immediately.
- ▶ Quick initialization requires a small amount of I/O to complete, and limits the number of volumes that can be initialized at the same time. Some volume actions, such as moving, expanding, shrinking, or adding a volume copy, are disabled when the specified volume is initializing. Those actions become available after the initialization process completes.
- ▶ The quick initialization process can be disabled in circumstances where it is not necessary. For example, if the volume is the target of a Copy Services function, the Copy Services operation formats the volume. The quick initialization process can also be disabled for performance testing so that the measurements of the raw system capabilities can take place without waiting for the process to complete.

## Creating mirrored volumes

To create a mirrored volume, complete the following steps:

1. In the Create Volumes window, click **Mirrored** and choose the **Pool** for **Copy1** and **Copy2** by using the drop-down menus. Although the mirrored volume can be created in the same pool, this setup is not typical. Generally, keep volumes copies on separate set of physical disks (Pools).
2. Enter the following Volume Details:
  - Quantity
  - Capacity
  - Capacity savings
  - Name

Leave the I/O group option at its default setting of Automatic (see Figure 6-16 on page 283).

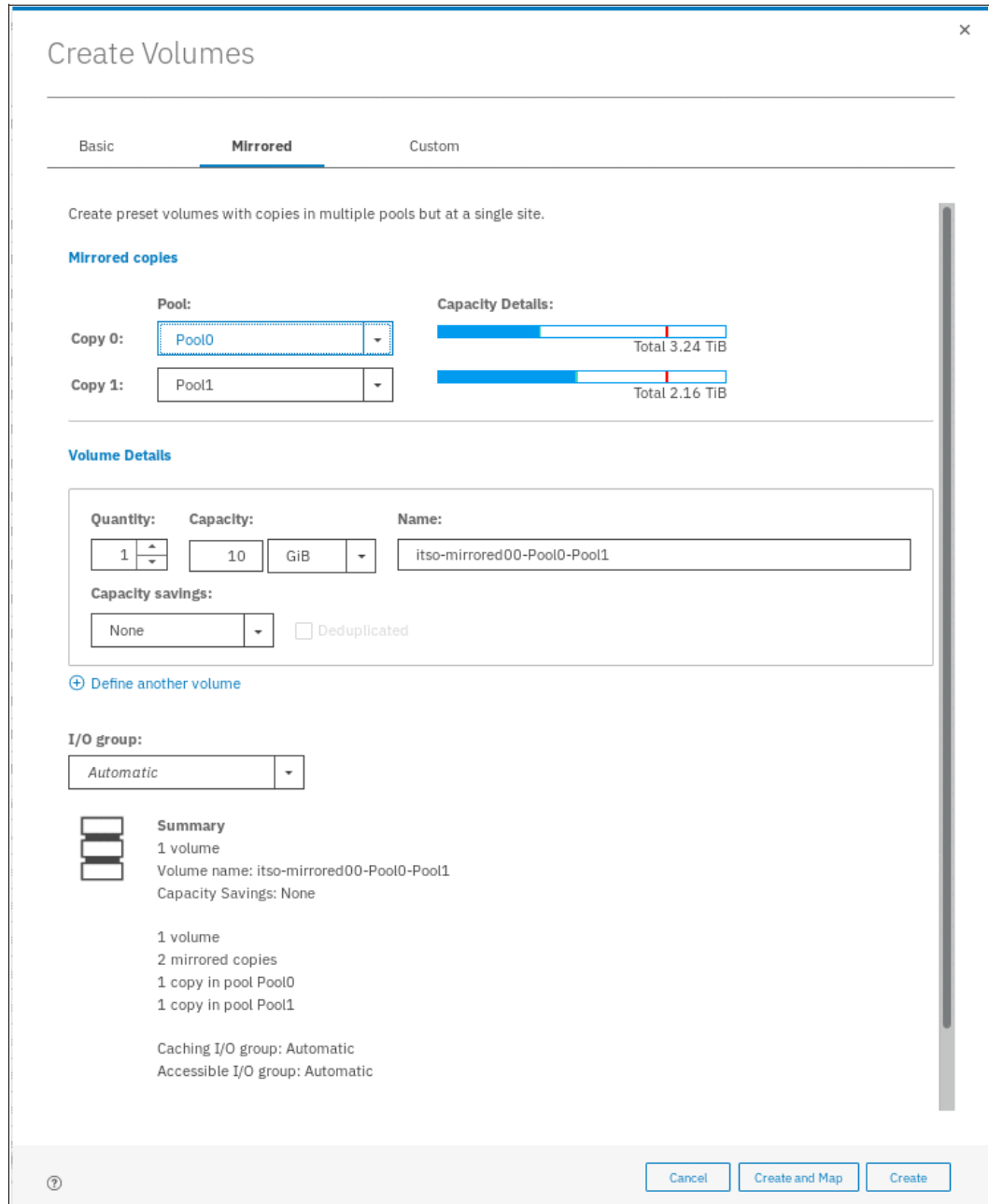


Figure 6-16 Mirrored Volume creation

3. Click **Create** (or **Create and Map**).

When the operation completes, the volume is seen in the Volumes pane in state `Online` (formatting), as shown in Figure 6-17 on page 284.

Name	State	Synchronized	Capacity	Pool
its0-basic00-Pool0	Online		10.00 GiB	Pool0
its0-mirrored00-Pool0-Pool1	Online (formatting)		10.00 GiB	Pool0
Copy 0*	Online (formatting)	Yes	10.00 GiB	Pool0
Copy 1	Online (formatting)	No	10.00 GiB	Pool1

Figure 6-17 Mirrored volume formatting

A mirrored volume is displayed in the GUI as configured in the pool in which it has its primary. In the example above volume `its0-mirrored00-Pool0-Pool1` is displayed as configured in Pool0 because it has its primary copy in Pool0.

**Note:** When creating a Mirrored volume by using this menu, you are not required to specify the Mirrored Sync rate (it defaults to 2 MBps). The synchronization rate can be customized by using the Custom menu.

### Capacity savings option

When the Basic or Mirrored method of volume creation is used, the GUI provides a Capacity Savings option, which enables altering the volume provisioning parameters without the use of the Custom volume provisioning method. You can select **Thin-provisioned** or **Compressed** from the drop-down menu, as shown in Figure 6-18, to create thin-provisioned or compressed volumes.

The screenshot shows the 'Volume Details' section of a GUI. It includes fields for 'Quantity' (set to 1), 'Capacity' (with a red 'x' icon and 'GiB' unit), and 'Name'. Below these is the 'Capacity savings:' section, which has a dropdown menu currently showing 'Thin-provisioned'. The dropdown menu is open, showing options: 'None', 'Thin-provisioned', and 'Compressed'. There is also a 'Deduplicated' checkbox which is currently unchecked. At the bottom, there is an 'I/O' section with a dropdown menu set to 'Automatic'.

Figure 6-18 Volume Creation with Capacity Saving option

**Note:** Consider the compression guidelines in Chapter 9, “Advanced features for storage efficiency” on page 449 before creating the first compressed volume copy on a system.

When a thin-provisioned or compressed volume is defined in a DRP, the Deduplicated option becomes available. Select this option to enable deduplication of the volume.

Thin provisioned and compressed volumes have a special icon in the **Capacity** column of the **Volumes** menu that makes it easy to distinguish them, as shown in Figure 6-19 on page 285.





Name	State	Synchronized	Capacity	Pool
itso-thin00-Pool0	✓ Online		 10.00 GiB	Pool0
itso-thin01-Pool1	✓ Online		 10.00 GiB	Pool1
itso-thin02-Pool1	✓ Online		 10.00 GiB	Pool1
itso-thin03-Pool1	✓ Online		 10.00 GiB	Pool1

Figure 6-19 Space-efficient volumes icon

Volume iso-thin00-Pool0 is thin-provisioned. Volume iso-thin01-Pool1 is compressed. There is no icon that indicates whether a volume is deduplicated.

## 6.5.2 Creating custom volumes

The Create Volumes window also enables Custom volume creation that expands the set of options for volume creation available to the administrator.

The Custom menu consists of the following panes:

- ▶ Volume Location: Mandatory; defines the number of volume copies, Pools to be used, and I/O group preferences.
- ▶ Volume Details: Mandatory; defines the Capacity savings option.
- ▶ Thin Provisioning: Enables configuration of Thin Provisioning settings if this capacity saving option is selected.
- ▶ Compressed: Enables configuration of Compression settings if this capacity saving option is selected.
- ▶ General: For configuring for Cache mode and Formatting.
- ▶ Summary

Work through these panes to customize your Custom volume as wanted, and then commit these changes by clicking **Create**.

You can mix and match settings on different panes to achieve the final volume configuration that meets your requirements

### Volume Location pane

The Volume Location pane is shown in Figure 6-20.

**Volume Location**

Volume copy type:  
 -

Pool:  
 -

-     
 Preferred node:  -     
 Accessible I/O groups:  -

Figure 6-20 Volume Location pane

This pane gives the following options:

- ▶ Volume copy type: You can choose between None (single volume copy) and Mirrored (two volume copies).
- ▶ Pool: Specifies storage pool to use for each of volume copies.
- ▶ Mirror sync rate: You can set the mirror sync rate for the volume copies. This option is displayed only for Mirrored volume copy type, and allows you to set the volume copy synchronization rate to a value 128 KiBps - 64 MiBps.
- ▶ Caching I/O group: You can choose between Automatic (allocated by the system) and manually specifying I/O group.
- ▶ Preferred node: You can choose between Automatic (allocated by the system) and manually specifying the preferred node for the volume.
- ▶ Accessible I/O groups: You can choose between Only the caching I/O group and All.

## Volume Details pane

The Volume Details pane is shown in Figure 6-21.



The screenshot shows the 'Volume Details' pane with the following fields and options:

- Quantity:** A spinner control set to 1.
- Capacity:** A field with a red border containing a red circle icon, followed by a dropdown menu set to 'GiB'.
- Name:** An empty text input field.
- Capacity savings:** A dropdown menu set to 'None'.
- Deduplicated:** An unchecked checkbox.
- Define another volume:** A blue link with a plus icon.

Figure 6-21 Volume Details pane

This pane gives the following options:

This pane features the following options:

- ▶ Quantity: Specify how many volumes are created.
- ▶ Capacity: Set the capacity of the volume.
- ▶ Name: Define the volume name.
- ▶ Capacity savings: Choose between None (standard-provisioned volume), Thin-provisioned and Compressed.
- ▶ Deduplicated: Thin-provisioned and Compressed volumes that are created in a DRP can be deduplicated.



If you click **Define another volume**, the GUI displays a subpane in which you can define the configuration of another volume, as shown in Figure 6-22.

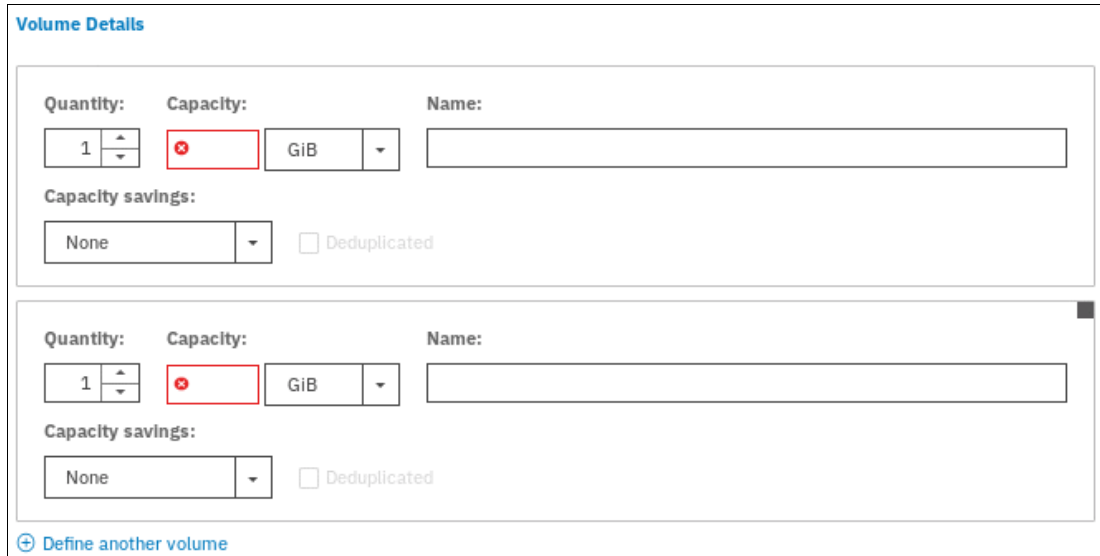


Figure 6-22 Volume Details pane with two volume subpanes

In this way, you can create volumes with different characteristics in a single invocation of the volume creation wizard.

### Thin Provisioning pane

If you choose to create a thin-provisioned volume, a Thin Provisioning pane is displayed, as shown in Figure 6-23.

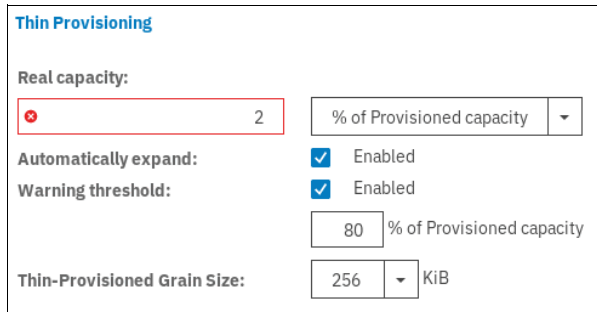


Figure 6-23 Custom volume creation – Thin Provisioning pane

This pane includes the following options:

- ▶ **Real capacity:** Real capacity of the volume, which is specified as percentage of the virtual capacity or in bytes.
- ▶ **Automatically expand:** Whether to automatically expand the real capacity of the volume if needed; defaults to Enabled.
- ▶ **Warning threshold:** Whether a warning message is sent and at what percentage of filled virtual capacity; defaults to Enabled, with a warning threshold set at 80%.
- ▶ **Thin-Provisioned Grain Size:** Allows you to define the grain size for the thin-provisioned volume; defaults to 256 KiB.

**Important:** If you do not use the **autoexpand** feature, the volume goes offline if it receives a write request after all real capacity is allocated.

The default grain size is 256 KiB. The optimum choice of grain size is dependent upon volume use type. Consider the following points:

- ▶ If you are *not* going to use the thin-provisioned volume as a FlashCopy source or target volume, use 256 KiB to maximize performance.
- ▶ If you *are* going to use the thin-provisioned volume as a FlashCopy source or target volume, specify the same grain size for the volume and for the FlashCopy function.
- ▶ If you plan to use Easy Tier with thin-provisioned volumes, see the IBM Support article [Performance Problem When Using Easy Tier With Thin Provisioned Volumes](#).

### Compressed pane

If you choose to create a compressed volume, a Compressed pane is displayed, as shown in Figure 6-24.

The screenshot shows a 'Compressed' configuration pane. It has a title bar 'Compressed' in blue. Below the title, there are three sections:

- Real capacity:** A text input field containing the number '2'. To its right is a dropdown menu labeled '% of Provisioned capacity'.
- Automatically expand:** A checkbox that is checked, followed by the text 'Enabled'.
- Warning threshold:** A checkbox that is checked, followed by the text 'Enabled'. Below this is a text input field containing '80' and a dropdown menu labeled '% of Provisioned capacity'.

Figure 6-24 Compressed pane

This pane gives the following options:

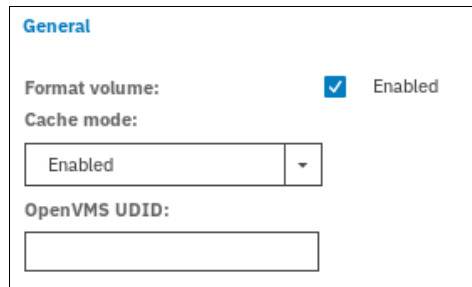
- ▶ **Real capacity:** Real capacity of the volume, which is specified as percentage of the virtual capacity or in bytes.
- ▶ **Automatically expand:** Whether to automatically expand the real capacity of the volume if needed; defaults to Enabled.
- ▶ **Warning threshold:** Whether a warning message is sent, and at what percentage of filled virtual capacity; defaults to Enabled, with a warning threshold set at 80%.

You cannot specify grain size for a compressed volume.

**Note:** Consider the compression guidelines in Chapter 9, “Advanced features for storage efficiency” on page 449 before creating the first compressed volume copy on a system.

## General pane

The General pane is shown in Figure 6-25.



General

Format volume:  Enabled

Cache mode:

OpenVMS UDID:

Figure 6-25 Custom volume creation – General pane

This pane includes the following options:

- ▶ **Format volume:** Controls whether the volume is formatted before being made available; defaults to Enabled.
- ▶ **Cache mode:** Controls volume caching; defaults to Enabled. Other available options are Read-only and Disabled.
- ▶ **OpenVMS unit device identifier (UDID):** Each OpenVMS Fibre Channel (FC)-attached volume requires a user-defined identifier or UDID. A UDID is a nonnegative integer that is used in the creation of the OpenVMS device name.

### 6.5.3 Creating volumes in multi-site topologies

This section describes creation of volumes in multi-site topologies.

**Note:** For volumes in multi-site topologies the asterisk (\*) does not indicate the primary copy, but the local volume copy that is used for data reads.

#### HyperSwap volumes

To create a HyperSwap volume, complete the following steps:

1. In the IBM Spectrum Virtualize GUI, select **Volumes** → **Volumes**, as shown in Figure 6-26 on page 290.

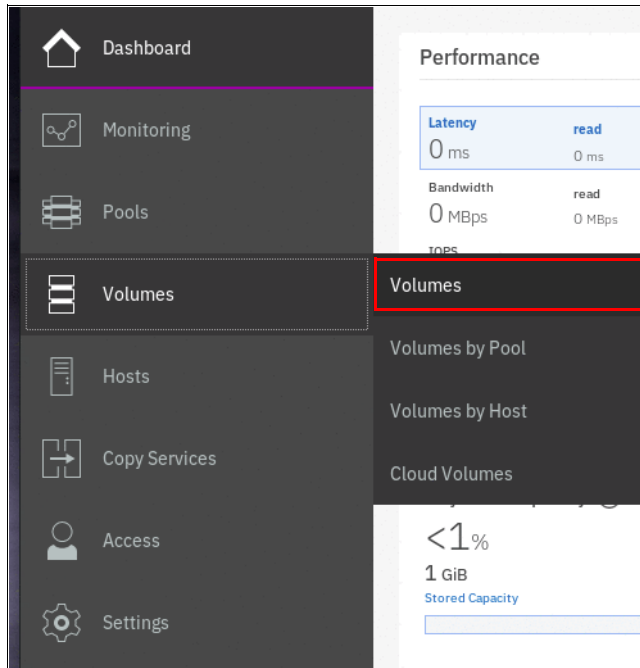


Figure 6-26 Volumes menu

A list of volumes, their state, capacity, and associated storage pools, is displayed.

2. To create a volume, click **Create Volumes**, as shown in Figure 6-27.

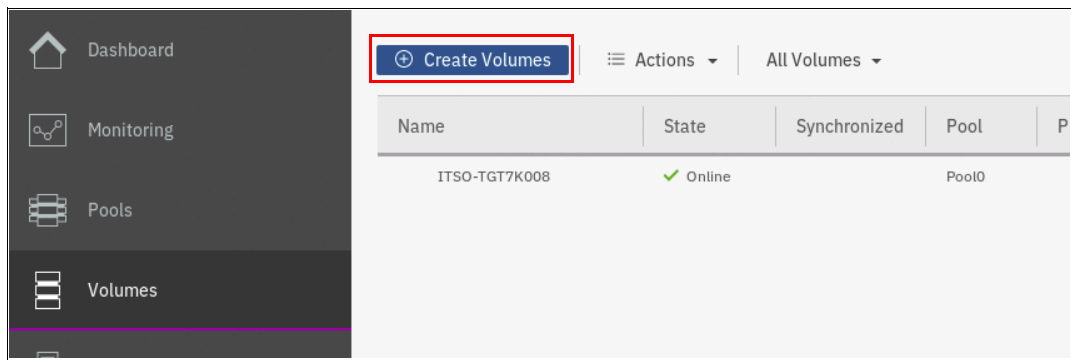


Figure 6-27 Create Volumes button

The Create Volumes tab opens the Create Volumes window, which displays available creation methods.

**Note:** The volume classes that are displayed in the Create Volumes window depend on the topology of the system.

The Create Volumes window for standard topology is shown in Figure 6-28.

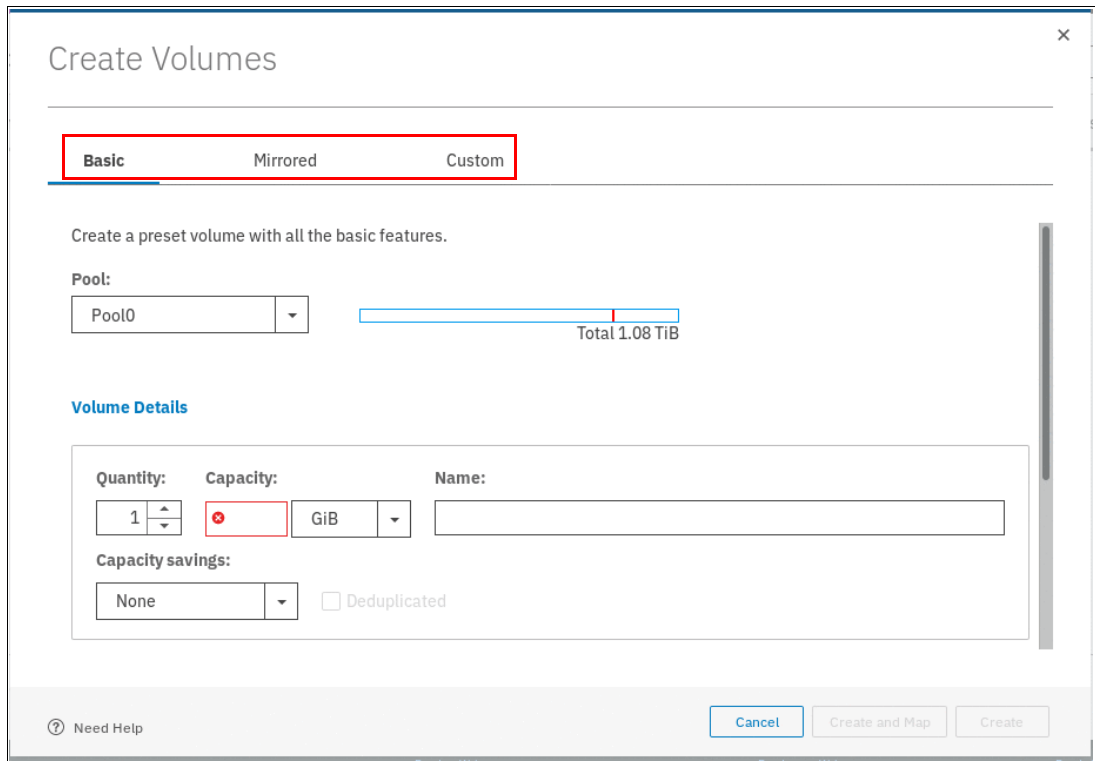


Figure 6-28 Basic, Mirrored, and Custom Volume Creation options

**Note:** Consider the following points:

- ▶ A *Basic volume* is a volume that has only one physical copy, uses storage that is allocated from a single pool on one site, and uses read/write cache mode.
- ▶ A *Mirrored volume* is a volume with two physical copies, where each volume copy can belong to a different storage pool.
- ▶ A *Custom volume* (in the context of this menu) is a Basic or Mirrored volume with values of some of its parameters that are changed from the defaults.
- ▶ The Create Volumes window also provides (by way of the Capacity Savings parameter) the ability to change the default provisioning of a Basic or Mirrored Volume to Thin-provisioned or Compressed. For more information, see “Capacity savings option” on page 284.

The notable difference between a HyperSwap volume and a basic volume creation is that in the HyperSwap Details the system asks for storage pool names at each site. The system uses its topology awareness to map storage pools to sites, which ensures that the data is correctly mirrored across locations.

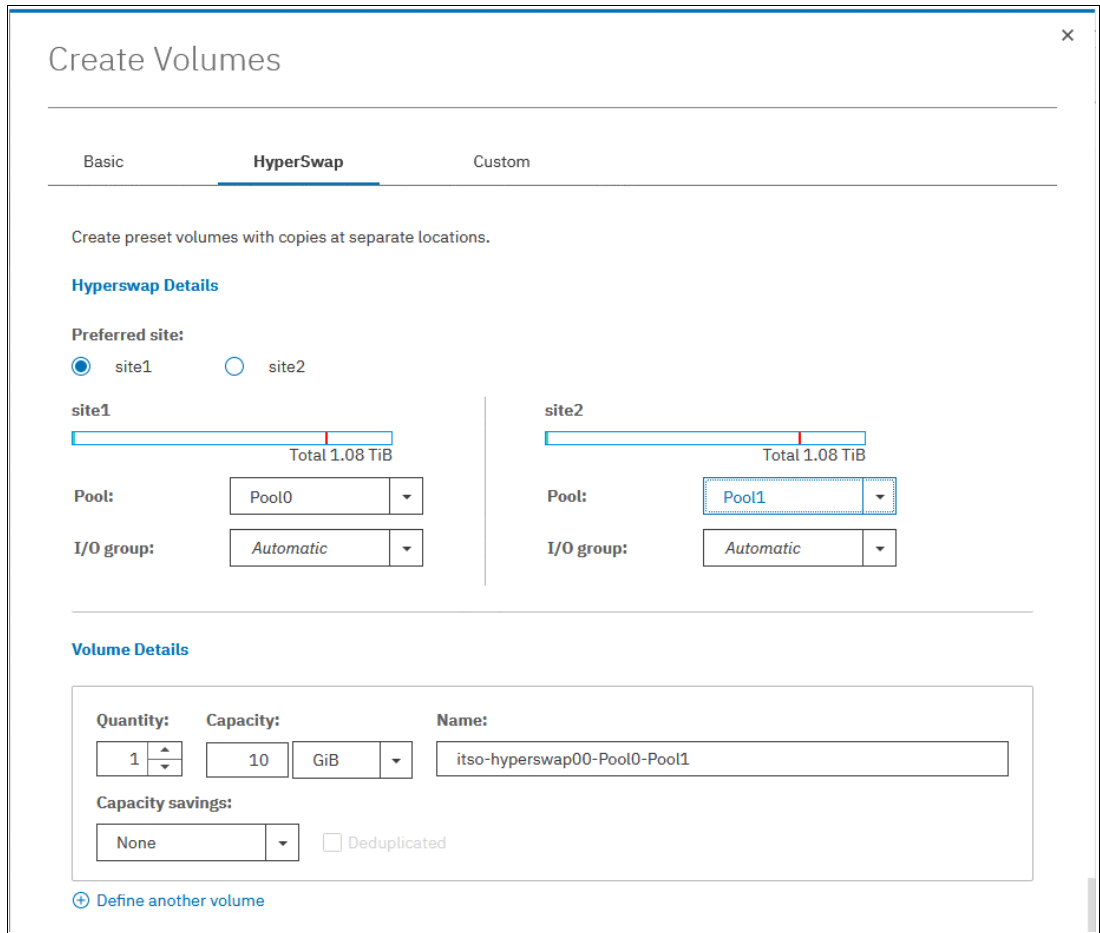


Figure 6-29 HyperSwap Volume creation window

As shown in Figure 6-29, a single volume is created, with volume copies in sites `site1` and `site2`. This volume is in an active-active (MM) relationship with extra resilience that is provided by two change volumes.

After the volume is created, it is visible in the volumes list, as shown in Figure 6-30.

Create Volumes		Actions	All Volumes	
Name	State	Synchronized	Pool	
itso-hyperswap00-Pool0-Pool1	Online (formatting)		Multiple	
itso-hyperswap00-Pool0-Pool1 (site1)	Online (formatting)	Yes	Pool0	
itso-hyperswap00-Pool0-Pool1 (site2)	Online (formatting)	Yes	Pool1	

Figure 6-30 HyperSwap volume visible in the Volumes list

The Pool column shows the value `Multiple`, which indicates that a volume is a HyperSwap volume. A volume copy at each site is visible and the change volumes that are used by the technology are not displayed in this GUI view.

A single `mkvolume` command allows creation of a HyperSwap volume. Up to IBM Spectrum Virtualize V7.5, this process required careful planning and the use of the following sequence of commands:

- ▶ `mkvdisk master_vdisk`
- ▶ `mkvdisk aux_vdisk`
- ▶ `mkvdisk master_change_volume`
- ▶ `mkvdisk aux_change_volume`
- ▶ `mkrcrelationship -activeactive`
- ▶ `chrcrelationship -masterchange`
- ▶ `chrcrelationship -auxchange`
- ▶ `addvdiskaccess`

### Stretched volumes

If the IBM Spectrum Virtualize system is set up in stretched topology, the **Create Volumes** menu shows the Stretched option. Click **Stretched** to create a stretched volume and define its attributes as shown in Figure 6-31 on page 294.

**Note:** For more information about Stretched topology, see IBM Knowledge Center and select **IBM SAN Volume Controller** → **Product overview** → **Technical Overview** → **Volumes** → **Stretched volumes**.

Creation options use site awareness of controllers and back-end storage, as shown in Figure 6-31.

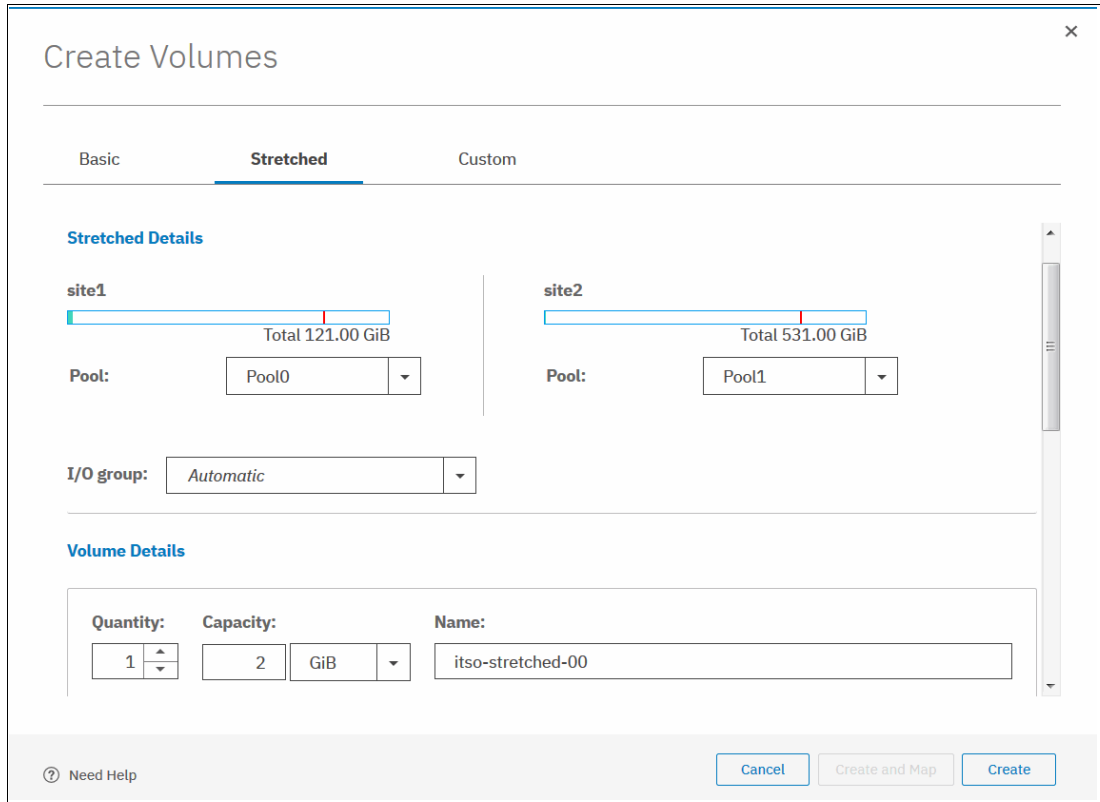


Figure 6-31 Creating a Stretched Volume

After you click **Create**, a volume is created and you are returned to the Volumes view. The newly created volume is visible, as shown in Figure 6-32.

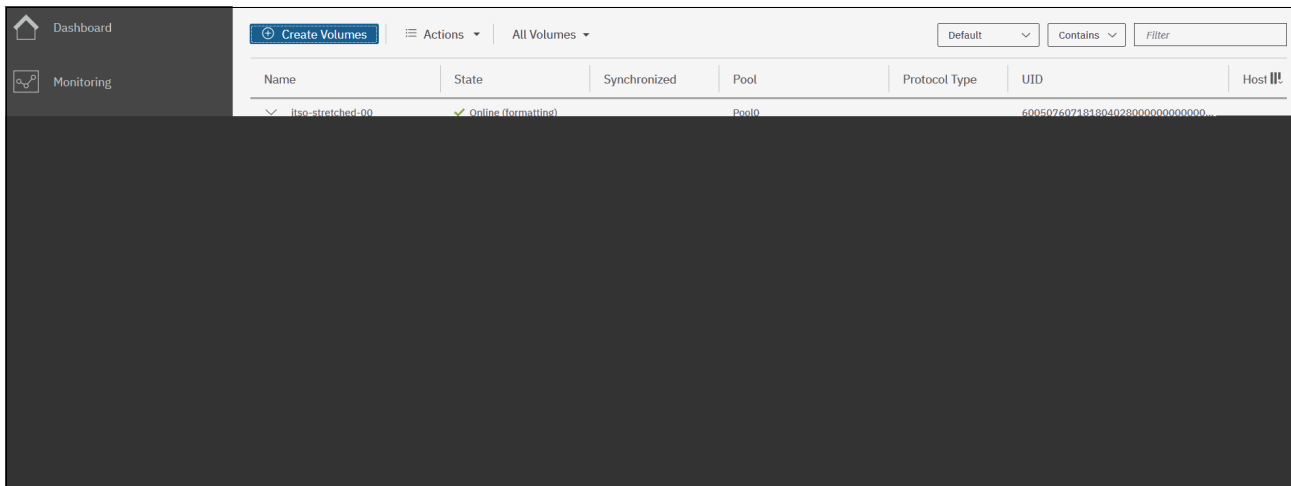


Figure 6-32 Stretched volume formatting after creation

**Tip:** There is no immediate indication that the volume is stretched. You can include a character string, such as `-stretch-` in volume names to indicate volumes that are stretched.



## 6.5.4 I/O throttling

This section describes how to work with I/O throttling on a volume.

### Defining a volume throttle

To set a volume throttle, complete the following steps:

1. Select **Volumes** → **Volumes**, and select the wanted volume to throttle. From the **Actions** menu, select **Edit Throttle**, as shown in Figure 6-33.

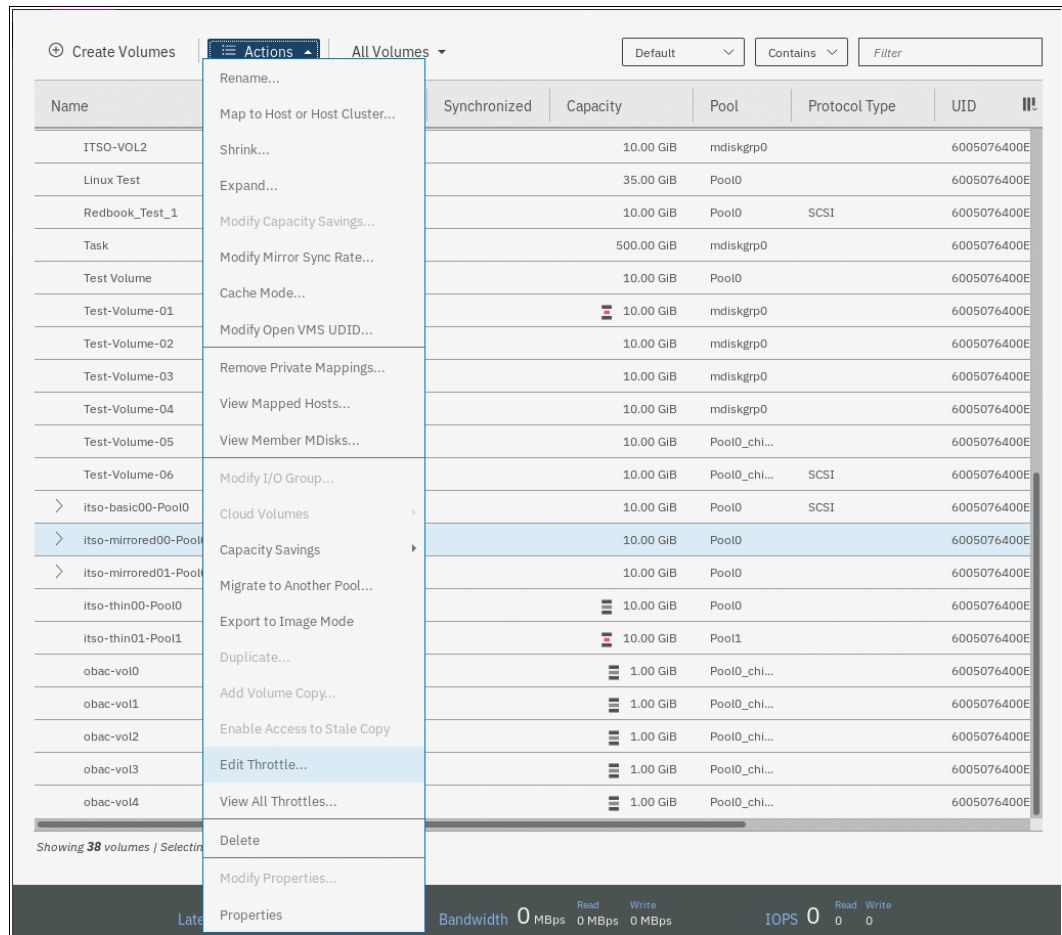


Figure 6-33 Edit throttle menu item

2. In the Edit Throttle window, define the throttle in terms of number of IOPS or bandwidth. In our example, we set an IOPS throttle of 10,000, as shown in Figure 6-34. Click **Create**.

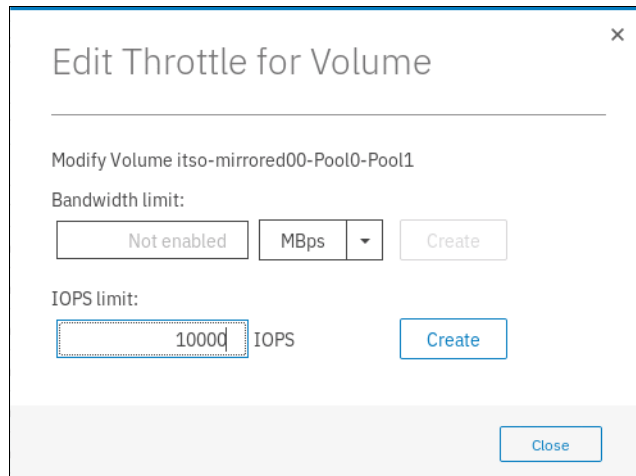


Figure 6-34 IOPS throttle on a volume

3. After the Edit Throttle task completes successfully, the Edit Throttle window is shown again. You can now set throttle based on the different metrics, modify throttle, or close the window without performing further actions by clicking **Close**.

## Listing volume throttles

To view volume throttles, select **Volumes** → **Volumes**, click the **Actions** menu and then, select **View All Throttles**, as shown in Figure 6-35.

The screenshot shows a storage management interface with a table of volumes and an 'Actions' menu open. The 'View All Throttles...' option is highlighted in the menu. The table below shows various volume configurations.

Name	Synchronized	Capacity	Pool	Protocol Type	UID
ITSO-VOL2		10.00 GiB	mdiskgrp0		6005076400E
Linux Test		35.00 GiB	Pool0		6005076400E
Redbook_Test_1		10.00 GiB	Pool0	SCSI	6005076400E
Task		500.00 GiB	mdiskgrp0		6005076400E
Test Volume		10.00 GiB	Pool0		6005076400E
Test-Volume-01		10.00 GiB	mdiskgrp0		6005076400E
Test-Volume-02		10.00 GiB	mdiskgrp0		6005076400E
Test-Volume-03		10.00 GiB	mdiskgrp0		6005076400E
Test-Volume-04		10.00 GiB	mdiskgrp0		6005076400E
Test-Volume-05		10.00 GiB	Pool0_chi...		6005076400E
Test-Volume-06		10.00 GiB	Pool0_chi...	SCSI	6005076400E
> itso-basic00-Pool0		10.00 GiB	Pool0	SCSI	6005076400E
> itso-mirrored00-Pool0		10.00 GiB	Pool0		6005076400E
> itso-mirrored01-Pool0		10.00 GiB	Pool0		6005076400E
its0-thin00-Pool0		10.00 GiB	Pool0		6005076400E
its0-thin01-Pool1		10.00 GiB	Pool1		6005076400E
obac-vol0		1.00 GiB	Pool0_chi...		6005076400E
obac-vol1		1.00 GiB	Pool0_chi...		6005076400E
obac-vol2		1.00 GiB	Pool0_chi...		6005076400E
obac-vol3		1.00 GiB	Pool0_chi...		6005076400E
obac-vol4		1.00 GiB	Pool0_chi...		6005076400E

Figure 6-35 View all throttles menu item

The **View All Throttles** menu shows all volume throttles defined in the system, as shown in Figure 6-36 on page 298.

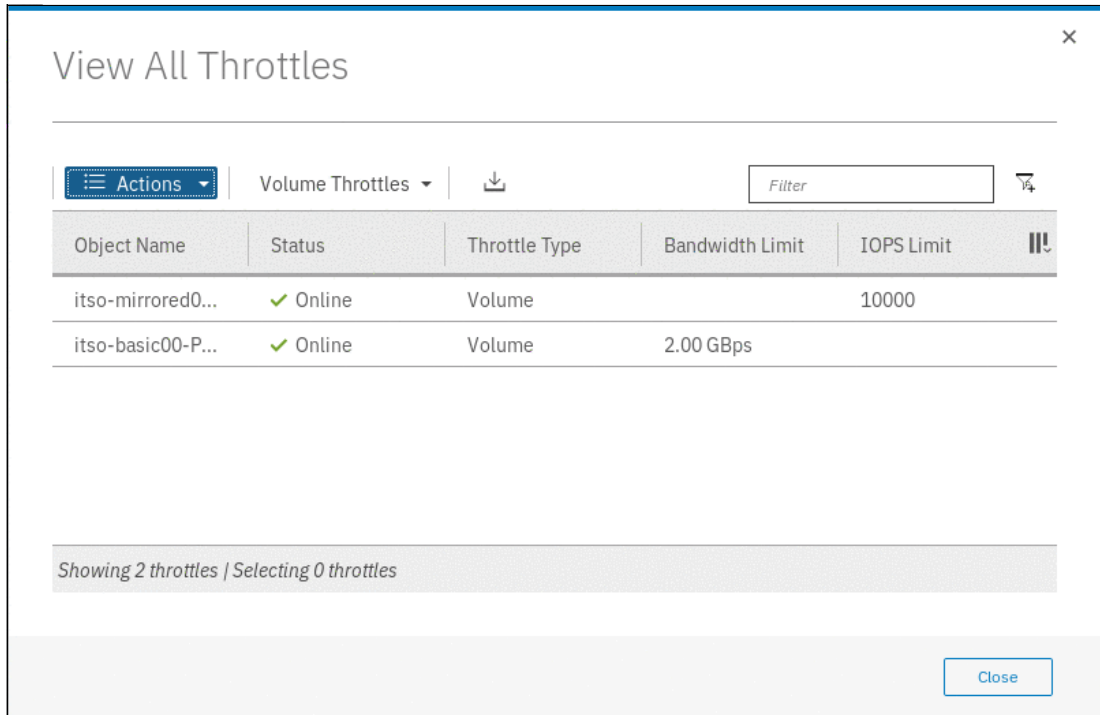


Figure 6-36 View volume throttles

You can view other throttles by selecting a different throttle type in the drop-down menu, as shown in Figure 6-37.

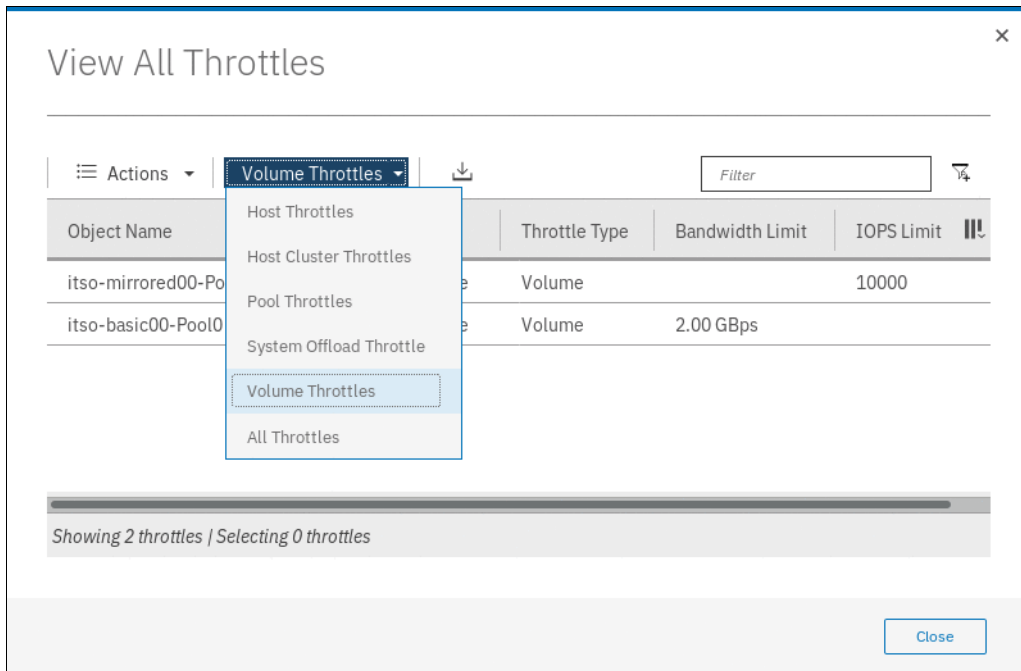


Figure 6-37 Filter throttle type

## Modifying or removing a volume throttle

To remove a volume throttle, complete the following steps:

1. From the **Volumes** menu, select the volume to which is attached the throttle that you want to remove. Select **Actions** → **Edit Throttle**, as shown in Figure 6-38.

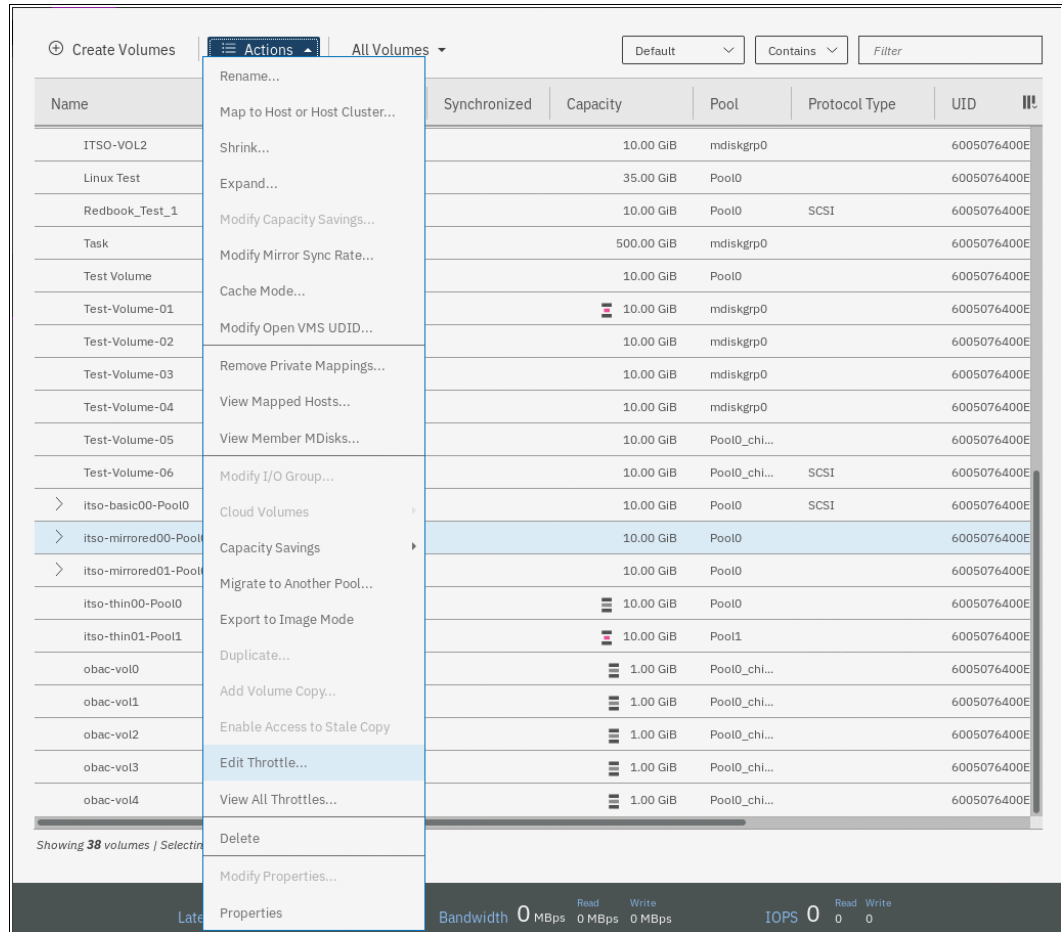


Figure 6-38 Edit throttle menu item

2. In the Edit Throttle window, click **Remove** for the throttle you want to remove. As shown in Figure 6-39, we remove the IOPS throttle from the volume.

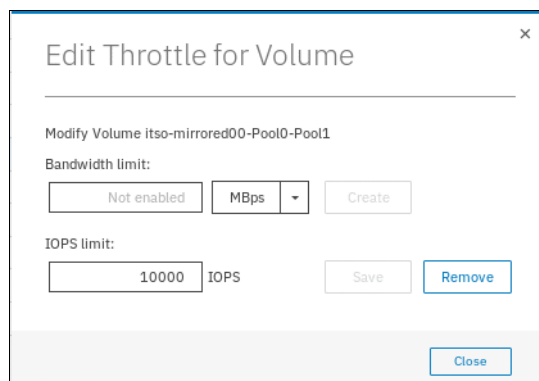


Figure 6-39 Remove throttle

After the Edit Throttle task completes successfully, the Edit Throttle window is shown again. You can now set throttle based on the different metrics, modify throttle, or close the window without performing any action by clicking **Close**.

## 6.5.5 Volume protection

To configure volume protection, select **Settings** → **System** → **Volume Protection**, as shown in Figure 6-40.

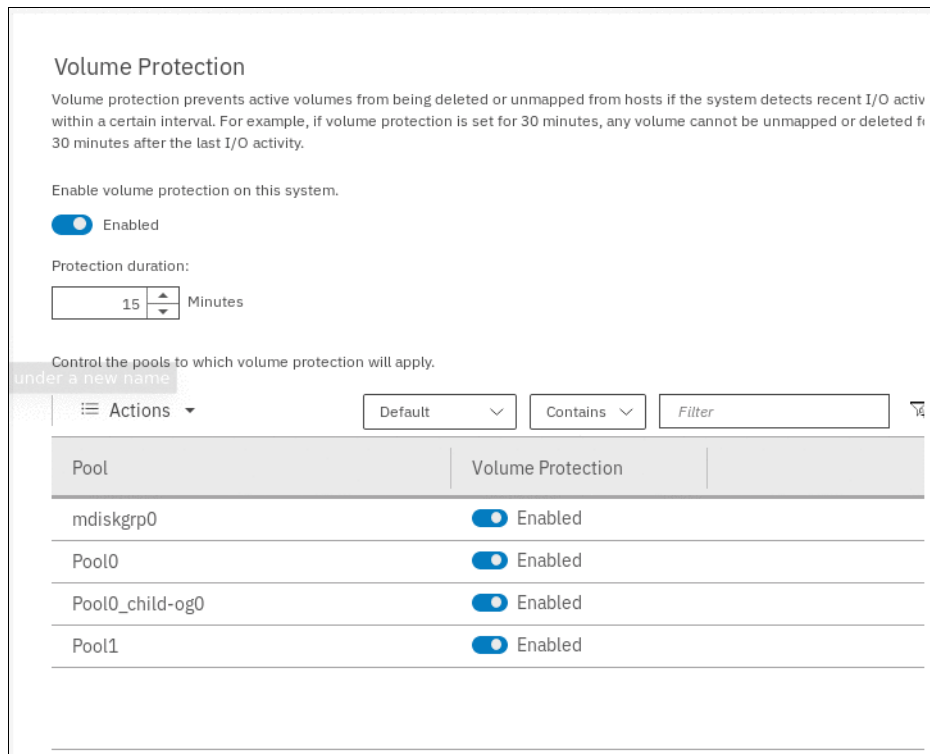


Figure 6-40 Volume protection configuration

In this view, you can configure system-wide volume protection (enabled by default), set the minimum inactivity period required to allow volume deletion (Protection duration) and configure volume protection for each configured pool (enabled by default). In the example, volume protection is enabled with 15 minutes minimum inactivity period and is turned on for all configured pools.

## 6.5.6 Modifying a volume

After a volume is created, many of its characteristics can be modified, as described next.

## Shrinking

To shrink a volume, complete the following steps:

1. Ensure that you have a current and verified backup of any in-use data that is stored on the volume that you intend to shrink.
2. From the **Volumes** menu, select the volume that you want to shrink. Select **Actions** → **Shrink...**, as shown in Figure 6-41.

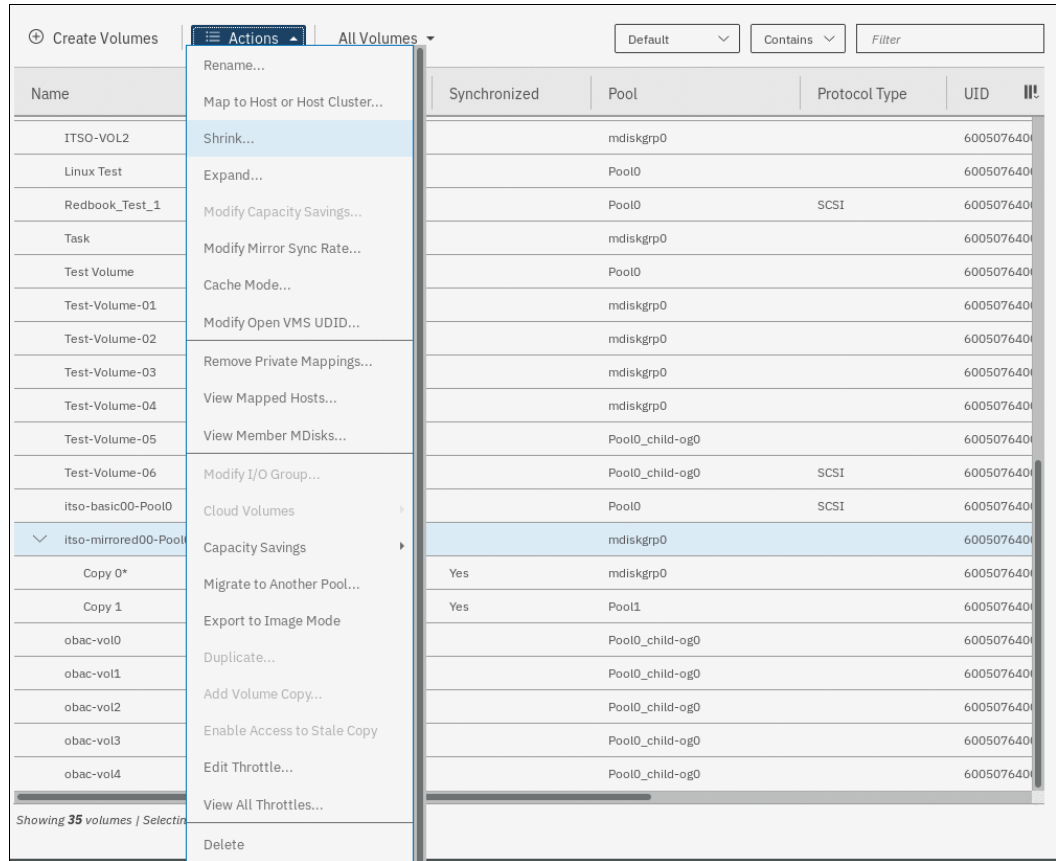


Figure 6-41 Volume shrink menu item

3. Specify **Shrink by** or **Final size** value (the other is calculated automatically), as shown in Figure 6-42.

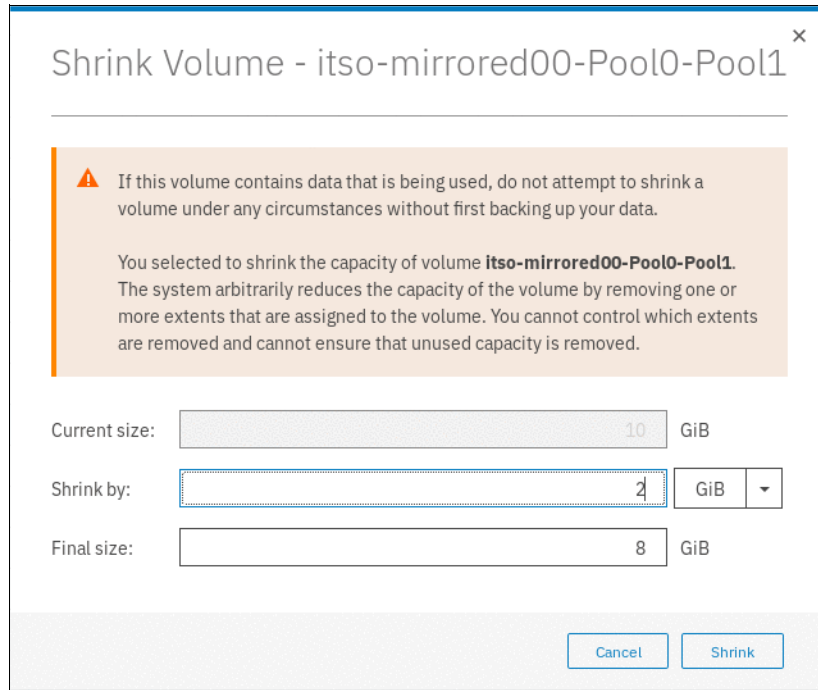


Figure 6-42 Specify the size of the shrunk volume

**Note:** The storage system reduces the volume capacity by removing one or more arbitrarily selected extents. Do not shrink a volume that contains data that is being used unless you have a current and verified backup of the data.

4. Click **OK** to confirm the action, as shown in Figure 6-43.

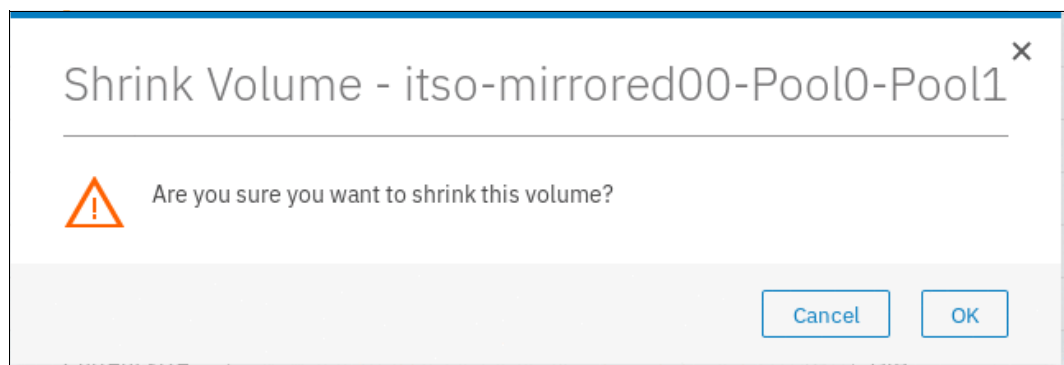


Figure 6-43 Confirm volume shrinking.



- After the operation completes, you can see the volume with the new size in the **Volumes** → **Volumes** view, as shown in Figure 6-44.

Name	State	Synchronized	Capacity	Pool
> its0-basic00-Pool0	✓ Online		10.00 GiB	Pool0
> its0-mirrored00-Pool0-Pool1	✓ Online		5.00 GiB	Pool0
> its0-mirrored01-Pool0-Pool1	✓ Online		10.00 GiB	Pool0

Figure 6-44 Shrank volume size

## Expanding

To expand a volume, complete the following steps:

- From the **Volumes** menu, select the volume that you want to expand. Select **Actions** → **Expand...**, as shown in Figure 6-45.

The screenshot shows a web-based interface for managing volumes. At the top, there are buttons for 'Create Volumes', 'Actions', and 'All Volumes'. Below these is a table of volumes with columns for Name, Synchronized, Capacity, Pool, and Protocol Type. The 'Expand...' option is highlighted in the 'Actions' dropdown menu for the selected volume 'its0-mirrored00-Pool0'. The table shows various volumes with different capacities and pool configurations.

Name	Synchronized	Capacity	Pool	Protocol Type
ITSO-VOL0_01		10.00 GiB	mdiskgrp0	
ITSO-VOL0_02		10.00 GiB	mdiskgrp0	
ITSO-VOL1		10.00 GiB	mdiskgrp0	
ITSO-VOL2		10.00 GiB	mdiskgrp0	
Linux Test		35.00 GiB	Pool0	
Redbook_Test_1		10.00 GiB	Pool0	SCSI
Task		500.00 GiB	mdiskgrp0	
Test Volume		10.00 GiB	Pool0	
Test-Volume-01		10.00 GiB	mdiskgrp0	
Test-Volume-02		10.00 GiB	mdiskgrp0	
Test-Volume-03		10.00 GiB	mdiskgrp0	
Test-Volume-04		10.00 GiB	mdiskgrp0	
Test-Volume-05		10.00 GiB	Pool0_child-og0	
Test-Volume-06		10.00 GiB	Pool0_child-og0	SCSI
its0-basic00-Pool0		10.00 GiB	Pool0	SCSI
its0-mirrored00-Pool0		5.00 GiB	mdiskgrp0	
Copy 0*		5.00 GiB	mdiskgrp0	
Copy 1		5.00 GiB	Pool1	
obac-vol0		1.00 GiB	Pool0_child-og0	
obac-vol1		1.00 GiB	Pool0_child-og0	
obac-vol2		1.00 GiB	Pool0_child-og0	
obac-vol3		1.00 GiB	Pool0_child-og0	
obac-vol4		1.00 GiB	Pool0_child-og0	

Figure 6-45 Volume expand menu item

- Specify **Expand by:** or **Final size:** (the other value is calculated automatically) and click **Expand**, as shown in Figure 6-46.

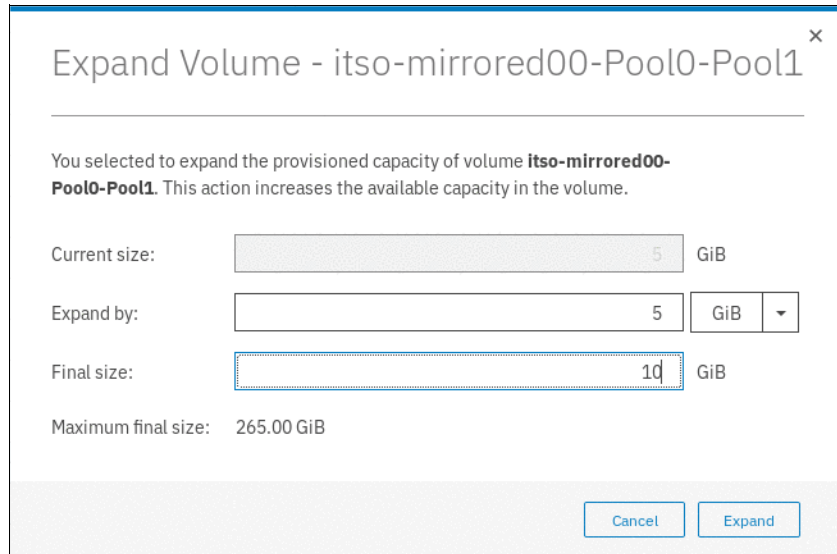


Figure 6-46 Specifying the expanded volume size

- After the operation completes (including formatting more space), you can see the volume with the new size in the **Volumes** → **Volumes** view, as shown in Figure 6-47.

Name	State	Synchronized	Capacity	Pool
Test-Volume-06	Online		10.00 GiB	Pool0_child-og0
itso-basic00-Pool0	Online		10.00 GiB	Pool0
itso-mirrored00-Pool0-Pool1	Online		10.00 GiB	mdiskgrp0
Copy 0*	Online	Yes	10.00 GiB	mdiskgrp0
Copy 1	Online	Yes	10.00 GiB	Pool1
obac-vol0	Online		1.00 GiB	Pool0_child-og0
obac-vol1	Online		1.00 GiB	Pool0_child-og0
obac-vol2	Online		1.00 GiB	Pool0_child-og0

Showing 35 volumes / Selecting 1 volume (10.00 GiB)

Figure 6-47 Expanded volume size

**Note:** Expanding a volume is not sufficient to increase available space that is visible to the host. The host must become aware of the changed volume size at the operating system level; for example, by way of a bus rescan. More operations at the Logical Volume Manager (LVM) or file system levels might be needed before more space is visible to applications that are running on the host.

## Modify capacity savings

This action is available for space-efficient volumes only. To modify capacity savings options for a volume, complete the following steps:

1. From the **Volumes** menu, select the volume that you want to modify. Select **Actions** → **Modify Capacity Savings...**, as shown in Figure 6-48.

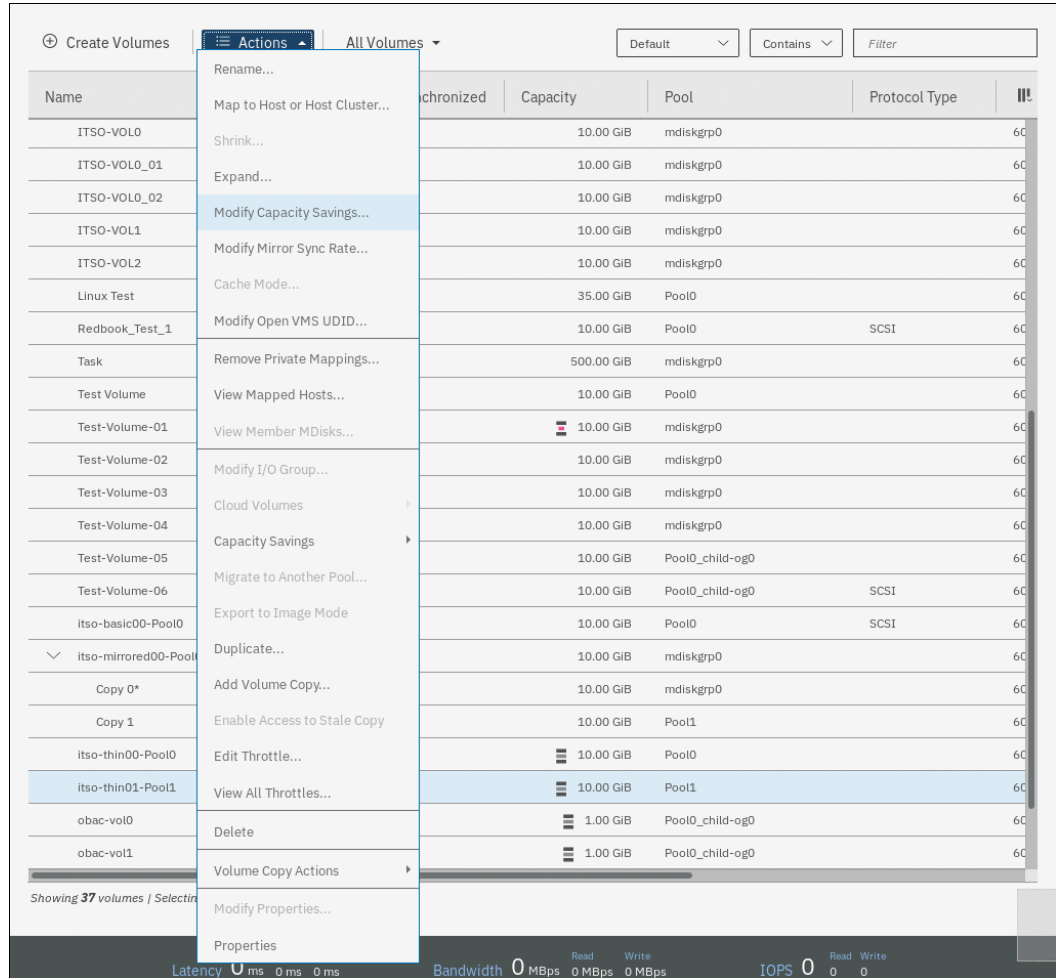


Figure 6-48 Modify capacity savings menu item

2. Select the wanted capacity savings option, as shown in Figure 6-49.

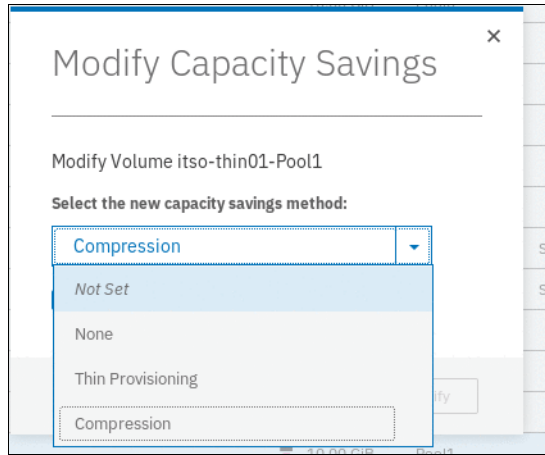


Figure 6-49 Capacity savings option for a volume

For volumes that are configured in a DRP, deduplication can be enabled, as shown in Figure 6-50.

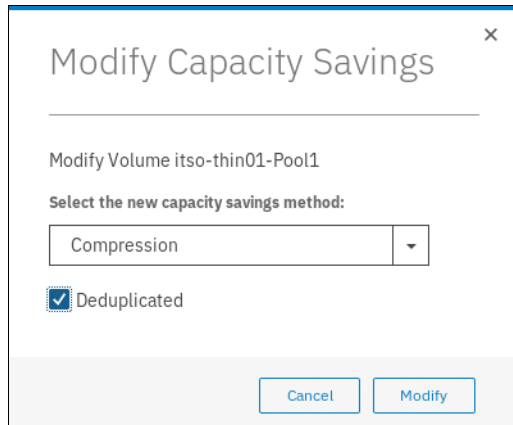


Figure 6-50 Enabling deduplication on a volume

After you configure the wanted capacity savings options of a volume, click **Modify** to apply them. When the operation completes, you are returned to the Volumes view.

### Modifying mirror sync rate

This action is available for mirrored volumes only. To modify mirror sync rate of a volume, complete the following steps:

1. From the **Volumes** menu, select the volume that you want to modify. Select **Actions** → **Modify Mirror Sync Rate...**, as shown in Figure 6-51 on page 307.

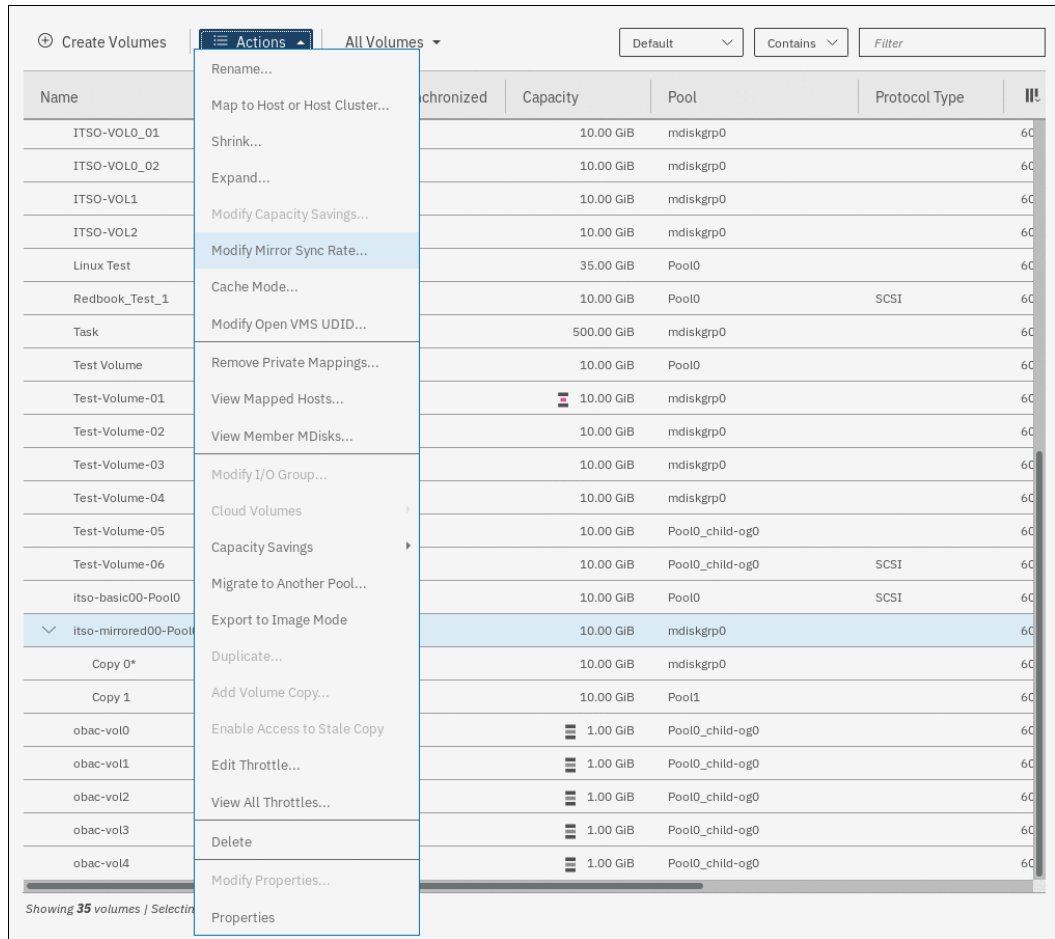


Figure 6-51 Modify mirror sync rate menu item

2. Select the mirror sync rate from the list. Available values are 0 KBps and 65 MBps. Click **Modify** to set the picked mirror sync rate for the volume, as shown in Figure 6-52.

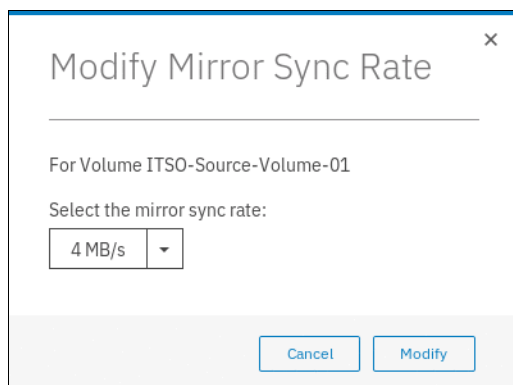


Figure 6-52 Set the mirror sync rate

When the operation completes, you are returned to the Volumes view.

## Changing cache mode

To change volume cache mode, complete the following steps:

1. From the **Volumes** menu, select the volume that you want to modify. Select **Actions** → **Cache Mode...**, as shown in Figure 6-53.

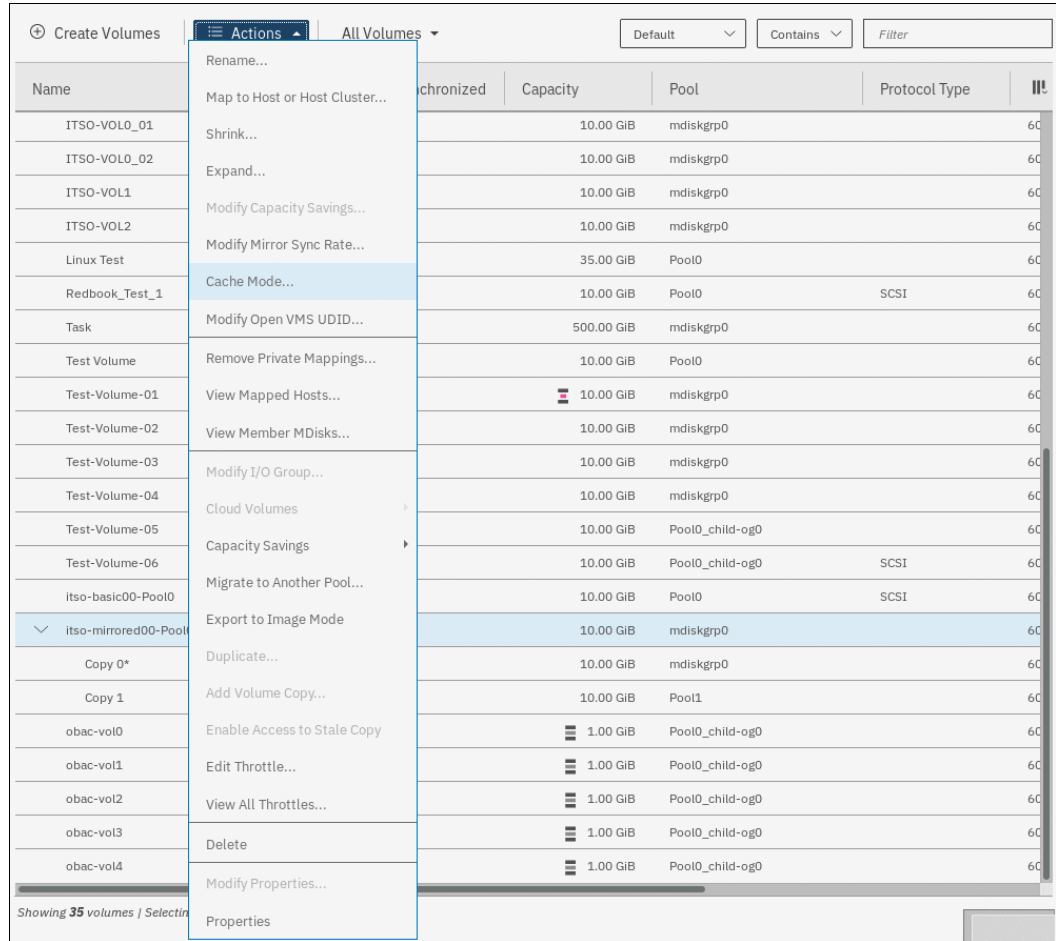


Figure 6-53 Modify volume cache mode menu item

2. Select the wanted cache mode for the volume from the drop-down list and click **OK**, as shown in Figure 6-54.

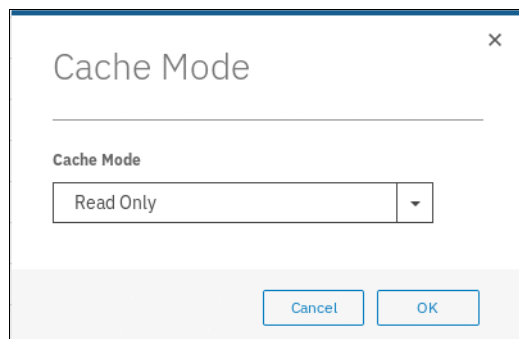


Figure 6-54 Setting volume cache mode

When the operation completes, you are returned to the Volumes view.

## Modifying OpenVMS UDID

To change the OpenVMS UDID, use the **Modify OpenVMS UDID** menu.

A *UDID* is a nonnegative integer that is used in the creation of the OpenVMS device name. All fibre-attached volumes have an allocation class of \$1\$, followed by the letters DGA, followed by the UDID. All storage unit LUNs that you assign to an OpenVMS system need an UDID so that the operating system can detect and name the device. LUN 0 must also have a UDID; however, the system displays LUN 0 as \$1\$GGA<UDID>, not as \$1\$DGA<UDID>.

For more information about fibre-attached storage devices, see, [Guidelines for OpenVMS Cluster Configurations](#).

To change volume OpenVMS UDID, complete the following steps:

1. From the **Volumes** menu, select the volume that you want to modify. Select **Actions** → **Modify Open VMS UDID...**, as shown in Figure 6-55.

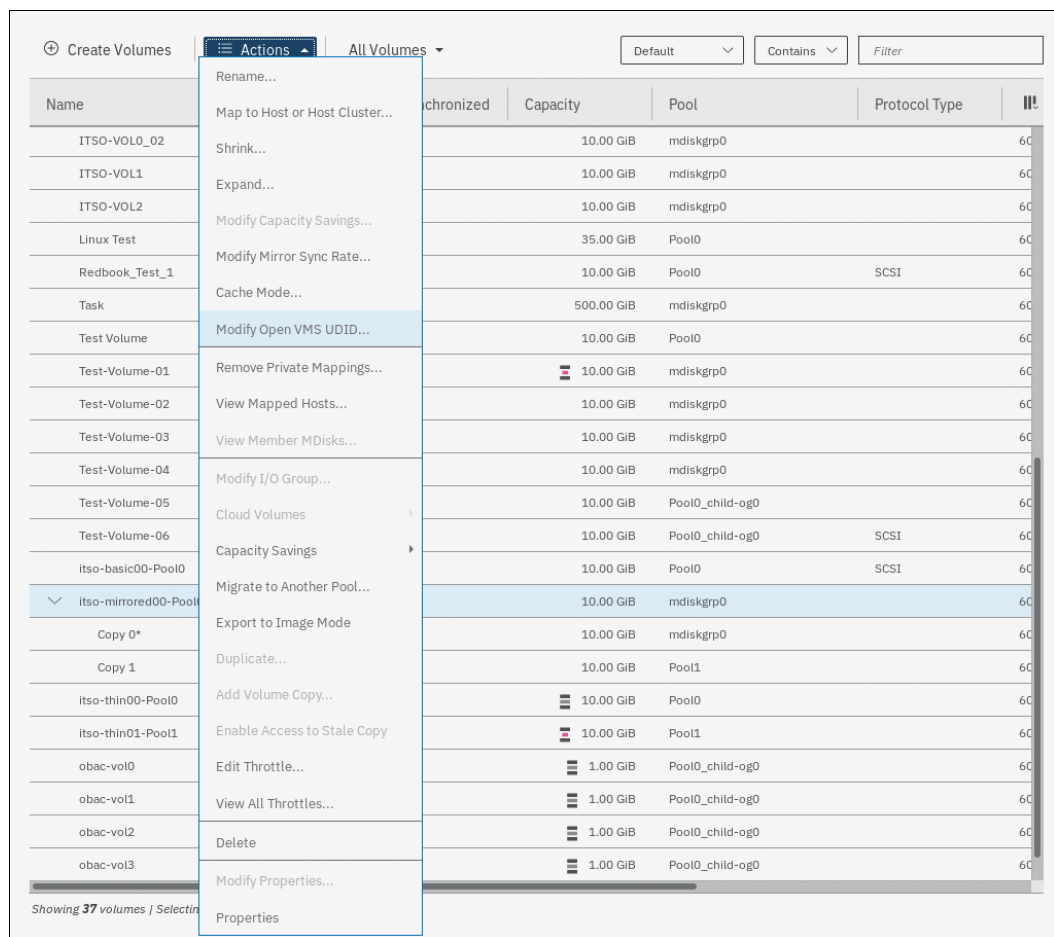


Figure 6-55 Modify Open VMS UDID menu item

2. Specify UDID for the volume and click **Modify**, as shown in Figure 6-56.

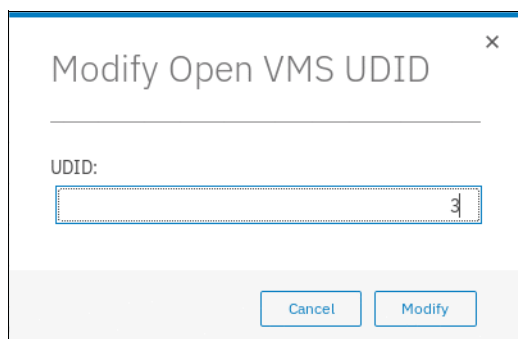


Figure 6-56 Set volume UDID

When the operation completes, you are returned to the Volumes view.

## 6.5.7 Deleting a volume

When you attempt to delete a volume, the system verifies whether it is a part of a host mapping, FlashCopy mapping, or remote-copy relationship. If any of these mappings exist, the delete attempt fails unless the **-force** parameter is specified in the corresponding remove commands.

If volume protection is enabled, a delete fails if the system finds a recent I/O activity to the volume, even if the **-force** parameter is specified. The **-force** parameter overrides the volume dependencies, not the volume protection setting.

To delete a volume, complete the following steps:

1. From the **Volumes** menu, select the volume that you want to modify. Select **Actions** → **Delete**, as shown in Figure 6-57 on page 311.



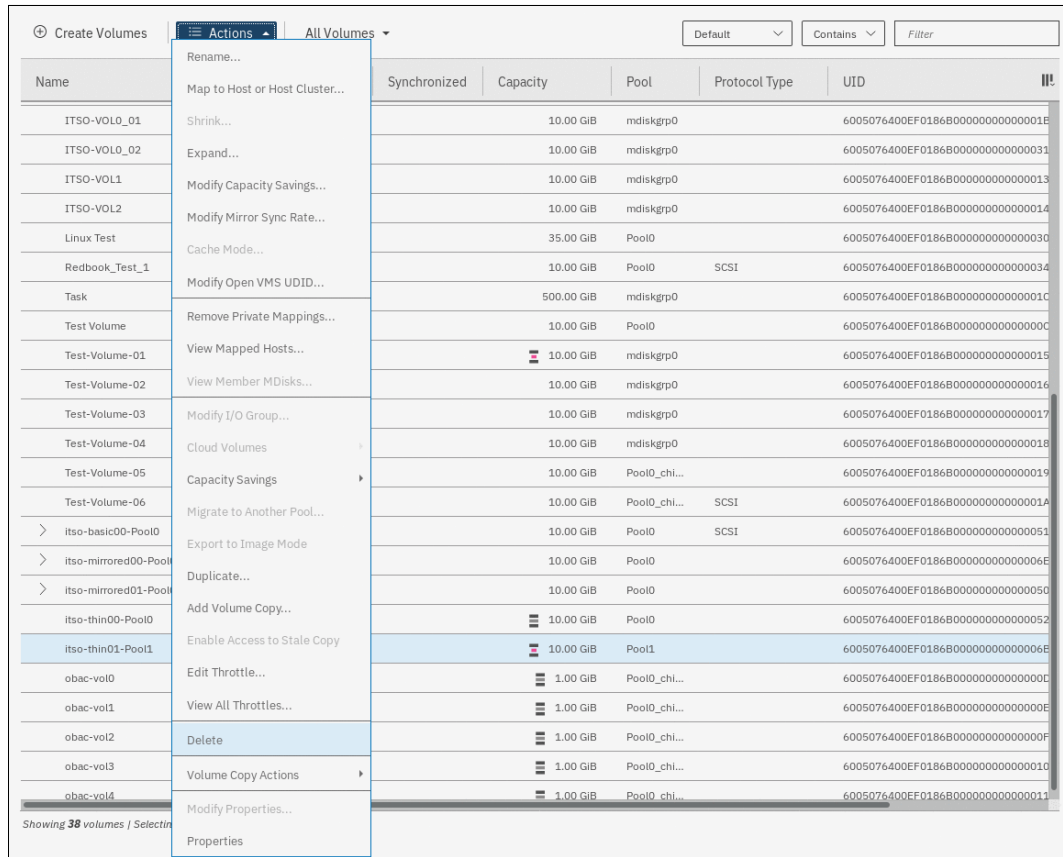


Figure 6-57 Volume delete menu item

- Review the list of volumes selected for deletion and enter the number of volumes you intend to delete, as shown in Figure 6-58. Click **Delete** to remove the volume from the system configuration.

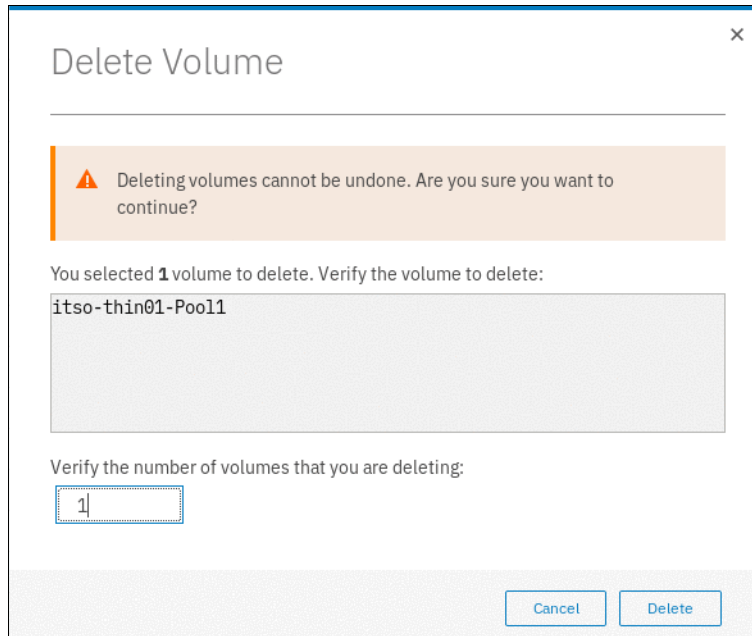


Figure 6-58 Confirm volume deletion

When the operation completes, you are returned to the Volumes view.

## 6.5.8 Mapping a volume to a host

To make a volume available to a host or cluster of hosts, it must be mapped. A volume can be mapped to the host at creation time, or later.

To map a volume to a host or cluster, complete the following steps:

1. From the **Volumes** menu, select the volume that you want to modify. From the **Actions** menu, select **Map to Host or Host Cluster...**, as shown in Figure 6-59.

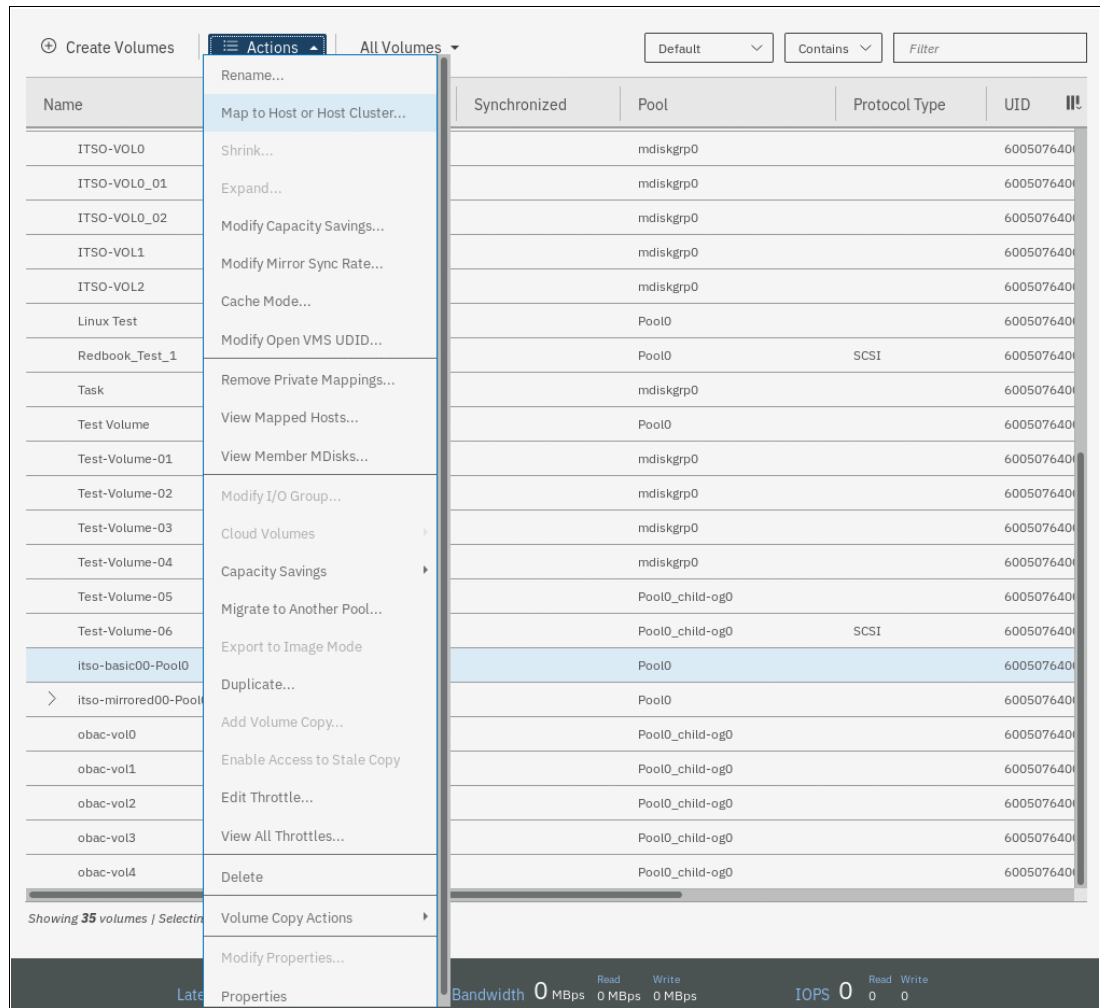


Figure 6-59 Volume mapping menu item

**Tip:** An alternative way of opening the Actions menu is to highlight (select) a volume and right-click.

2. A Create Mapping window opens. In this window, select whether to create a mapping to a Host or Host Cluster. The list of objects of the suitable type is displayed. Select to which hosts or host clusters the volume should be mapped.

You can allow the storage system to assign the SCIS LUN ID to the volume by selecting **System Assign**, or select **Self Assign** and enter the LUN ID by yourself. Click **Next**.  
 In the example that is shown in Figure 6-60, a single volume is mapped to a host and the storage system assigns the SCSI LUN IDs.

## Create Mapping ×

---

Create Mappings to:

Hosts  
 Host Clusters

Select hosts to map to itso-basic00-Pool0

🔍

Name	Status	Host Type	Host Mappings	Owners
ITSO-VMHOST-02	✔ Online	Generic	Yes	
ITSO-host2	✘ Offline	Generic	No	

Showing 2 hosts / Selecting 1 host

Would you like the system to assign SCSI LUN IDs or manually assign these IDs?

System Assign  
 Self Assign

Cancel
◀ Back
Next ▶

Figure 6-60 Mapping a volume to a host

3. A summary window is displayed all volume mappings for the selected host. The new mapping is highlighted, as shown in Figure 6-61. Review the future configuration state and click **Map Volumes** to map the volume.

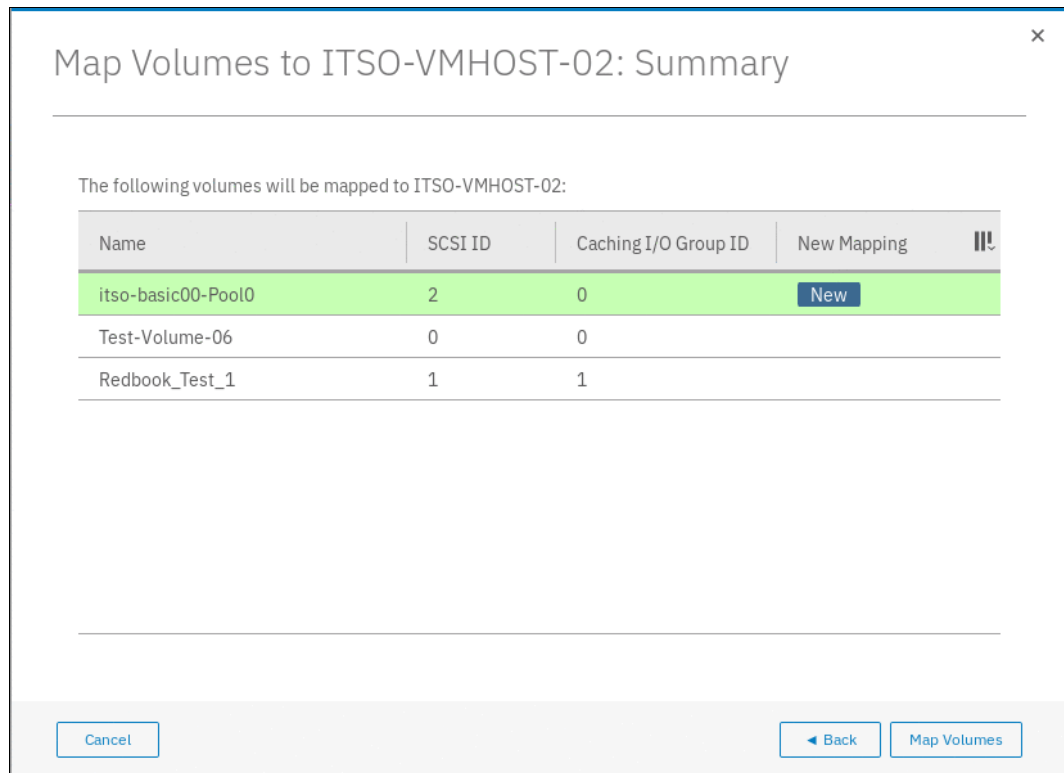


Figure 6-61 Map volume to host cluster summary

4. After the task completes, the wizard returns to the Volumes window. You can list volumes mapped to the host by selecting **Hosts** → **Mappings**, as shown in Figure 6-62.

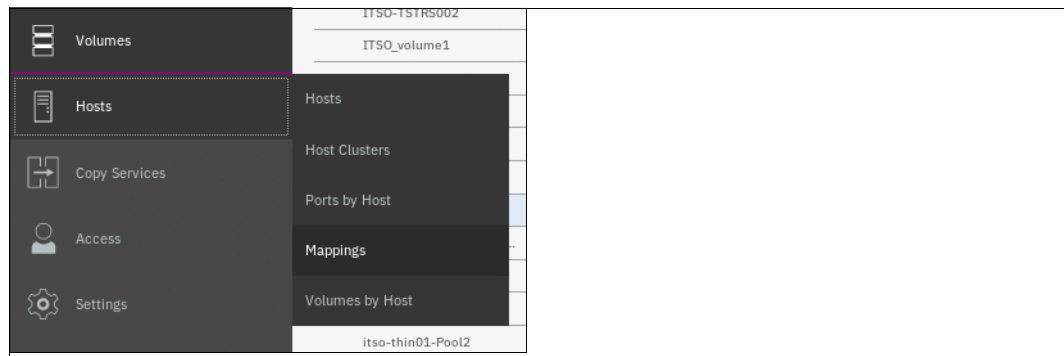


Figure 6-62 Accessing the Hosts Mapping menu

5. List of volumes mapped to all hosts is displayed, as shown in Figure 6-63.

Host Name	SCSI ID	Volume Name	UID	I/O Group ID	I/O Group Name
ITSO-VMHOST-02	2	itso-basic00-Pool0	6005076400EF0186B000000000000051	0	io_grp0
ITSO-VMHOST-02	1	Redbook_Test_1	6005076400EF0186B000000000000034	1	io_grp1
ITSO-VMHOST-02	0	Test-Volume-06	6005076400EF0186B00000000000001A	0	io_grp0

Figure 6-63 List of volume to host mappings

To see volumes that are mapped to clusters instead of hosts, change the value that is shown in the upper left corner (see Figure 6-63) from **Private Mappings** to **Shared Mappings**.

**Note:** You can use the filter to display only the wanted hosts or volumes.

The host can now access the mapped volume. For more information about discovering the volumes on the host, see Chapter 7, “Hosts” on page 351.

### 6.5.9 Migrating a volume to another storage pool

IBM Spectrum Virtualize enables online volume migration with no applications downtime. Volumes can be moved between storage pools without affecting business workloads that are running on these volumes.

There are two ways to perform volume migration: by using the volume migration feature and by creating a volume copy. These methods are described next.

#### Volume migration using migration feature

The migration process is a low priority process that does not affect the performance of the IBM Spectrum Virtualize system. However, as subsequent volume extents are moved to the new storage pool, the performance of the volume is determined more by the characteristics of the new storage pool.

**Note:** You cannot move a volume copy that is compressed to an I/O group that contains at least one node that does not support compressed volumes.

To migrate a volume to another storage pool, complete the following steps:

1. In the **Volumes** menu, highlight the volume that you want to migrate. Select **Actions** → **Migrate to Another Pool...**, as shown in Figure 6-64.

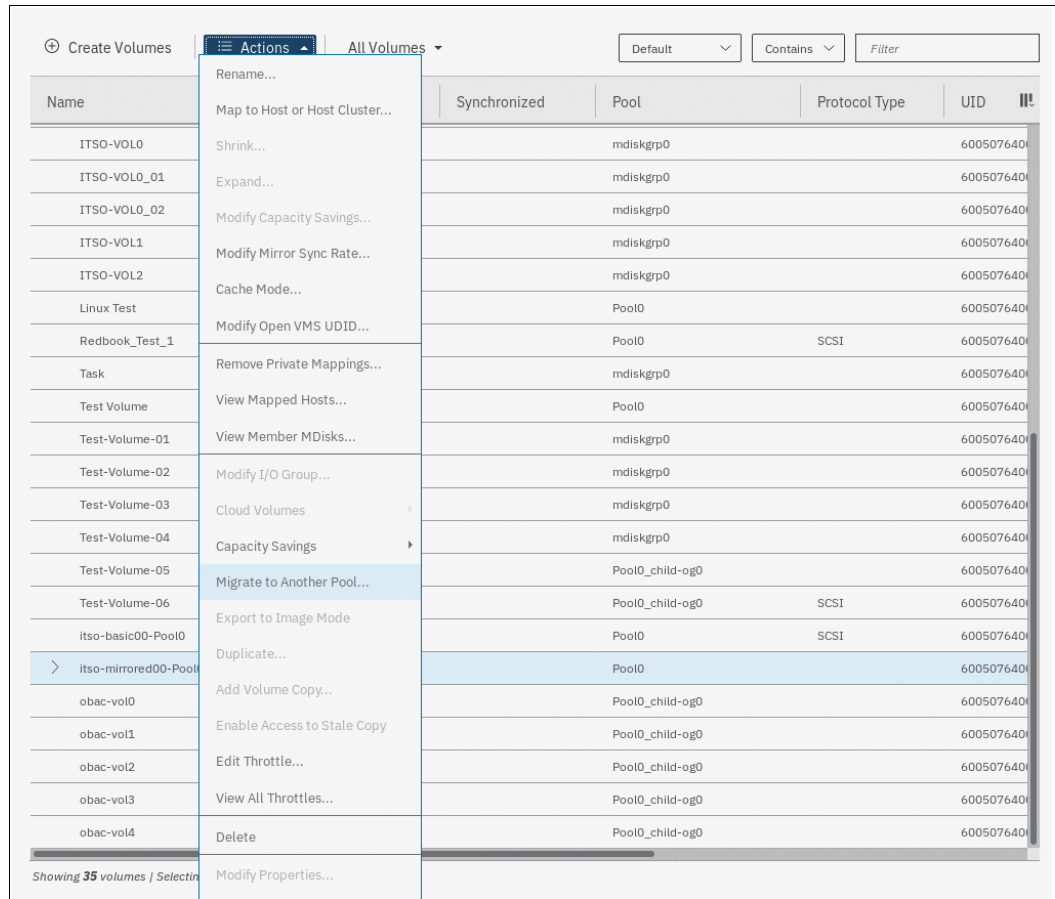


Figure 6-64 Migrate Volume Copy window: Select menu option

- The Migrate Volume Copy window opens. If your volume consists of more than one copy, select the copy that you want to migrate to another storage pool, as shown in Figure 6-65. If the selected volume consists of one copy, the volume copy selection pane is not displayed.

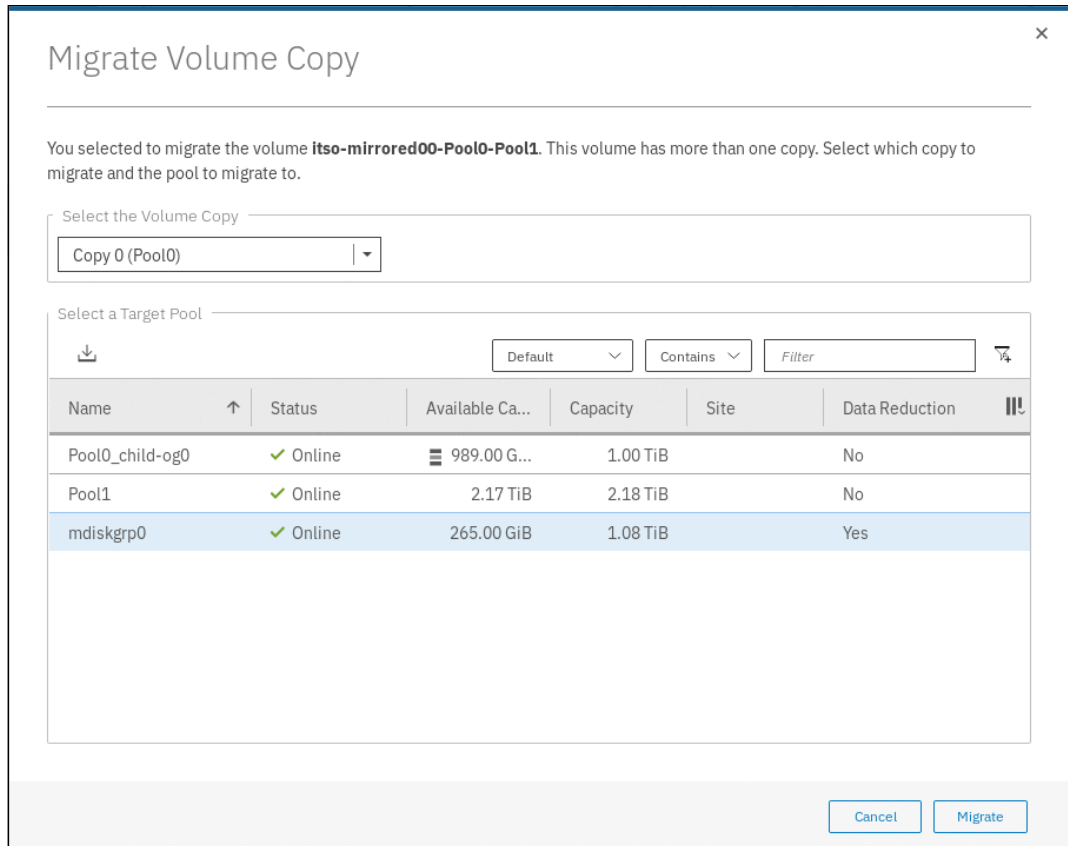


Figure 6-65 Migrate Volume Copy: Selecting the volume copy

Select the new target storage pool and click **Migrate**, as shown in Figure 6-65. The Select a Target Pool pane displays the list of all pools that are a valid migration copy target for the selected volume copy.

- You are returned to the Volumes view. The time that it takes for the migration process to complete depends on the size of the volume. The status of the migration can be monitored by selecting **Monitoring** → **Background Tasks**, as shown in Figure 6-66.

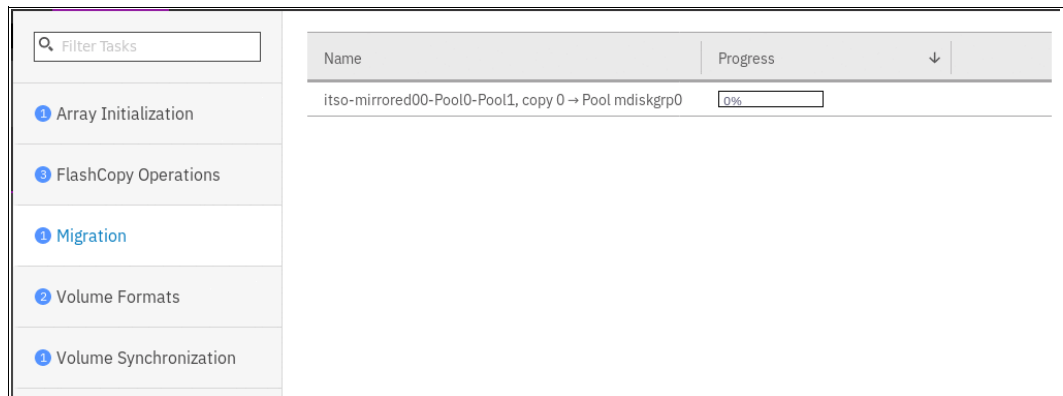


Figure 6-66 Migration progress

After the migration task completes, the completed migration task is visible in Recently Completed Task pane of the Background Tasks menu, as shown in Figure 6-67.

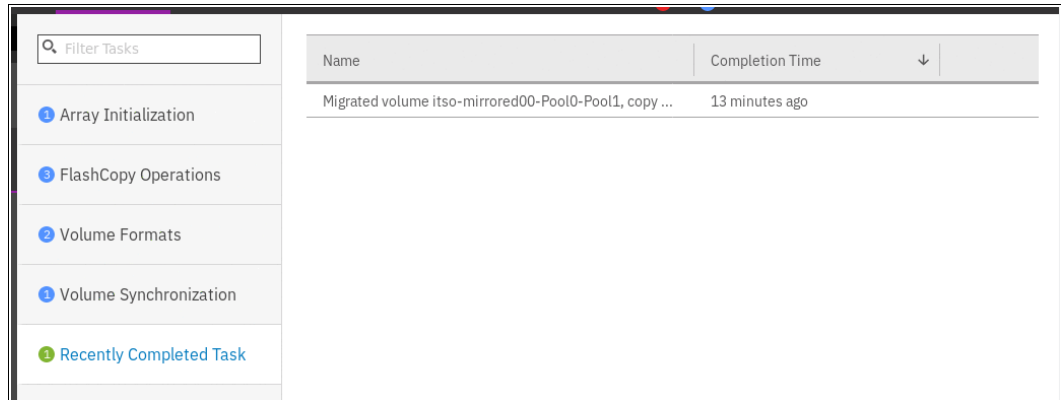


Figure 6-67 Migration complete

In the **Volumes** → **Volumes** menu, the volume copy is now displayed in the target storage pool, as shown in Figure 6-68.

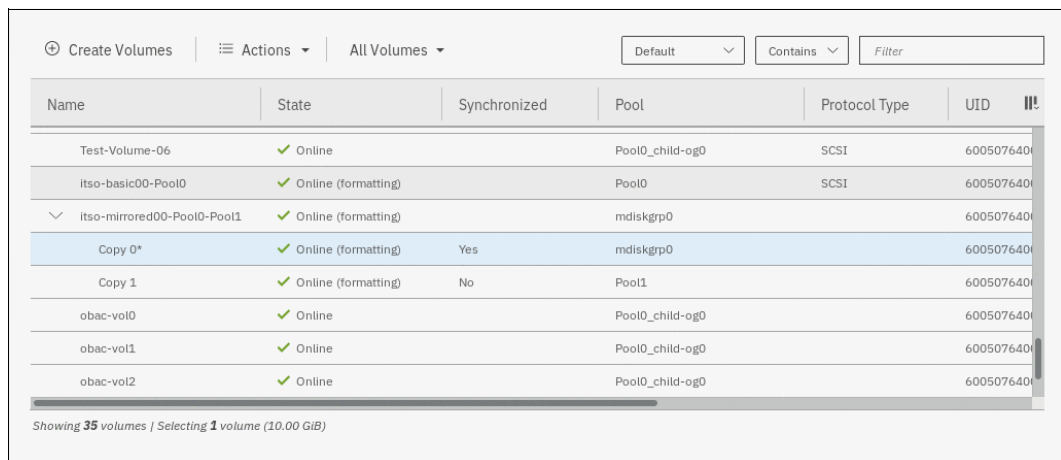


Figure 6-68 Volume copy after migration

The volume copy is now migrated without any host or application downtime to the new storage pool.

Another way to migrate single-copy volumes to another pool is to use the volume copy feature, as described in “Volume migration by adding a volume copy” on page 319.

**Note:** Migrating a volume between storage pools with different extent sizes is *not* supported. If you must migrate a volume to a storage pool with a different extent size, use the volume migration by adding a volume copy method.



## Volume migration by adding a volume copy

IBM Spectrum Virtualize supports creating, synchronizing, splitting, and deleting volume copies. A combination of these tasks can be used to migrate volumes to other storage pools.

The easiest way to migrate volume copies is to use the migration feature that is described in 6.5.9, “Migrating a volume to another storage pool” on page 315. However, in some use cases, the preferred or only method of volume migration is to create a copy of the volume in the target storage pool and then remove the old copy.

**Note:** You can specify storage efficiency characteristics of the new volume copy differently than these of the primary copy. For example, you can make a thin-provisioned copy of a standard-provisioned volume.

This volume migration option can be used only for single-copy volumes. If you must move a copy of a mirrored volume by using this method, you must delete one of the volume copies first, and then create a copy in the target storage pool. This process causes a temporary loss of redundancy while the volume copies synchronize.

To migrate a volume by using the volume copy feature, complete the following steps:

1. Select the volume that you want to move, and click **Actions** → **Add Volume Copy**, as shown in Figure 6-69.

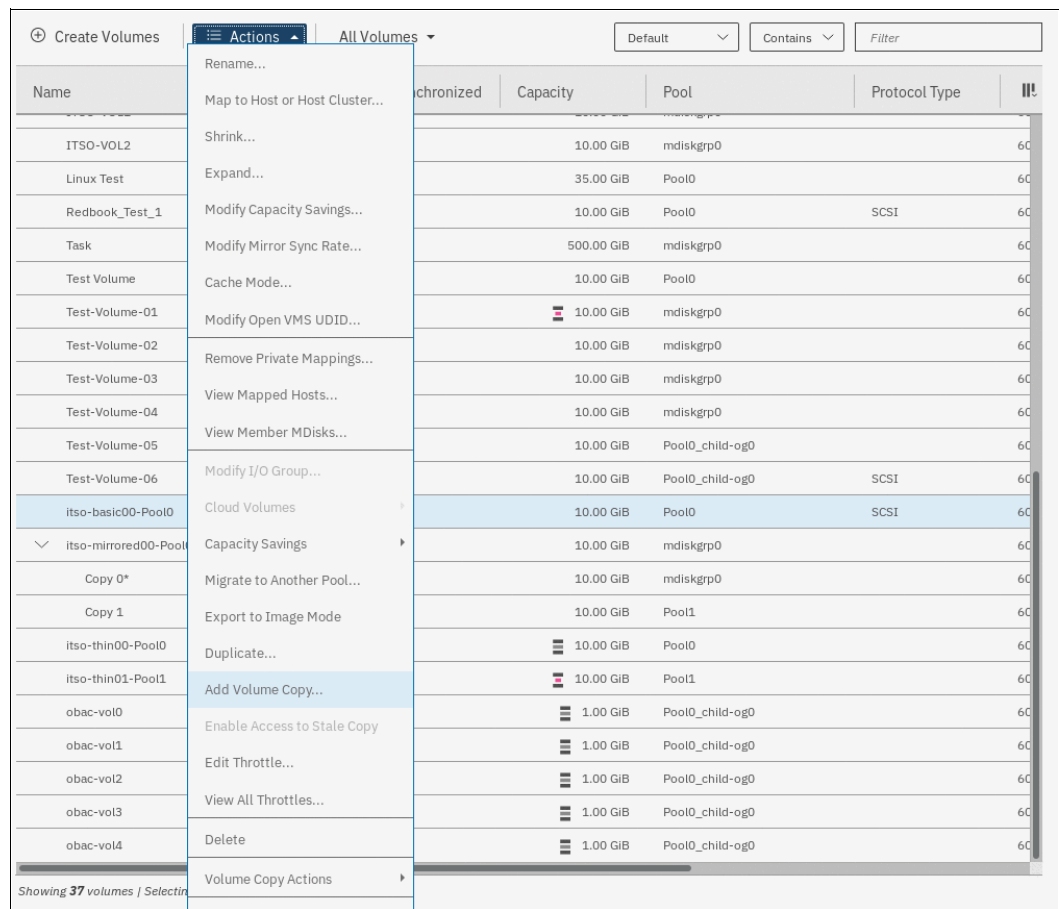


Figure 6-69 Adding the volume copy

2. Create a second copy of your volume in the target storage pool, as shown in Figure 6-70. In our example, a compressed copy of the volume is created in target pool Pool12. Click **Add**. The Deduplication option is not available if either of the volume copies is not in a DRP.

**Add Volume Copy**

Create preset volumes with copies in multiple pools but at a single site.

**Copy 0:** Pool: Pool0 Capacity Details: Total 3.24 TiB

**Copy 1:** Pool: Pool1 Capacity Details: Total 2.16 TiB

Capacity savings: None  Deduplicated

**Summary**  
 1 volume  
 2 mirrored copies  
 1 copy in pool Pool0  
 1 copy in pool Pool1

Cancel Add

Figure 6-70 Defining the new volume copy

Wait until the copies are synchronized, as shown in Figure 6-71.

Name	State	Synchronized	Capacity	Pool
itso-basic00-Pool0	✓ Online		10.00 GiB	Pool0
Copy 0*	✓ Online	Yes	10.00 GiB	Pool0
Copy 1	✓ Online	Yes	10.00 GiB	Pool1

Figure 6-71 Verifying that the volume copies are synchronized

- Change the roles of the volume copies by making the new copy the primary copy, as shown in Figure 6-72. The current primary copy is displayed with an asterisk next to its name.

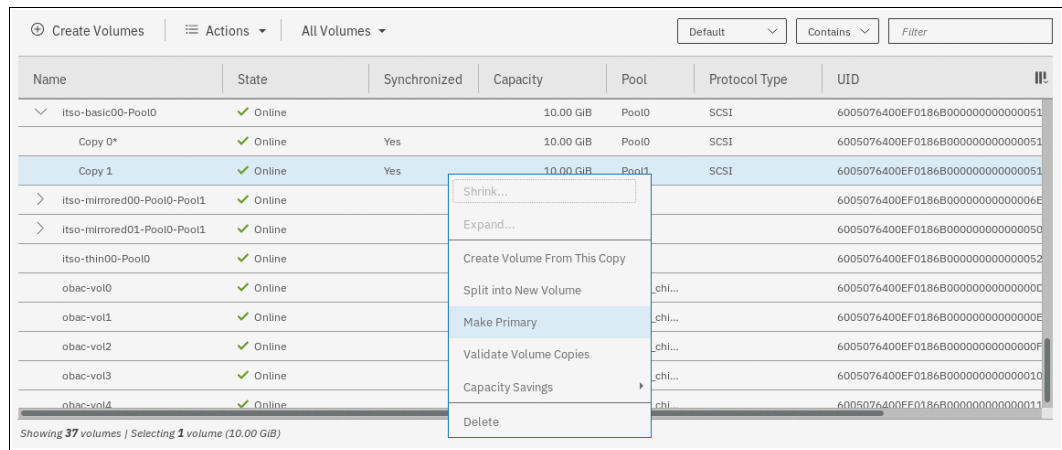


Figure 6-72 Setting the volume the copy in the target storage pool as the primary

- Split or delete the volume copy in the source pool, as shown in Figure 6-73.

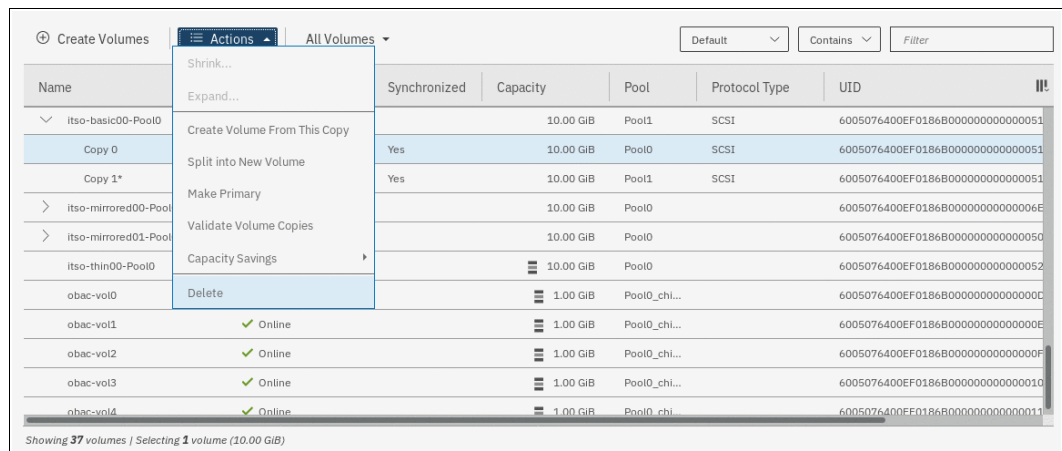


Figure 6-73 Deleting the volume copy in the source pool

- Confirm the removal of the volume copy as shown in Figure 6-74.

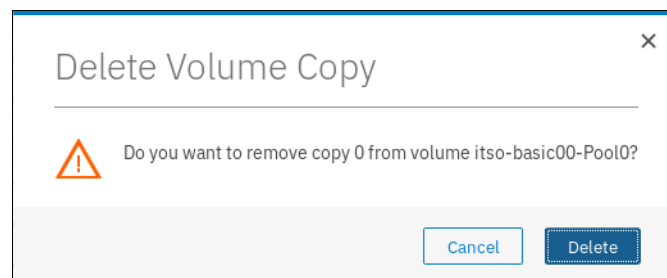


Figure 6-74 Confirm deletion of a volume copy

The Volumes view now shows that the volume has a single copy in the target pool, as shown in Figure 6-75.

Name	State	Synchronized	Capacity	Pool
test-volume-05	✓ Online		10.00 GiB	Pool0_chi...
Test-Volume-06	✓ Online		10.00 GiB	Pool0_chi...
its0-basic00-Pool0	✓ Online		10.00 GiB	Pool1

Figure 6-75 Volume copy in the target storage pool

The migration of volumes by using the volume copy feature requires more user interaction, but might be a preferred option for specific use cases. One such example is migrating a volume from a tier 1 storage pool to a lower performance tier 2 storage pool.

First, the volume copy feature can be used to create a copy in the tier 2 pool (steps 1 and 2). All reads are still performed in the tier 1 pool to the primary copy. After the volume copies are synchronized (step 3), all writes are destaged to both pools, but the reads are still only done from the primary copy.

To test the performance of the volume in the new pool switch the roles of the volume copies making the new copy the primary (step 4). If the performance is acceptable, the volume copy in tier 1 can be split or deleted. If the testing of the tier 2 pool shows unsatisfactory performance, the old copy can be made the primary again to switch volume reads back to the tier 1 copy.

## 6.6 Volume operations by using the CLI

This section describes how to perform various volume configuration and administrative tasks by using the CLI.

For more information about how to set up CLI access, see Appendix B, “CLI setup” on page 849.

### 6.6.1 Displaying volume information

Use the `lsvdisk` command to display information about all volumes that are defined within the IBM Spectrum Virtualize environment. To display more information about a specific volume, run the command again and provide the volume name or the volume ID as the command parameter, as shown in Example 6-1.

*Example 6-1 The lsvdisk command*

```
IBM_Storage:ITS0:superuser>lsvdisk -delim ' '
id name IO_group_id IO_group_name status mdisk_grp_id mdisk_grp_name capacity type FC_id
FC_name RC_id RC_name vdisk_UID fc_map_count copy_count fast_write_state se_copy_count
RC_change compressed_copy_count parent_mdisk_grp_id parent_mdisk_grp_name formatting
encrypt volume_id volume_name function
0 A_MIRRORED_VOL_1 0 io_grp0 online many many 10.00GB many
6005076400F580049800000000000002 0 2 empty 0 no 0 many many no yes 0 A_MIRRORED_VOL_1
1 COMPRESSED_VOL_1 0 io_grp0 online 1 Pool1 15.00GB striped
6005076400F580049800000000000003 0 1 empty 0 no 1 1 Pool1 no yes 1 COMPRESSED_VOL_1
2 vdisk0 0 io_grp0 online 0 Pool0 10.00GB striped 6005076400F580049800000000000004 0 1
empty 0 no 0 0 Pool0 no yes 2 vdisk0
```

```

3 THIN_PROVISION_VOL_1 0 io_grp0 online 0 Poo10 100.00GB striped
6005076400F580049800000000000005 0 1 empty 1 no 0 0 Poo10 no yes 3 THIN_PROVISION_VOL_1
4 COMPRESSED_VOL_2 0 io_grp0 online 1 Poo11 30.00GB striped
6005076400F580049800000000000006 0 1 empty 0 no 1 1 Poo11 no yes 4 COMPRESSED_VOL_2
5 COMPRESS_VOL_3 0 io_grp0 online 1 Poo11 30.00GB striped
6005076400F580049800000000000007 0 1 empty 0 no 1 1 Poo11 no yes 5 COMPRESS_VOL_3
6 MIRRORED_SYNC_RATE_16 0 io_grp0 online many many 10.00GB many
6005076400F580049800000000000008 0 2 empty 0 no 0 many many no yes 6 MIRRORED_SYNC_RATE_16
7 THIN_PROVISION_MIRRORED_VOL 0 io_grp0 online many many 10.00GB many
6005076400F580049800000000000009 0 2 empty 2 no 0 many many no yes 7
THIN_PROVISION_MIRRORED_VOL
8 Tiger 0 io_grp0 online 0 Poo10 10.00GB striped      6005076400F580049800000000000010 0 1
not_empty 0 no 0 0 Poo10 yes yes 8 Tiger
12 vdisk0_restore 0 io_grp0 online 0 Poo10 10.00GB striped
6005076400F58004980000000000000E 0 1 empty 0 no 0 0 Poo10 no yes 12 vdisk0_restore
13 vdisk0_restore1 0 io_grp0 online 0 Poo10 10.00GB striped
6005076400F58004980000000000000F 0 1 empty 0 no 0 0 Poo10 no yes 13 vdisk0_restore1

```

---

## 6.6.2 Creating a volume

Running the `mkvdisk` command creates sequential, striped, or image mode volumes. When they are mapped to a host object, these objects are seen as disk drives on which the host can perform I/O operations.

**Creating an image mode disk:** If you do not specify the `-size` parameter when you create an image mode disk, the entire MDisk capacity is used.

You must know the following information before you start to create the volume:

- ▶ In which storage pool the volume has its extents
- ▶ From which I/O Group the volume is accessed
- ▶ Which IBM Spectrum Virtualize node is the preferred node for the volume
- ▶ Size of the volume
- ▶ Name of the volume
- ▶ Type of the volume
- ▶ Whether this volume is to be managed by IBM Easy Tier to optimize its performance

When you are ready to create your striped volume, run the `mkvdisk` command. The command that is shown in Example 6-2 creates a 10 GB striped volume with volume ID 8 within the storage pool `Poo10` and assigns it to the I/O group `io_grp0`. Its preferred node is node 1.

*Example 6-2 The `mkvdisk` command*

```

IBM_Storwize:ITS0:superuser>mkvdisk -mdiskgrp Poo10 -iogrp io_grp0 -size 10 -unit gb -name
Tiger
Virtual Disk, id [8], successfully created

```

---

To verify the results, run the `lsvdisk` command and provide the volume ID as the command parameter, as shown in Example 6-3.

*Example 6-3 The `lsvdisk` command*

```

IBM_Storwize:ITS0:superuser>lsvdisk 8
id 8
name Tiger
IO_group_id 0
IO_group_name io_grp0

```

```

status online
mdisk_grp_id 0
mdisk_grp_name Pool0
capacity 10.00GB
type striped
formatted no
formatting yes
mdisk_id
mdisk_name
FC_id
FC_name
RC_id
RC_name
vdisk_UID 6005076400F580049800000000000010
preferred_node_id 2
fast_write_state not_empty
cache readwrite
udid
fc_map_count 0
sync_rate 50
copy_count 1
se_copy_count 0
File system
mirror_write_priority latency
RC_change no
compressed_copy_count 0
access_IO_group_count 1
last_access_time
parent_mdisk_grp_id 0
parent_mdisk_grp_name Pool0
owner_type none
owner_id
owner_name
encrypt yes
volume_id 8
volume_name Tiger
function
throttle_id
throttle_name
IOPs_limit
bandwidth_limit_MB
volume_group_id
volume_group_name
cloud_backup_enabled no
cloud_account_id
cloud_account_name
backup_status off
last_backup_time
restore_status none
backup_grain_size
deduplicated_copy_count 0

copy_id 0
status online
sync yes
auto_delete no
primary yes
mdisk_grp_id 0
mdisk_grp_name Pool0
type striped

```

```
mdisk_id
mdisk_name
fast_write_state not_empty
used_capacity 10.00GB
real_capacity 10.00GB
free_capacity 0.00MB
overallocation 100
autoexpand
warning
grainsize
se_copy no
easy_tier on
easy_tier_status balanced
tier tier0_flash
tier_capacity 0.00MB
tier tier1_flash
tier_capacity 0.00MB
tier tier_enterprise
tier_capacity 0.00MB
tier tier_nearline
tier_capacity 10.00GB
compressed_copy no
uncompressed_used_capacity 10.00GB
parent_mdisk_grp_id 0
parent_mdisk_grp_name Poo10
encrypt yes
deduplicated_copy no
used_capacity_before_reduction
0.00MB
```

---

The required tasks to create a volume are complete.

### 6.6.3 Creating a thin-provisioned volume

Example 6-4 shows an example of creating a thin-provisioned volume. The following parameters must be specified:

- |                    |   |
|--------------------|---|
| <b>-rsize</b>      | This parameter makes the volume a thin-provisioned volume. If this parameter is missing, the volume is created as standard-provisioned.                               |
| <b>-autoexpand</b> | This parameter specifies that thin-provisioned volume copies automatically expand their real capacities by allocating new extents from their storage pool (optional). |
| <b>-grainsize</b>  | This parameter sets the grain size in kilobytes (KB) for a thin-provisioned volume (optional).  |

*Example 6-4 Using the command mkvdisk*

---

```
IBM_Storwize:ITS0:superuser>mkvdisk -mdiskgrp Poo10 -iogrp 0 -vtype striped -size 10 -unit
gb -rsize 50% -autoexpand -grainsize 256
Virtual Disk, id [9], successfully created
```

---

This command creates a thin-provisioned volume with 10 GB of virtual capacity. The volume belongs to the storage pool that is named Site1\_Pool and is owned by input/output (I/O) Group io\_grp0. The real capacity automatically expands until the real volume size of 10 GB is reached. The grain size is set to 256 KB, which is the default.

**Disk size:** When the `-rsize` parameter is used to specify the real physical capacity of a thin-provisioned volume, the following options are available to specify the physical capacity: `disk_size`, `disk_size_percentage`, and `auto`.

Use the `disk_size_percentage` option to define initial real capacity by using a percentage of the disk's virtual capacity that is defined by the `-size` parameter. This option takes as a parameter an integer, or an integer that is immediately followed by the percent (%) symbol.

Use `disk_size` option to directly specify the real physical capacity by specifying its size in the units defined by using the `-unit` parameter (the default unit is MB). The `-rsize` value can be greater than, equal to, or less than the size of the volume.

The `auto` option creates a volume copy that uses the entire size of the MDisk. If you specify the `-rsize auto` option, you must also specify the `-vtype image` option.

An entry of 1 GB uses 1,024 MB.

## 6.6.4 Creating a volume in image mode

Use image mode volume to bring a non-virtualized disk (for example, from a pre-virtualization environment) under the control of the IBM Spectrum Virtualize system. After it is managed by the system, you can migrate the volume to the standard MDisk.

When an image mode volume is created, it directly maps to the thus far unmanaged MDisk from which it is created. Therefore, except for a thin-provisioned image mode volume, the volume's LBA  $x$  equals MDisk LBA  $x$ .

**Size:** An image mode volume must be at least 512 bytes (the capacity cannot be 0) and always occupies at least one extent.

You must use the `-mdisk` parameter to specify an MDisk that has a mode of unmanaged. The `-fntdisk` parameter cannot be used to create an image mode volume.

**Capacity:** If you create a mirrored volume from two image mode MDisk without specifying a `-capacity` value, the capacity of the resulting volume is the smaller of the two MDisk. The remaining space on the larger MDisk is inaccessible.

If you do not specify the `-size` parameter when you create an image mode disk, the entire MDisk capacity is used.

Run the `mkvdisk` command to create an image mode volume, as shown in Example 6-5.

*Example 6-5 The `mkvdisk` (image mode) command*

---

```
IBM_2145:ITS0_CLUSTER:superuser>mkvdisk -mdiskgrp ITS0_Pool1 -iogrp 0 -mdisk mdisk25 -vtype
image -name Image_Volume_A
Virtual Disk, id [6], successfully created
```

---

As shown in Example 6-5, an image mode volume that is named `Image_Volume_A` is created that uses the `mdisk25` MDisk. The MDisk is moved to the storage pool `ITS0_Pool1`, and the volume is owned by the I/O group `io_grp0`.



If you run the `lsvdisk` command, it shows volume named `Image_Volume_A` with type `image`, as shown in Example 6-6.

*Example 6-6 The lsvdisk command*

---

```

IBM_2145:ITSO_CLUSTER:superuser>lsvdisk -filtervalue type=image
id name          IO_group_id IO_group_name status mdisk_grp_id mdisk_grp_name capacity
type FC_id FC_name RC_id RC_name vdisk_UID          fc_map_count copy_count
fast_write_state se_copy_count RC_change compressed_copy_count parent_mdisk_grp_id
parent_mdisk_grp_name formatting encrypt volume_id volume_name function
6 Image_Volume_A 0          io_grp0      online 5          ITSO_Pool1    1.00GB
image
          6005076801FE80840800000000000021 0          1
empty          0          no          no          0          5
ITSO_Pool1          no          no          6          Image_Volume_A

```

---

### 6.6.5 Adding a volume copy

You can add a copy to a volume. If volume copies are defined on different MDisks, the volume remains accessible, even when the MDisk on which one of its copies depends becomes unavailable. You can also create a copy of a volume on a dedicated MDisk by creating an image mode copy of the volume. Although volume copies can increase the availability of data, they are not separate objects.

Volume mirroring can be also used as an alternative method of migrating volumes between storage pools.

To create a copy of a volume, run the `addvdiskcopy` command. This command creates a copy of the chosen volume in the specified storage pool, which changes a non-mirrored volume into a mirrored one.

The following scenario shows the how to create a copy of a volume in a different storage pool. As shown in Example 6-7, the volume initially has a single copy with `copy_id 0` that is provisioned in pool `Pool0`.

*Example 6-7 The lsvdisk command*

---

```

IBM_Storwize:ITSO:superuser>lsvdisk 2
id 2
name vdisk0
IO_group_id 0
IO_group_name io_grp0
status online
mdisk_grp_id 0
mdisk_grp_name Pool0
capacity 10.00GB
type striped
formatted yes
formatting no
mdisk_id
mdisk_name
FC_id
FC_name
RC_id
RC_name
vdisk_UID 6005076400F580049800000000000004
preferred_node_id 2
fast_write_state empty
cache readonly
udid

```

```

fc_map_count 0
sync_rate 50
copy_count 1
se_copy_count 0
File system
mirror_write_priority latency
RC_change no
compressed_copy_count 0
access_IO_group_count 1
last_access_time
parent_mdisk_grp_id 0
parent_mdisk_grp_name Pool0
owner_type none
owner_id
owner_name
encrypt yes
volume_id 2
volume_name vdisk0
function
throttle_id
throttle_name
IOPs_limit
bandwidth_limit_MB
volume_group_id
volume_group_name
cloud_backup_enabled no
cloud_account_id
cloud_account_name
backup_status off
last_backup_time
restore_status none
backup_grain_size
deduplicated_copy_count 0

copy_id 0
status online
sync yes
auto_delete no
primary yes
mdisk_grp_id 0
mdisk_grp_name Pool0
type striped
mdisk_id
mdisk_name
fast_write_state empty
used_capacity 10.00GB
real_capacity 10.00GB
free_capacity 0.00MB
overallocation 100
autoexpand
warning
grainsize
se_copy no
easy_tier on
easy_tier_status balanced
tier tier0_flash
tier_capacity 0.00MB
tier tier1_flash
tier_capacity 0.00MB
tier tier_enterprise

```

```
tier_capacity 0.00MB
tier tier_nearline
tier_capacity 10.00GB
compressed_copy no
uncompressed_used_capacity 10.00GB
parent_mdisk_grp_id 0
parent_mdisk_grp_name Pool0
encrypt yes
deduplicated_copy no
used_capacity_before_reduction 0.00MB
```

---

Example 6-8 shows adding the second volume copy by running the **addvdiskcopy** command.

*Example 6-8 The addvdiskcopy command*

---

```
IBM_Storwize:ITS0:superuser>addvdiskcopy -mdiskgrp Pool1 -vtype striped -unit gb vdisk0
Vdisk [2] copy [1] successfully created
```

---

During the synchronization process, you can see the status by running the **lsvdiskssyncprogress** command.

As shown in Example 6-9, the first time that the status is checked, the synchronization progress is at 48%, and the estimated completion time is 161026203918. The estimated completion time is displayed in the YYMMDDHHMMSS format. In our example, it is 2016, Oct-26 20:39:18. When the command is run again, the progress status is at 100%, and the synchronization is complete.

*Example 6-9 Synchronization*

---

```
IBM_Storwize:ITS0:superuser>lsvdiskssyncprogress
vdisk_id vdisk_name copy_id progress estimated_completion_time
2        vdisk0      1        0        171018232305
IBM_Storwize:ITS0:superuser>lsvdiskssyncprogress
vdisk_id vdisk_name copy_id progress estimated_completion_time
2        vdisk0      1        100
```

---

As shown in Example 6-10, the new volume copy (copy\_id 1) was added and appears in the output of the **lsvdisk** command.

*Example 6-10 The lsvdisk command*

---

```
IBM_Storwize:ITS0:superuser>lsvdisk vdisk0
id 2
name vdisk0
IO_group_id 0
IO_group_name io_grp0
status online
mdisk_grp_id many
mdisk_grp_name many
capacity 10.00GB
type many
formatted yes
formatting no
mdisk_id many
mdisk_name many
FC_id
FC_name
RC_id
RC_name
```

```

vdisk_UID 6005076400F580049800000000000004
preferred_node_id 2
fast_write_state empty
cache readonly
udid
fc_map_count 0
sync_rate 50
copy_count 2
se_copy_count 0
File system
mirror_write_priority latency
RC_change no
compressed_copy_count 0
access_IO_group_count 1
last_access_time
parent_mdisk_grp_id many
parent_mdisk_grp_name many
owner_type none
owner_id
owner_name
encrypt yes
volume_id 2
volume_name vdisk0
function
throttle_id
throttle_name
IOPs_limit
bandwidth_limit_MB
volume_group_id
volume_group_name
cloud_backup_enabled no
cloud_account_id
cloud_account_name
backup_status off
last_backup_time
restore_status none
backup_grain_size
deduplicated_copy_count 0

copy_id 0
status online
sync yes
auto_delete no
primary yes
mdisk_grp_id 0
mdisk_grp_name Pool0
type striped
mdisk_id
mdisk_name
fast_write_state empty
used_capacity 10.00GB
real_capacity 10.00GB
free_capacity 0.00MB
overallocation 100
autoexpand
warning
grainsize
se_copy no
easy_tier on
easy_tier_status balanced

```

```
tier tier0_flash
tier_capacity 0.00MB
tier tier1_flash
tier_capacity 0.00MB
tier tier_enterprise
tier_capacity 0.00MB
tier tier_nearline
tier_capacity 10.00GB
compressed_copy no
uncompressed_used_capacity 10.00GB
parent_mdisk_grp_id 0
parent_mdisk_grp_name Pool0
encrypt yes
deduplicated_copy no
used_capacity_before_reduction 0.00MB
```

```
copy_id 1
status online
sync yes
auto_delete no
primary no
mdisk_grp_id 1
mdisk_grp_name Pool1
type striped
mdisk_id
mdisk_name
fast_write_state empty
used_capacity 10.00GB
real_capacity 10.00GB
free_capacity 0.00MB
overallocation 100
autoexpand
warning
grainsize
se_copy no
easy_tier on
easy_tier_status balanced
tier tier0_flash
tier_capacity 0.00MB
tier tier1_flash
tier_capacity 0.00MB
tier tier_enterprise
tier_capacity 0.00MB
tier tier_nearline
tier_capacity 10.00GB
compressed_copy no
uncompressed_used_capacity 10.00GB
parent_mdisk_grp_id 1
parent_mdisk_grp_name Pool1
encrypt yes
deduplicated_copy no
used_capacity_before_reduction 0.00MB
```

---

When adding a volume copy, you can define it with different parameters than the original volume copy. For example, you can create a thin-provisioned copy of a standard-provisioned volume to migrate a thick-provisioned volume to a thin-provisioned volume. The migration also can be done in the opposite direction.

**Volume copy mirror parameters:** To change the parameters of a volume copy, you must delete the volume copy and redefine it with the new values.

In Example 6-11, the volume name is changed from VOL\_NO\_MIRROR to VOL\_WITH\_MIRROR.

*Example 6-11 Volume name changes*

---

```
IBM_Storwize:ITS0:superuser>chvdisk -name VOL_WITH_MIRROR VOL_NO_MIRROR
IBM_Storwize:ITS0:superuser>
```

---

## Using -autodelete flag to migrate a volume

This section shows the use of the **addvdiskcopy** command with the **-autodelete** flag set. The **-autodelete** flag causes the primary copy to be deleted after the secondary copy is synchronized.

Example 6-12 shows a shortened **lsvdisk** output for an uncompressed volume with a single volume copy.

*Example 6-12 An uncompressed volume*

---

```
IBM_Storwize:ITS0:superuser>lsvdisk UNCOMPRESSED_VOL
id 9
name UNCOMPRESSED_VOL
IO_group_id 0
IO_group_name io_grp0
status online
...

copy_id 0
status online
sync yes
auto_delete no
primary yes
...
compressed_copy no
...
```

---

Example 6-13 adds a compressed copy with the **-autodelete** flag set.

*Example 6-13 Compressed copy*

---

```
IBM_Storwize:ITS0:superuser>addvdiskcopy -autodelete -rsize 2 -mdiskgrp 0 -compressed
UNCOMPRESSED_VOL
Vdisk [9] copy [1] successfully created
```

---

Example 6-14 shows the **lsvdisk** output with another compressed volume (copy 1) and volume copy 0 being set to **auto\_delete yes**.

*Example 6-14 The lsvdisk command output*

---

```
IBM_Storwize:ITS0:superuser>lsvdisk UNCOMPRESSED_VOL
id 9
name UNCOMPRESSED_VOL
IO_group_id 0
IO_group_name io_grp0
status online
...
compressed_copy_count 2
```

---

```
...  
  
copy_id 0  
status online  
sync yes  
auto_delete yes  
primary yes  
...  
  
copy_id 1  
status online  
sync no  
auto_delete no  
primary no  
...  
...
```

---

When copy 1 is synchronized, copy 0 is deleted. You can monitor the progress of volume copy synchronization by running the `lsvdisk syncprogress` command.

## 6.6.6 Splitting a mirrored volume

Running the `splitvdiskcopy` command creates an independent volume in the specified I/O group from a volume copy of the specified mirrored volume. In effect, the command changes a volume with two copies into two independent volumes, each with a single copy.

If the copy that you are splitting is not synchronized, you must use the `-force` parameter. If you are attempting to remove the only synchronized copy of the source volume, the command fails. However, you can run the command when either copy of the source volume is offline.

Example 6-15 shows the `splitvdiskcopy` command, which is used to split a mirrored volume. It creates a volume named `SPLIT_VOL` from copy with ID 1 of the volume named `VOLUME_WITH_MIRRORED_COPY`.

### *Example 6-15 Split volume*

---

```
IBM_Storwize:ITS0:superuser>splitvdiskcopy -copy 1 -iogrp 0 -name SPLIT_VOL  
VOLUME_WITH_MIRRORED_COPY  
Virtual Disk, id [1], successfully created
```

---

As you can see in Example 6-16, the new volume is created as an independent volume.

### *Example 6-16 The lsvdisk command*

---

```
IBM_Storwize:ITS0:superuser>lsvdisk SPLIT_VOL  
id 1  
name SPLIT_VOL  
IO_group_id 0  
IO_group_name io_grp0  
status online  
mdisk_grp_id 1  
mdisk_grp_name Pool1  
capacity 10.00GB  
type striped  
formatted yes  
formatting no  
mdisk_id  
mdisk_name  
FC_id
```

```

FC_name
RC_id
RC_name
vdisk_UID 6005076400F580049800000000000012
preferred_node_id 1
fast_write_state empty
cache readwrite
udid
fc_map_count 0
sync_rate 50
copy_count 1
se_copy_count 0
File system
mirror_write_priority latency
RC_change no
compressed_copy_count 0
access_IO_group_count 1
last_access_time
parent_mdisk_grp_id 1
parent_mdisk_grp_name Pool1
owner_type none
owner_id
owner_name
encrypt yes
volume_id 1
volume_name SPLIT_VOL
function
throttle_id
throttle_name
IOPs_limit
bandwidth_limit_MB
volume_group_id
volume_group_name
cloud_backup_enabled no
cloud_account_id
cloud_account_name
backup_status off
last_backup_time
restore_status none
backup_grain_size
deduplicated_copy_count 0

copy_id 0
status online
sync yes
auto_delete no
primary yes
mdisk_grp_id 1
mdisk_grp_name Pool1
type striped
mdisk_id
mdisk_name
fast_write_state empty
used_capacity 10.00GB
real_capacity 10.00GB
free_capacity 0.00MB
overallocation 100
autoexpand
warning
grainsize

```



```
se_copy no
easy_tier on
easy_tier_status balanced
tier tier0_flash
tier_capacity 0.00MB
tier tier1_flash
tier_capacity 0.00MB
tier tier_enterprise
tier_capacity 0.00MB
tier tier_nearline
tier_capacity 10.00GB
compressed_copy no
uncompressed_used_capacity 10.00GB
parent_mdisk_grp_id 1
parent_mdisk_grp_name Pool1
encrypt yes
deduplicated_copy no
used_capacity_before_reduction 0.00MB
```

---

## 6.6.7 Modifying a volume

Running the **chvdisk** command modifies a single property of a volume. Only one property can be modified at a time. Therefore, changing the volume name and modifying its I/O group requires two invocations of the command.

**Tips:** Changing the I/O group with which this volume is associated requires a flush of the cache within the nodes in the current I/O group to ensure that all data is written to disk. I/O must be suspended at the host level before you perform this operation.

If the volume has a mapping to any hosts, it is impossible to move the volume to an I/O group that does not include any of those hosts.

This operation fails if insufficient space exists to allocate bitmaps for a mirrored volume in the target I/O group.

If the **-force** parameter is used and the system is unable to destage all write data from the cache, the contents of the volume are corrupted by the loss of the cached data.

If the **-force** parameter is used to move a volume that has out-of-sync copies, a full resynchronization is required.

## 6.6.8 Deleting a volume

Run the **rmvdisk** command to delete a volume. When this command is run on a managed mode volume, any data on the volume is lost, and the extents that made up this volume are returned to the pool of free extents in the storage pool.

If any RC, IBM FlashCopy, or host mappings still exist for the target of **rmvdisk** command, the delete fails unless the **-force** flag is specified. This flag causes the deletion of the volume and any volume to host mappings and copy mappings.

If the volume is migrated to image mode, the delete fails unless the **-force** flag is specified. Use of the **-force** flag halts the migration and then deletes the volume.

If the command succeeds (without the **-force** flag) for an image mode volume, the write cache data is flushed to the storage before the volume is removed. Therefore, the underlying LU is consistent with the disk state from the point of view of the host by using the image-mode volume (crash-consistent file system). If the **-force** flag is used, consistency is not ensured; that is, data the host believes to be written might not be present on the LU.

If any non-destaged data exists in the fast write cache for the target of **rmvdisk** command, the deletion of the volume fails unless the **-force** flag is specified, in which case, any non-destaged data in the fast write cache is deleted.

Example 6-17 shows how to use the **rmvdisk** command to delete a volume from your IBM Spectrum Virtualize configuration.

*Example 6-17 The rmvdisk command*

---

```
IBM_2145:ITS0_CLUSTER:superuser>rmvdisk volume_A
```

---

This command deletes the `volume_A` volume from the IBM Spectrum Virtualize configuration. If the volume is assigned to a host, you must use the **-force** flag to delete the volume, as shown in Example 6-18.

*Example 6-18 The rmvdisk -force command*

---

```
IBM_2145:ITS0_CLUSTER:superuser>rmvdisk -force volume_A
```

---

## 6.6.9 Volume protection

To prevent active volumes or host mappings from being deleted inadvertently, the system supports a global setting that prevents these objects from being deleted if the system detects recent I/O activity to these objects.

Run the **chsystem** command to set the time interval for which the volume must be idle before it can be deleted from the system. This setting affects the following commands:

- ▶ **rmvdisk**
- ▶ **rmvolume**
- ▶ **rmvdiskcopy**
- ▶ **rmvdiskhostmap**
- ▶ **rmmdiskgrp**
- ▶ **rmhostiogr**
- ▶ **rmhost**
- ▶ **rmhostport**

These commands fail unless the volume was idle for the specified interval or the **-force** parameter was used.

To enable volume protection by setting the required inactivity interval, run the following command:

```
svctask chsystem -vdiskprotectionenabled yes -vdiskprotectiontime 60
```

The **-vdiskprotectionenabled yes** parameter enables volume protection and the **-vdiskprotectiontime** parameter specifies for how long a volume must be inactive (in minutes) before it can be deleted. In this example, volumes can be deleted only if they were inactive for over 60 minutes.

To disable volume protection, run the following command:

```
svctask chsystem -vdiskprotectionenabled no
```

## 6.6.10 Expanding a volume

Expanding a volume presents a larger capacity disk to your operating system. Although this expansion can be easily performed by using IBM Spectrum Virtualize, you must ensure that your operating system supports expansion before this function is used.

Assuming that your operating system supports expansion, you can run the **expandvdisksize** command to increase the capacity of a volume, as shown in Example 6-19.

*Example 6-19 The expandvdisksize command*

---

```
IBM_2145:ITS0_CLUSTER:superuser>expandvdisksize -size 5 -unit gb volume_C
```

---

This command expands the `volume_C` volume (which was 35 GB) by another 5 GB to give it a total size of 40 GB.

To expand a thin-provisioned volume, you can use the **-rsize** option, as shown in Example 6-20. This command changes the real size of the `volume_B` volume to a real capacity of 55 GB. The capacity of the volume is unchanged.

*Example 6-20 The lsvdisk command*

---

```
IBM_Storwize:ITS0:superuser>lsvdisk volume_B
id 26
capacity 100.00GB
type striped
.
.
copy_id 0
status online
used_capacity 0.41MB
real_capacity 50.02GB
free_capacity 50.02GB
overallocation 199
autoexpand on
warning 80
grainsize 32
se_copy yes
```

```
IBM_Storwize:ITS0:superuser>expandvdisksize -rsize 5 -unit gb volume_B
```

```
IBM_Storwize:ITS0:superuser>lsvdisk volume_B
id 26
name volume_B
capacity 100.00GB
type striped
.
.
copy_id 0
status online
used_capacity 0.41MB
real_capacity 55.02GB
free_capacity 55.02GB
overallocation 181
autoexpand on
warning 80
```

```
grainsize 32
se_copy yes
```

---

**Important:** If a volume is expanded, its type becomes striped, even if it was previously sequential or in image mode.

If not enough extents are available to expand your volume to the specified size, the following error message is displayed:

```
CMMVC5860E The action failed because there were not enough extents in the
storage pool.
```

### 6.6.11 HyperSwap volume modification with CLI

The following new CLI commands for administering volumes were released in IBM Spectrum Virtualize V7.6. However, the GUI uses the new commands only for HyperSwap volume creation (**mkvolume**) and deletion (**rmvolume**):

- ▶ **mkvolume**
- ▶ **mkimagevolume**
- ▶ **addvolumecopy\***
- ▶ **rmvolumecopy\***
- ▶ **rmvolume**

In addition, the **lsvdisk** output shows more fields: **volume\_id**, **volume\_name**, and **function**, which help to identify the individual VDisks that make up a HyperSwap volume. This information is used by the GUI to provide views that reflect the client's view of the *HyperSwap* volume and its site-dependent copies, as opposed to the "low-level" VDisks and VDisk Change Volumes.

The following individual commands are related to HyperSwap:

- ▶ **mkvolume**

Creates an empty volume by using storage from a storage pool. The type of volume that is created is determined by the system topology and the number of storage pools specified. Volume is always formatted (zeroed). The **mkvolume** command can be used to create the following objects:

- Basic volume: Any topology
- Mirrored volume: Standard topology
- Stretched volume: Stretched topology
- HyperSwap volume: HyperSwap topology

- ▶ **rmvolume**

Removes a volume. For a HyperSwap volume, this process includes deleting the active-active relationship and the change volumes.

The **-force** parameter that is used by **rmvdisk** is replaced by a set of override parameters, one for each operation-stopping condition, which makes it clearer to the user exactly what protection they are bypassing.

- ▶ **mkimagevolume**

Creates an image mode volume. This command can be used to import a volume, which preserves data. It can be implemented as a separate command to provide greater differentiation between the action of creating an empty volume and creating a volume by importing data on an MDisk.

► `addvolumecopy`

Adds a copy to a volume. The new copy is always synchronized from the existing copy. For stretched and HyperSwap topology systems, this command creates an HA volume. This command can be used to create the following volume types:

- Mirrored volume: Standard topology
- Stretched volume: Stretched topology
- HyperSwap volume: HyperSwap topology

► `rmvolumecopy`

Removes a copy of a volume. This command leaves the volume intact. It also converts a Mirrored, Stretched, or HyperSwap volume to a basic volume. For a HyperSwap volume, this command includes deleting the active-active relationship and the change volumes.

This command enables a copy to be identified by its site.

The `-force` parameter that is used by `rmvdi skcopy` is replaced by a set of override parameters, one for each operation-stopping condition, which makes it clearer to the user exactly what protection they are bypassing.

## 6.6.12 Mapping a volume to a host

Run the `mkvdi skhostmap` command to map a volume to a host. This mapping makes the volume available to the host for I/O operations. A host can perform I/O operations only on volumes that are mapped to it.

When the host bus adapter (HBA) on the host scans for devices that are attached to it, the HBA discovers all of the volumes that are mapped to its FC ports and their SCSI identifiers (SCSI LUN IDs).

For example, the first disk that is found is generally SCSI LUN 1. You can control the order in which the HBA discovers volumes by assigning the SCSI LUN ID as required. If you do not specify a SCSI LUN ID when mapping a volume to the host, the storage system automatically assigns the next available SCSI LUN ID, based on any mappings that exist with that host.

Example 6-21 shows how to map volumes `volume_B` and `volume_C` to defined host `Almaden` by running the `mkvdi skhostmap` command.

*Example 6-21 The `mkvdiskhostmap` command*

---

```
IBM_Storwize:ITS0:superuser>mkvdiskhostmap -host Almaden volume_B
Virtual Disk to Host map, id [0], successfully created
IBM_Storwize:ITS0:superuser>mkvdiskhostmap -host Almaden volume_C
Virtual Disk to Host map, id [1], successfully created
```

---

Example 6-22 shows the output of the `lshostvdiskmap` command, which shows that the volumes are mapped to the host.

*Example 6-22 The `lshostvdiskmap -delim` command*

---

```
IBM_2145:ITS0_CLUSTER:superuser>lshostvdiskmap -delim :
id:name:SCSI_id:vdisk_id:vdisk_name:vdisk_UID
2:Almaden:0:26:volume_B:6005076801AF813F1000000000000020
2:Almaden:1:27:volume_C:6005076801AF813F1000000000000021
```

---

**Assigning a specific LUN ID to a volume:** The optional `-scsi scsi_lun_id` parameter can help assign a specific LUN ID to a volume that is to be associated with a host. The default (if nothing is specified) is to assign the next available ID based on current volume mapped to the host.

Certain HBA device drivers stop when they find a gap in the sequence of SCSI LUN IDs, as shown in the following examples:

- ▶ Volume 1 is mapped to Host 1 with SCSI LUN ID 1.
- ▶ Volume 2 is mapped to Host 1 with SCSI LUN ID 2.
- ▶ Volume 3 is mapped to Host 1 with SCSI LUN ID 4.

When the device driver scans the HBA, it might stop after discovering volumes 1 and 2 because no SCSI LUN is mapped with ID 3.

**Important:** Ensure that the SCSI LUN ID allocation is contiguous.

If you use host clusters, run the `mkvolumehostclustermap` command to map a volume to a host cluster instead (see Example 6-23).

*Example 6-23 The `mkvolumehostclustermap` command*

---

```
BM_Storwize:ITS0:superuser>mkvolumehostclustermap -hostcluster vmware_cluster
UNCOMPRESSED_VOL
Volume to Host Cluster map, id [0], successfully created
```

---

### 6.6.13 Listing volumes mapped to the host

Run the `lshostvdiskmap` command to show the volumes that are mapped to the specific host, as shown in Example 6-24.

*Example 6-24 The `lshostvdiskmap` command*

---

```
IBM_2145:ITS0_CLUSTER:superuser>lshostvdiskmap -delim , Siam
id,name,SCSI_id,vdisk_id,vdisk_name,wwpn,vdisk_UID
3,Siam,0,0,volume_A,210000E08B18FF8A,60050768018301BF280000000000000C
```

---

In the output of the command, you can see that only one volume (`volume_A`) is mapped to the host `Siam`. The volume is mapped with SCSI LUN ID 0.

If no host name is specified as `lshostvdiskmap` command, it returns all defined host-to-volume mappings.

**Specifying the flag before the host name:** Although the `-delim` flag normally comes at the end of the command string, you must specify this flag before the host name in this case. Otherwise, it returns the following message:

```
CMMVC6070E An invalid or duplicated parameter, unaccompanied argument, or
incorrect argument sequence has been detected. Ensure that the input is as per
the help.
```

You can also run the `lshostclustermap` command to show the volumes that are mapped to a specific host cluster, as shown in Example 6-25 on page 341.

*Example 6-25 The lshostclustervolumemap command*

---

```
IBM_Storwize:ITS0:superuser>lshostclustervolumemap
id name          SCSI_id volume_id volume_name      volume_UID
IO_group_id IO_group_name
0 vmware_cluster 0          9          UNCOMPRESSED_VOL 6005076400F580049800000000000011 0
io_grp0
```

---

## 6.6.14 Listing hosts mapped to the volume

To identify the hosts to which a specific volume was mapped, run the `lsvdiskhostmap` command, as shown in Example 6-26.

*Example 6-26 The lsvdiskhostmap command*

---

```
IBM_2145:ITS0_CLUSTER:superuser>lsvdiskhostmap -delim , volume_B
id,name,SCSI_id,host_id,host_name,vdisk_UID
26,volume_B,0,2,Almaden,6005076801AF813F1000000000000020
```

---

This command shows the list of hosts to which the volume `volume_B` is mapped.

**Specifying the `-delim` flag:** Although the optional `-delim` flag normally comes at the end of the command string, you must specify this flag before the volume name in this case. Otherwise, the command does not return any data.

## 6.6.15 Deleting a volume to host mapping

Deleting a volume mapping does not affect the volume. Instead, it removes only the host's ability to use the volume. Use the `rmvdiskhostmap` command to unmap a volume from a host, as shown in Example 6-27.

*Example 6-27 The rmvdiskhostmap command*

---

```
IBM_2145:ITS0_CLUSTER:superuser>rmvdiskhostmap -host Tiger volume_D
```

---

This command unmaps the volume that is called `volume_D` from the host that is called `Tiger`.

You can also run the `rmvolumehostclustermap` command to delete a volume mapping from a host cluster, as shown in Example 6-28.

*Example 6-28 The rmvolumehostclustermap command*

---

```
IBM_Storwize:ITS0:superuser>rmvolumehostclustermap -hostcluster vmware_cluster
UNCOMPRESSED_VOL
```

---

This command unmaps the volume called `UNCOMPRESSED_VOL` from the host cluster called `vmware_cluster`.

**Note:** Removing a volume mapped to the host makes the volume unavailable for I/O operations. Ensure that the host is prepared for this before removing a volume mapping.

## 6.6.16 Migrating a volume

You might want to migrate volumes from one set of MDisks to another set of MDisks to decommission an old disk subsystem to better distribute load across your virtualized environment, or to migrate data into the IBM Spectrum Virtualize environment by using image mode. For more information about migration, see Chapter 8, “Storage migration” on page 429.

**Important:** After migration is started, it continues until completion unless it is stopped or suspended by an error condition or the volume that is being migrated is deleted.

As you can see from the parameters that are shown in Example 6-29, before you can migrate your volume, you must determine the name of the volume that you want to migrate and the name of the storage pool to which you want to migrate it. To list the names of volumes and storage pools, run the `lsvdisk` and `lsmdiskgrp` commands.

The command that is shown in Example 6-29 moves `volume_C` to the storage pool named `STGPoo1_DS5000-1`.

*Example 6-29 The migratevdisk command*

---

```
IBM_2145:ITSO_CLUSTER:superuser>migratevdisk -mdiskgrp STGPoo1_DS5000-1 -vdisk volume_C
```

---

**Note:** If insufficient extents are available within your target storage pool, you receive an error message. Ensure that the source MDG and target MDG have the same extent size.

You can use the optional `threads` parameter to control priority of the migration process. The default is 4, which is the highest priority setting. However, if you want the process to take a lower priority over other types of I/O, you can specify 3, 2, or 1.

You can run the `lsmigrate` command at any time to see the status of the migration process, as shown in Example 6-30.

*Example 6-30 The lsmigrate command*

---

```
IBM_2145:ITSO_CLUSTER:superuser>lsmigrate
migrate_type MDisk_Group_Migration
progress 0
migrate_source_vdisk_index 27
migrate_target_mdisk_grp 2
max_thread_count 4
migrate_source_vdisk_copy_id 0
```

```
IBM_2145:ITSO_CLUSTER:superuser>lsmigrate
migrate_type MDisk_Group_Migration
progress 76
migrate_source_vdisk_index 27
migrate_target_mdisk_grp 2
max_thread_count 4
migrate_source_vdisk_copy_id 0
```

---

**Progress:** The progress is shown in terms of percentage complete. If no output is displayed when running the command, all volume migrations finished.



## 6.6.17 Migrating a fully managed volume to an image mode volume

Migrating a fully managed volume to an image mode volume enables the IBM Spectrum Virtualize system to be removed from the data path. This feature might be useful when the IBM Spectrum Virtualize system is used as a data mover.

To migrate a fully managed volume to an image mode volume, the following rules apply:

- ▶ Cloud snapshots must not be enabled on the source volume.
- ▶ The destination MDisk must be greater than or equal to the size of the volume.
- ▶ The MDisk that is specified as the target must be in an unmanaged state.
- ▶ Regardless of the mode in which the volume starts, it is reported as a managed mode during the migration.
- ▶ If the migration is interrupted by a system recovery or cache problem, the migration resumes after the recovery completes.

Example 6-31 shows running the `migratetoimage` command to migrate the data from `volume_A` onto `mdisk10`, and to put the MDisk `mdisk10` into the `STGPool_IMAGE` storage pool.

*Example 6-31 The `migratetoimage` command*

---

```
IBM_2145:ITSO_CLUSTER:superuser>migratetoimage -vdisk volume_A -mdisk mdisk10 -mdiskgrp
STGPool_IMAGE
```

---

## 6.6.18 Shrinking a volume

Running the `shrinkvdisksize` command reduces the capacity that is allocated to the particular volume by the specified amount. You cannot shrink the real size of a thin-provisioned volume to less than its used size. All capacities (including changes) must be in multiples of 512 bytes. An entire extent is reserved even if it is only partially used. The default capacity units are MBs.

You can use this command to shrink the physical capacity of a volume or to reduce the virtual capacity of a thin-provisioned volume without altering the physical capacity that is assigned to the volume. Use the following parameters to change volume size:

- ▶ For a standard-provisioned volume, use the `-size` parameter.
- ▶ For a thin-provisioned volume's real capacity, use the `-rsize` parameter.
- ▶ For a thin-provisioned volume's virtual capacity, use the `-size` parameter.

When the virtual capacity of a thin-provisioned volume is changed, the warning threshold is automatically scaled.

If the volume contains data that is being used, do not shrink the volume without backing up the data first. The system reduces the capacity of the volume by removing arbitrarily chosen extent, or extents from those sets allocated to the volume. You cannot control which extents are removed. Therefore, you cannot assume that it is unused space that is removed.

Image mode volumes cannot be reduced in size. To reduce their size they must first be migrated to fully managed mode.

Before the `shrinkvdisksize` command is used on a mirrored volume, all copies of the volume must be synchronized.

**Important:** Consider the following guidelines when you are shrinking a disk:

- ▶ If the volume contains data or host-accessible metadata (for example, an empty physical volume of an LVM), do not shrink the disk.
- ▶ This command can shrink a FlashCopy target volume to the same capacity as the source.
- ▶ Before you shrink a volume, validate that the volume is not mapped to any host objects.
- ▶ You can determine the exact capacity of the source or master volume by running the `svcinfolsvdisk -bytes vdiskname` command.

Shrink the volume by the required amount by running the `shrinkvdisksize -size disk_size -unit b | kb | mb | gb | tb | pb vdisk_name | vdisk_id` command.

Example 6-32 shows running `shrinkvdisksize` command to reduce the size of volume `volume_D` from a total size of 80 GB by 44 GB to the new total size of 36 GB.

*Example 6-32 The shrinkvdisksize command*

---

```
IBM_2145:ITSO_CLUSTER:superuser>shrinkvdisksize -size 44 -unit gb volume_D
```

---

## 6.6.19 Listing volumes using the MDisk

Run the `lsmdiskmember` command to identify which volumes use space on the specified MDisk. Example 6-33 displays list of volume IDs of all volume copies that use `mdisk8`. To correlate the IDs that are displayed in this output to volume names, run the `lsvdisk` command.

*Example 6-33 The lsmdiskmember command*

---

```
IBM_2145:ITSO_CLUSTER:superuser>lsmdiskmember mdisk8
id copy_id
24 0
27 0
```

---

## 6.6.20 Listing MDisks used by the volume

Run the `lsvdiskmember` command to list MDisks that supply space used by the specified volume. Example 6-34 lists MDisk IDs of all MDisks used by volume with ID 0.

*Example 6-34 The lsvdiskmember command*

---

```
IBM_2145:ITSO_CLUSTER:superuser>lsvdiskmember 0
id
4
5
6
7
```

---

If you want to know more about these MDisks, you can run the `lsmdisk` command, providing as a parameter the MDisk ID that is listed in the output of the `lsvdiskmember` command.

## 6.6.21 Listing volumes defined in the storage pool

Run the `lsvdisk -filtervalue` command to list volumes that are defined in the specified storage pool. Example 6-35 shows how to use the `lsvdisk -filtervalue` command to list all volumes that are defined in the storage pool named Pool0.

*Example 6-35 The lsvdisk -filtervalue command: volumes in the pool*

---

```
IBM_Storwize:ITS0:superuser>lsvdisk -filtervalue mdisk_grp_name=Pool0 -delim ,
id,name,IO_group_id,IO_group_name,status,mdisk_grp_id,mdisk_grp_name,capacity,type,FC_id,FC
_name,RC_id,RC_name,vdisk_UID,fc_map_count,copy_count,fast_write_state,se_copy_count,RC_cha
nge,compressed_copy_count,parent_mdisk_grp_id,parent_mdisk_grp_name,formatting,encrypt,volu
me_id,volume_name,function
0,A_MIRRORED_VOL_1,0,io_grp0,online,0,Pool0,10.00GB,striped,,,,,6005076400F5800498000000000
00002,0,1,empty,0,no,0,0,Pool0,no,yes,0,A_MIRRORED_VOL_1,
2,VOLUME_WITH_MIRRORED_COPY,0,io_grp0,online,0,Pool0,10.00GB,striped,,,,,6005076400F5800498
00000000000004,0,1,empty,0,no,0,0,Pool0,no,yes,2,VOLUME_WITH_MIRRORED_COPY,
3,THIN_PROVISION_VOL_1,0,io_grp0,online,0,Pool0,100.00GB,striped,,,,,6005076400F58004980000
0000000005,0,1,empty,1,no,0,0,Pool0,no,yes,3,THIN_PROVISION_VOL_1,
6,MIRRORED_SYNC_RATE_16,0,io_grp0,online,0,Pool0,10.00GB,striped,,,,,6005076400F58004980000
0000000008,0,1,empty,0,no,0,0,Pool0,no,yes,6,MIRRORED_SYNC_RATE_16,
7,THIN_PROVISION_MIRRORED_VOL,0,io_grp0,online,0,Pool0,10.00GB,striped,,,,,6005076400F58004
980000000000009,0,1,empty,1,no,0,0,Pool0,no,yes,7,THIN_PROVISION_MIRRORED_VOL,
8,Tiger,0,io_grp0,online,0,Pool0,10.00GB,striped,,,,,6005076400F580049800000000000010,0,1,e
mpty,0,no,0,0,Pool0,no,yes,8,Tiger,
9,UNCOMPRESSED_VOL,0,io_grp0,online,0,Pool0,10.00GB,striped,,,,,6005076400F58004980000000000
00011,0,1,empty,0,no,1,0,Pool0,no,yes,9,UNCOMPRESSED_VOL,
12,vdisk0_restore,0,io_grp0,online,0,Pool0,10.00GB,striped,,,,,6005076400F5800498000000000000
000E,0,1,empty,0,no,0,0,Pool0,no,yes,12,vdisk0_restore,
13,vdisk0_restore1,0,io_grp0,online,0,Pool0,10.00GB,striped,,,,,6005076400F5800498000000000000
0000F,0,1,empty,0,no,0,0,Pool0,no,yes,13,vdisk0_restore1,
```

---

## 6.6.22 Listing storage pools in which a volume has its extents

Run the `lsvdisk` command to show to which storage pool a specific volume belongs, as shown in Example 6-36.

*Example 6-36 The lsvdisk command: Storage pool ID and name*

---

```
IBM_Storwize:ITS0:superuser>lsvdisk 0
id 0
name A_MIRRORED_VOL_1
IO_group_id 0
IO_group_name io_grp0
status online
mdisk_grp_id 0
mdisk_grp_name Pool0
capacity 10.00GB
type striped
formatted yes
formatting no
mdisk_id
mdisk_name
FC_id
FC_name
RC_id
RC_name
vdisk_UID 6005076400F580049800000000000002
```

```

preferred_node_id 2
fast_write_state empty
cache readwrite
udid 4660
fc_map_count 0
sync_rate 50
copy_count 1
se_copy_count 0
File system
mirror_write_priority latency
RC_change no
compressed_copy_count 0
access_IO_group_count 1
last_access_time
parent_mdisk_grp_id 0
parent_mdisk_grp_name Pool0
owner_type none
owner_id
owner_name
encrypt yes
volume_id 0
volume_name A_MIRRORED_VOL_1
function
throttle_id 1
throttle_name throttle1
IOPs_limit 233
bandwidth_limit_MB 122
volume_group_id
volume_group_name
cloud_backup_enabled no
cloud_account_id
cloud_account_name
backup_status off
last_backup_time
restore_status none
backup_grain_size
deduplicated_copy_count 0

copy_id 0
status online
sync yes
auto_delete no
primary yes
mdisk_grp_id 0
mdisk_grp_name Pool0
type striped
mdisk_id
mdisk_name
fast_write_state empty
used_capacity 10.00GB
real_capacity 10.00GB
free_capacity 0.00MB
overallocation 100
autoexpand
warning
grainsize
se_copy no
easy_tier on
easy_tier_status measured
tier tier0_flash

```

```

tier_capacity 0.00MB
tier tier1_flash
tier_capacity 0.00MB
tier tier_enterprise
tier_capacity 0.00MB
tier tier_nearline
tier_capacity 10.00GB
compressed_copy no
uncompressed_used_capacity 10.00GB
parent_mdisk_grp_id 0
parent_mdisk_grp_name Pool0
encrypt yes
deduplicated_copy no
used_capacity_before_reduction0.00MB

```

---

For more information about these storage pools, run the `lsmdiskgrp` command, as described in Chapter 5, “Storage pools” on page 199.

### 6.6.23 Tracing a volume from a host back to its physical disks

In some cases, you might need to verify exactly which physical disks are used to store data of a volume. This information is not directly available to the host; however, it might be obtained by using a sequence of queries.

The first step is to unequivocally map a logical device that is seen by the host to a volume that is presented by the storage system. The best volume characteristics for this purpose is the volume ID. This ID is available to the operating system in the Vendor Specified Identifier field of page 0x80 or 0x83 (vital product data, VPD), which the storage device sends in response to SCSI INQUIRY command from the host.

In practice, the ID can be obtained from the multipath driver in the operating system. After the volume ID is known, it can be used to identify physical location of data.

**Note:** For sequential and image mode volumes, a volume copy is mapped to one MDisk. This configuration often is not the case for striped volumes, unless volume size is not greater than an extent size. Therefore, a single striped volume uses multiple mDisks in a typical case.

On hosts that are running IBM System Storage Multipath Subsystem Device Driver (SDD), you can obtain the volume ID from the output of the `datapath query device` command. You see a long disk serial number for each vpath device, as shown in Example 6-37.

*Example 6-37 The datapath query device command*

```

DEV#: 0 DEVICE NAME: Disk1 Part0 TYPE: 2145 POLICY: OPTIMIZED
SERIAL: 60050768018301BF2800000000000005
=====
Path# Adapter/Hard Disk State Mode Select Errors
  0 Scsi Port2 Bus0/Disk1 Part0 OPEN NORMAL 20 0
  1 Scsi Port3 Bus0/Disk1 Part0 OPEN NORMAL 2343 0

DEV#: 1 DEVICE NAME: Disk2 Part0 TYPE: 2145 POLICY: OPTIMIZED
SERIAL: 60050768018301BF2800000000000004
=====
Path# Adapter/Hard Disk State Mode Select Errors
  0 Scsi Port2 Bus0/Disk2 Part0 OPEN NORMAL 2335 0

```

```

1      Scsi Port3 Bus0/Disk2 Part0      OPEN  NORMAL          0          0

DEV#:  2  DEVICE NAME: Disk3 Part0  TYPE: 2145          POLICY: OPTIMIZED
SERIAL: 60050768018301BF280000000000006
=====
Path#      Adapter/Hard Disk      State Mode      Select      Errors
  0      Scsi Port2 Bus0/Disk3 Part0      OPEN  NORMAL      2331        0
  1      Scsi Port3 Bus0/Disk3 Part0      OPEN  NORMAL        0          0

```

**State:** In Example 6-38, the state of each path is OPEN. Sometimes, the state is CLOSED. This state does not necessarily indicate a problem because it might be a result of the path's processing stage.

On a Linux host running native multipath driver, you can use the output of the `multipath -ll` command, as shown in Example 6-38.

*Example 6-38 Volume ID as returned by multipath -ll*

```

mpath1 (360050768018301BF2800000000000004)  IBM,2145
[size=2.0G][features=0][hwhandler=0]
\_ round-robin 0 [prio=200][ enabled]
\_ 4:0:0:1 sdd 8:48 [active][ready]
\_ 5:0:0:1 sdt 65:48 [active][ready]
\_ round-robin 0 [prio=40][ active]
\_ 4:0:2:1 sdak 66:64 [active][ready]
\_ 5:0:2:1 sda1 66:80 [active][ready]

```

**Note:** Volume ID shown in the output of `multipath -ll` is generated by Linux `scsi_id`. For systems that provide the VPD by way of page 0x83 (such as IBM Spectrum Virtualize devices), the ID that is obtained from the VPD page is prefixed with number 3, which is the Network Address Authority (NAA) type identifier. Therefore, the volume NAA identifier (that is, volume ID obtained by way of the SCSI INQUIRY command) starts at the second displayed digit. In Example 6-38, the volume ID starts with digit 6.

After you know the volume ID, complete the following steps:

1. Run the `lshostvdiskmap` command to list volumes mapped to the host. Example 6-39 lists volumes mapped to host Almaden.

*Example 6-39 The lshostvdiskmap command*

```

IBM_2145:ITSO_CLUSTER:superuser>lshostvdiskmap -delim , Almaden
id,name,SCSI_id,vdisk_id,vdisk_name,vdisk_UID
2,Almaden,0,26,volume_B,60050768018301BF28000000000000005
2,Almaden,1,27,volume_A,60050768018301BF28000000000000004
2,Almaden,2,28,volume_C,60050768018301BF28000000000000006

```

Look for the VDisk unique identifier (UID) that matches volume UID that was identified and note the volume name (or ID) for a volume with this UID.

2. Run the `lsvdiskmember vdiskname` command to list of the MDisks that contain extents allocated to the specified volume, as shown in Example 6-40.

*Example 6-40 The lsvdiskmember command*

```

IBM_2145:ITSO_CLUSTER:superuser>lsvdiskmember volume_A
id

```

0  
1  
2  
3  
4  
10  
11  
13  
15  
16  
17

---

3. For each of MDisk IDs that were obtained in the previous step, run the `lsmdisk mdiskID` command to discover MDisk controller and LUN information. Example 6-41 shows output for MDisk 0. The output displays the back-end storage controller name and the controller LUN ID to help you to track back to a LUN within the disk subsystem.

*Example 6-41 The lsmdisk command*

---

```
IBM_2145:ITS0_CLUSTER:superuser>lsmdisk 0
id 0
name mdisk0
status online
mode managed
mdisk_grp_id 0
mdisk_grp_name STGPool_DS3500-1
capacity 128.0GB
quorum_index 1
block_size 512
controller_name ITS0-DS3500
ctrl_type 4
ctrl_WWNN 20080080E51B09E8
controller_id 2
path_count 4
max_path_count 4
ctrl_LUN_# 0000000000000000
UID 60080e50001b0b62000007b04e731e4d00000000000000000000000000000000
preferred_WWPN 20580080E51B09E8
active_WWPN 20580080E51B09E8
fast_write_state empty
raid_status
raid_level
redundancy
strip_size
spare_goal
spare_protection_min
balanced
tier generic_hdd
```

---

You can identify the back-end storage that is presenting the LUN by using the value of the `controller_name` field that was returned for the MDisk.

On the back-end storage, you can identify which physical disks make up the LUN that was presented to the Storage Virtualize system by using the volume ID that is displayed in the UID field.







# Hosts

This chapter describes the host configuration procedures that are required to attach supported hosts to the systems. It also describes Host Clustering and N-Port Virtualization ID (NPV) support from a host's perspective.

This chapter includes the following topics:

- ▶ 7.1, “Host attachment overview” on page 352
- ▶ 7.2, “Host clusters” on page 353
- ▶ 7.3, “NVMe over Fibre Channel” on page 353
- ▶ 7.4, “N\_Port ID Virtualization support” on page 354
- ▶ 7.5, “Hosts operations by using the GUI” on page 365
- ▶ 7.6, “Performing hosts operations by using the command-line interface” on page 415

## 7.1 Host attachment overview

The storage systems support a wide range of host types (from both IBM and non-IBM). This feature makes it possible to consolidate storage in an open systems environment into a common pool of storage. Then, you can use and manage the storage pool more efficiently as a single entity from a central point on the storage area network (SAN).

The ability to consolidate storage for attached open systems hosts provides the following benefits:

- ▶ Easier storage management
- ▶ Increased utilization rate of the installed storage capacity
- ▶ Advanced Copy Services functions that are offered across storage systems from separate vendors
- ▶ Only one multipath driver is required for attached hosts

Hosts can be connected to the IBM Spectrum Virtualize system by using any of the following protocols:

- ▶ Fibre Channel (FC)
- ▶ Fibre Channel over Ethernet (FCoE)
- ▶ Internet Small Computer Systems Interface (iSCSI)
- ▶ iSCSI Extensions for RDMA (iSER)
- ▶ Non-Volatile Memory Express (NVMe)

Hosts that connect to the IBM Spectrum Virtualize system by using the fabric switches that use the FC or FCoE protocol must be zoned appropriately, as indicated in Chapter 2, “Planning” on page 53.

Hosts that connect to the IBM Spectrum Virtualize system with iSCSI or iSER protocol must be configured appropriately, as indicated in Chapter 2, “Planning” on page 53.

**Note:** Certain host operating systems can be directly connected to the IBM Spectrum Virtualize system without the need for FC fabric switches. For more information, go to the [IBM SSIC](#):

For load balancing and access redundancy on the host side, you must use a host multipathing driver. A host multipathing I/O driver is required in the following situations:

- ▶ Protection from fabric link failures, including port failures on the IBM Spectrum Virtualize system nodes
- ▶ Protection from a host bus adapter (HBA) failure (if two HBAs are in use)
- ▶ Protection from fabric failures if the host is connected through two HBAs to two separate fabrics
- ▶ To provide load balancing across the host HBAs

To learn about various host operating systems and versions that are supported by IBM SAN Volume Controller, go to the [IBM SSIC](#).

For more information about how to attach various supported host operating systems to the systems, see the “Host Attachment” section of [IBM Knowledge Center](#).

If your host OS is not in SSIC, you can ask an IBM representative to submit a special request for support by contacting your business partner, account manager, or IBM Support.

## 7.2 Host clusters

IBM Spectrum Virtualize software supports host clusters starting with V7.7.1. The host cluster allows a user to create a group of hosts to form a cluster. A cluster is treated as one single entity, which allows multiple hosts to have access to the same volumes.

Volumes that are mapped to a host cluster are assigned to all members of the host cluster with the same Small Computer System Interface (SCSI) ID.

A typical use case is to define a host cluster that contains all of the worldwide port names (WWPNs) that belong to the hosts that are participating in a host OS-based cluster, such as IBM PowerHA® or Microsoft Cluster Server (MSCS). Before the host clusters, as an example, an ESX cluster can be created as a single host object, containing up to 32 ports (WWPN). Within the host cluster object, you can have up to 128 hosts in a single host cluster object. With that setup, managing host clusters becomes easier.

The following commands deal with host clusters:

- ▶ `lshostcluster`
- ▶ `lshostclustermember`
- ▶ `lshostclustervolumemap`
- ▶ `mkhost` (modified to put a host into a host cluster on creation)
- ▶ `rmhostclustermember`
- ▶ `rmhostcluster`
- ▶ `rmvolumehostclustermap`
- ▶ `chostcluster`

## 7.3 NVMe over Fibre Channel

Running V8.2 and above, IBM SAN Volume Controller can attach NVMe hosts by using FC-NVMe. FC-NVMe uses Fibre Channel Protocol (FCP) as its underlying transport, which puts the data transfer in control of the target and transfers data direct from host memory, similar to Remote Direct Memory Access (RDMA). In addition, FC-NVMe allows for a host to send commands and data together (first burst), which eliminates the first data “read” by the target and provides better performance at distances.

The limitation that existed in V8.2 of being able to attach only one of SCSI or NVMe per I/O group was removed in V8.3. You can now run SCSI and NVMe in parallel. The limit of 512 host objects per I/O group remains in place as described in Chapter 2, “Planning” on page 53. However, when running SCSI and NVMe in parallel, there are limits for each protocol, as listed in Table 7-1.

*Table 7-1 Defined Host Object Limits Per I/O group*

SCSI Host Objects	NVMe Host Objects	Total Host Objects
496	16	512

**Note:** Although these are the maximum amount of each type of host attachment you can have, if you do not have the maximum NVMe Host Objects defined, the amount of Total Host Objects that you can have is not reduced.

For example, if you had 10 NVMe Host Objects defined, you could have up to 502 SCSI Host Objects defined. The only hard limit is 16 NVMe Host Objects per I/O group. You should be diligent when planning a parallel SCSI and NVMe deployment because this can be resource intensive, especially with large deployments (many hosts). This is because of NVMe is more sensitive to delays than SCSI. It is always recommended to check the [IBM Support Configuration Limits we page](#) for your product and specifically the “NVMe over Fibre Channel Host Properties” section.

Ensure you select the correct product.

Although the protocol attachment limit was removed, a volume can still be mapped only to a host by using one protocol. This is to avoid any potential interoperability problems. Flashcopy, Volume Mirroring, Remote Copy (RC) and data reduction pools (DRPs) are all supported with NVMe over Fabric (NVMe-oF). New with V8.3.1 is full support for Stretched Cluster configurations.

**Note:** There is currently no support for Hyperswap or Non-Disruptive Volume Move (NDVM). This is because NVMe as a protocol does not allow accessing the same volume from more than one NVMe subsystem. In 8.2, each I/O group was its own NVMe subsystem, and so it was not possible to use NDVM or Hyperswap. However, in V8.3.1.0, we have made an improvement so that the entire cluster is one single NVMe subsystem.

## 7.4 N\_Port ID Virtualization support

The usage model for all IBM Spectrum Virtualize products is based around a two-way active/active node model. A pair of distinct control modules shares active/active access for any specific volume. These nodes each have their own FC worldwide node name (WWNN). Therefore, all ports that are presented from each node have a set of WWPNs that is presented to the fabric.

Traditionally, if one node fails or is removed for some reason, the paths that are presented for volumes from that node go offline. In this case, it is up to the native operating system multipathing software to fail over from using both sets of WWPN to just those that remain online. Although this process is what multipathing software is designed to do, occasionally it can be problematic, particularly if paths are not seen as coming back online for some reason.

Starting with IBM Spectrum Virtualize V7.7, the system can be enabled in NPIV mode. When NPIV mode is enabled on the IBM Spectrum Virtualize system, ports do not come online until they are ready to service I/O, which improves host behavior around node unpendes. In addition, path failures because of an offline node are masked from hosts and their multipathing driver do not need to do any path recovery.

From IBM Spectrum Virtualize V8.2 and later, the IBM SAN Volume Controller system can now attach to NVMe hosts by using FC-NVMe. FC-NVMe uses the FCP as its underlying transport, which puts the data transfer in control of the target and transfers data directly from host memory, similar to RDMA. In addition, FC-NVMe allows for a host to send commands and data together (first burst), which eliminates the first data “read” by the target and providing better performance at distances.

For more information about NVMe, see *IBM Storage and the NVM Express Revolution*, REDP-5437.

When NPIV is enabled on IBM Spectrum Virtualize system nodes, each physical WWPN reports up to four virtual WWPNs in addition to the physical one, as listed in Table 7-2.

Table 7-2 IBM Spectrum Virtualize NPIV ports’s

NPIV port	Port description
Primary Port	This is the WWPN that communicates with back-end storage. It can be used for node to node traffic (local or remote).
Primary SCSI Host Attach Port	This is the WWPN that communicates with hosts. It is a target port only. This is the primary port, so it is based on this local node’s WWNN.
Failover SCSI Host Attach Port	This is a standby WWPN that communicates with hosts and is brought online only if the partner node within the I/O group goes offline. This is the same as the Primary Host Attach WWPN of the partner node.
Primary NVMe Host Attach Port	This is the WWPN that communicates with hosts. It is a target port only. This is the primary port, so it is based on this local node’s WWNN.
Failover NVMe Host Attach Port	This is a standby WWPN that communicates with hosts and is brought online only if the partner node within the I/O group goes offline. This is the same as the Primary Host Attach WWPN of the partner node.

Figure 7-1 shows the five WWPNs that are associated with an IBM SAN Volume Controller port when NPIV is enabled.

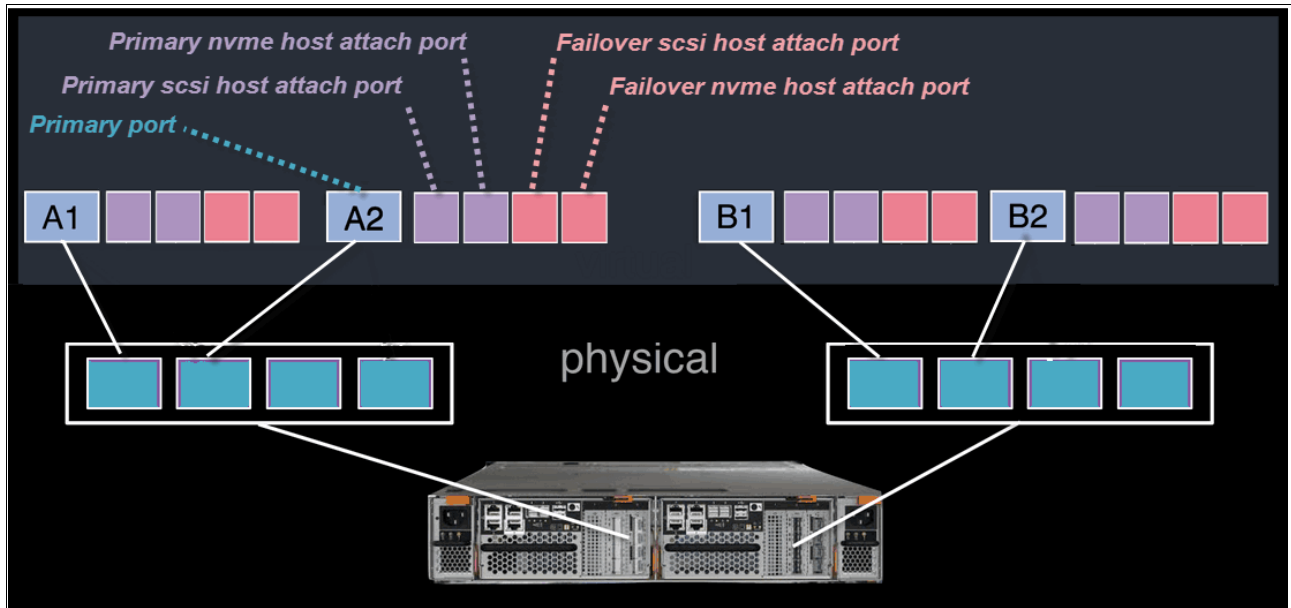


Figure 7-1 Allocation of NPIV virtual WWPN ports per physical port

The *failover host attach ports* are not currently active. Figure 7-2 shows what happens when the partner node fails. After the node failure, the failover host attach ports on the remaining node become active and take on the WWPN of the failed node's primary host attach port.

**Note:** Figure 7-2 shows only two ports per node in detail, but the same details apply to all physical ports. The effect is the same for NVMe ports because they use the same NPIV structure, but with the topology NVMe instead of regular SCSI.

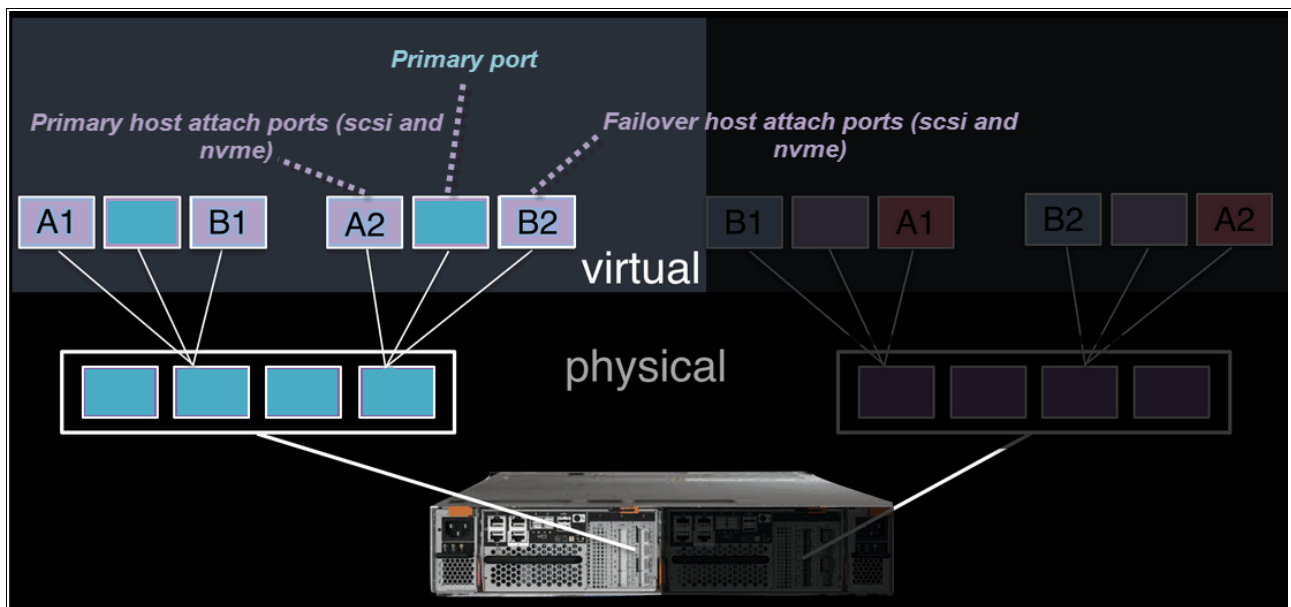


Figure 7-2 Allocation of NPIV virtual WWPN ports per physical port after a node failure

With Version 7.7 onwards, this process happens automatically when NPIV is enabled at a system level in the SAN Volume Controller system. This failover happens only between the two nodes in the same I/O group. Similar NPIV capabilities were introduced with Version 8.1, allowing for a hot spare node to swap into an I/O group.

**Note:** When the Hot Spare Node function is used, the NPIV ports move to the partner node on the I/O group in a node failure. This action occurs to make the ports available immediately. After the hot spare node finishes this process, the NPIV ports automatically move to the hot spare node.

A transitional mode allows migration of hosts from previous non-NPIV-enabled systems to enabled NPIV systems, which allows for a transition period as hosts are rezoned to the *primary host attach* WWPNs.

The process for enabling NPIV on a new system is slightly different than on an existing system. For more information, see [IBM Knowledge Center](#).

**Note:** NPIV is supported only for FCP. It is not supported for the FCoE or iSCSI protocols.

## 7.4.1 NPIV prerequisites

Consider the following key points for NPIV enablement:

- ▶ The IBM Spectrum Virtualize system must be running Version 7.7 or later.
- ▶ A Version 7.7 or later system with NPIV enabled as back-end storage for a system that is earlier than Version 7.7 is not supported.
- ▶ Both nodes in an IO group should have identical hardware to allow failover to work as expected.
- ▶ The FC switches that the SAN Volume Controller ports are attached to must support NPIV and have this feature enabled.

## 7.4.2 Enabling NPIV on a new system

New IBM SAN Volume Controller systems that are shipped with Version 7.7 and later should have NPIV enabled by default. If your new IBM Spectrum Virtualize system does not have NPIV enabled, it can be enabled by completing the following steps:

1. Run the `lsvg` command to determine the <id> of I/O groups that are present in the system, as shown in Example 7-1. In our example, we have a single I/O group with ID 0.

*Example 7-1 Listing the I/O groups in a system*

---

```
IBM_2145:ITS0-SV1:superuser>lsvg
id name          node_count vdisk_count host_count site_id site_name
0  io_grp0        2           8           0           0
1  io_grp1        0           0           0           0
2  io_grp2        0           0           0           0
3  io_grp3        0           0           0           0
4  recovery_io_grp 0           0           0           0
IBM_2145:ITS0-SV1:superuser>
```

---

- Run the `lsgigrp <id> | grep fctargetportmode` command for the specific I/O group ID to display the `fctargetportmode` setting. If this is enabled, as shown in Example 7-2, NPIV host target port mode is enabled.

*Example 7-2 Checking the NPIV mode with the `fctargetportmode` field*

```
IBM_2145:ITS0-SV1:superuser>lsgigrp 0 | grep fctargetportmode
fctargetportmode enabled
IBM_2145:ITS0-SV1:superuser>
```

- The virtual WWPNs can be listed by running the `lstargetportfc` command, as shown in Example 7-3. Look for the `host_io_permitted` and `virtualized` columns to be `yes`, meaning the WWPN in those lines is a primary host attach port and should be used when zoning the hosts to the SAN Volume Controller system.

*Example 7-3 Listing the virtual WWPNs*

```
IBM_2145:ITS0-SV1:superuser>lstargetportfc
id WWPN          WWNN          port_id owning_node_id current_node_id nportid host_io_permitted virtualized protocol
1 500507680140A288 500507680100A288 1 1 1 010A00 no no scsi
2 500507680142A288 500507680100A288 1 1 1 010A02 yes yes scsi
3 500507680144A288 500507680100A288 1 1 1 010A01 yes yes nvme
4 500507680130A288 500507680100A288 2 1 1 010400 no no scsi
5 500507680132A288 500507680100A288 2 1 1 010401 yes yes scsi
6 500507680134A288 500507680100A288 2 1 1 010402 yes yes nvme
7 500507680110A288 500507680100A288 3 1 1 010500 no no scsi
8 500507680112A288 500507680100A288 3 1 1 010501 yes yes scsi
9 500507680114A288 500507680100A288 3 1 1 010502 yes yes nvme
10 500507680120A288 500507680100A288 4 1 1 010A00 no no scsi
11 500507680122A288 500507680100A288 4 1 1 010A02 yes yes scsi
12 500507680124A288 500507680100A288 4 1 1 010A01 yes yes nvme
49 500507680C110009 500507680C000009 1 2 2 010500 no no scsi
50 500507680C150009 500507680C000009 1 2 2 010502 yes yes scsi
51 500507680C190009 500507680C000009 1 2 2 010501 yes yes nvme
52 500507680C120009 500507680C000009 2 2 2 010400 no no scsi
53 500507680C160009 500507680C000009 2 2 2 010401 yes yes scsi
54 500507680C1A0009 500507680C000009 2 2 2 010402 yes yes nvme
55 500507680C130009 500507680C000009 3 2 2 010900 no no scsi
56 500507680C170009 500507680C000009 3 2 2 010902 yes yes scsi
57 500507680C1B0009 500507680C000009 3 2 2 010901 yes yes nvme
58 500507680C140009 500507680C000009 4 2 2 010900 no no scsi
59 500507680C180009 500507680C000009 4 2 2 010901 yes yes scsi
60 500507680C1C0009 500507680C000009 4 2 2 010902 yes yes nvme
IBM_2145:ITS0-SV1:superuser>
```

- At this point, you can zone the hosts using the primary host attach ports (virtual WWPNs) of the IBM SAN Volume Controller ports, as shown in `bo1d` in the output in Example 7-3.

**Note:** If supported on the host, you can use NVMe ports in the zoning. Make sure your host supports NVMe before creating the zones by using `nvme` ports. At the IBM SAN Volume Controller system, a host can have ports of only one topology; you cannot have NVMe and SCSI ports on the same host object.

- If the status of `fctargetportmode` is disabled and this is a new installation, run the `chiogrp` command to set enabled NPIV mode, as shown in Example 7-4.

*Example 7-4 Changing the NPIV mode to enabled*

```
IBM_2145:ITS0-SV1:superuser>chiogrp -fctargetportmode enabled 0
```



6. NPIV enablement can be verified by checking the `fctargetportmode` field, as shown in Example 7-5.

*Example 7-5 NPIV enablement verification*

---

```

IBM_2145:ITS0-SV1:superuser>lsiogrp 0
id 0
name io_grp0
node_count 2
vdisk_count 8
host_count 0
flash_copy_total_memory 20.0MB
flash_copy_free_memory 19.9MB
remote_copy_total_memory 20.0MB
remote_copy_free_memory 20.0MB
mirroring_total_memory 20.0MB
mirroring_free_memory 19.9MB
raid_total_memory 40.0MB
raid_free_memory 40.0MB
maintenance no
compression_active no
accessible_vdisk_count 8
compression_supported no
max_enclosures 20
encryption_supported yes
flash_copy_maximum_memory 2048.0MB
site_id
site_name
fctargetportmode enabled
compression_total_memory 0.0MB
deduplication_supported yes
deduplication_active no
nqn nqn.1986-03.com.ibm:nvme:2145.0000020067214511.iogroup0
IBM_2145:ITS0-SV1:superuser>

```

---

You can now configure zones for hosts by using the primary host attach ports (virtual WWPNs) of the IBM Spectrum Virtualize ports, as shown in **bold** in the output that is shown in Example 7-3 on page 358.

### 7.4.3 Enabling NPIV on a system

If your IBM Spectrum Virtualize system was running before upgrading to Version 7.7.1 or later, the NPIV feature is not turned on by default because it results in all hosts losing access to the SAN Volume Controller disk. Check your system by running the `lsiogrp` command and looking for the `fctargetportmode` setting, as shown in Example 7-6.

*Example 7-6 Checking whether fctargetportmode is disabled*

---

```

IBM_2145:ITS0-SV1:superuser>lsiogrp
id name          node_count vdisk_count host_count site_id site_name
0 io_grp0        2          8           0          0
1 io_grp1        0          0           0          0
2 io_grp2        0          0           0          0
3 io_grp3        0          0           0          0
4 recovery_io_grp 0          0           0          0
IBM_2145:ITS0-SV1:superuser>lsiogrp 0 | grep fctargetportmode

```

```
fctargetportmode disabled
IBM_2145:ITS0-SV1:superuser>
```

---

If your system is not running with NPIV enabled for host attachment, enable NPIV by completing the following steps after ensuring that you meet the prerequisites:

1. Audit your SAN fabric layout and zoning rules because NPIV has stricter requirements. Ensure that equivalent ports are on the same fabric and in the same zone.
2. Check the path count between your hosts and the IBM Spectrum Virtualize system to ensure that the number of paths is half of the usual supported maximum.

For more information, see the topic about zoning considerations for NPIV in [IBM Knowledge Center](#).

3. Run the `lstargetportfc` command to discover the primary host attach WWPNs (virtual WWPNs), as shown in **bold** in Example 7-7. You can identify these virtual WWPNs because they currently do not allow host I/O or have a nportid assigned because NPIV is not yet enabled.

*Example 7-7 Using the `lstargetportfc` command to get primary host WWPNs (virtual WWPNs)*

---

```
IBM_2145:ITS0-SV1:superuser>lstargetportfc
id WWPN          WWN          port_id owning_node_id current_node_id nportid host_io_permitted virtualized protocol
1  500507680140A288 500507680100A288 1      1              1              010A00 yes             no             scsi
2  500507680142A288 500507680100A288 1      1              1              000000 no             yes            scsi
3  500507680144A288 500507680100A288 1      1              1              000000 no             yes            nvme
4  500507680130A288 500507680100A288 2      1              1              010400 yes             no             scsi
5  500507680132A288 500507680100A288 2      1              1              000000 no             yes            scsi
6  500507680134A288 500507680100A288 2      1              1              000000 no             yes            nvme
7  500507680110A288 500507680100A288 3      1              1              010500 yes             no             scsi
8  500507680112A288 500507680100A288 3      1              1              000000 no             yes            scsi
9  500507680114A288 500507680100A288 3      1              1              000000 no             yes            nvme
10 500507680120A288 500507680100A288 4      1              1              010A00 yes             no             scsi
11 500507680122A288 500507680100A288 4      1              1              000000 no             yes            scsi
12 500507680124A288 500507680100A288 4      1              1              000000 no             yes            nvme
49 500507680C110009 500507680C000009 1      2              2              010500 yes             no             scsi
50 500507680C150009 500507680C000009 1      2              2              000000 no             yes            scsi
51 500507680C190009 500507680C000009 1      2              2              000000 no             yes            nvme
52 500507680C120009 500507680C000009 2      2              2              010400 yes             no             scsi
53 500507680C160009 500507680C000009 2      2              2              000000 no             yes            scsi
54 500507680C1A0009 500507680C000009 2      2              2              000000 no             yes            nvme
55 500507680C130009 500507680C000009 3      2              2              010900 yes             no             scsi
56 500507680C170009 500507680C000009 3      2              2              000000 no             yes            scsi
57 500507680C1B0009 500507680C000009 3      2              2              000000 no             yes            nvme
58 500507680C140009 500507680C000009 4      2              2              010900 yes             no             scsi
59 500507680C180009 500507680C000009 4      2              2              000000 no             yes            scsi
60 500507680C1C0009 500507680C000009 4      2              2              000000 no             yes            nvme
IBM_2145:ITS0-SV1:superuser>
```

---

4. Enable transitional mode for NPIV on IBM Spectrum Virtualize system (see Example 7-8).

*Example 7-8 NPIV in transitional mode*

---

```
IBM_2145:ITS0-SV1:superuser>chiogrp -fctargetportmode transitional 0
IBM_2145:ITS0-SV1:superuser>lsiogrp 0 | grep fctargetportmode
fctargetportmode transitional
IBM_2145:ITS0-SV1:superuser>
```

---

Alternatively, to activate NPIV in transitional mode by using the GUI, select **Settings** → **System** → **I/O Groups** as shown in Figure 7-3.

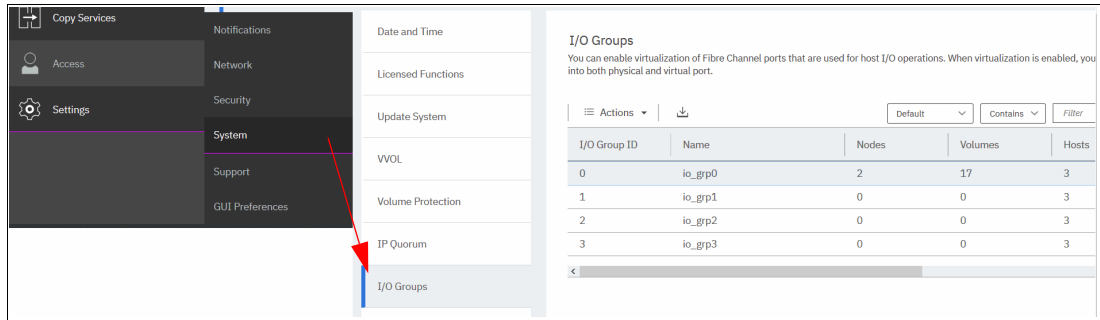


Figure 7-3 I/O Groups Menu

You can then check the current NPIV setting by viewing the **NPIV** column, which shows disabled if NPIV is not yet enabled. Select the I/O group that you want to enable NPIV on and select **Actions** → **Change NPIV Settings**, as shown in Figure 7-4.

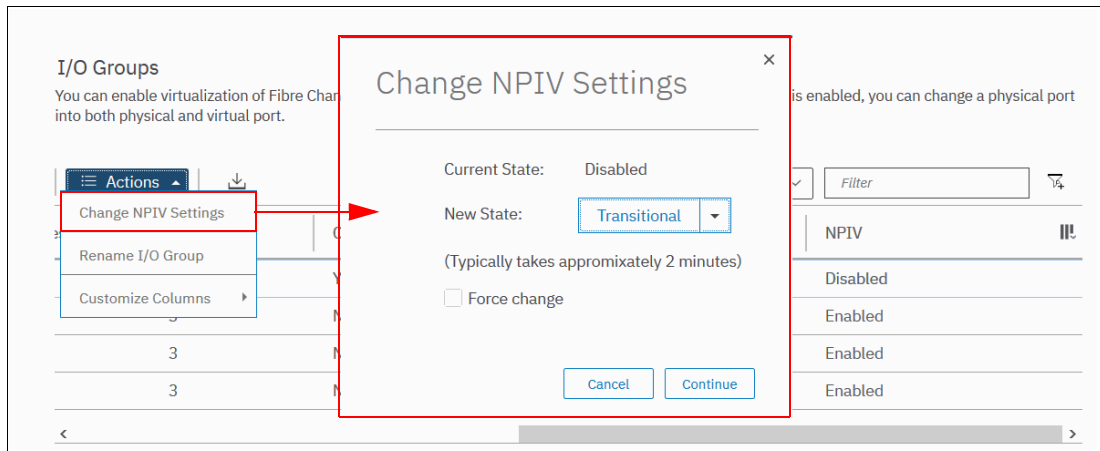


Figure 7-4 Change NPIV Settings

After you select **Continue**, NPIV is enabled in Transitional Mode.

Ensure that the primary host attach WWPNNs (virtual WWPNNs) now allow host traffic, as shown in **bold** in Example 7-9 on page 362.

- Ensure that the primary host attach WWPNs (virtual WWPNs) now allow host traffic, as shown in **bold** in Example 7-9.

*Example 7-9 Host attach WWPNs (virtual WWPNs) permitting host traffic*

```

IBM_2145:ITS0-SV1:superuser>lstargetportfc
id WWPN WWN port_id owning_node_id current_node_id nportid host_io_permitted virtualized protocol
1 500507680140A288 500507680100A288 1 1 1 010A00 yes no scsi
2 500507680142A288 500507680100A288 1 1 1 010A02 yes yes scsi
3 500507680144A288 500507680100A288 1 1 1 010A01 yes yes nvme
4 500507680130A288 500507680100A288 2 1 1 010400 yes no scsi
5 500507680132A288 500507680100A288 2 1 1 010401 yes yes scsi
6 500507680134A288 500507680100A288 2 1 1 010402 yes yes nvme
7 500507680110A288 500507680100A288 3 1 1 010500 yes no scsi
8 500507680112A288 500507680100A288 3 1 1 010501 yes yes scsi
9 500507680114A288 500507680100A288 3 1 1 010502 yes yes nvme
10 500507680120A288 500507680100A288 4 1 1 010A00 yes no scsi
11 500507680122A288 500507680100A288 4 1 1 010A02 yes yes scsi
12 500507680124A288 500507680100A288 4 1 1 010A01 yes yes nvme
49 500507680C110009 500507680C000009 1 2 2 010500 yes no scsi
50 500507680C150009 500507680C000009 1 2 2 010502 yes yes scsi
51 500507680C190009 500507680C000009 1 2 2 010501 yes yes nvme
52 500507680C120009 500507680C000009 2 2 2 010400 yes no scsi
53 500507680C160009 500507680C000009 2 2 2 010401 yes yes scsi
54 500507680C1A0009 500507680C000009 2 2 2 010402 yes yes nvme
55 500507680C130009 500507680C000009 3 2 2 010900 yes no scsi
56 500507680C170009 500507680C000009 3 2 2 010902 yes yes scsi
57 500507680C1B0009 500507680C000009 3 2 2 010901 yes yes nvme
58 500507680C140009 500507680C000009 4 2 2 010900 yes no scsi
59 500507680C180009 500507680C000009 4 2 2 010901 yes yes scsi
60 500507680C1C0009 500507680C000009 4 2 2 010902 yes yes nvme
IBM_2145:ITS0-SV1:superuser>

```

- Add the primary host attach ports (virtual WWPNs) to your host zones, but do not remove the current IBM SAN Volume Controller WWPNs that are in the zones. Example 7-10 shows a host zone to the Primary Port WWPNs of the IBM SAN Volume Controller nodes.

*Example 7-10 Legacy host zone*

```

zone: WINDOWS_HOST_01_IBM_ITSOSV1
      10:00:00:05:1e:0f:81:cc
      50:05:07:68:01:40:A2:88
      50:05:07:68:01:10:A2:88
      50:05:07:68:0C:11:00:09
      50:05:07:68:0C:13:00:09

```

Example 7-11 shows that we added the primary host attach ports (virtual WWPNs) to our example host zone to allow us to change the host without disrupting its availability.

*Example 7-11 Transitional host zone*

```

zone: WINDOWS_HOST_01_IBM_ITSOSV1
      10:00:00:05:1e:0f:81:cc
      50:05:07:68:01:40:A2:88
      50:05:07:68:01:10:A2:88
      50:05:07:68:0C:11:00:09
      50:05:07:68:0C:13:00:09
      50:05:07:68:01:42:A2:88
      50:05:07:68:01:12:A2:88
      50:05:07:68:0C:15:00:09
      50:05:07:68:0C:17:00:09

```

7. With the transitional zoning active in your fabrics, ensure that the host uses the new NPIV ports for host I/O. Example 7-12 shows the before and after pathing for our host. Notice that the select count now increases on the new paths and stopped on the old paths.

*Example 7-12 Host device pathing: Before and after*

C:\Program Files\IBM\SDDDSM>datapath query device

Total Devices : 1

DEV#: 0 DEVICE NAME: Disk2 Part0 TYPE: 2145 POLICY: OPTIMIZED  
 SERIAL: 600507680C838020E800000000000002 LUN SIZE: 20.0GB

```
=====
```

Path#	Adapter/Hard Disk	State	Mode	Select	Errors
0	Scsi Port2 Bus0/Disk2 Part0	OPEN	NORMAL	10626	0
1	Scsi Port2 Bus0/Disk2 Part0	OPEN	NORMAL	10425	0
2 *	Scsi Port2 Bus0/Disk2 Part0	OPEN	NORMAL	0	0
3 *	Scsi Port2 Bus0/Disk2 Part0	OPEN	NORMAL	0	0
4	Scsi Port3 Bus0/Disk2 Part0	OPEN	NORMAL	1128804	0
5	Scsi Port3 Bus0/Disk2 Part0	OPEN	NORMAL	1129439	0
6 *	Scsi Port3 Bus0/Disk2 Part0	OPEN	NORMAL	0	0
7 *	Scsi Port3 Bus0/Disk2 Part0	OPEN	NORMAL	0	0

C:\Program Files\IBM\SDDDSM>datapath query device

Total Devices : 1

DEV#: 0 DEVICE NAME: Disk2 Part0 TYPE: 2145 POLICY: OPTIMIZED  
 SERIAL: 600507680C838020E800000000000002 LUN SIZE: 20.0GB

```
=====
```

Path#	Adapter/Hard Disk	State	Mode	Select	Errors
0 *	Scsi Port2 Bus0/Disk2 Part0	OPEN	NORMAL	10630	1
1 *	Scsi Port2 Bus0/Disk2 Part0	OPEN	NORMAL	10427	1
2 *	Scsi Port2 Bus0/Disk2 Part0	OPEN	NORMAL	0	0
3 *	Scsi Port2 Bus0/Disk2 Part0	OPEN	NORMAL	0	0
4 *	Scsi Port3 Bus0/Disk2 Part0	OPEN	NORMAL	1128809	2
5 *	Scsi Port3 Bus0/Disk2 Part0	OPEN	NORMAL	1129445	1
6 *	Scsi Port3 Bus0/Disk2 Part0	OPEN	NORMAL	0	0
7 *	Scsi Port3 Bus0/Disk2 Part0	OPEN	NORMAL	0	0
<b>8</b>	<b>Scsi Port3 Bus0/Disk2 Part0</b>	<b>OPEN</b>	<b>NORMAL</b>	<b>76312</b>	<b>0</b>
9	Scsi Port3 Bus0/Disk2 Part0	OPEN	NORMAL	76123	0
<b>10 *</b>	<b>Scsi Port3 Bus0/Disk2 Part0</b>	<b>OPEN</b>	<b>NORMAL</b>	<b>0</b>	<b>0</b>
11 *	Scsi Port3 Bus0/Disk2 Part0	OPEN	NORMAL	0	0
<b>12</b>	<b>Scsi Port2 Bus0/Disk2 Part0</b>	<b>OPEN</b>	<b>NORMAL</b>	<b>623</b>	<b>0</b>
13	Scsi Port2 Bus0/Disk2 Part0	OPEN	NORMAL	610	0
<b>14 *</b>	<b>Scsi Port2 Bus0/Disk2 Part0</b>	<b>OPEN</b>	<b>NORMAL</b>	<b>0</b>	<b>0</b>
15 *	Scsi Port2 Bus0/Disk2 Part0	OPEN	NORMAL	0	0

```
=====
```

**Remember:** Consider the following points:

- ▶ You can verify that you are logged in to the NPIV ports by running the `lsfabric -host host_id_or_name` command. If I/O activity is occurring, each host has at least one line in the command output that corresponds to a host port and shows `active` in the activity field:
  - Hosts where no I/O occurred in the past 5 minutes show `inactive` for any login.
  - Hosts that do not adhere to preferred paths might still be processing I/O to primary ports.
- ▶ Depending on the host operating system, rescanning for storage might be required on some hosts to recognize extra paths that are now provided by using primary host attach ports (virtual WWPNS).

8. After all hosts are rezoned and the pathing is validated, change the system NPIV to enabled mode by running the command that is shown in Example 7-13.

*Example 7-13 Enabling NPIV*

```
IBM_2145:ITS0-SV1:superuser>chiogrp -fctargetportmode enabled 0
```

Alternatively, to Enable NPIV by way of the Storwize GUI, browse to the I/O Group window as shown in Step 4, select the wanted I/O group, click on **Actions** and then **Change NPIV Settings**, which opens the NPIV Settings window, as shown in Figure 7-5.

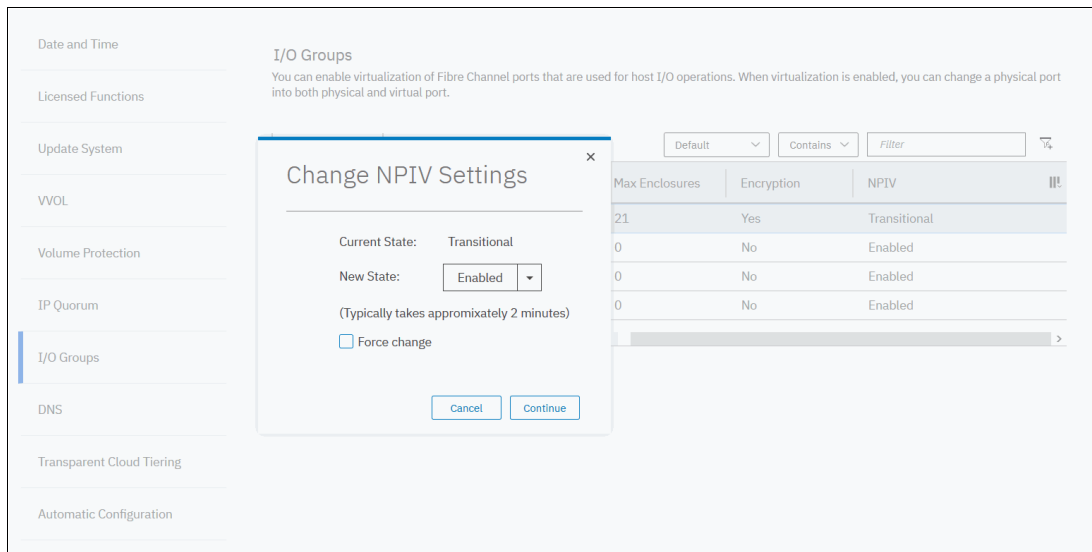


Figure 7-5 Change NPIV Settings - Enabling

NPIV is enabled on the IBM Spectrum Virtualize system, and you confirmed that the hosts use the virtualized WWPNS for I/O. To complete the NPIV implementation, the host zones can be amended to remove the old primary attach port WWPNS. Example 7-14 shows our final zone with the host HBA and the IBM SAN Volume Controller virtual WWPNS.

*Example 7-14 Final host zone*

```
zone: WINDOWS_HOST_01_IBM_ITSOSV1
      10:00:00:05:1e:0f:81:cc
      50:05:07:68:01:42:A2:88
      50:05:07:68:01:12:A2:88
```

50:05:07:68:0C:15:00:09  
50:05:07:68:0C:17:00:09

**Note:** If there are still hosts that are configured to use the physical ports on the IBM SAN Volume Controller system, the system prevents you from changing `fctargetportmode` from `transitional` to `enabled`, and shows the following error:

CMMVC8019E Task could interrupt IO and force flag not set.

## 7.5 Hosts operations by using the GUI

This section describes performing the following host operations by using the IBM Spectrum Virtualize GUI:

- ▶ Creating hosts
- ▶ Advanced host administration
- ▶ Adding and deleting host ports
- ▶ Host mappings overview

### 7.5.1 Creating hosts

This section describes how to create FC and iSCSI hosts by using the IBM Spectrum Virtualize GUI. It is assumed that hosts are prepared for attachment, as described in IBM Knowledge Center, and that the host WWPNs or their iSCSI initiator names are known.

For more information, see the Host Attachment section of [IBM Knowledge Center](#).

To create a host, complete the following steps:

1. Open the host configuration window by clicking **Hosts** (see Figure 7-6).

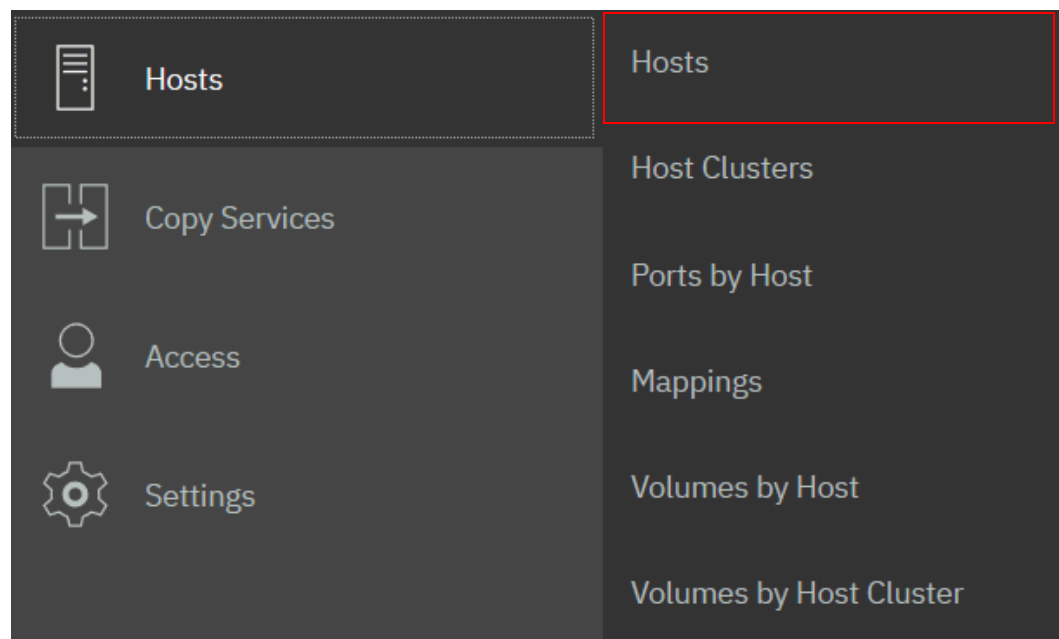


Figure 7-6 Opening the host window

2. To create a host, click **Add Host**. If you want to create an FC host, continue with “Creating Fibre Channel hosts” on page 366. To create an iSCSI host, see “Creating iSCSI hosts” on page 370.

### Creating Fibre Channel hosts

To create FC hosts, complete the following steps:

1. Select **Fibre Channel**. The FC host configuration window opens (see Figure 7-7).

**Add Host**

Because NPIV is enabled on this system, host traffic is only allowed over the storage system's virtual ports. Ensure that SAN zoning allows connectivity between virtual ports and the host.

**Required Fields**

Name:

Host connections:

Host port (WWPN):

**Optional Fields**

Host type:

I/O groups:

Host cluster:

Ownership Group:

Figure 7-7 Fibre Channel host configuration



2. Enter a name for your host and click the **Host Port (WWPN)** menu to get a list of all discovered WWPNs (see Figure 7-8).

The screenshot shows a dialog box titled "Add Host" with a close button (X) in the top right corner. Below the title is a horizontal line. An information icon (i) is followed by a text block: "Because NPIV is enabled on this system, host traffic is only allowed over the storage system's virtual ports. Ensure that SAN zoning allows connectivity between virtual ports and the host." Below this is a section titled "Required Fields" with the following fields: "Name:" with the value "Windows-Host-01"; "Host connections:" with a dropdown menu showing "Fibre Channel (SCSI)"; and "Host port (WWPN):" with a dropdown menu showing a list of WWPNs: "10000090FADD245B" and "10000090FADD245C". To the right of the WWPN list are plus (+) and minus (-) icons. Below the required fields is a section titled "Optional Fields" with the following fields: "Host type:" with a dropdown menu showing "Generic"; "I/O groups:" with a dropdown menu showing "All"; "Host cluster:" with a dropdown menu showing "No Host Cluster Selected"; and "Ownership Group:" with a dropdown menu showing "No ownership group selected". At the bottom right of the dialog box are "Cancel" and "Add" buttons.

Figure 7-8 Available WWPNs

3. Select one or more WWPNs for your host. The IBM Spectrum Virtualize system should have the host port WWPNs available if the host is prepared. If they do not appear in the list, scan for new disks as required on the respective operating system and click the **Rescan** icon in the WWPN box. If they still do not appear, check the SAN zoning and repeat the scanning.

**Creating offline hosts:** If you want to create a host that is offline or not connected at the moment, it is also possible to enter manually the WWPNs. enter them into the Host Ports field to add them to the host.

4. If you want to add ports to your host, click the Plus sign (+) to add all the ports that belong to the specific host.

5. If you are creating a Hewlett-Packard UNIX (HP-UX) or Target Port Group Support (TPGS) host, click the **Host Type** menu (see Figure 7-9). Select your host type. If your specific host type is not listed, select **Generic**.

The screenshot shows a window titled "Add Host" with a close button (X) in the top right corner. Below the title is a horizontal line. An information icon (i) is followed by a text block: "Because NPIV is enabled on this system, host traffic is only allowed over the storage system's virtual ports. Ensure that SAN zoning allows connectivity between virtual ports and the host." Below this is a section titled "Required Fields" containing three input fields: "Name:" with the value "Windows-Host-01", "Host connections:" with a dropdown menu showing "Fibre Channel (SCSI)", and "Host port (WWPN):" with the value "1000090FADD245B" and refresh (+) and minus (-) icons. Below the required fields is a section titled "Optional Fields" containing four labels: "Host type:", "I/O groups:", "Host cluster:", and "Ownership Group:". The "Host type:" dropdown menu is open, showing a list of options: "Generic", "HP/UX", "OpenVMS", "TPGS", and "VVOL". The "Generic" option is highlighted. At the bottom right of the form is an "Add" button.

Figure 7-9 Host type selection

6. If you set up object-based access control (OBAC) as described in Chapter 11, "Ownership groups" on page 673, select the ownership group that you want the host to be a part of from the **Ownership Group** drop-down menu, as shown in Figure 7-10 on page 369.

## Add Host ✕

---

**i** Because NPIV is enabled on this system, host traffic is only allowed over the storage system's virtual ports. Ensure that SAN zoning allows connectivity between virtual ports and the host.

**Required Fields**

Name:

Host connections:  ▼

Host port (WWPN):

---

**Optional Fields**

Host type:  ▼

I/O groups:  ▼

Host cluster:  ▼

Ownership Group:  ▼

**No ownership group selected**

ownershipgroup0

Figure 7-10 Adding the new host to an ownership group

7. Click **Add** to create the host object.
8. Click **Close** to return to the host window. Repeat these steps for all of your FC hosts. Figure 7-11 shows the Hosts window after creating a host.

Name	Status	Host Type	# of Ports	Host Mappings	Host Cluster ID	Host Cluster Name
Windows-Host-01	✔ Online	Generic	2	No		

Figure 7-11 Hosts view after creating a host

After you add FC hosts, see Chapter 6, “Volumes” on page 255 to create volumes and map them to the created hosts.

## Creating iSCSI hosts

When creating an iSCSI attached host, consider the following points:

- ▶ iSCSI IP addresses can fail over to the partner node in the I/O group if a node fails. This design reduces the need for multipathing support in the iSCSI host.
- ▶ The iSCSI Qualified Name (IQN) of the host is added to an IBM Spectrum Virtualize host object in the same way that you add FC WWPNs.
- ▶ Host objects can have WWPNs and IQNs.
- ▶ Standard iSCSI host connection procedures can be used to discover and configure IBM Spectrum Virtualize as an iSCSI target.
- ▶ IBM Spectrum Virtualize supports the Challenge Handshake Authentication Protocol (CHAP) authentication methods for iSCSI.
- ▶ The name `iqn.1986-03.com.ibm:2145.<cluster_name>.<node_name>` is the IQN for an IBM Spectrum Virtualize node. Because the IQN contains the clustered system name and the node name, it is important not to change these names after iSCSI is deployed.
- ▶ Each node can be given an iSCSI alias as an alternative to the IQN.

To create iSCSI hosts, complete the following steps:

1. Click **iSCSI** and the iSCSI configuration window opens (Figure 7-12).

**Add Host**

**i** Because NPIV is enabled on this system, host traffic is only allowed over the storage system's virtual ports. Ensure that SAN zoning allows connectivity between virtual ports and the host.

**Required Fields**

Name: VMware-Host-01

Host connections: iSCSI (SCSI)

Host IQN: iqn.1998-01.com.vmware:esx6-8hq

**Optional Fields**

CHAP authentication:

CHAP secret: Enter 1 to 79 characters

CHAP username: Enter 1 to 31 characters

Host type: Generic

I/O groups: All

Host cluster: No Host Cluster Selected

Ownership Group: No ownership group selected

Cancel Add

Figure 7-12 Adding an iSCSI host

2. Enter a host name and the iSCSI initiator name into the **iSCSI host IQN** field. Click the plus sign (+) if you want to add initiator names to one host.
3. If you are connecting an HP-UX or TPGS host, click the **Host type** menu and then select the correct host type. For our ESX host, we selected **VVOL**. However, **generic** is good if you are not using VMware vSphere Virtual Volumes (VVOLs).
4. Click **Add** and then, click **Close** to complete the host object definition.
5. Repeat these steps for every iSCSI host that you want to create. Figure 7-13 shows the Hosts window after creating two FC hosts and one iSCSI host.

Name	Status	Host Type	# of Ports	Host Mappings	Host Cluster ID	Host Cluster Name
RHEL-Host-01	Online	Generic	2	No		
VMware-Host-01	Offline	Generic	1	No		
Windows-Host-01	Online	Generic	2	No		

Figure 7-13 Defined Hosts list

Although the iSCSI host is now configured to provide connectivity, the iSCSI Ethernet ports must also be configured.

Complete the following steps to enable iSCSI connectivity:

1. Select **Settings** → **Network** and select the iSCSI tab (see Figure 7-14).

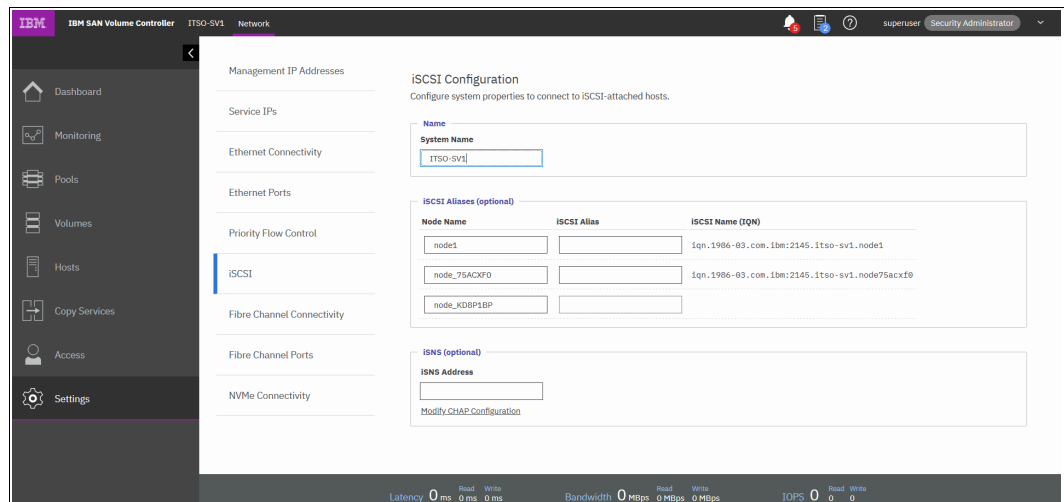


Figure 7-14 Network: iSCSI settings

- In the iSCSI Configuration window, you can modify the system name, node names, and provide optional iSCSI Alias for each node if you want (see Figure 7-15).

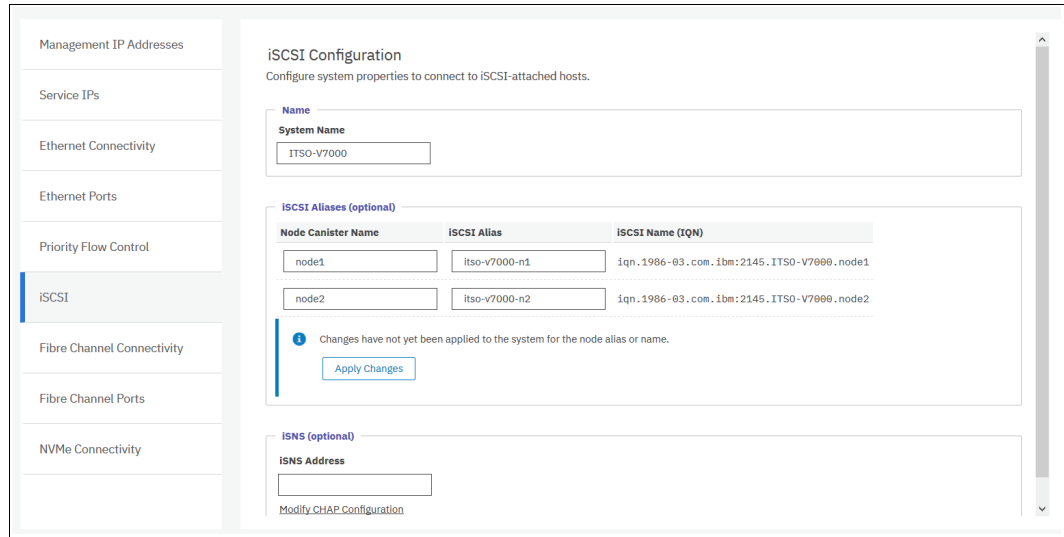


Figure 7-15 iSCSI Configuration window

- The window displays an Apply Changes prompt to apply any changes that you made before continuing.
- In the lower part of the configuration window, you can configure internet Storage Name Service (iSNS) addresses and CHAP if you need them in your environment.

**Note:** The authentication of hosts is optional. By default, it is disabled. The user can choose to enable CHAP or *CHAP authentication*, which involves sharing a CHAP secret between the cluster and the host. If the correct key is not provided by the host, IBM Spectrum Virtualize does not allow it to perform I/O to volumes. Also, you can assign a CHAP secret to the cluster.

- Click the **Ethernet Ports** tab to see the list of ports to configure iSCSI IPs (see Figure 7-16).

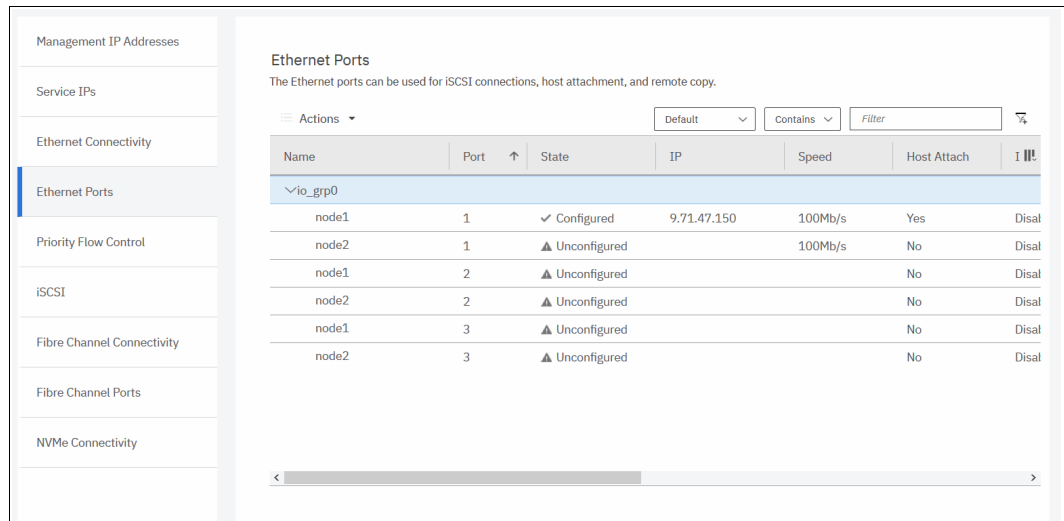


Figure 7-16 Ethernet port list

6. Select the port to set the iSCSI IP information, click **Actions** → **Modify IP Settings** (see Figure 7-17).

Modify Port 1 of Node node2 ×

---

**IPv4 address:**

**Subnet mask:**

**Gateway:**

▶ **IPv6**

Figure 7-17 Modifying the IP settings

7. After you enter the IP address for a port, click **Modify** to enable the configuration. After the changes are successfully applied, click **Close**.

- You can see that iSCSI is enabled for host I/O on the required interfaces (Yes) under the Host Attach column, as shown in Figure 7-18.

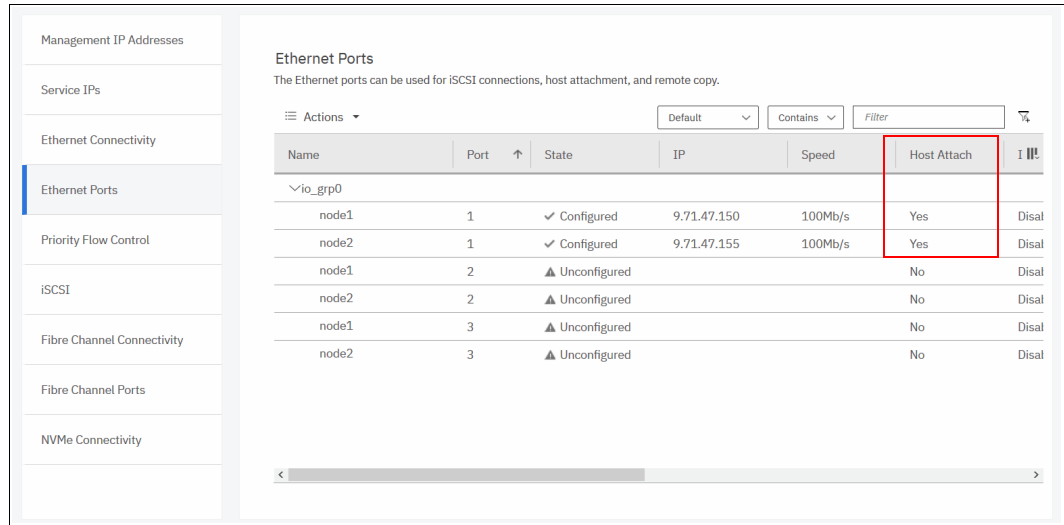


Figure 7-18 Host attach permitted on port

- By default, iSCSI host connection is enabled after setting the IP address. You can enable or disable Internet Protocol Version 4 (IPv4) or Internet Protocol Version 6 (IPv6) iSCSI host to use the port by clicking **Actions** → **Modify iSCSI Hosts** (see Figure 7-19).

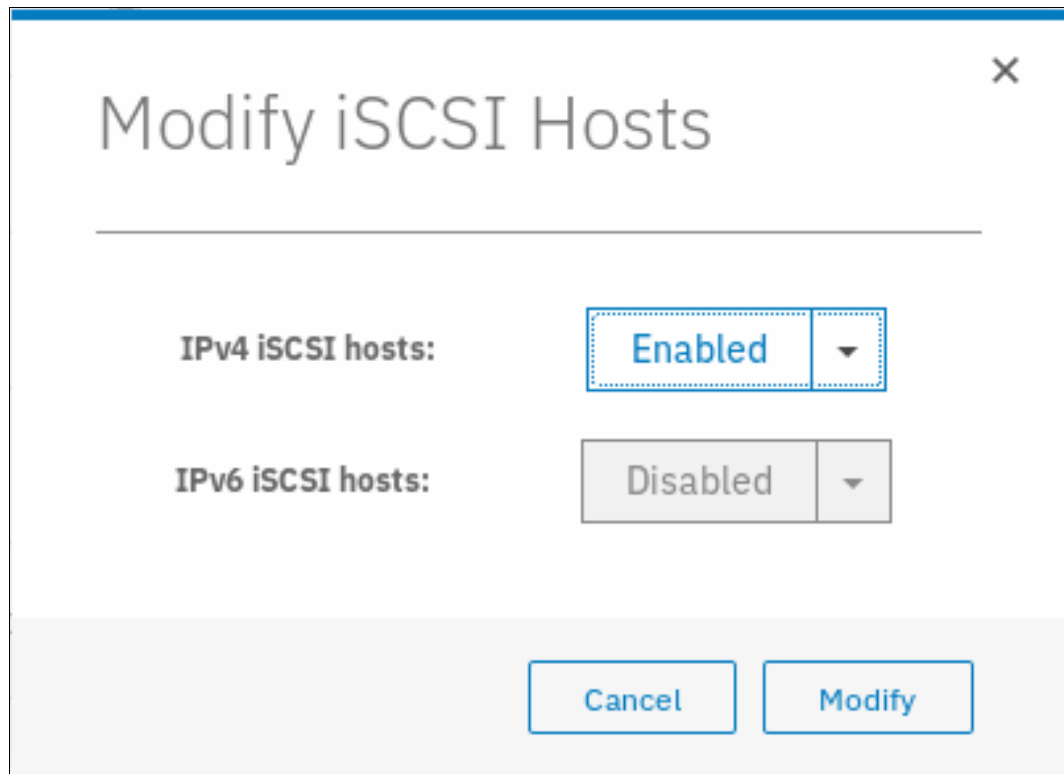


Figure 7-19 Modify iSCSI host connections



10. Use the iSCSI network in a separate subnet. It is also possible to set a virtual local area network (VLAN) for the iSCSI traffic. To enable the VLAN, click **Actions** → **Modify VLAN**, as shown in Figure 7-20.

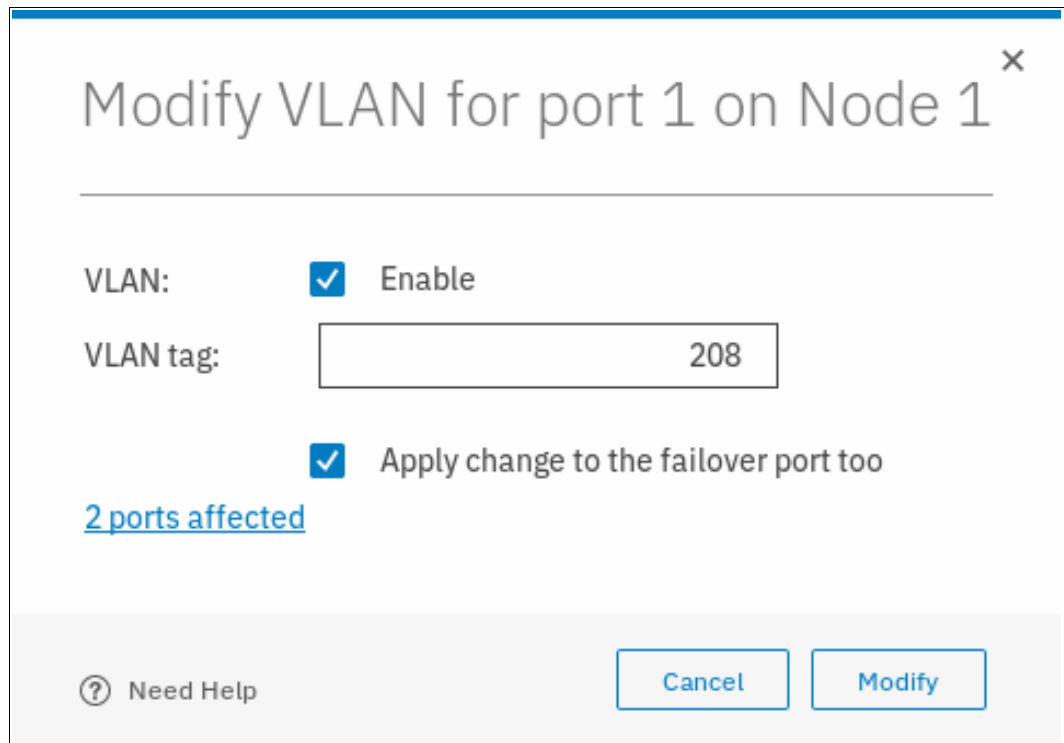


Figure 7-20 Modifying the VLAN

The IBM SAN Volume Controller system is now configured and ready for iSCSI host use. Note the initiator IQN names of your IBM SAN Volume Controller nodes (see Figure 7-15 on page 372) because you need them when adding storage on your host.

For more information about how to create volumes and map them to a host, see Chapter 6, “Volumes” on page 255.

## Creating NVMe hosts

To configure a NVMe host, complete the following steps:

1. Go to the host window, and click **Add Host**. In **Host connections**, select **NVMe**, as shown in Figure 7-21.

**Add Host** [X]

**Info** Because NPIV is enabled on this system, host traffic is only allowed over the storage system's virtual ports. Ensure that SAN zoning allows connectivity between virtual ports and the host.

**Required Fields**

Name:

Host connections:  ▼

Host port (NQN):  (+) (-)

**Optional Fields**

Host type:  ▼

I/O groups:  ▼

Host cluster:  ▼

Ownership Group:  ▼

Figure 7-21 Host connection: NVMe

2. Enter the host name and the NVMe Qualified Name (NQN) of the host, as shown in Figure 7-22.

Figure 7-22 Defining NQN

3. Click **Add**. Your host is shown in the defined host list.
4. The I/O group NQN must be configured on the host. To find the I/O group NQN, run the **lsiogrp** command, as shown in Example 7-15.

*Example 7-15 The lsiogrp command*

```

IBM_2145:ITS0-SV1:superuser>lsiogrp 0
id 0
name io_grp0
node_count 2
vdisk_count 8
host_count 3
flash_copy_total_memory 20.0MB
flash_copy_free_memory 19.9MB
remote_copy_total_memory 20.0MB
remote_copy_free_memory 20.0MB
mirroring_total_memory 20.0MB

```

```
mirroring_free_memory 19.9MB
raid_total_memory 40.0MB
raid_free_memory 40.0MB
maintenance no
compression_active no
accessible_vdisk_count 8
compression_supported no
max_enclosures 20
encryption_supported yes
flash_copy_maximum_memory 2048.0MB
site_id
site_name
fctargetportmode disabled
compression_total_memory 0.0MB
deduplication_supported yes
deduplication_active no
nqn nqn.1986-03.com.ibm:nvme:2145.000020067214511.iogroup0
IBM_2145:ITS0-SV1:superuser>
```

---

You can now configure your NVMe host to use the SAN Volume Controller system as a target.

**Note:** For more information about a compatibility matrix and supported hardware, see [IBM Knowledge Center](#) and the [IBM SSIC](#).

## 7.5.2 Host clusters

IBM Spectrum Virtualize V7.7 introduced the concept of a *host cluster*. A host cluster allows a user to group individual hosts to form a cluster, which is treated as one entity instead of dealing with all of the hosts individually in the cluster.

The host cluster is useful for hosts that are participating in a cluster at host operating system levels. Examples are Microsoft Clustering Server, IBM PowerHA, Red Hat Cluster Suite, and VMware ESX. By defining a host cluster, a user can map one or more volumes to the host cluster object.

As a result, the volume or set of volumes is mapped to each individual host object that is part of the host cluster. Each of the volumes is mapped by using a single command with the same SCSI ID to each host that is part of the host cluster.

Although a host is part of a host cluster, volumes can still be assigned to an individual host in a non-shared manner. A policy can be devised that can pre-assign a standard set of SCSI IDs for volumes to be assigned to the host cluster, and devise another set of SCSI IDs to be used for individual assignments to hosts.

**Note:** For example, SCSI IDs 0 - 100 for individual host assignment and SCSI IDs above 100 can be used for a host cluster. By using such a policy, wanted volumes cannot be shared, and others can be shared. For example, the boot volume of each host can be kept private while data and application volumes can be shared.

## Creating a host cluster

This section describes how to create a host cluster. It is assumed that individual hosts are already created, as described in 7.5.1, “Creating hosts” on page 365.

Complete the following steps:

1. From the menu on the left, select **Hosts** → **Host Clusters** (see Figure 7-23).

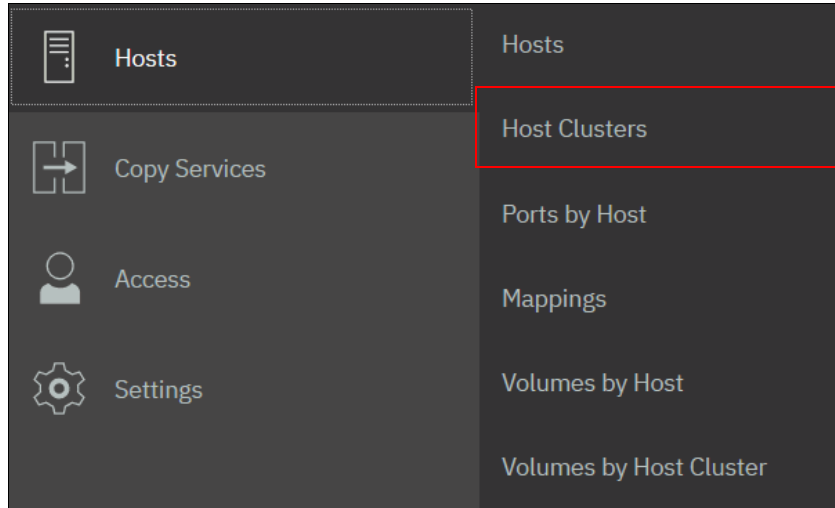


Figure 7-23 Host clusters

2. Click **Create Host Cluster**, as shown in Figure 7-24.

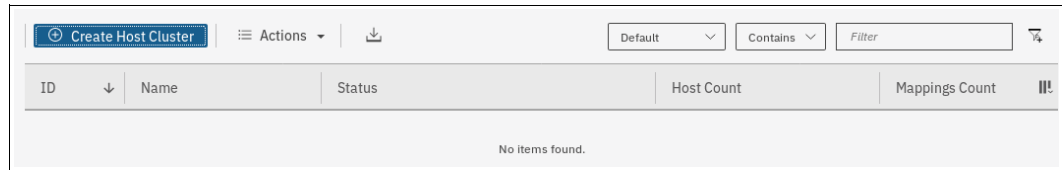


Figure 7-24 Creating a host cluster

3. Enter a cluster name and select the individual hosts that you want in the cluster object, as shown in Figure 7-25.

## Create Host Cluster ✕

---

Name:

**Optional:** Select hosts to assign to a new host cluster. Any current volume mappings become the shared mappings for all the hosts in the host cluster.

i It is recommended that all hosts in a host cluster have access to the same I/O Groups.

Ownership Group:  ▼

↓

Default ▼

Contains ▼

🔍

Name	Status	Host Type	Host Mappings	C
Windows-Host-01	✓ Online	Generic	No	owne
Windows-Host-02	✓ Online	Generic	No	owne

<
>

Showing 2 hosts | Selecting 2 hosts

? Cancel
◀ Back
Next ▶

Figure 7-25 Creating a host cluster: Details

4. Click **Next**. A summary window opens, in which you confirm that you selected the correct hosts. Click **Make Host Cluster** (see Figure 7-26).

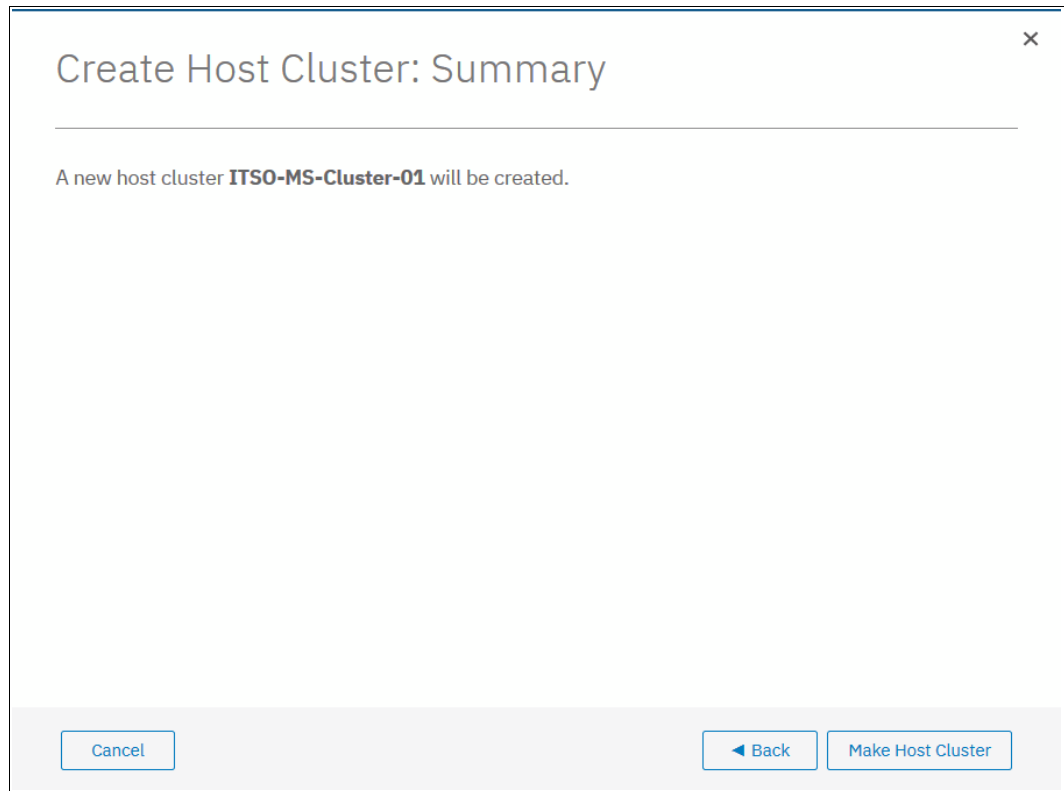


Figure 7-26 Create host cluster summary

5. After the task completes, click **Close** to return to the Host Cluster view, where you can see the cluster that you created (see Figure 7-27).

ID	Name	Status	Host Count	Mappings Count	Ports Count
0	ITSO-MS-Cluster-01	Online	2	0	2

Figure 7-27 Host Cluster view

**Note:** The host cluster status depends on its member hosts. One offline or degraded host sets the host cluster status as Degraded.

From the Host Clusters view, you have many options to manage and configure the host cluster. These options are accessed by selecting a cluster and clicking **Actions** (see Figure 7-28).

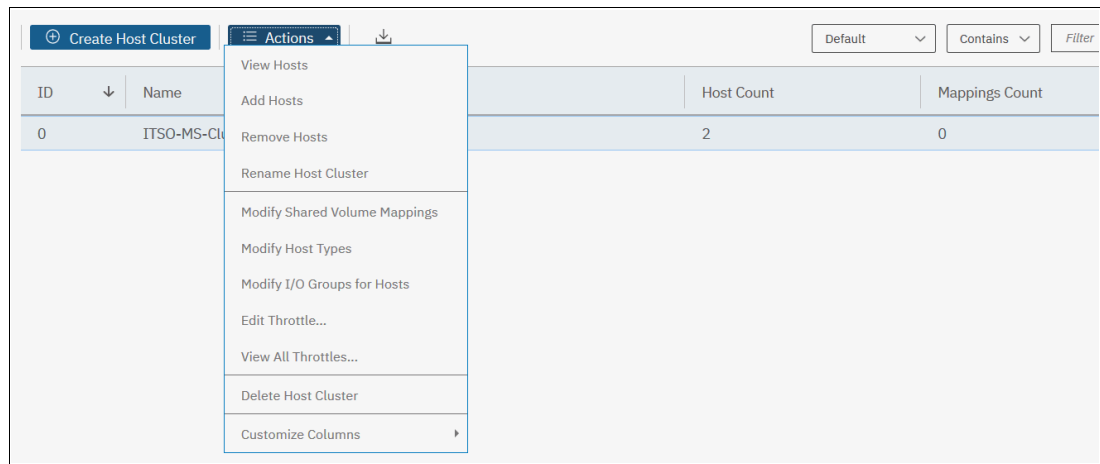


Figure 7-28 Host Clusters Actions menu

From the **Actions** menu, the following tasks can be performed:

- ▶ View Hosts status within the cluster.
- ▶ Add or Remove Hosts from the cluster.
- ▶ Rename the host cluster.
- ▶ Modify Shared Volume mappings allows you to add or remove volumes that are mapped to all hosts in the cluster while maintaining the same SCSI ID for all hosts.
- ▶ Modify Host Type can be used to change from generic to VVOLs, as an example.
- ▶ Modify I/O Groups for Hosts is used to assign or restrict volume access to specific I/O groups.
- ▶ Edit Throttle is used to restrict MBps or input/output operations per second (IOPS) bandwidth for the host cluster.
- ▶ View All Throttles displays any throttling settings and allows for changing, deleting, or refining Throttle settings.
- ▶ Delete Host Cluster to delete the cluster entity. When deleting the host cluster, you can choose to keep the mappings on the hosts.
- ▶ Customize Columns modifies which columns are displayed that show the properties of the host cluster

### 7.5.3 Advanced host administration

This section covers host administration, including topics such as host modification, host mappings, and deleting hosts. Basic host creation by using FC and iSCSI connectivity is described in 7.5.1, “Creating hosts” on page 365.



It is assumed that a few hosts are created by using the IBM Spectrum Virtualize GUI and that some volumes are already mapped to them. This section describes three functions that are covered in the Hosts section of the IBM Spectrum Virtualize GUI (Figure 7-29):

- ▶ Hosts (“Modifying mappings” on page 384)
- ▶ Ports by Host (7.5.4, “Adding and deleting host ports” on page 400)
- ▶ Host Mappings (7.5.5, “Host mappings overview” on page 410)

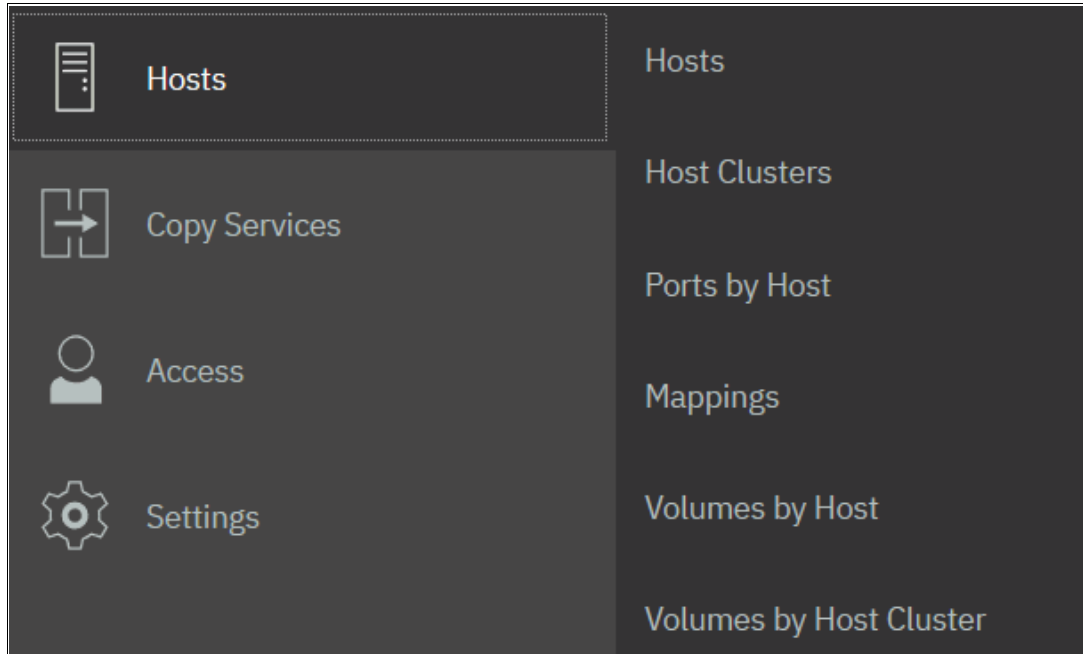


Figure 7-29 IBM Spectrum Virtualize Hosts menu

In the **Hosts** → **Hosts** view, three hosts are created and the volumes are mapped to them in our example. If needed, you can now modify these hosts by selecting a host and clicking **Actions** or right-clicking the host to see the available tasks (see Figure 7-30).

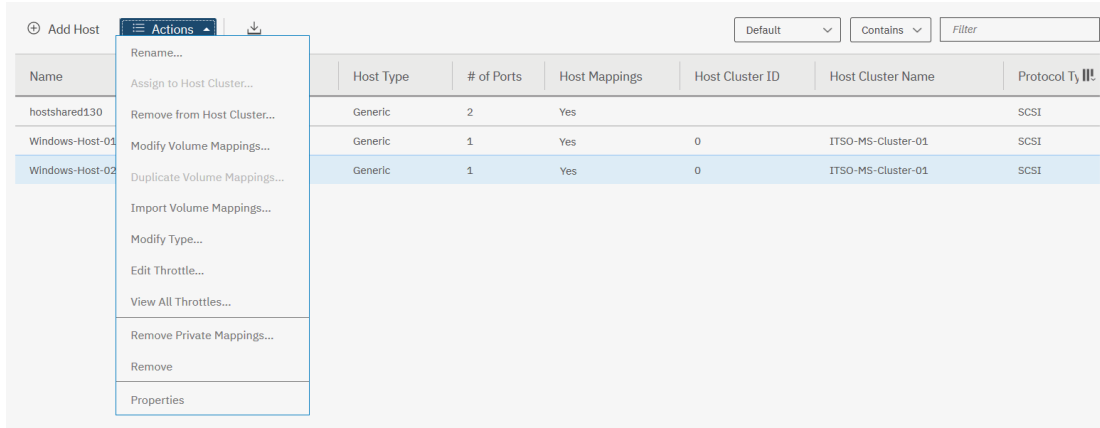


Figure 7-30 Host actions

## Modifying mappings

To modify what volumes are mapped to a specific host, complete the following steps:

1. Select **Actions** → **Modify Volume Mappings** (Figure 7-30). The window that is shown in Figure 7-31 opens. At the upper left, you can confirm that the correct host is targeted. The list shows all volumes that are mapped to the selected host. In our example, one volume with SCSI ID 0 is mapped to the host ITSO-VMHOST-01.

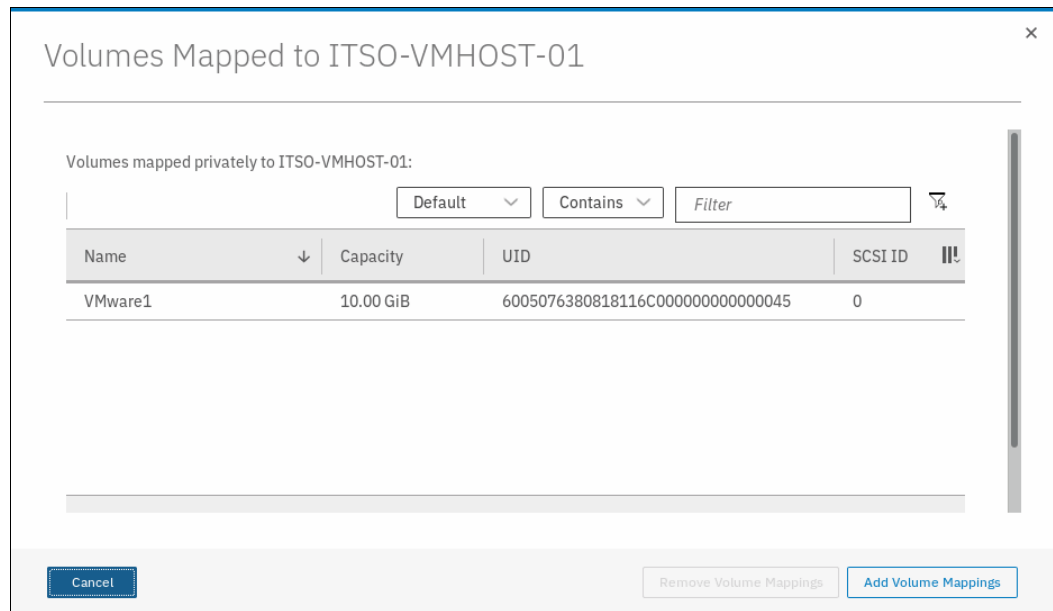


Figure 7-31 Modifying the host volume mappings

**Note:** When you change volume mappings in a **host cluster**, the changes apply to the shared mappings only. For example, when a volume mapping is added to the host cluster, it becomes a shared mapping among all the hosts within the cluster. When a mapping is removed, the volume is removed from the shared mappings for the host cluster. However, you can select specific hosts that retain access to that volume as a private mapping.

2. By selecting a listed volume, you can remove that volume map from the host. However, in our case we want to add a volume to our host. Continue by clicking **Add Volume Mapping** (see Figure 7-32).

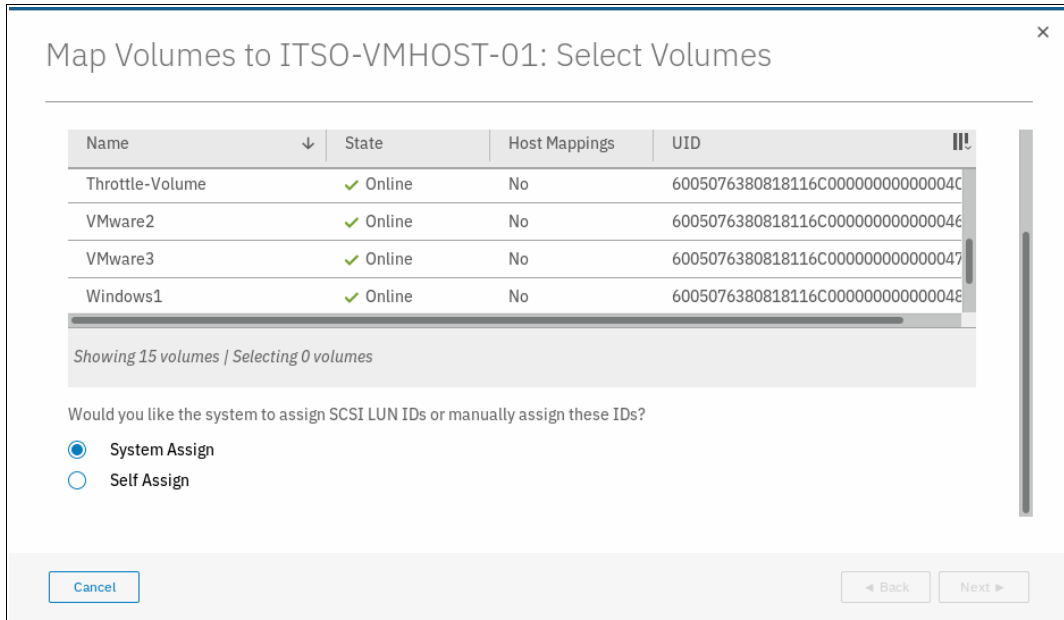


Figure 7-32 Volumes selection list

A new list opens that shows all volumes. You can easily identify whether a volume you want to map is mapped to another host, as shown in Figure 7-33.

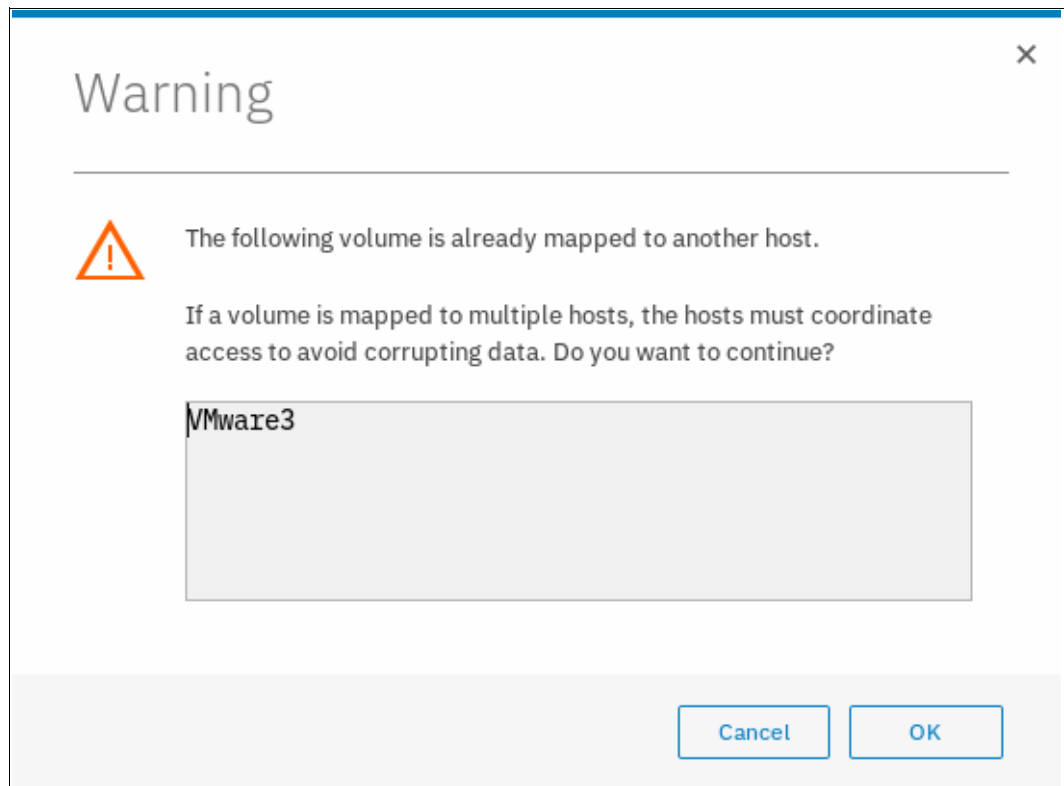


Figure 7-33 Volume mapped to another host warning

- To map a volume, select it and click **Next** to map it to the host. The volume is assigned the next available SCSI ID if you leave **System Assign** selected. However, by selecting **Self Assign**, you can manually set the SCSI IDs (see Figure 7-34).

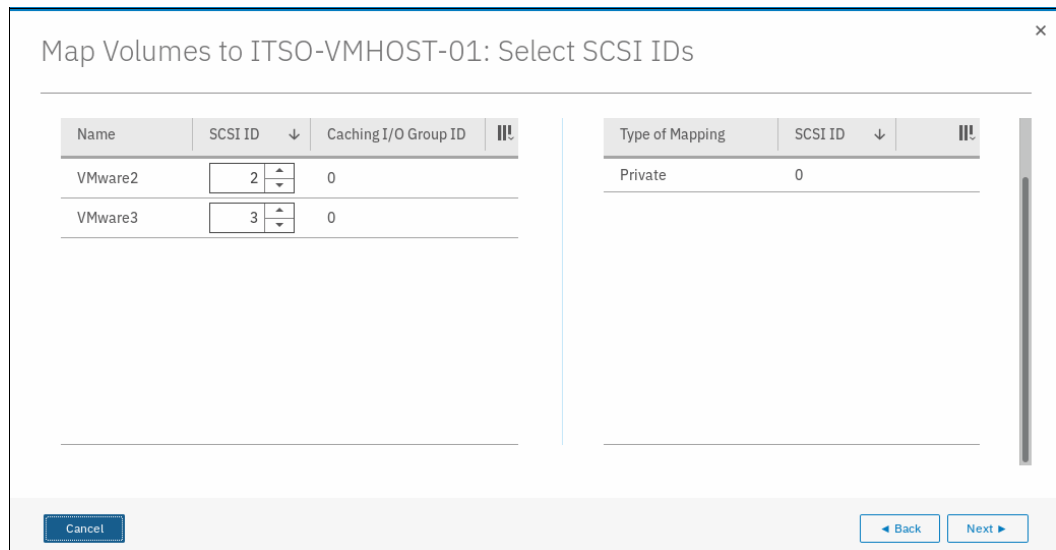


Figure 7-34 Modify Host Volume Mappings: Assigning a SCSI ID

If you select a SCSI ID that is in use for the host, you cannot proceed. In Figure 7-34, we selected SCSI ID 0. However, in the right column you can see SCSI ID 0 is already allocated. By changing to SCSI ID 1, we can click **Next**.

- A summary window opens that shows the new mapping details (see Figure 7-35). After confirming that this mapping is what you planned, click **Map Volumes** and then, click **Close**.

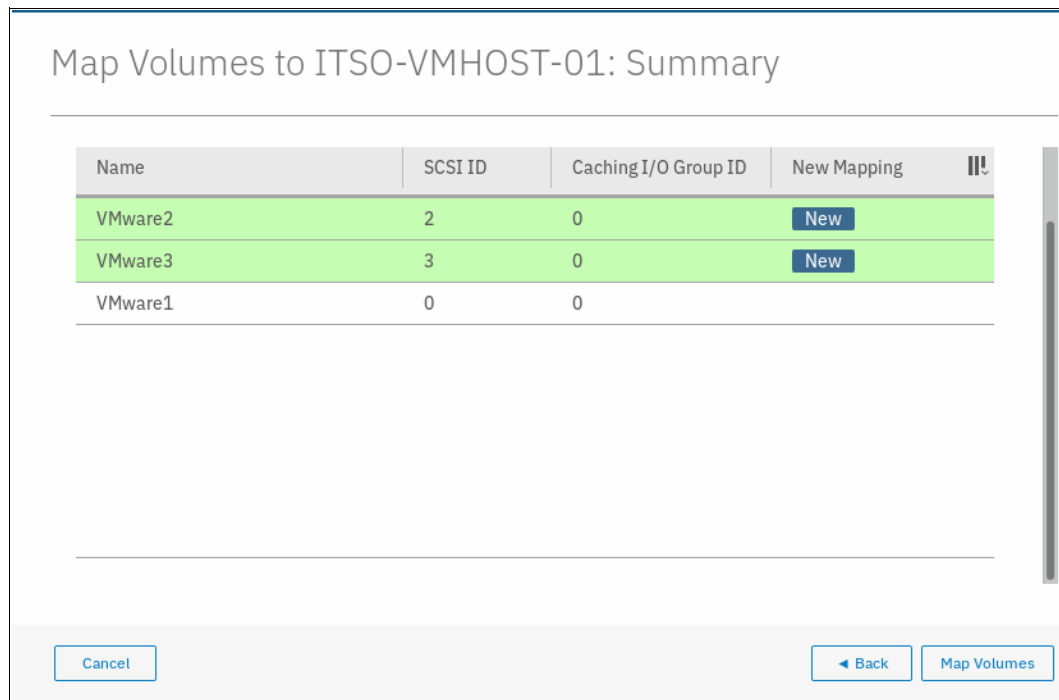


Figure 7-35 Confirming the modified mappings

**Note:** The SCSI ID of the volume can be changed only before it is mapped to a host. Changing it later is not possible unless the volume is unmapped again.

## Removing private mappings from a host

A host can access only those volumes on an IBM SAN Volume Controller system that are mapped to it. If you want to remove access to all volumes for one host regardless of how many volumes are mapped to it, complete the following steps:

1. From the Hosts pane, select the host and click **Actions** → **Remove Private Mappings** to remove all access that the selected host has to its volumes (see Figure 7-36).

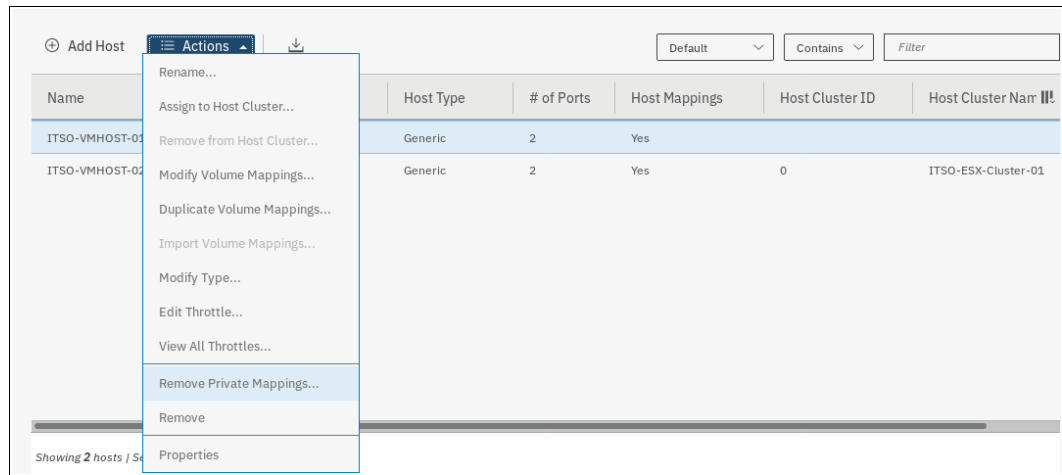


Figure 7-36 Unmapping all volumes action

- You are prompted to confirm the number of mappings to be removed. To confirm your action, enter the number of volumes to be removed and click **Remove** (see Figure 7-37). In this example, we remove three volume mappings.

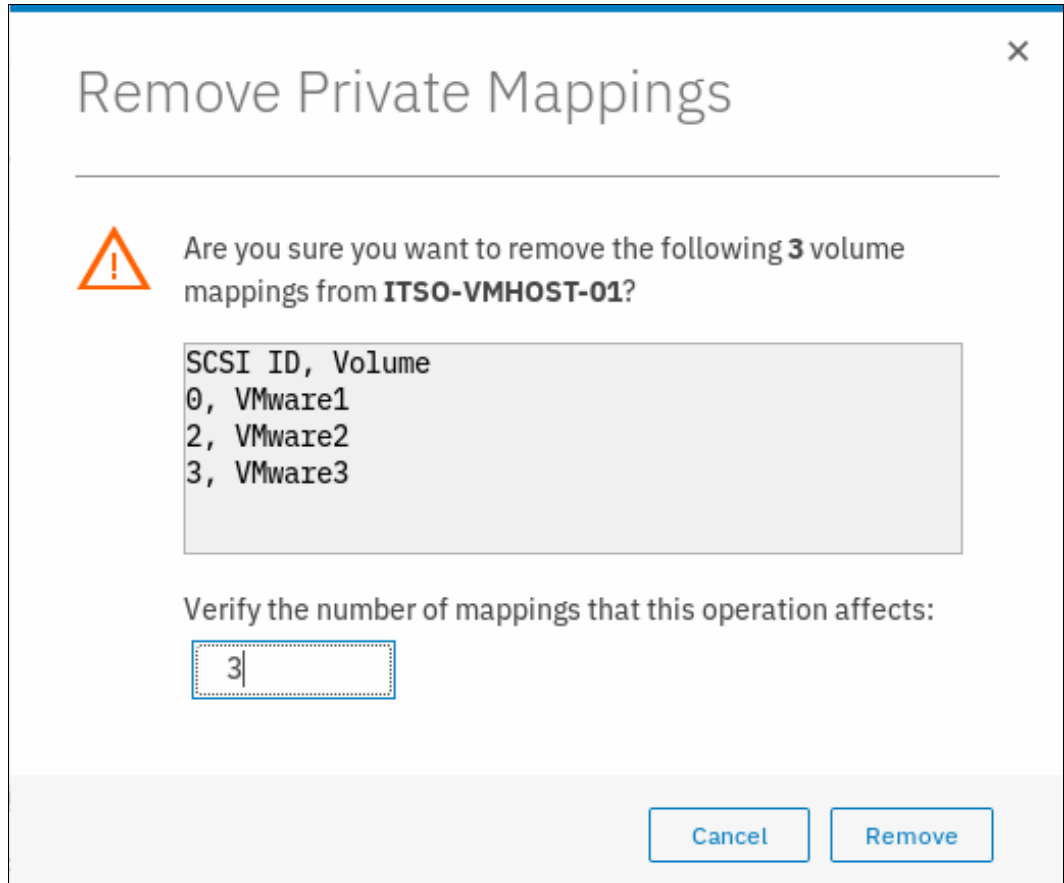


Figure 7-37 Confirming the number of mappings to be removed

**Unmapping:** If you click **Remove**, all access for this host to volumes that are controlled by the IBM SAN Volume Controller system is removed. Ensure that you run the required procedures on your host operating system, such as unmounting the file system, taking disks offline, or disabling the volume group, before removing the volume mappings from your host object by using the IBM Spectrum Virtualize GUI.

- The changes are applied to the system. Click **Close**. Figure 7-38 shows that the selected host no longer has any host mappings.

Name	Status	Host Type	# of Ports	Host Mappings	Host Cluster ID	Host Cluster Name
iscsihost	Offline	Generic	1	No		
ITSO-VMHOST-01	Offline	Generic	1	No		

Figure 7-38 All mappings for host ITSO-VMHOST-01 were removed

## Duplicating and importing mappings

Volumes that are assigned to one host can be quickly and simply mapped to another host object. You might do this, for example, when replacing an aging host's hardware and want to ensure that the replacement host node can access the same set of volumes as the old host.

You can accomplish this task in two ways: By duplicating the mappings on the existing host object to the new host object, or by importing the host mappings to the new host. To duplicate the mappings, complete the following steps:

1. To duplicate an existing host mapping, select the host that you want to duplicate and select **Actions** → **Duplicate Volume Mappings** (see Figure 7-39). In our example, we duplicate the volumes that are mapped to host ITS0-VMHOST-01 to the new host ITS0-VMHOST-02.

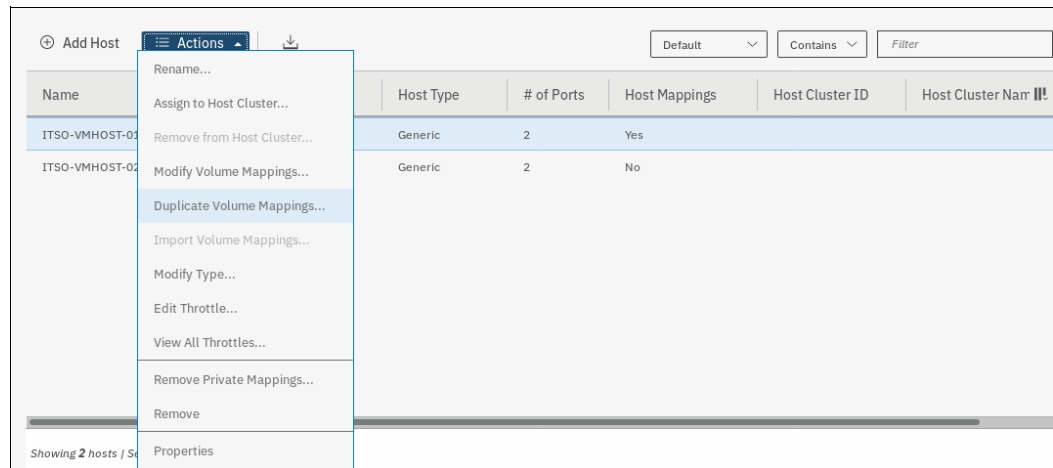


Figure 7-39 Duplicating host volume mappings

2. The Duplicate Mappings window opens. Select a listed target host object to which you want to map all the existing source host volumes and click **Duplicate** (see Figure 7-40).

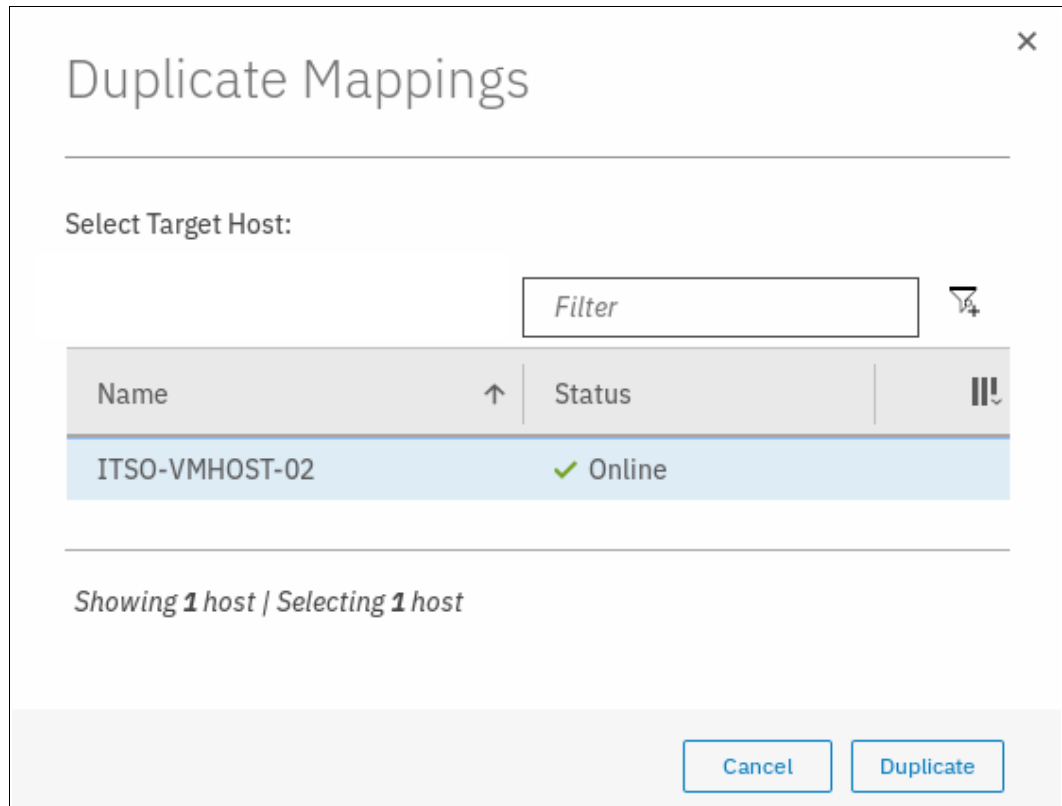


Figure 7-40 Duplicate mappings window

**Note:** You can duplicate mappings only to a host that has no volumes that are mapped.



- After the task completion is displayed, verify the new mappings on the new host object. From the **Hosts** menu (see Figure 7-39 on page 389), right-click the target host and select **Properties**.
- Click the **Mapped Volumes** tab and verify that the required volumes are mapped to the new host (see Figure 7-41).

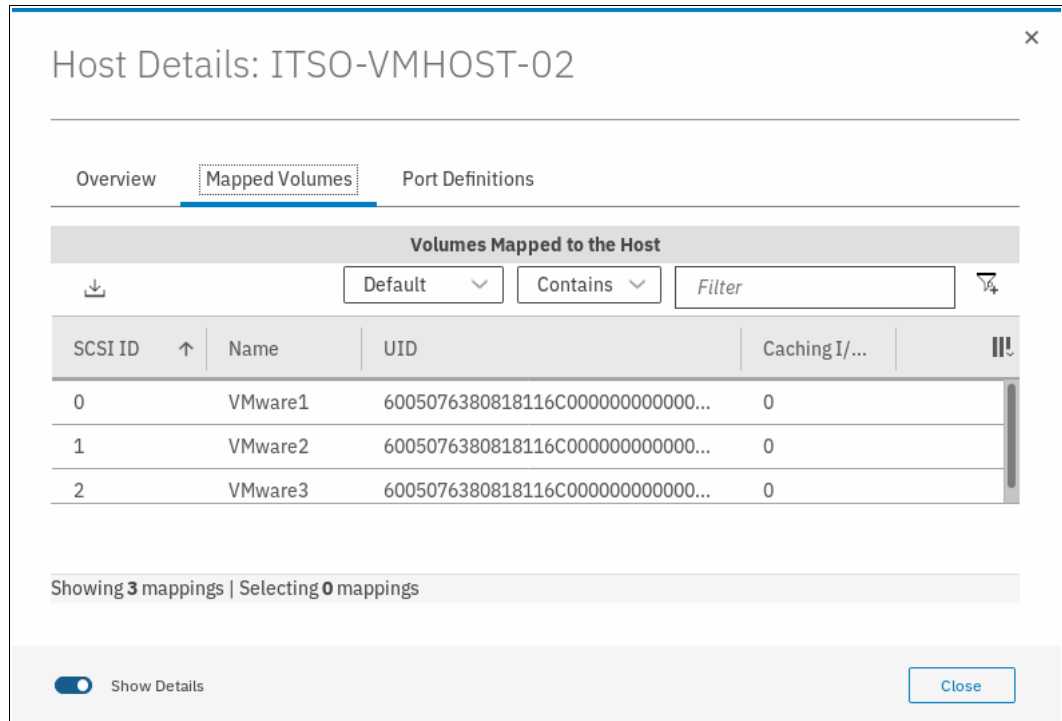


Figure 7-41 Host Details: New mappings on the target host

The same mapping process as the GUI can be accomplished by importing existing hosts mappings to the new host:

- Select the new host without any mapped volumes and click **Actions** → **Import Volume Mappings** (see Figure 7-42).

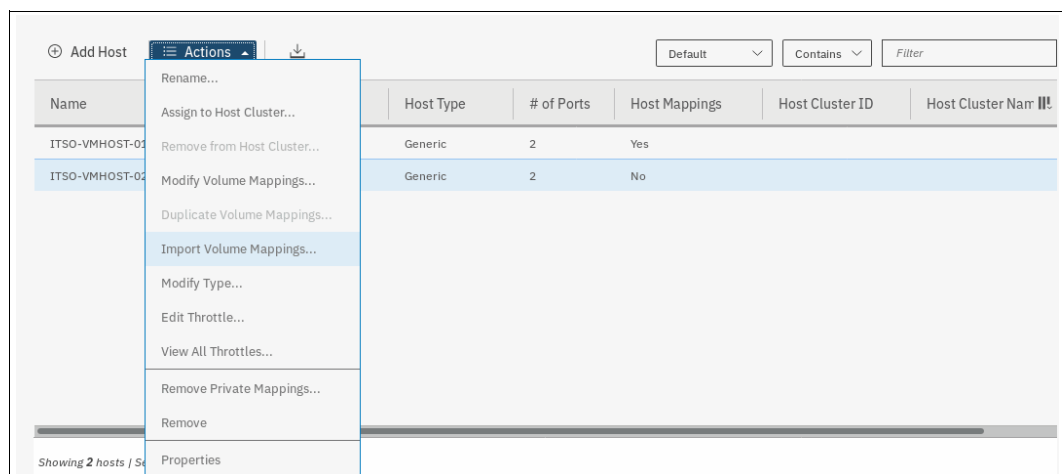


Figure 7-42 Hosts Actions: Importing volume mappings

2. The Import Mappings window opens. Select the source host from which you want to import the volume mappings. As shown in Figure 7-43, we select the host ITS0-VMHOST-01 and click **Import**.

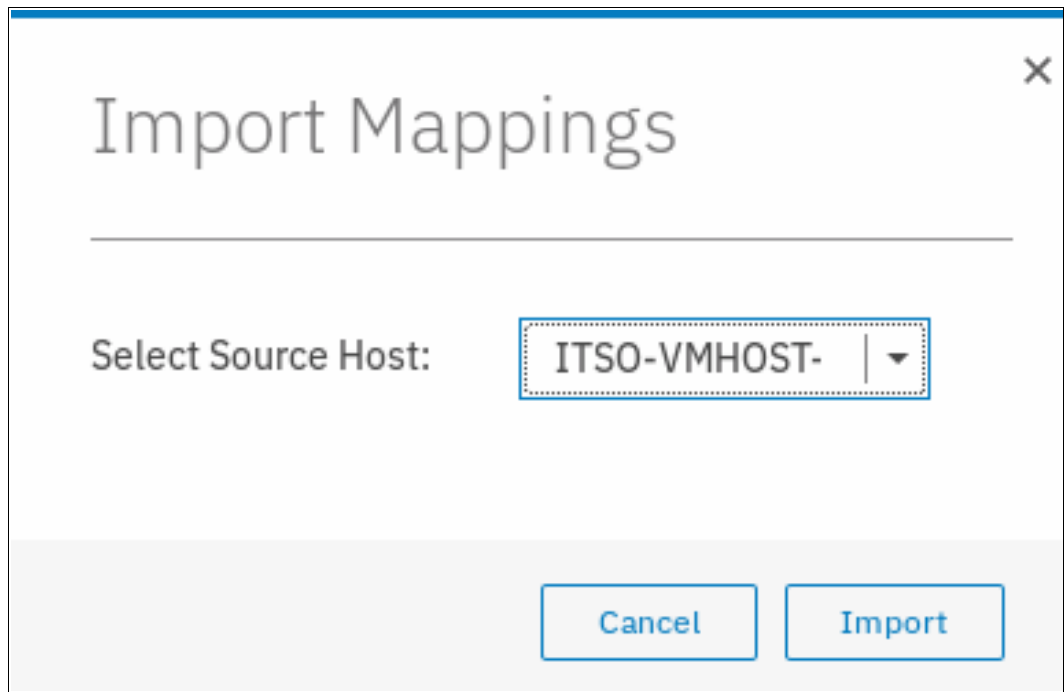


Figure 7-43 Selecting the source host

3. After the task completes, verify that the mappings are as expected by opening the **Hosts** menu (see Figure 7-30 on page 384), right-clicking the target host, and selecting **Properties**. Then, click the **Mapped Volumes** tab and verify that the required volumes are mapped to the new host (see Figure 7-41 on page 391).

**Note:** You can import only mappings from a source host that is in the same ownership group as your target host. If they are not, the import fails with the following message:  
The command failed because the objects are in different ownership groups.

## Renaming a host

To rename a host, complete the following steps:

1. Select the host, and then right-click and select **Rename** (see Figure 7-44).

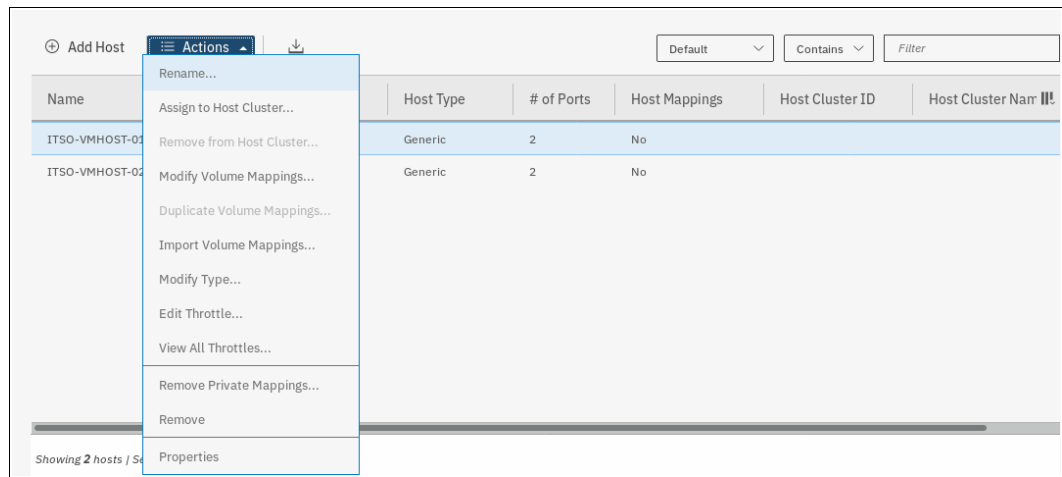


Figure 7-44 Renaming a host

2. Enter a new name and click **Rename** (see Figure 7-45). If you click **Reset**, the changes are reset to the original host name.

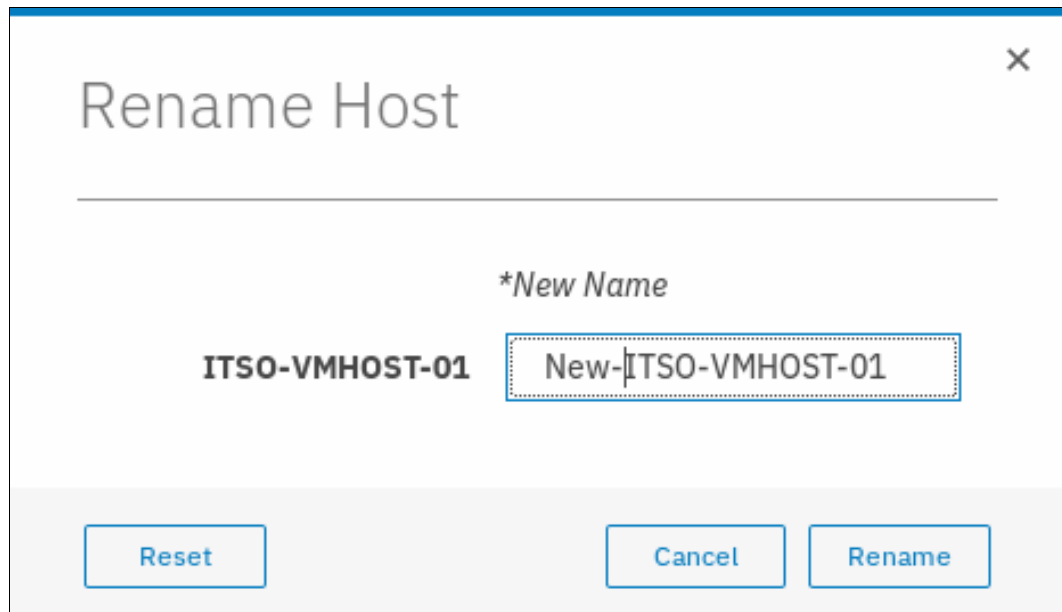


Figure 7-45 Rename Host window

3. After the changes are applied to the system, click **Close**.

## Removing a host

To remove a host object definition, complete the following steps:

1. From the Hosts pane, select the host and right-click it or click **Actions** → **Remove** (see Figure 7-46).

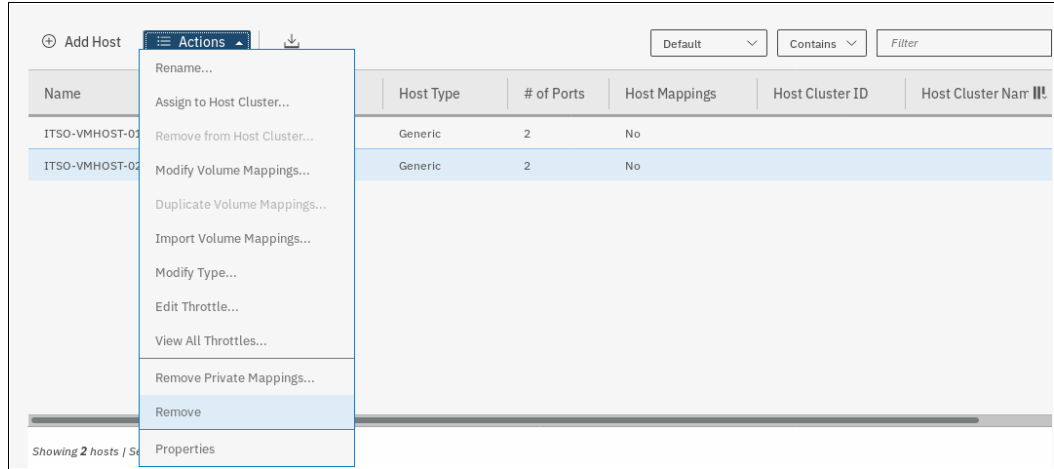


Figure 7-46 Removing a host

2. Confirm that the window displays the correct list of hosts that you want to remove by entering the number of hosts to remove and clicking **Delete** (see Figure 7-47).

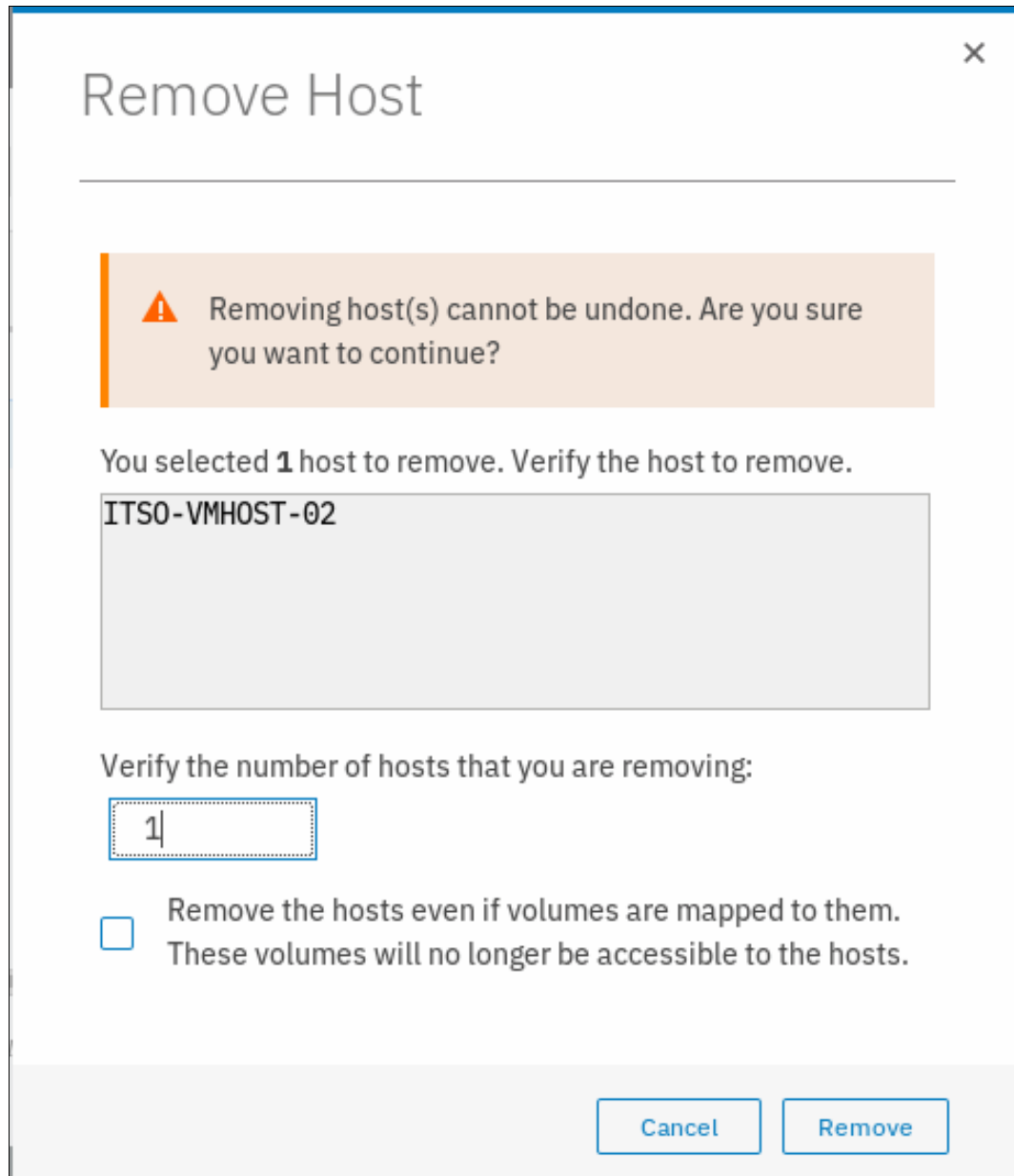


Figure 7-47 Confirm the removal of the host

3. If the host that you are removing has volumes that are mapped to it, force the removal by selecting the check box in the lower part of the window. By selecting this check box, the host is removed and it no longer has access to any volumes on this system.
4. After the task completes, click **Close**.

## Host properties

To view a host object's properties, complete the following steps:

1. From the IBM Spectrum Virtualize GUI Hosts pane, select a host, and right-click it or click **Actions** → **Properties** (see Figure 7-48).

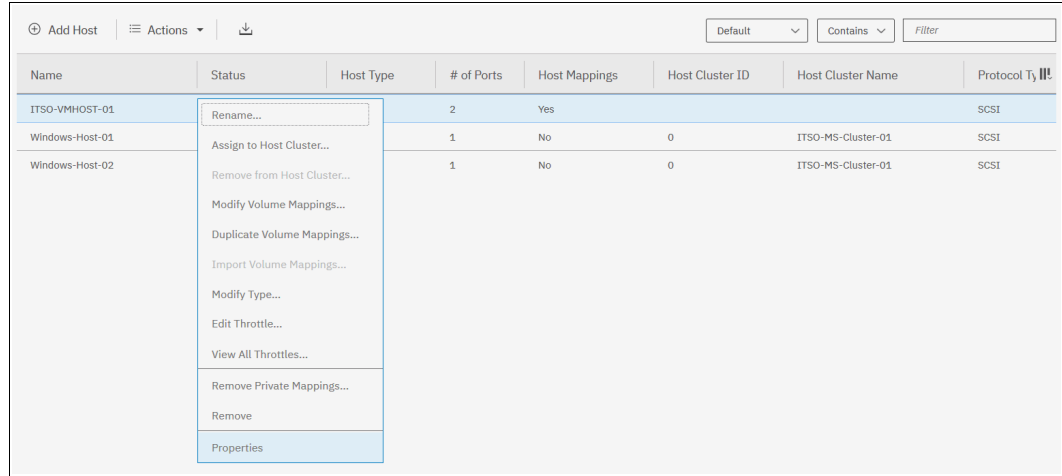


Figure 7-48 Host properties

The Host Details window opens (see Figure 7-49).

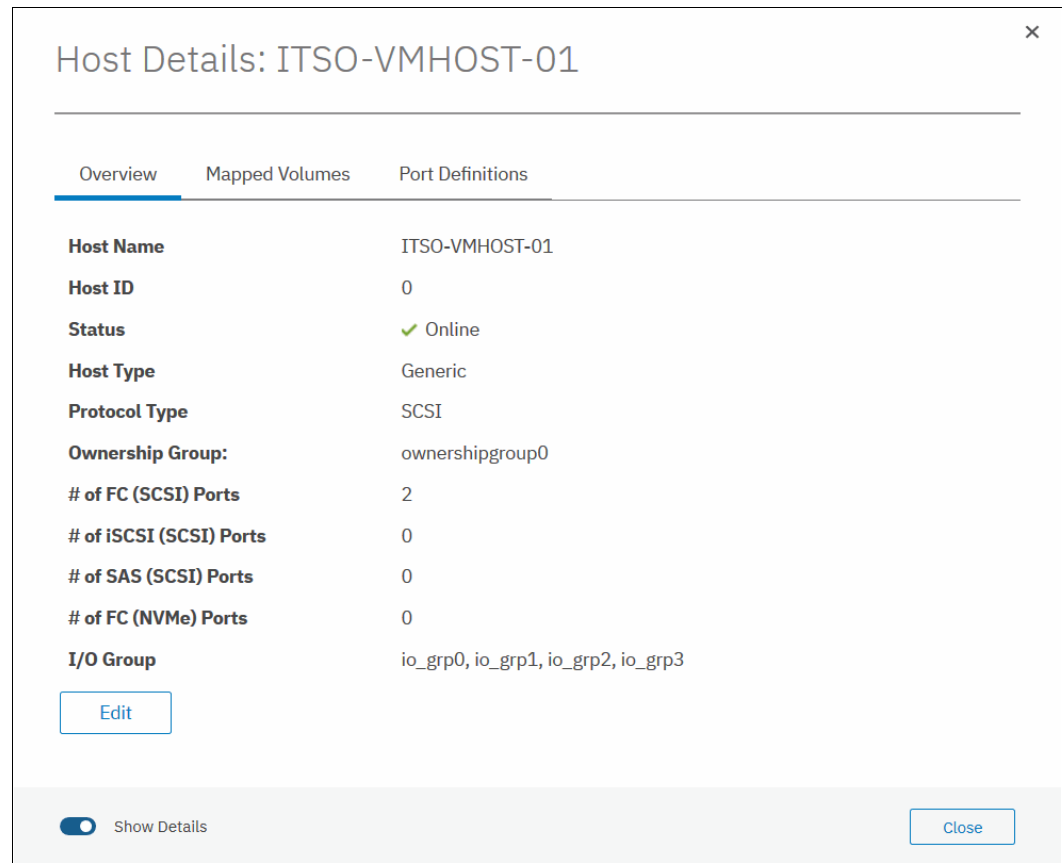


Figure 7-49 Host properties overview

The Host Details window shows an overview of the selected host properties. It has three tabs: Overview, Mapped Volumes, and Port Definitions. The Overview tab is shown in Figure 7-49 on page 396.

2. Select the **Show Details** slider to see more details about the host.
3. Click **Edit** to change the host properties (see Figure 7-50).

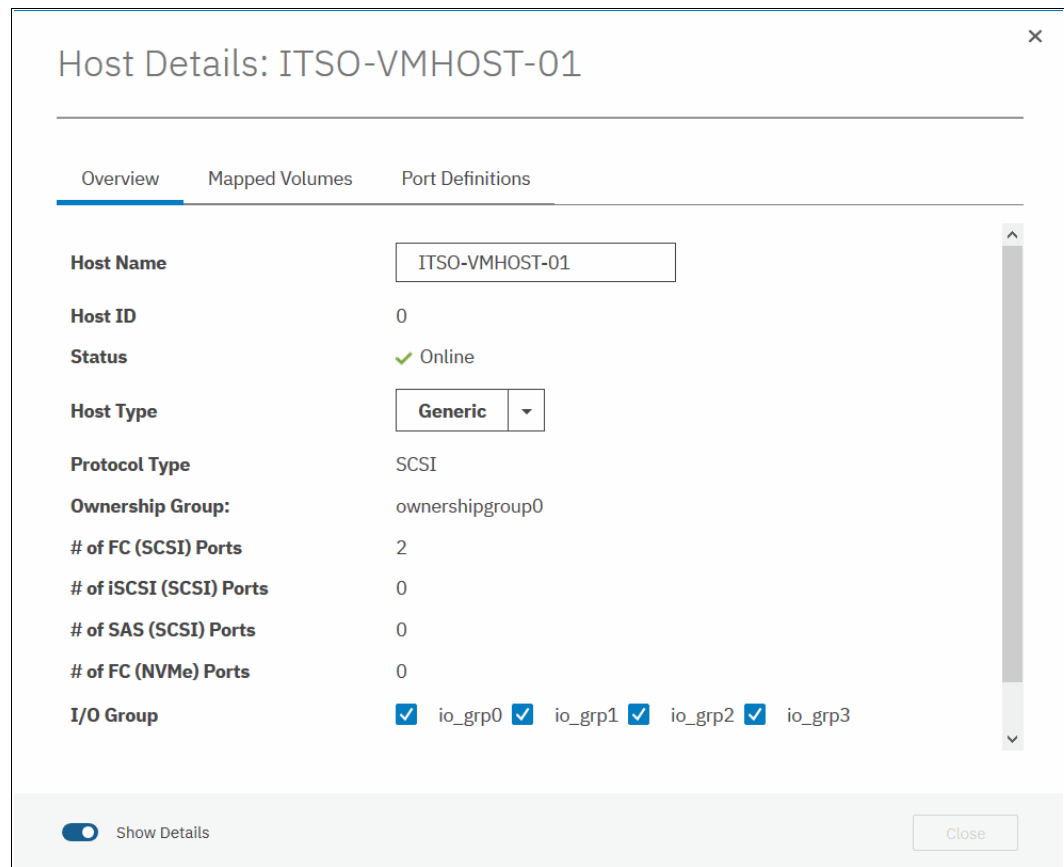


Figure 7-50 Editing the host properties

In the window that is shown in Figure 7-50, you can modify the following properties:

- Host Name: Change the host name.
- Host Type: If you are going to attach HP-UX, OpenVMS, or TPGS hosts, change this setting.
- Protocol Type: If you want to change a host protocol type between SCSI and NVMe, use this setting.
- I/O Group: The host has access to volumes that are mapped from selected I/O groups.
- iSCSI CHAP Secret: Enter or change the iSCSI CHAP secret if this host is using iSCSI.

**Note:** You can change the protocol type of a host only if the host has no ports that are configured.

4. When you are finished making changes, click **Save** to apply them. The editing window closes.

The Mapped Volumes tab shows a summary of which volumes are mapped with which SCSI ID and unique identifier (UID) to this host (see Figure 7-51). The Show Details slider does not show any more information for this list.

Host Details: ITSO-VMHOST-01

Overview **Mapped Volumes** Port Definitions

Volumes Mapped to the Host

Default Contains Filter

SCSI ID	Name	UID	Caching I/...
0	VMware1	60050768019C851444000000000000...	0
1	VMware2	60050768019C851444000000000000...	0
2	VMware3	60050768019C851444000000000000...	0

Showing 3 mappings | Selecting 0 mappings

Show Details  Close

Figure 7-51 Mapped volumes tab



The Port Definitions tab shows the configured host ports of a host and provides status information about them (see Figure 7-52). This is also where you can find the NQN of your NVMe host.

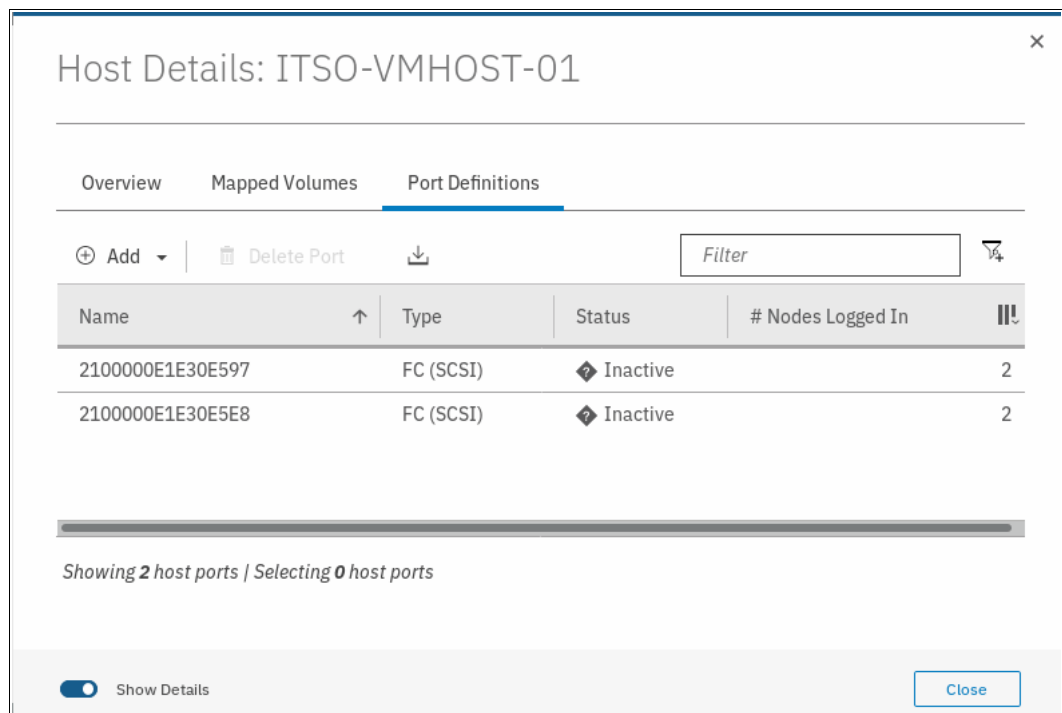


Figure 7-52 Port definitions

This window offers the option to **Add** or **Delete Port** on the host, as described in 7.5.4, “Adding and deleting host ports” on page 400.

5. Click **Close** to close the Host Details window.

## 7.5.4 Adding and deleting host ports

To configure host ports, from the left menu, select **Hosts** → **Ports by Host** to open the associated pane (see Figure 7-53).

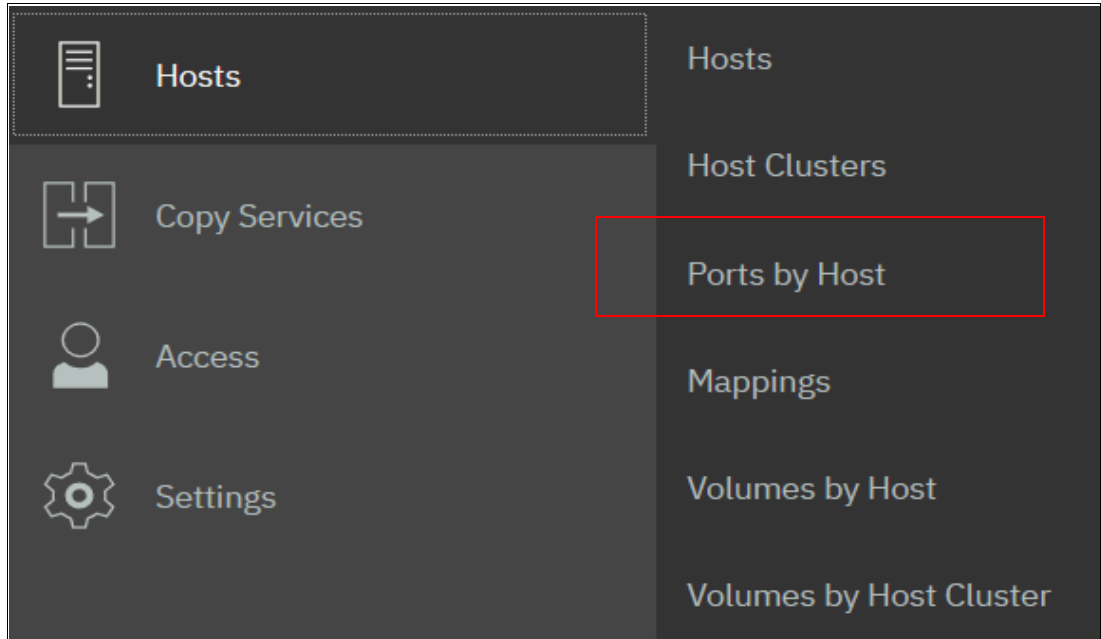


Figure 7-53 Ports by Host pane

A list of all the hosts is displayed. The function icons indicate whether the host is FC-, iSCSI-, or serial-attached SCSI (SAS)-attached. The port details of the selected host are shown to the right, as shown in Figure 7-54.

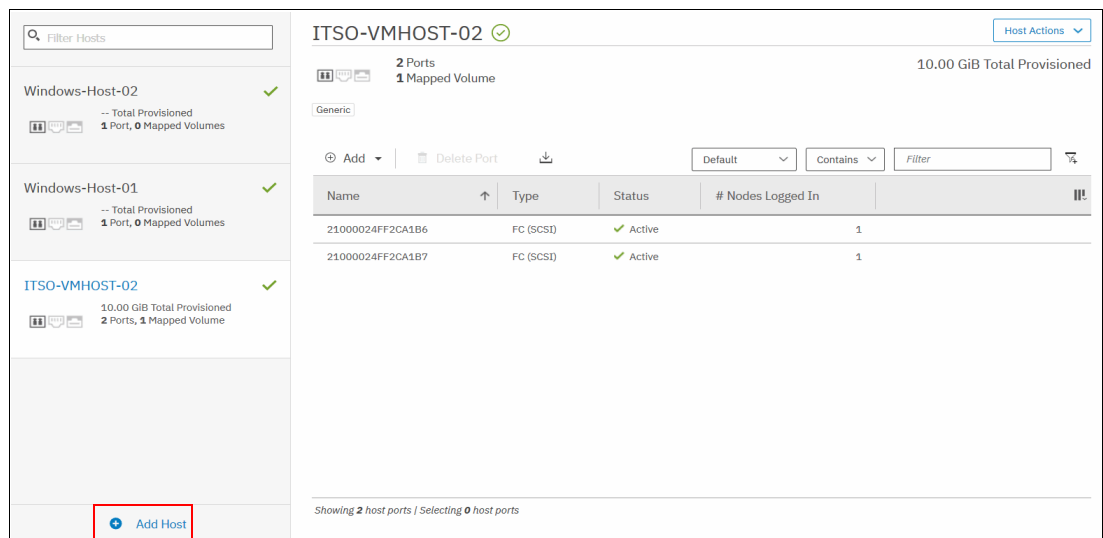


Figure 7-54 Ports by Host actions

## Adding a Fibre Channel or iSCSI host port

To add a host port, complete the following steps:

1. Select the host.
2. Click **Add** (see Figure 7-55) and select one of the following options:
  - **Fibre Channel (SCSI) Port** (see “Adding a Fibre Channel port”).
  - **iSCSI (SCSI) Port** (see “Adding an iSCSI host port” on page 405).

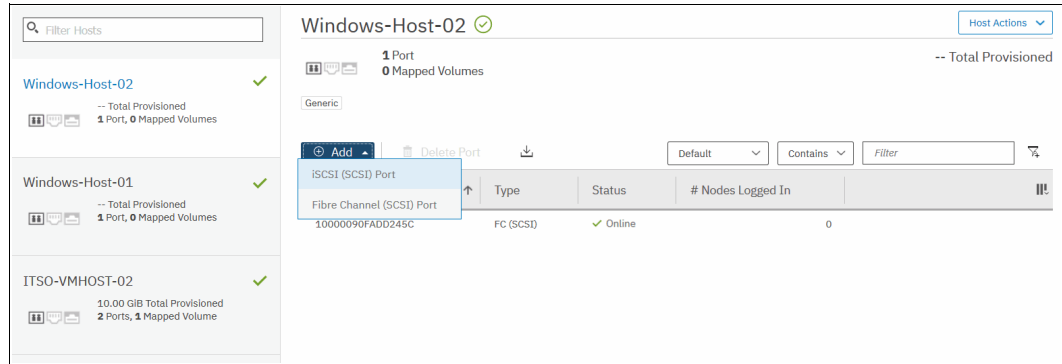


Figure 7-55 Adding host ports

## Adding a Fibre Channel port

To add an FC port, complete the following steps:

1. Click **Fibre Channel Port** (see Figure 7-55 on page 401). The Add Fibre Channel Ports window opens (see Figure 7-56).

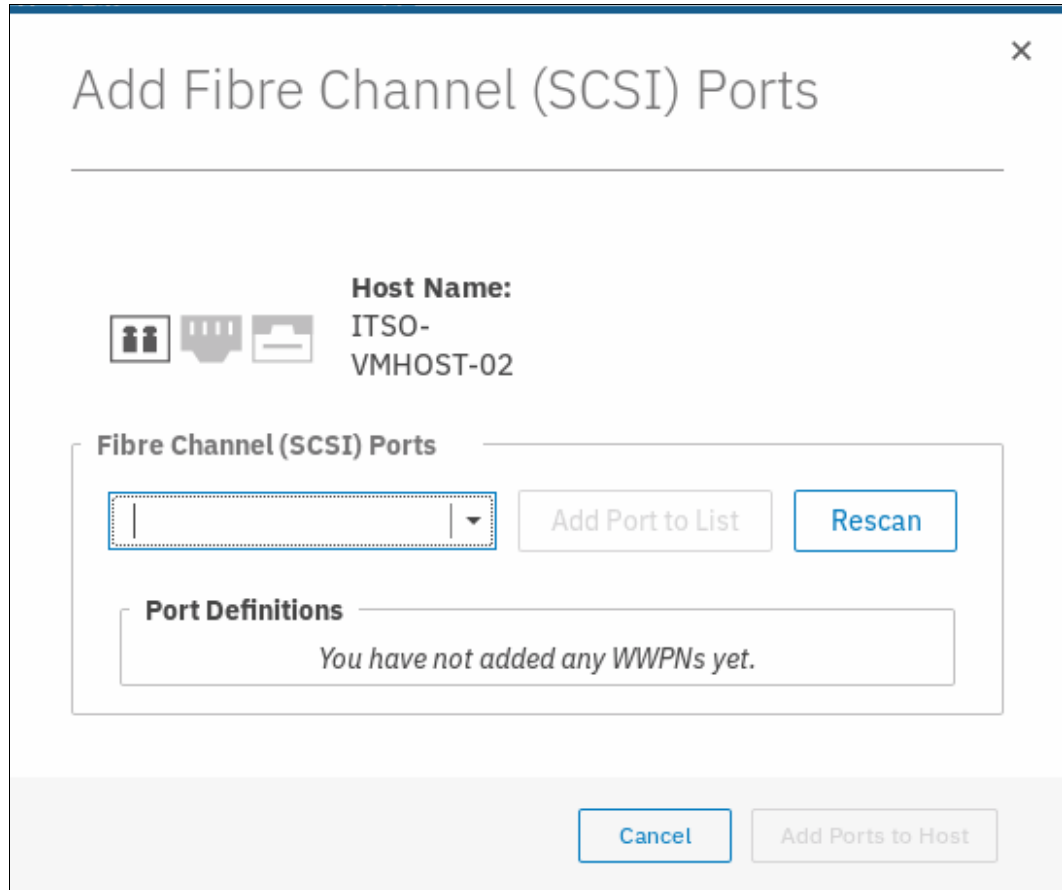


Figure 7-56 Add Fibre Channel Ports window

2. Click the **Fibre Channel (SCSI) Ports** drop-down menu to display a list of all discovered FC WWPNs. If the WWPN of your host is not available in the menu, enter it manually or check the SAN zoning to ensure that connectivity is configured, and then, rescan storage from the host.

3. Select the WWPN that you want to add and click **Add Port to List** (see Figure 7-57).

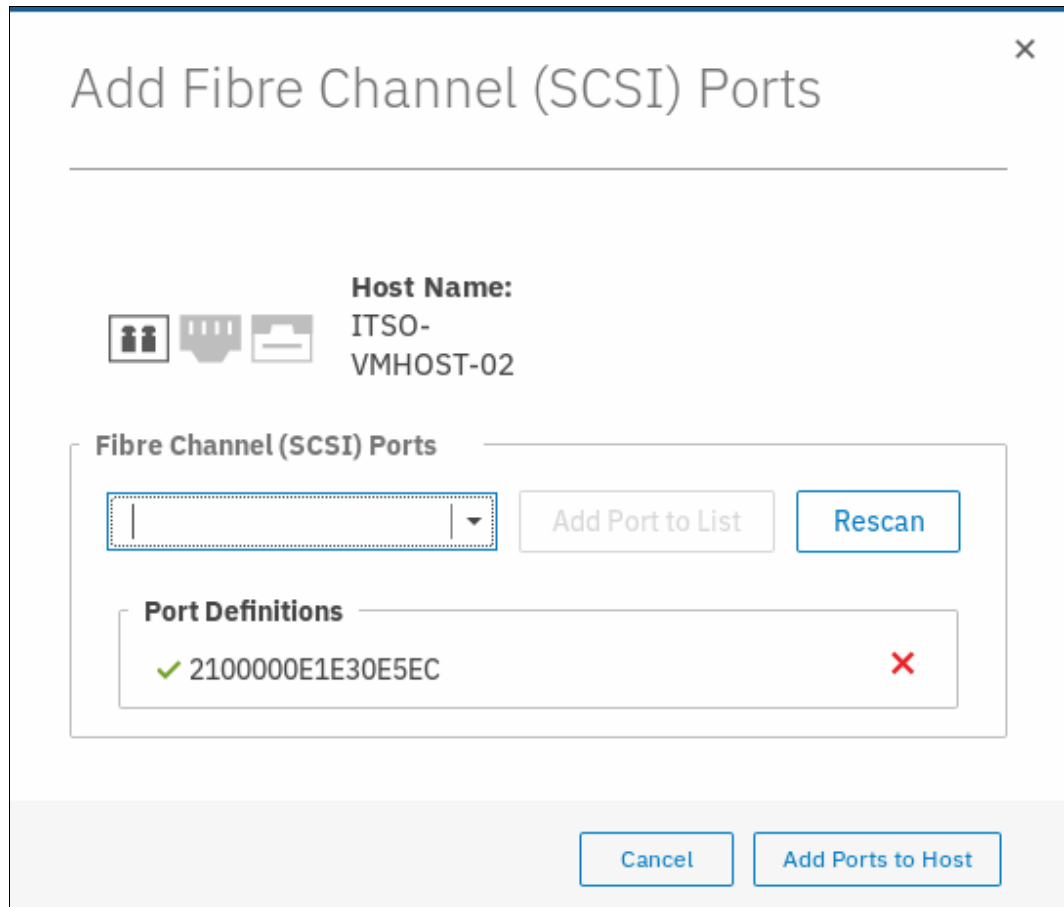


Figure 7-57 Adding a port to a list

This step can be repeated to add ports to the host.

4. To add an offline port (if the WWPN of your host is not available in the drop-down menu), manually enter the WWPN of the port in the **Fibre Channel Ports** field and click **Add Port to List**.

The port is unverified (see Figure 7-58) because it is not logged on to the IBM SAN Volume Controller system. The first time that it logs on, its state is automatically changed to online, and the mapping is applied to this port.

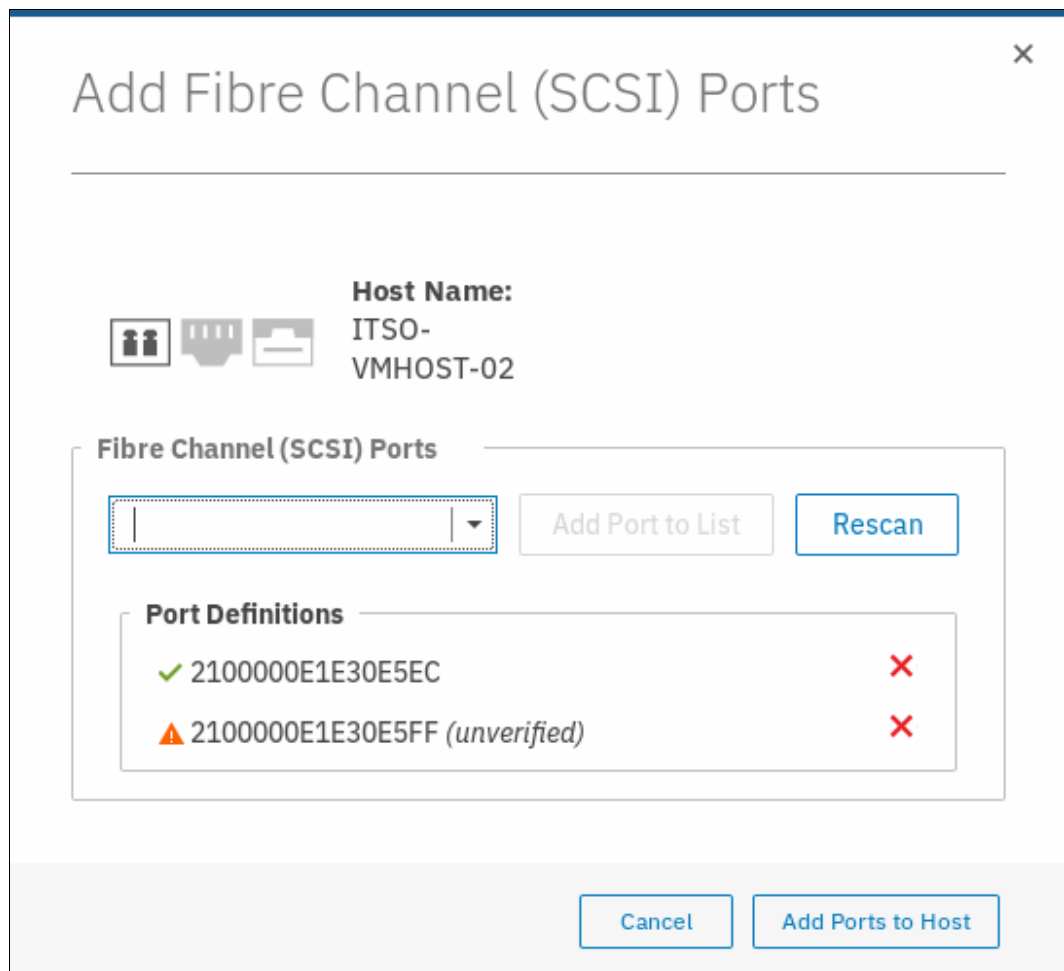


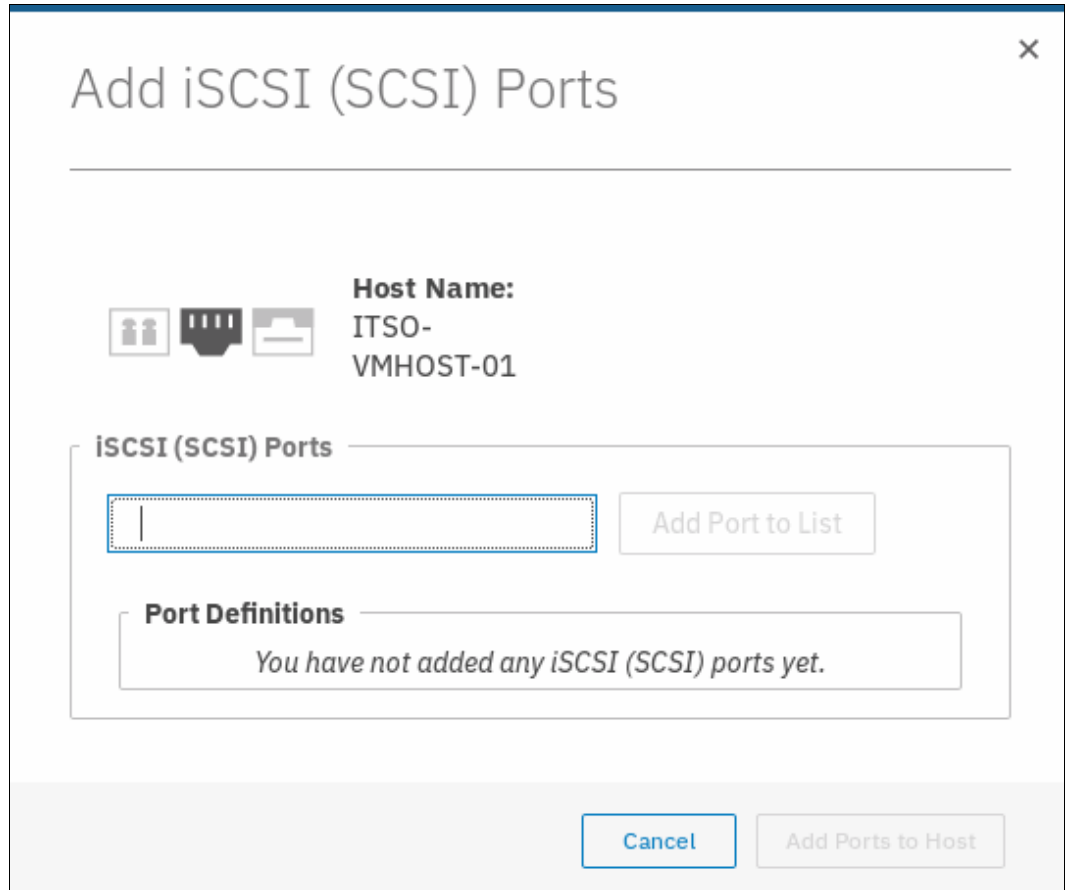
Figure 7-58 Unverified port

5. To remove a port from the list, click the red X next to the port. In this example, we delete the manually added FC port so only the detected port remains.
6. Click **Add Ports to Host** to apply the changes and click **Close**.

## Adding an iSCSI host port




To add an iSCSI host port, complete the following steps:

1. Click **iSCSI Port** (see Figure 7-55 on page 401). The Add iSCSI Ports window opens (see Figure 7-59).



**Add iSCSI (SCSI) Ports** [X]

---

   **Host Name:**  
ITSO-  
VMHOST-01

**iSCSI (SCSI) Ports**

**Port Definitions**

*You have not added any iSCSI (SCSI) ports yet.*

Figure 7-59 Adding iSCSI host ports

2. Enter the initiator name of your host (see Figure 7-60) and click **Add Port to List**.

The screenshot shows a dialog box titled "Add iSCSI (SCSI) Ports" with a close button (X) in the top right corner. Below the title bar, there are three icons (two people, a network port, and a document) followed by the text "Host Name: ITSO-VMHOST-01". Below this, there is a section titled "iSCSI (SCSI) Ports" containing a text input field and a disabled "Add Port to List" button. Underneath, there is a "Port Definitions" section with a list box containing the entry "iqn.2003-01.com.vmware:00.fcd0ab21.vmhost01" and a red "X" icon to its right. At the bottom of the dialog, there are two buttons: "Cancel" and "Add Ports to Host".

Figure 7-60 Entering the initiator name

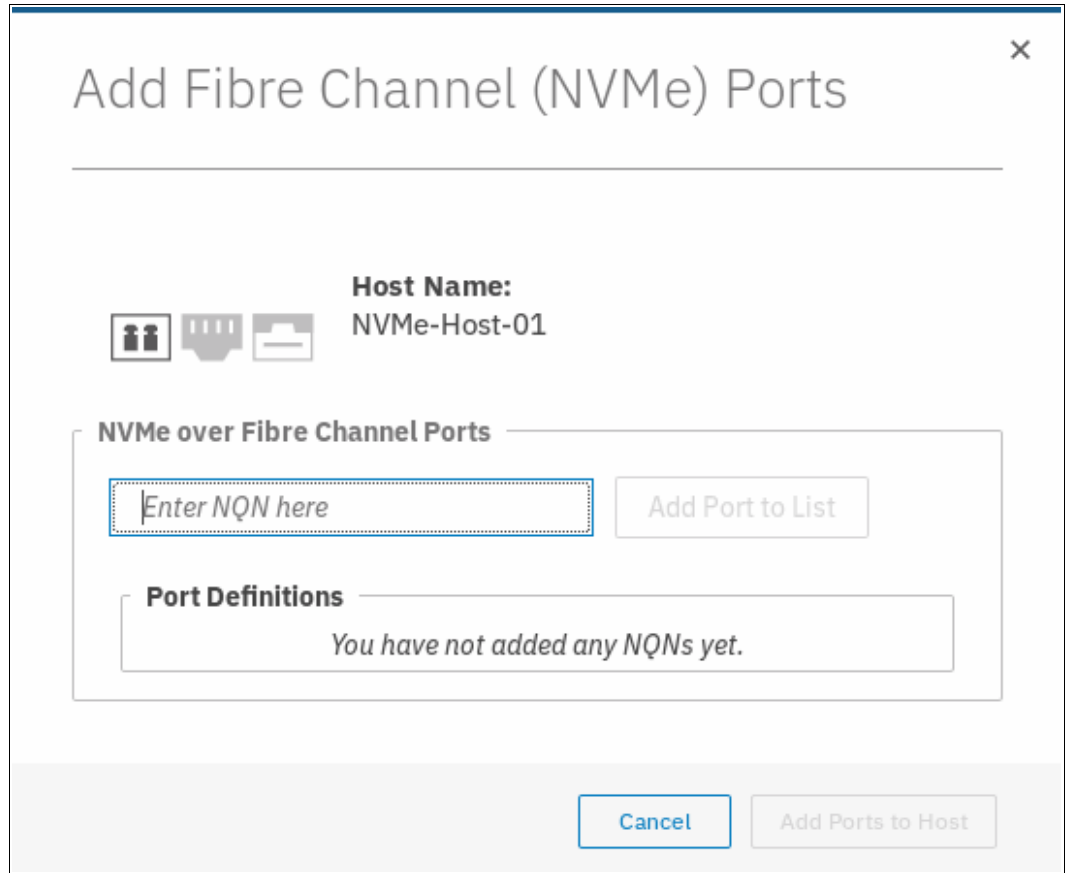
3. Click **Add Ports to Host** to apply the changes to the system and click **Close**.



## Adding a NVMe host port

To add an NVMe host port, complete the following steps:

1. Click **Add** → **Fibre Channel (NVMe) Port**. The Add Fibre Channel (NVMe) Ports window opens, as shown in Figure 7-61.



The screenshot shows a window titled "Add Fibre Channel (NVMe) Ports" with a close button (X) in the top right corner. Below the title bar, there are three icons (two people, a server rack, and a document) to the left of the text "Host Name: NVMe-Host-01". Below this, there is a section titled "NVMe over Fibre Channel Ports" which contains a text input field with the placeholder text "Enter NQN here" and a button labeled "Add Port to List". Below the input field and button is a section titled "Port Definitions" which contains the text "You have not added any NQNs yet.". At the bottom of the window, there are two buttons: "Cancel" and "Add Ports to Host".

Figure 7-61 Add Fibre Channel (NVMe) Ports

2. Enter the NQN and then click **Add Ports to Host**, as shown in Figure 7-62.

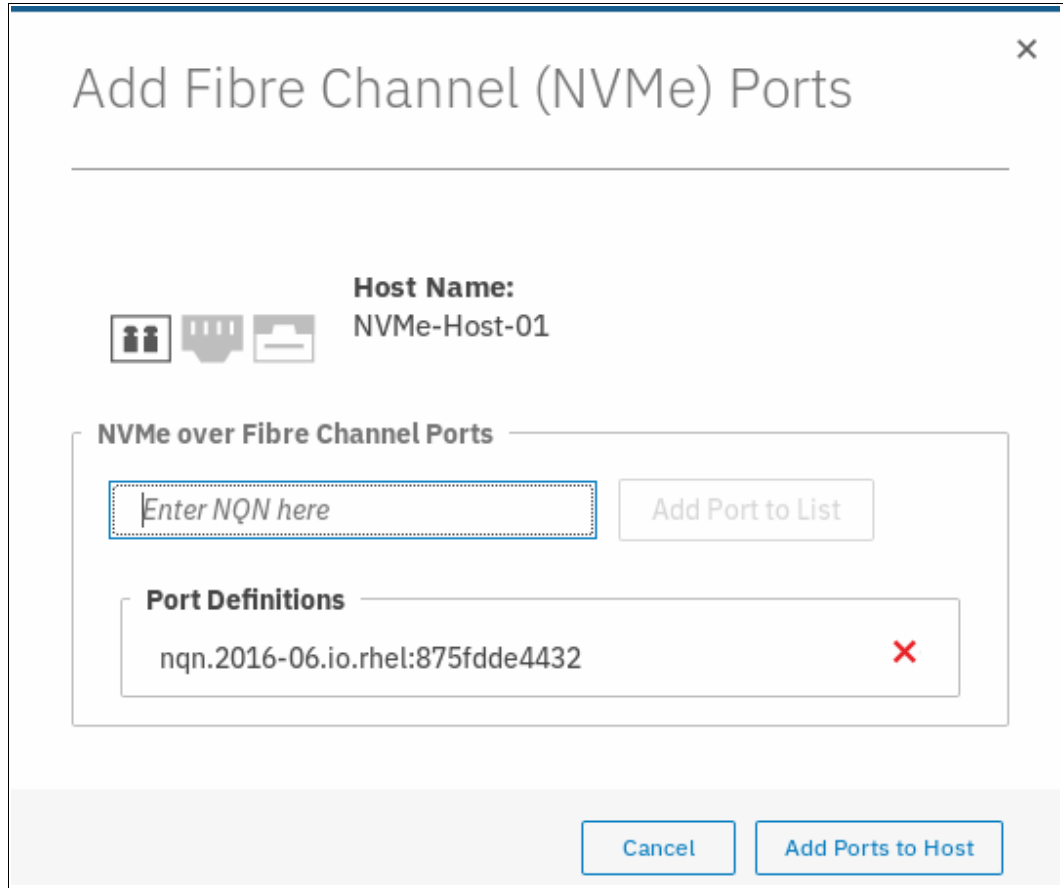


Figure 7-62 Entering the NVMe Qualified Name

### Deleting a host port

To delete a host port, complete the following steps:

1. Highlight the host port and right-click it or click **Delete Port** (see Figure 7-63).

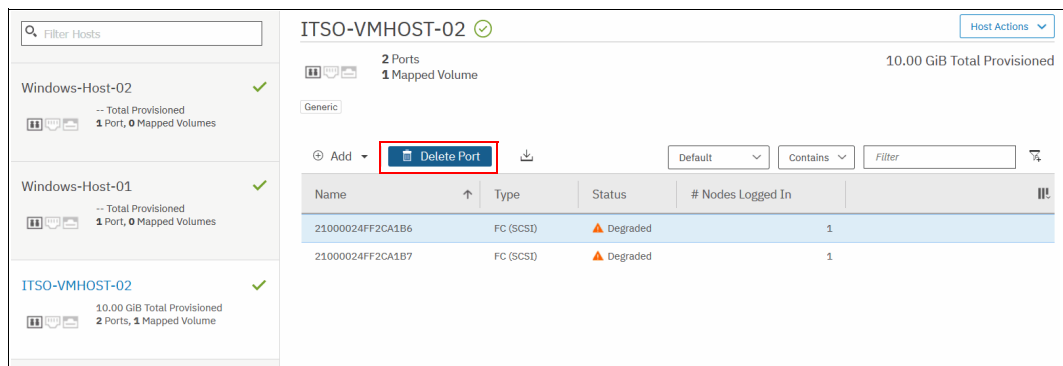


Figure 7-63 Deleting a host port

You can also press the **Ctrl** key to select several host ports to delete (see Figure 7-64).

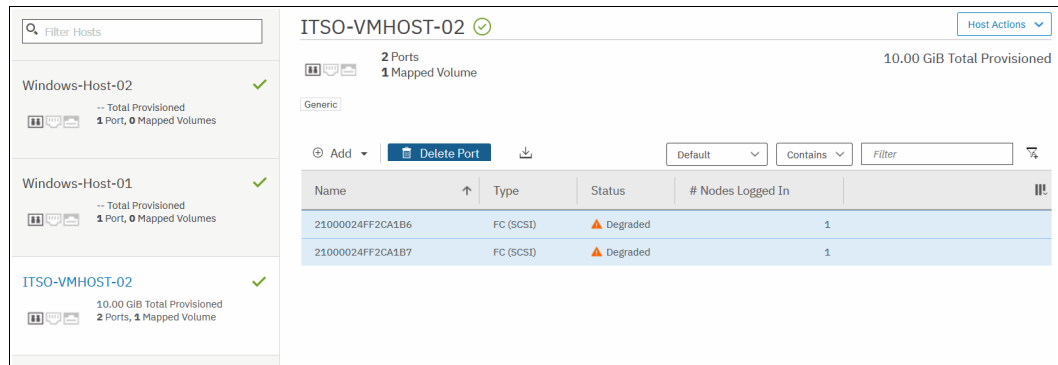


Figure 7-64 Deleting several host ports

2. Click **Delete** and confirm the number of host ports that you want to remove by entering that number in the **Verify** field (see Figure 7-65).

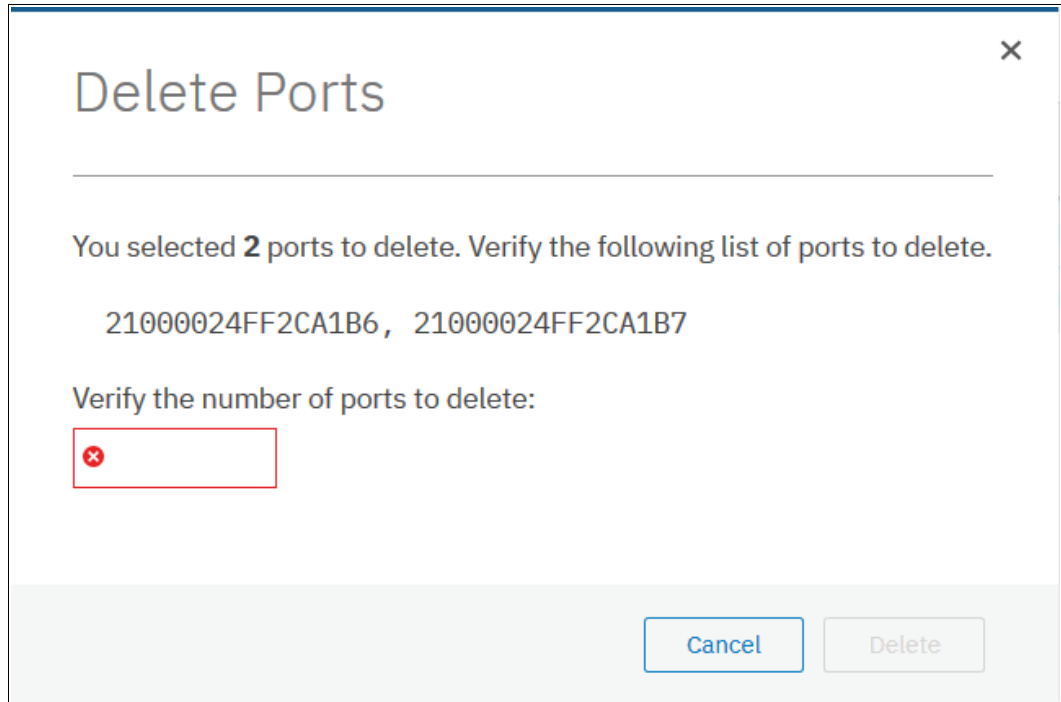


Figure 7-65 Entering the number of host ports to delete

3. Click **Delete** to apply the changes and then click **Close**.

**Note:** The same process is used to delete FC (including NVMe) and iSCSI ports.

## 7.5.5 Host mappings overview

Click **Hosts** → **Mappings**, as shown in Figure 7-66.

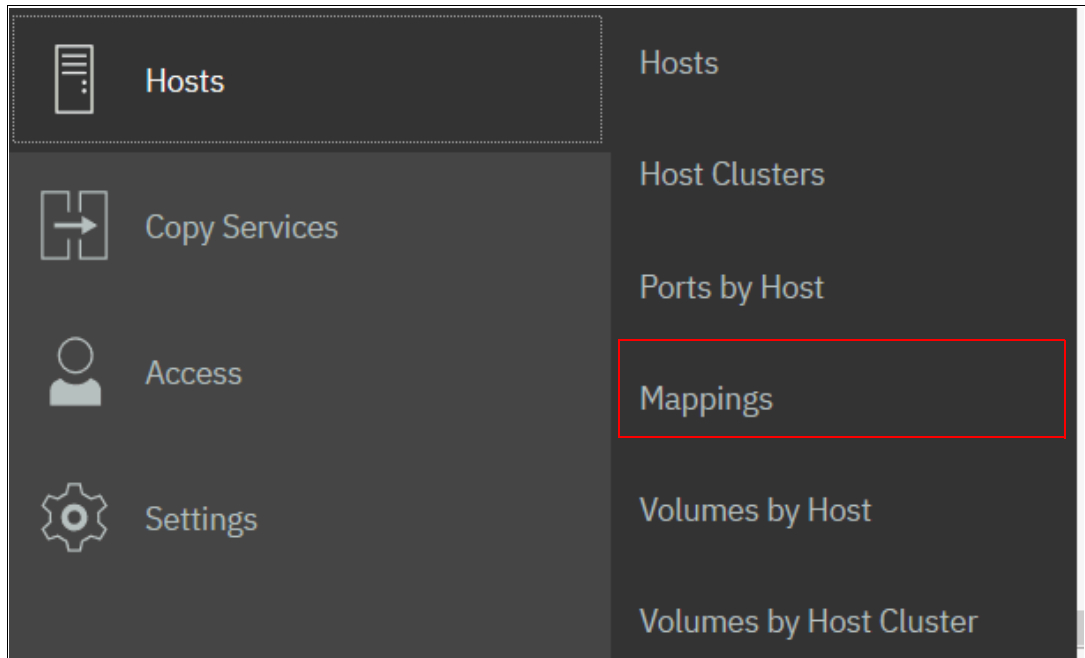


Figure 7-66 Host mappings pane

A list of all volume mappings is shown (see Figure 7-67).

Host Name	SCSI ID	Volume Name	Mapping Type	UID	I/O
ITSO-VMHOST-01	0	VMware1	Private	6005076380818116C000000000000045	0
ITSO-VMHOST-01	1	VMware2	Private	6005076380818116C000000000000046	0
ITSO-VMHOST-02	0	VMware3	Private	6005076380818116C000000000000047	0
RHEL-Host-01	0	Linux1	Private	6005076380818116C00000000000003F	0
RHEL-Host-01	5	Linux6	Private	6005076380818116C000000000000044	0
RHEL-Host-01	4	Linux5	Private	6005076380818116C000000000000043	0
RHEL-Host-01	2	Linux3	Private	6005076380818116C000000000000041	0
RHEL-Host-01	3	Linux4	Private	6005076380818116C000000000000042	0

Showing 9 mappings | Selecting 0 mappings

Figure 7-67 Mappings list

This window lists all hosts and volumes. This example shows that the host ITSO-VMHOST-01 has two mapped volumes, and their associated SCSI IDs, Volume Names, and Volume Unique Identifiers. If you have more than one caching I/O group, you also see which volume is handled by which I/O group.

If you select one line and click **Actions** (see Figure 7-68 on page 411), the following tasks are available:

- Unmap Volumes

- ▶ Properties (Host)
- ▶ Properties (Volume)

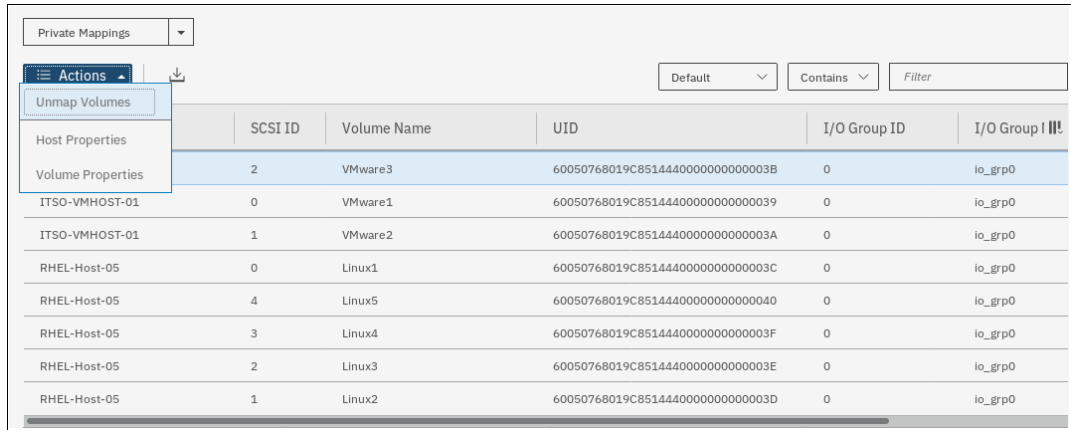


Figure 7-68 Host Mappings Actions menu

## Unmapping a volume

This action removes the mappings for all selected entries. From the **Actions** menu that is shown in Figure 7-68, select one or more lines (while pressing the **Ctrl** key), and click **Unmap Volumes**. Confirm how many volumes are to be unmapped by entering that number in the **Verify** field (see Figure 7-69), and then, click **Unmap**.

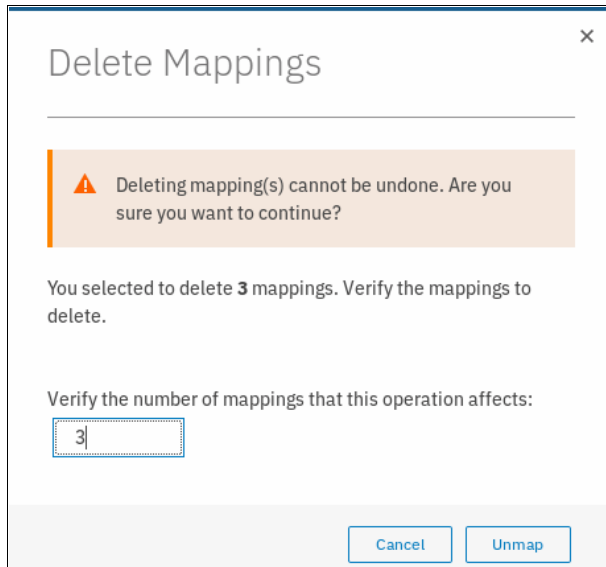


Figure 7-69 Unmapping selected volumes

## Properties (Host)

Selecting an entry and clicking **Properties (Host)**, as shown in Figure 7-68, opens the Host Properties window. The contents of this window are described in “Host properties” on page 396.

## Properties (Volume)

Selecting an entry and clicking **Properties (Volume)**, as shown in Figure 7-68, opens the Volume Properties view. The contents of this window are described in Chapter 6, “Volumes” on page 255.

## 7.5.6 Listing Volumes by Host

To see an overview of the host mappings, click **Hosts** → **Volumes by Host**, as shown in Figure 7-70.

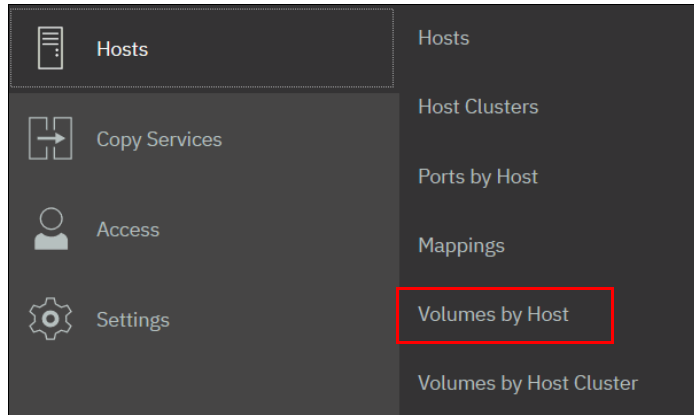


Figure 7-70 Volumes by Host

A list of Volumes mapped per host is displayed, as shown in Figure 7-71. This window differs from the Host Mappings window as you see which volumes are mapped by host, rather than by type of mapping.

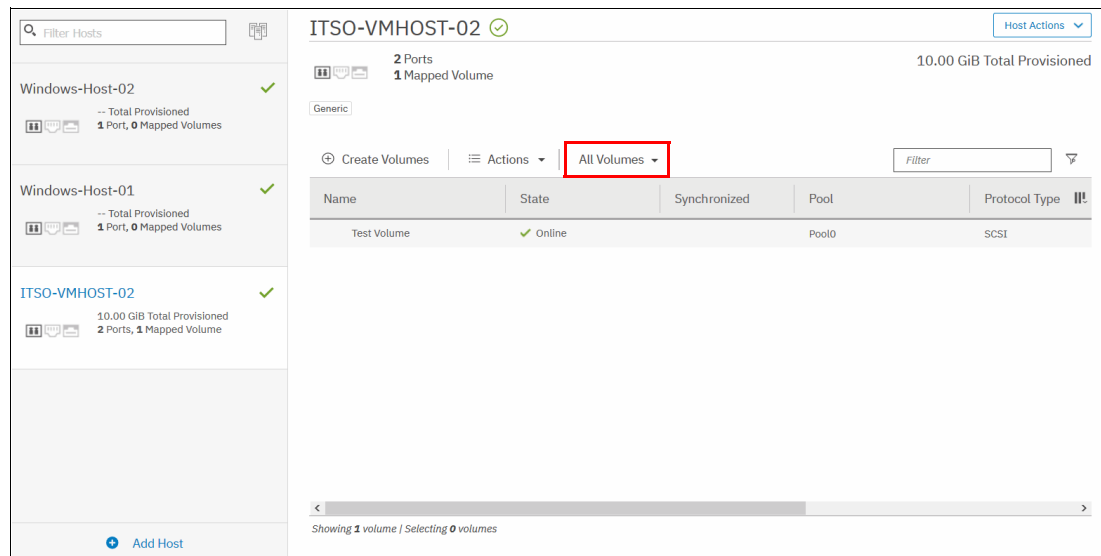


Figure 7-71 Volumes by Host Actions Menu

You can also filter by type of volume by selecting from the following options in the volumes drop down menu:

- ▶ All Volumes
- ▶ Thin-Provisioned Volumes
- ▶ Compressed Volumes
- ▶ Deduplicated Volumes

Finally, you can create and map a new Volume by using this Window. Select **Create Volumes** and the volume creation window appears, as shown in Figure 7-72.

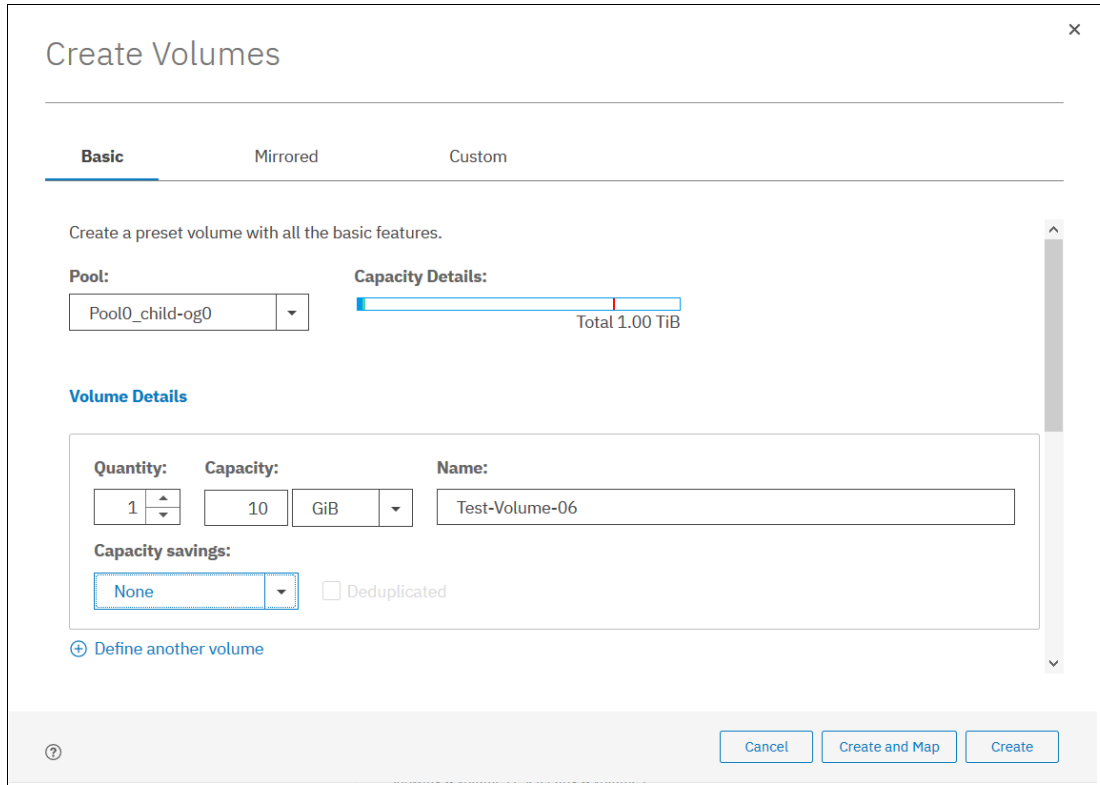


Figure 7-72 Volume Creation Window

After you chose the wanted volume settings, click **Create and Map**, which brings you to the **Create Mapping** window, as shown in Figure 7-73.

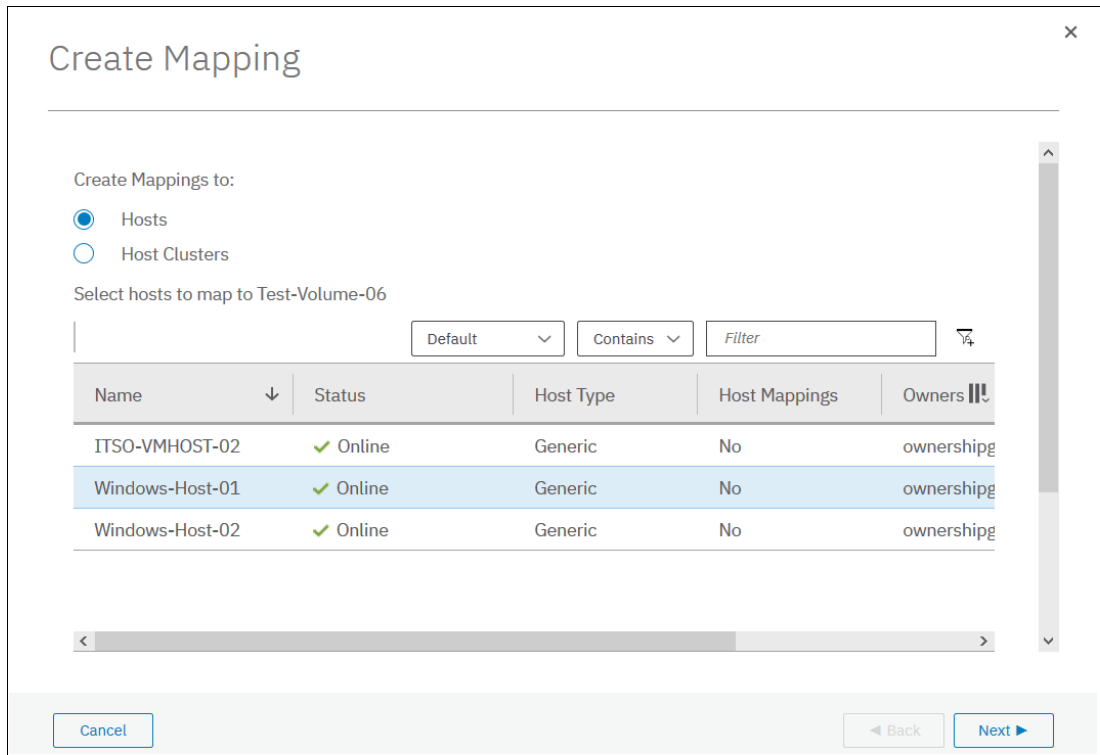


Figure 7-73 Create Mapping Window

Select the host to which you want to map the new volume. If you want to map the volume to multiple hosts, hold down **SHIFT** or **CTRL** and select the wanted hosts. Then, click **Next**, which brings you to a confirmation window, at which point you can click **Map Volumes**, which completes the process. You can then see your new volume listed under the host that you mapped it to, as shown in Figure 7-74.

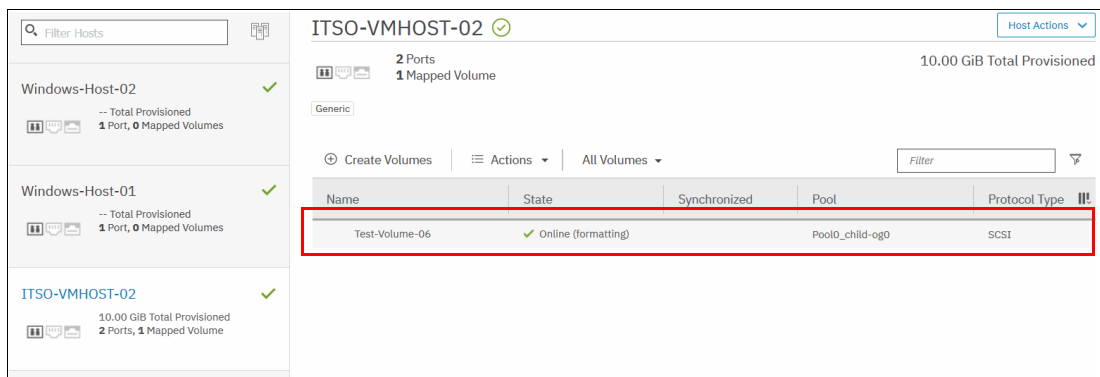


Figure 7-74 New Volume Mapped to Host



**Note:** If you configured object-based access control on your system, the following components must be in the same ownership group when creating and mapping a volume in this manner:

- ▶ User you are logged in as
- ▶ Pool that you are creating the volume on
- ▶ Host that you are mapping the volume to

If not, when you come to map the new volume to the host as shown in Figure 7-73 on page 414, the host does not appear in the list of hosts that are available to create the mapping.

## 7.5.7 Listing Volumes by Host Cluster

To see an overview of the host mappings, click **Hosts** → **Volumes by Host Cluster** as shown in Figure 7-75.

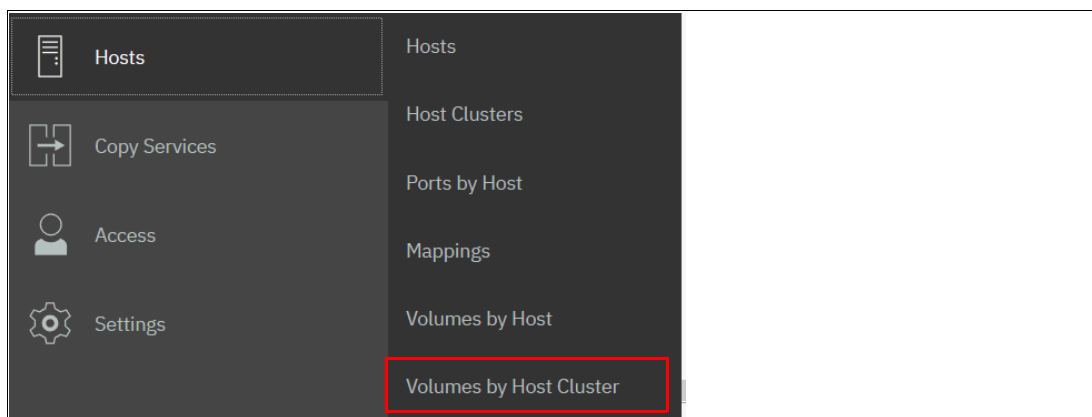


Figure 7-75 Volumes by Host Cluster

This window has the same options as the Volumes by Host option (see 7.5.6, “Listing Volumes by Host” on page 412) except it is for Host Clusters rather than single hosts. If you have no Host Clusters defined in your system, no Host Objects are displayed.

## 7.6 Performing hosts operations by using the command-line interface

This section describes some of the host-related actions that can be taken within the SAN Volume Controller system by using the command-line interface (CLI).

### 7.6.1 Creating a host by using the CLI

This section describes how to create FC and iSCSI hosts by using the IBM Spectrum Virtualize CLI. It is assumed that hosts are prepared for attachment as per the guidelines in the “Host Attachment” section of [IBM Knowledge Center](#).

## Creating Fibre Channel hosts

To create an FC host, complete the following steps:

1. Rescan the SAN on the SAN Volume Controller system by using the **detectmdisk** command, as shown in Example 7-16.

*Example 7-16 Rescanning the SAN on the IBM SAN Volume Controller system*

---

```
IBM_2145:ITS0:superuser>detectmdisk
```

---

**Note:** The **detectmdisk** command does not return any response.

If the zoning is implemented correctly, any new WWPNs should be discovered by the SAN Volume Controller system after running the **detectmdisk** command.

2. List the candidate WWPNs and identify the WWPNs belonging to the new host, as shown in Example 7-17.

*Example 7-17 Available WWPNs*

---

```
IBM_2145:ITS0-SV1:superuser>lsfcportcandidate
fc_WWPN
2100000E1E09E3E9
2100000E1E30E5E8
2100000E1E30E60F
2100000E1EC2E5A2
2100000E1E30E597
2100000E1E30E5EC
```

---

3. Run the **mkhost** command with the required parameters, as shown in Example 7-18.

*Example 7-18 Host creation*

---

```
IBM_2145:ITS0-SV1:superuser>mkhost -name ITS0-VMHOST-03 -fcwwpn
2100000E1E30E597:2100000E1E30E5EC
Host, id [3], successfully created
IBM_2145:ITS0-SV1:superuser>
```

---

## Creating iSCSI hosts

Before you create an iSCSI host in the SAN Volume Controller system, you must know the IQN address of the host. To find the IQN of the host, see your host operating system-specific documentation.

Create a host by completing the following steps:

1. Create the iSCSI host by using the **mkhost** command, as shown in Example 7-19.

*Example 7-19 Creating an iSCSI host by using mkhost*

---

```
IBM_2145:ITS0-SV1:superuser>mkhost -iscsiname iqn.1994-05.com.redhat:e6ff477b58 -name RHEL-Host-06
Host, id [4], successfully created
IBM_2145:ITS0-SV1:superuser>
```

---

2. The iSCSI host can be verified by using the **lshost** command, as shown in Example 7-20.

*Example 7-20 Verifying iSCSI host by using the lshost command*

---

```
IBM_2145:ITS0-SV1:superuser>lshost 4
id 4
name RHEL-Host-06
```

---



```

status offline
site_id
site_name
host_cluster_id
host_cluster_name
protocol nvme
status_policy redundant
status_site all
nqn nqn.2014-08.com.redhat:nvme:nvm-nvmehost01-edf223876
node_logged_in_count 0
state offline
owner_id 0
IBM_2145:ITS0-SV1:superuser>

```

---

**Note:** If you have OBAC set up, you can use the `-ownershipgroup` parameter when creating a host to add the new host to a pre-configured ownership group. You can use either the ownership group name or ID. An example command can be found below:

```

svctask mkhost -name NVMe-Host-01 -nqn
nqn.2014-08.com.redhat:nvme:nvm-nvmehost01-edf223876 -protocol nvme -type
generic -ownershipgroup ownershipgroup0

```

## 7.6.2 Performing advanced host administration by using the CLI

This section describes the following advanced host operations that can be carried out by using the CLI:

- ▶ Mapping a volume to a host
- ▶ Mapping a volume that is already mapped to a different host
- ▶ Unmapping a volume from a host
- ▶ Renaming a host
- ▶ Host properties

### Mapping a volume to a host

To map a volume, complete the following steps:

1. To map an existing volume to a host, run the `mkvdiskhostmap` command, as shown in Example 7-23.

*Example 7-23 Mapping a volume*

```

IBM_2145:ITS0-SV1:superuser>mkvdiskhostmap -host RHEL_HOST -scsi 0 RHEL_VOLUME
Virtual Disk to Host map, id [0], successfully created

```

---

2. The volume mapping can then be checked by running the `lshostvdiskmap` command against that particular host, as shown in Example 7-24.

*Example 7-24 Checking the mapped volume*

```

IBM_2145:ITS0-SV1:superuser>lshostvdiskmap RHEL_HOST
id name      SCSI_id vdisk_id vdisk_name  vdisk_UID
IO_group_id IO_group_name mapping_type host_cluster_id host_cluster_name
7 RHEL_HOST 0      109      RHEL_VOLUME 600507680C81825B00000000000000154 0
io_grp0      private

```

---

## Mapping a volume that is already mapped to a different host

To map a volume to another host that already is mapped to one host, complete the following steps:

1. Run `mkvdiskhost -force` command, as shown in Example 7-25.

### Example 7-25 Mapping the same volume to a second host

```
IBM_2145:ITS0-SV1:superuser>svctask mkvdiskhostmap -force -host RHEL-Host-06
-scsi 0 Linux1
Virtual Disk to Host map, id [0], successfully created
IBM_2145:ITS0-SV1:superuser>
```

**Note:** The volume `Linux1` is mapped to both hosts by using the same SCSI ID. Typically, that is a requirement for most host-based clustering software, such as Microsoft Clustering Service (MSCS), IBM PowerHA, and so on.

2. The volume `Linux1` is mapped to two hosts (`RHEL-Host-05` and `RHEL-Host-06`), which can be seen by running `lsvdiskhostmap`, as shown in Example 7-26.

### Example 7-26 Ensuring that the same volume is mapped to multiple hosts

```
IBM_2145:ITS0-SV1:superuser>lsvdiskhostmap Linux1
id name   SCSI_id host_id host_name   vdisk_UID          IO_group_id IO_group_name mapping_type
host_cluster_id host_cluster_name protocol
11 Linux1 0      3      RHEL-Host-05 60050768019C8514440000000000003C 0      io_grp0     private
scsi
11 Linux1 0      4      RHEL-Host-06 60050768019C8514440000000000003C 0      io_grp0     private
scsi
IBM_2145:ITS0-SV1:superuser>
```

## Unmapping a volume from a host

To unmap a volume from the host, run the `rmvdiskhostmap` command, as shown in Example 7-27.

### Example 7-27 Unmapping a volume from a host

```
IBM_2145:ITS0-SV1:superuser>rmvdiskhostmap -host RHEL-Host-06 Linux1
IBM_2145:ITS0-SV1:superuser>
```

**Note:** Before unmapping a volume from a host on the SAN Volume Controller system, ensure that the host-side action is completed on that volume by using the respective host operating system platform commands, such as unmounting the file system or removing the volume or volume group. Otherwise, unmapping can potentially result in data corruption.

## Renaming a host

To rename an existing host definition, run `chhost -name`, as shown in Example 7-28.

### Example 7-28 Renaming a host

```
IBM_2145:ITS0-SV1:superuser>chhost -name RHEL-Host-07 RHEL-Host-06
IBM_2145:ITS0-SV1:superuser>
```

In Example 7-28, the host `RHEL-Host-06` is renamed to `RHEL-Host-07`.



These options can only be changed using the **chhost** command. When the host is created using **mkhost** the default policy of **redundant** is set.

### 7.6.3 Adding and deleting a host port by using the CLI

This section describes adding and deleting a host port to and from the SAN Volume Controller system.

#### Adding ports to a defined host

If an HBA is added or a network interface controller (NIC) are added to a server that is defined within the SAN Volume Controller system, run the **addhostport** command to add the new port definitions to the host configuration.

If the host is connected through SAN with FC, and if the WWPN is zoned to the SAN Volume Controller system, run the **lsfcportcandidate** command to compare your information with the information that is available from the server administrator, as shown in Example 7-31.

*Example 7-31 Listing the newly available WWPN*

---

```
IBM_2145:ITS0-SV1:superuser>lsfcportcandidate
fc_WWPN
2100000E1E09E3E9
2100000E1E30E5E8
2100000E1E30E60F
2100000E1EC2E5A2
IBM_2145:ITS0-SV1:superuser>
```

---

Use host or SAN switch utilities to verify whether the WWPN matches the information for the new WWPN. If the WWPN matches, run the **addhostport** command to add the port to the host, as shown in Example 7-32.

*Example 7-32 Adding the newly discovered WWPN to the host definition*

---

```
IBM_2145:ITS0-SV1:superuser>addhostport -hbawwpm 2100000E1E09E3E9:2100000E1E30E5E8
ITS0-VMHOST-01
IBM_2145:ITS0-SV1:superuser>
```

---

This command adds the WWPNs 2100000E1E09E3E9 and 2100000E1E30E5E8 to the ITS0-VMHOST-01 host.

If the new HBA is not connected or zoned, the **lshbaportcandidate** command does not display your WWPN. In this case, you can manually enter the WWPN of your HBA or HBAs and use the **-force** flag to create the host, as shown in Example 7-33.

*Example 7-33 Adding a WWPN to the host definition by using the -force option*

---

```
IBM_2145:ITS0-SV1:superuser>addhostport -hbawwpm 2100000000000001 -force ITS0-VMHOST-01
IBM_2145:ITS0-SV1:superuser>
```

---

This command forces the addition of the WWPN 2100000000000001 to the host that is called ITS0-VMHOST-01.

**Note:** WWPNs are not case-sensitive within the CLI.

The host port count can be verified by running the **lshost** command again. The host **ITS0-VMHOST-01** has an updated port count of 3, as shown in Example 7-34.

*Example 7-34 Host with updated port count*

---

```
IBM_2145:ITS0-SV1:superuser>lshost
id name          port_count iogrp_count status  site_id site_name host_cluster_id
host_cluster_name protocol
0 ITS0-VMHOST-02 0          2          offline scsi
1 ITS0-VMHOST-01 3          2          degraded scsi
2 RHEL-Host-01   1          2          offline  scsi
3 ITS0-VMHOST-03 2          2          online   scsi
IBM_2145:ITS0-SV1:superuser>
```

---

If the host uses iSCSI as a connection method, the new iSCSI IQN ID should be used to add the port. Unlike FC-attached hosts, with iSCSI, available candidate ports cannot be checked.

After getting the other iSCSI IQN, run the **addhostport** command, as shown in Example 7-35.

*Example 7-35 Adding an iSCSI port to the defined host*

---

```
IBM_2145:ITS0-SV1:superuser>addhostport -iscsiname iqn.1994-05.com.redhat:e6ddffaab567
RHEL-Host-05
IBM_2145:ITS0-SV1:superuser>
```

---





To remove the iSCSI IQN, run the command, as shown in Example 7-39.

*Example 7-39 Removing iSCSI port from the host*

---

```
IBM_2145:ITS0-SV1:superuser>rmhostport -iscsiname iqn.1994-05.com.redhat:e6ddffaab567
RHEL-Host-05
IBM_2145:ITS0-SV1:superuser>
```

---

To remove the NVMe NQN, use the command that is shown in Example 7-40.

*Example 7-40 Removing NQN port from the host*

---

```
IBM_2145:ITS0-SV1:superuser>rmhostport -nqn nqn.2016-06.io.rhel:875adad3345
RHEL-Host-08
IBM_2145:ITS0-SV1:superuser>
```

---

**Note:** Multiple ports can be removed concurrently by using separators or colons (:) between the port names, as shown in the following example:

```
rmhostport -hbawpnr 210000E08B054CAA:210000E08B892BCD Angola
```

## 7.6.4 Host cluster operations

This section describes the following host cluster operations that can be performed by using the CLI:

- ▶ Creating a host cluster (**mkhostcluster**)
- ▶ Adding a member to the host cluster (**addhostclustermember**)
- ▶ Listing a host cluster (**lshostcluster**)
- ▶ Listing a host cluster member (**lshostclustermember**)
- ▶ Assigning a volume to the host cluster (**mkvolumehostclustermap**)
- ▶ Listing a host cluster for mapped volumes (**lshostclustermap**)
- ▶ Unmapping a volume from the host cluster (**rmvolumehostclustermap**)
- ▶ Removing a host cluster member (**rmhostclustermember**)
- ▶ Removing the host cluster (**rmhostcluster**)

### Creating a host cluster

To create a host cluster, run the **mkhostcluster** command, as shown in Example 7-41.

*Example 7-41 Creating a host cluster by using mkhostcluster*

---

```
IBM_2145:ITS0-SV1:superuser>mkhostcluster -name ITS0-ESX-Cluster-01
Host cluster, id [0], successfully created.
IBM_2145:ITS0-SV1:superuser>
```

---

**Note:** While creating the host cluster, if you want it to inherit the volumes that are mapped to a particular host, use the **-seedfromhost** flag option. Any volume mapping that does not need to be shared can be kept private by using the **-ignoreseedvolume** flag option.

## Adding a host to a host cluster

After creating a host cluster, a host or a list of hosts can be added by running the `addhostclustermember` command, as shown in Example 7-42.

*Example 7-42 Adding a host or hosts to a host cluster*

```
IBM_2145:ITS0-SV1:superuser>addhostclustermember -host ITS0-VMHOST-01:ITS0-VMHOST-02
ITS0-ESX-Cluster-01
IBM_2145:ITS0-SV1:superuser>
```

In Example 7-42, the hosts ITS0-VMHOST-01 and ITS0-VMHOST-02 were added as part of host cluster ITS0-ESX-Cluster-01.

## Listing the host cluster member

To list the host members that are part of a particular host cluster, run the `lshostclustermember` command, as shown in Example 7-43.

*Example 7-43 Listing the host cluster members by running lshostclustermember*

```
IBM_2145:ITS0-SV1:superuser>lshostclustermember ITS0-ESX-Cluster-01
host_id host_name      status type   site_id site_name
0       ITS0-VMHOST-01 offline generic
4       ITS0-VMHOST-02 offline generic
IBM_2145:ITS0-SV1:superuser>
```

## Mapping a volume to a host cluster

To map a volume to a host cluster so that it automatically is mapped to member hosts, run the `mkvolumehostclustermap` command, as shown in Example 7-44.

*Example 7-44 Mapping the volume to the host cluster*

```
IBM_2145:ITS0-SV1:superuser>mkvolumehostclustermap -hostcluster ITS0-ESX-Cluster-01 VMware1
Volume to Host Cluster map, id [0], successfully created
IBM_2145:ITS0-SV1:superuser>
```

**Note:** When a volume is mapped to a host cluster, that volume is mapped to all the members of the host cluster with the same SCSI\_ID.

## Listing the volumes that are mapped to a host cluster

To list the volumes that are mapped to a host cluster, run `lshostclustervolumemap` command, as shown in Example 7-45.

*Example 7-45 Listing volumes that are mapped to a host cluster by using lshostclustervolumemap*

```
IBM_2145:ITS0-SV1:superuser>lshostclustervolumemap ITS0-ESX-Cluster-01
id name          SCSI_id volume_id volume_name volume_UID          IO_group_id
IO_group_name protocol
0 ITS0-ESX-Cluster-01 0      8      VMware1    60050768019C85144400000000000039 0      io_grp0
scsi
0 ITS0-ESX-Cluster-01 1      9      VMware2    60050768019C8514440000000000003A 0      io_grp0
scsi
0 ITS0-ESX-Cluster-01 2      10     VMware3    60050768019C8514440000000000003B 0      io_grp0
scsi
IBM_2145:ITS0-SV1:superuser>
```

**Note:** The `lshostvdiskmap` command can be run against each host that is part of a host cluster to ensure that the mapping type for the shared volume is shared and is private for the non-shared volume.

## Removing a volume mapping from a host cluster

To remove a volume mapping from a host cluster, run the `rmvolumehostclustermap` command, as shown in Example 7-46.

*Example 7-46 Removing a volume mapping*

---

```
IBM_2145:ITS0-SV1:superuser>rmvolumehostclustermap -hostcluster ITS0-ESX-Cluster-01 VMware3
IBM_2145:ITS0-SV1:superuser>
```

---

In Example 7-46, volume `VMware3` is unmapped from the host cluster `ITS0-ESX-Cluster-01`. The current volume mapping can be checked to ensure that it is unmapped, as shown in Example 7-45 on page 425.

**Note:** To specify the host or hosts that acquire private mappings from the volume that is being removed from the host cluster, specify the `-makeprivate` flag.

## Removing a host cluster member

To remove a host cluster member, run the `rmhostclustermember` command, as shown in Example 7-47.

*Example 7-47 Removing a host cluster member*

---

```
IBM_2145:ITS0-SV1:superuser>rmhostclustermember -host ITS0-VMHOST-02 -removemappings
ITS0-ESX-Cluster-01
IBM_2145:ITS0-SV1:superuser>
```

---

In Example 7-47, the host `ITS0-VMHOST-02` was removed as a member from the host cluster `ITS0-ESX-Cluster-01`, along with the associated volume mappings due to the `-removemappings` flag being specified.

## Removing a host cluster

To remove a host cluster, run the `rmhostcluster` command, as shown in Example 7-48.

*Example 7-48 Removing a host cluster*

---

```
IBM_2145:ITS0-SV1:superuser>rmhostcluster -removemappings ITS0-ESX-Cluster-01
IBM_2145:ITS0-SV1:superuser>
```

---

The `-removemappings` flag also causes the system to remove any host mappings to volumes that are shared. The mappings are deleted before the host cluster is deleted.

**Note:** To keep the volumes mapped to the host objects even after the host cluster is deleted, specify the `-keepmappings` flag instead of `-removemappings` for the `rmhostcluster` command. When `-keepmappings` is specified, the host cluster is deleted, but the volume mapping to the host becomes private instead of shared.

## Adding a host or host cluster to an ownership group

To add a host or a host cluster to an ownership group, use the **chhost** or **chhostcluster** command with the **-ownershipgroup** parameter, as shown in Example 7-49.

*Example 7-49 Adding a host cluster to an ownershipgroup*

---

```
IBM_Storwize:ITS0-V7000:JackUser>chhostcluster -ownershipgroup 1 0
IBM_Storwize:ITS0-V7000:JackUser>
```

---

**Note:** You must specify first the ID of the ownership group you wish to add the host to, and then the ID of the host or host cluster. So the above command in Example 7-49 would add host cluster ID 0 to ownership group ID 1.

## Removing a host or host cluster from an ownership group

To remove a host or a host cluster from an ownership group, use the **chhost** or **chhostcluster** command with the **-noownershipgroup** parameter, as shown in Example 7-50.

*Example 7-50 Removing a host cluster from an ownership group*

---

```
IBM_Storwize:ITS0-V7000:JackUser>chhostcluster -noownershipgroup 1
IBM_Storwize:ITS0-V7000:JackUser>
```

---

This command would remove host cluster 1 from the ownershipgroup assigned to it.





## Storage migration

This chapter describes the steps that are involved in migrating data from an existing external storage system to the capacity of the IBM SAN Volume Controller by using the storage migration wizard. Migrating data from other storage systems to the IBM SAN Volume Controller consolidates storage. It also allows for IBM Spectrum Virtualize features, such as IBM Easy Tier, thin provisioning, compression, encryption, storage replication, and the easy-to-use GUI to be used across all volumes.

Storage migration uses the volume mirroring functionality to allow reads and writes during the migration, and minimizing disruption and downtime. After the migration is complete, the existing system can be retired.

IBM SAN Volume Controller supports migration through Fibre Channel (FC) and internet Small Computer Systems Interface (iSCSI) connections.

This chapter includes the following topics:

- ▶ 8.1, “Storage migration overview” on page 430
- ▶ 8.2, “Storage migration wizard” on page 432

**Note:** This chapter does not cover migration outside of the storage migration wizard. To migrate data outside of the wizard, you must use the Import option in the GUI. This chapter also does not cover virtualization of external storage. For more information about these topics, see Chapter 5, “Storage pools” on page 199.

## 8.1 Storage migration overview

To migrate data from an existing storage system to the IBM SAN Volume Controller, it is necessary to use the built-in External Virtualization capability. This capability places external connected logical units (LUs) under the control of the IBM SAN Volume Controller, which acts as a proxy while hosts continue to access them. The volumes are then fully virtualized in the IBM SAN Volume Controller.

**Attention:** A license is required for external systems that are being virtualized based on Storage Capacity Units (SCUs). Data can be migrated from existing storage systems to your system by using the External Virtualization function within 45 days of purchase of the system without purchase of a license. After 45 days, any ongoing use of the External Virtualization function requires a license.

Set the license temporarily during the migration process to prevent messages that indicate that you are in violation of the license agreement from being sent. When the migration is complete, or after 45 days, reset the license to its original limit or purchase a new license.

Consider the following points about the storage migration process:

- ▶ Typically, storage systems divide storage into many Small Computer System Interface (SCSI) LUs that are presented to hosts.
- ▶ I/O to the LUs must be stopped and changes made to the mapping of the external storage system LUs and to the fabric or iSCSI configuration so that the original LUs are presented directly to the IBM SAN Volume Controller and not to the hosts anymore. The IBM SAN Volume Controller discovers the external LUs as *unmanaged* managed disks (MDisks).
- ▶ The unmanaged MDisks are *imported* to the IBM SAN Volume Controller as *image-mode volumes* and placed into a temporary storage pool. This storage pool is now a logical container for the LUs.
- ▶ Each MDisk has a one-to-one mapping with an image-mode volume. From a data perspective, the image-mode volumes represent the LUs exactly as they were before the import operation. The image-mode volumes are on the same physical drives of the external storage system and the data remains unchanged. The IBM SAN Volume Controller is presenting active images of the LUs and is acting as a proxy.
- ▶ The hosts must have the existing storage system multipath device driver removed, and are then configured for IBM SAN Volume Controller attachment. The IBM SAN Volume Controller hosts are defined with worldwide port names (WWPNs) or iSCSI Qualified Names (IQNs), and the volumes are mapped to the hosts. After the volumes are mapped, the hosts discover the IBM SAN Volume Controller volumes through a host rescan or reboot operation.
- ▶ IBM Spectrum Virtualize volume mirroring operations are then started. The image-mode volumes are mirrored to generic volumes. Volume mirroring is an online migration task, which means a host can still access and use the volumes during the mirror synchronization process.
- ▶ After the mirror operations are complete, the image-mode volumes are removed. The external storage system LUs are now migrated and the now redundant storage can be decommissioned or reused elsewhere.



**Important:** If you are migrating volumes from an IBM Storwize or FlashSystem family product, be aware that the source system must be configured in the *storage* layer. Otherwise, the IBM SAN Volume Controller does not discover the source as a backend controller.

The default layer setting for Storwize and FlashSystem family systems is *storage*. For more information about layers and how to change them, see Chapter 5, “Storage pools” on page 199.

## 8.1.1 Interoperability and compatibility

Interoperability is an important consideration when a new storage system is set up in an environment that contains existing storage infrastructure. Before attaching any external storage systems to the IBM SAN Volume Controller, see the [IBM System Storage Interoperation Center \(SSIC\)](#).

Select **IBM System Storage IBM SAN Volume Controller** in Storage Family, then **SVC Storage Controller Support** in Storage Model. You can then refine your search by selecting the external storage controller that you want to use in the **Storage Controller** menu.

The matrix results give you indications about the external storage that you want to attach to the IBM SAN Volume Controller, such as minimum firmware level or support for disks greater than 2 TB.

## 8.1.2 Prerequisites

Before the storage migration wizard can be started, the external storage system must be visible to the IBM SAN Volume Controller. You also need to confirm that the restrictions and prerequisites are met.

Administrators can migrate data from the external storage system to the system that uses iSCSI or FC or Fibre Channel over Ethernet (FCoE) connections. For more information about how to manage external storage, see Chapter 5, “Storage pools” on page 199.

### Common prerequisites

If you have VMware ESX server hosts, you must change settings on the VMware host so copies of the volumes can be recognized by the system after the migration is completed. To enable volume copies to be recognized by the system for VMware ESX hosts, you must complete one of the following actions:

- ▶ Enable the EnableResignature setting.
- ▶ Disable the DisallowSnapshotLUN setting.

To learn more about these settings, see the documentation for the VMware ESX host.

**Note:** Test the setting changes on a non-production server. The logical unit number (LUN) has a different unique identifier (UID) after it is imported. It resembles a mirrored volume to the VMware server.

## Prerequisites for Fibre Channel Connection

The following prerequisites for Fibre Channel Connection (FICON) must be met:

- ▶ Cable this system into the SAN of the external storage that you want to migrate. Ensure that your system is cabled into the same SAN as the external storage system that you are migrating.
- ▶ If you are using FC, connect the FC cables to the FC ports in *both* nodes of your system, and then to the FC network. If you are using FCoE, connect Ethernet cables to the 10 Gbps Ethernet ports.

For more information, see 2.6, “Fibre Channel SAN configuration planning” on page 58. Alternatively, directly attach the external storage system to the IBM SAN Volume Controller nodes instead of using a switched fabric.

## Prerequisites for iSCSI connections

The following prerequisites for iSCSI connections must be met:

- ▶ Cable this system to the external storage system with a redundant switched fabric. Migrating iSCSI external storage requires that the system and the storage system are connected through an Ethernet switch. Symmetric ports on *all* nodes of the system must be connected to the same switch and must be configured on the same subnet.
- ▶ Modify the Ethernet port attributes to enable the external storage on the Ethernet port to enable external storage connectivity. To modify the Ethernet port for external storage, click **Network** → **Ethernet Ports** and right-click a configured port. Select **Modify Storage Ports** to enable the port for external storage connections.
- ▶ Cable the Ethernet ports on the storage system to fabric in the same way as the system and ensure that they are configured in the same subnet. Optionally, you can use a virtual local area network (VLAN) to define network traffic for the system ports.
- ▶ For full redundancy, configure two Ethernet fabrics with separate Ethernet switches. If the source system nodes and the external storage system both have more than two Ethernet ports, extra redundant iSCSI connection can be established for increased throughput.

## 8.2 Storage migration wizard

The storage migration wizard simplifies the migration task. The wizard features easy-to-follow windows that guide users through the entire process. The wizard shows you which commands are being run so that you can see exactly what is being performed throughout the process.

**Attention:** The risk of losing data when the storage migration wizard is used correctly is low. However, it is prudent to avoid potential data loss by creating a backup of all the data that is stored on the hosts, the existing storage systems, and the IBM SAN Volume Controller before the wizard is used.

Complete the following steps to complete the migration by using the storage migration wizard:

1. Select **Pools** → **System Migration**, as shown in Figure 8-1 on page 433. The System Migration pane provides access to the storage migration wizard and displays information about the migration progress.

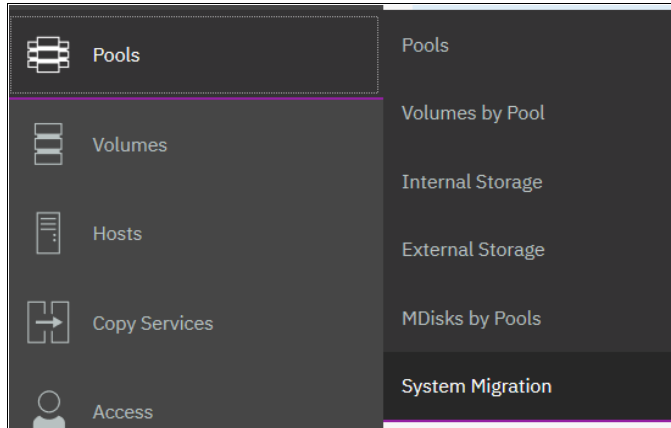


Figure 8-1 Navigating to System Migration

2. Click **Start New Migration** to begin the storage migration wizard, as shown in Figure 8-2.

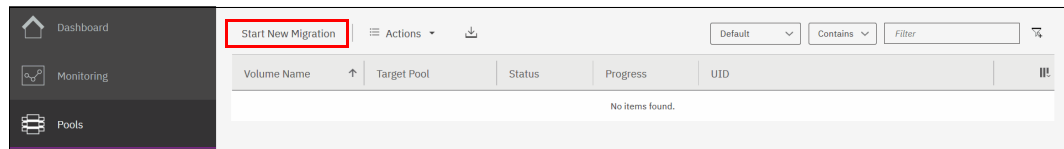


Figure 8-2 Starting a migration

**Note:** Starting a new migration adds the volume to be migrated to the list displayed in the pane. After a volume is migrated, it remains in the list until you finalize the migration.

3. If both FC and iSCSI external systems are detected, a dialogue is shown prompting you to decide which protocol should be used. Select the type of attachment between the IBM SAN Volume Controller and the external system from which you want to migrate volumes and click **Next**. If only one type of attachment is detected, this dialogue is not displayed.  
If the external storage system is not detected, the warning message that is shown in Figure 8-3 is displayed when you attempt to start the migration wizard. Click **Close** and correct the problem before you try to start the migration wizard again.

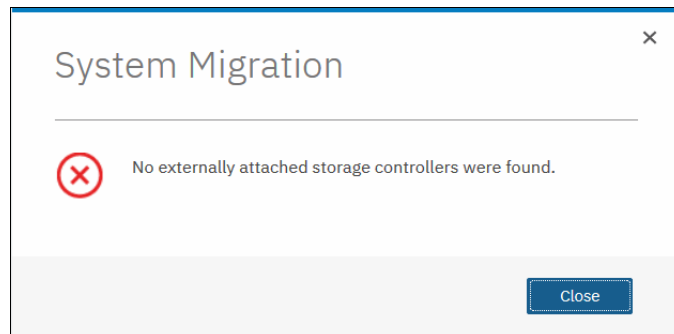


Figure 8-3 Error message if no external storage is detected

4. When the wizard starts, you are prompted to verify the restrictions and prerequisites that are listed in Figure 8-4 on page 435. Address the following restrictions and prerequisites:

– Restrictions:

- You are not using the storage migration wizard to migrate clustered hosts, including clusters of VMware hosts and Virtual I/O Servers (VIOS).
- You are not using the storage migration wizard to migrate SAN boot images.

If you have either of these two environments, the migration must be performed outside of the wizard because more steps are required.

The VMware vSphere Storage vMotion feature might be an alternative for migrating VMware clusters. For information, see this [web page](#).

– Prerequisites:

- IBM SAN Volume Controller nodes and the external storage system are connected to the same SAN fabric.
- If there are VMware ESX hosts involved in the data migration, the VMware ESX hosts are set to allow volume copies to be recognized.

For more information about the Storage Migration prerequisites, see 8.1.2, “Prerequisites” on page 431.

If all restrictions are satisfied and prerequisites are met, select all of the boxes and click **Next**, as shown in Figure 8-4.

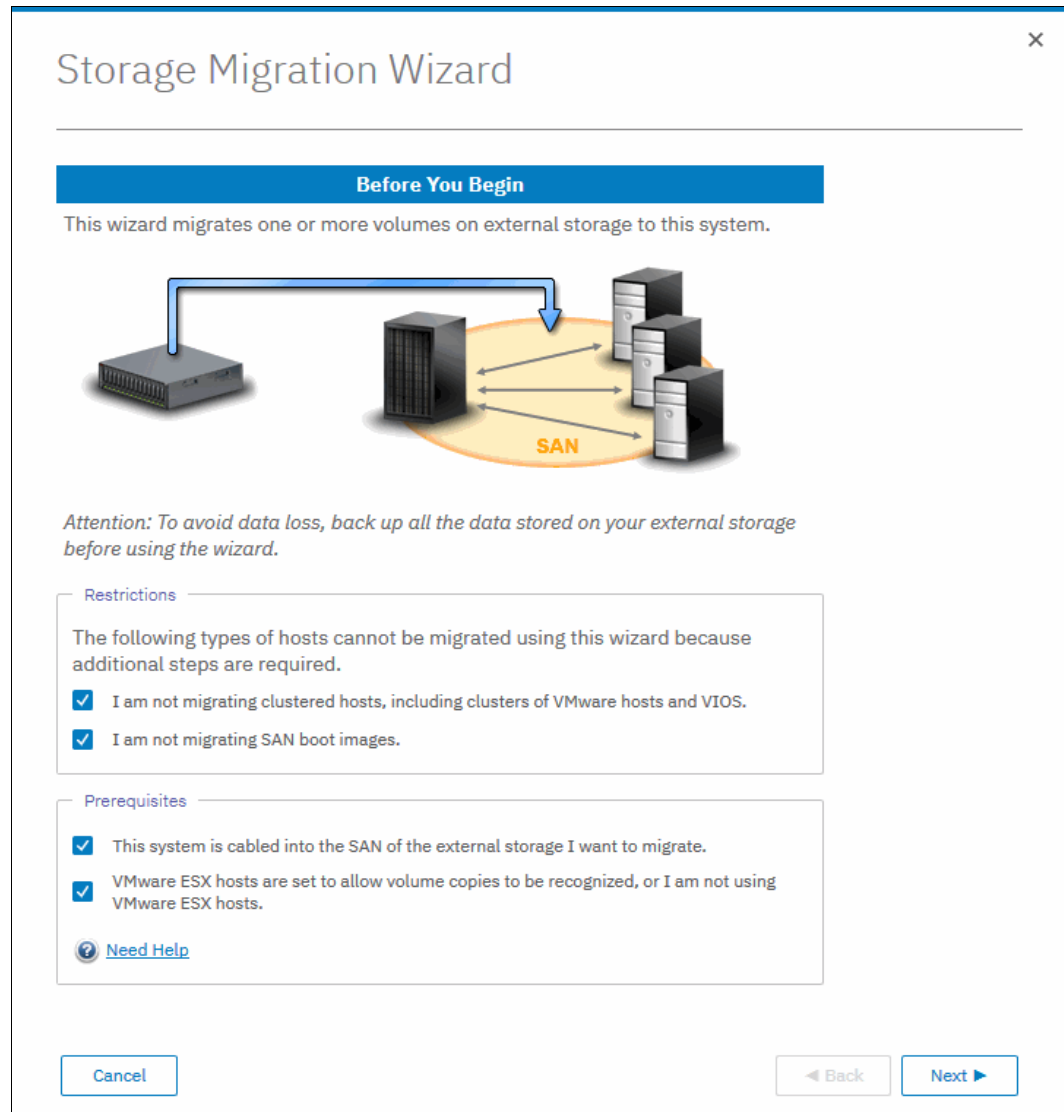


Figure 8-4 Restrictions and prerequisites confirmation

5. Prepare the environment migration by following the instructions that are shown in Figure 8-5.

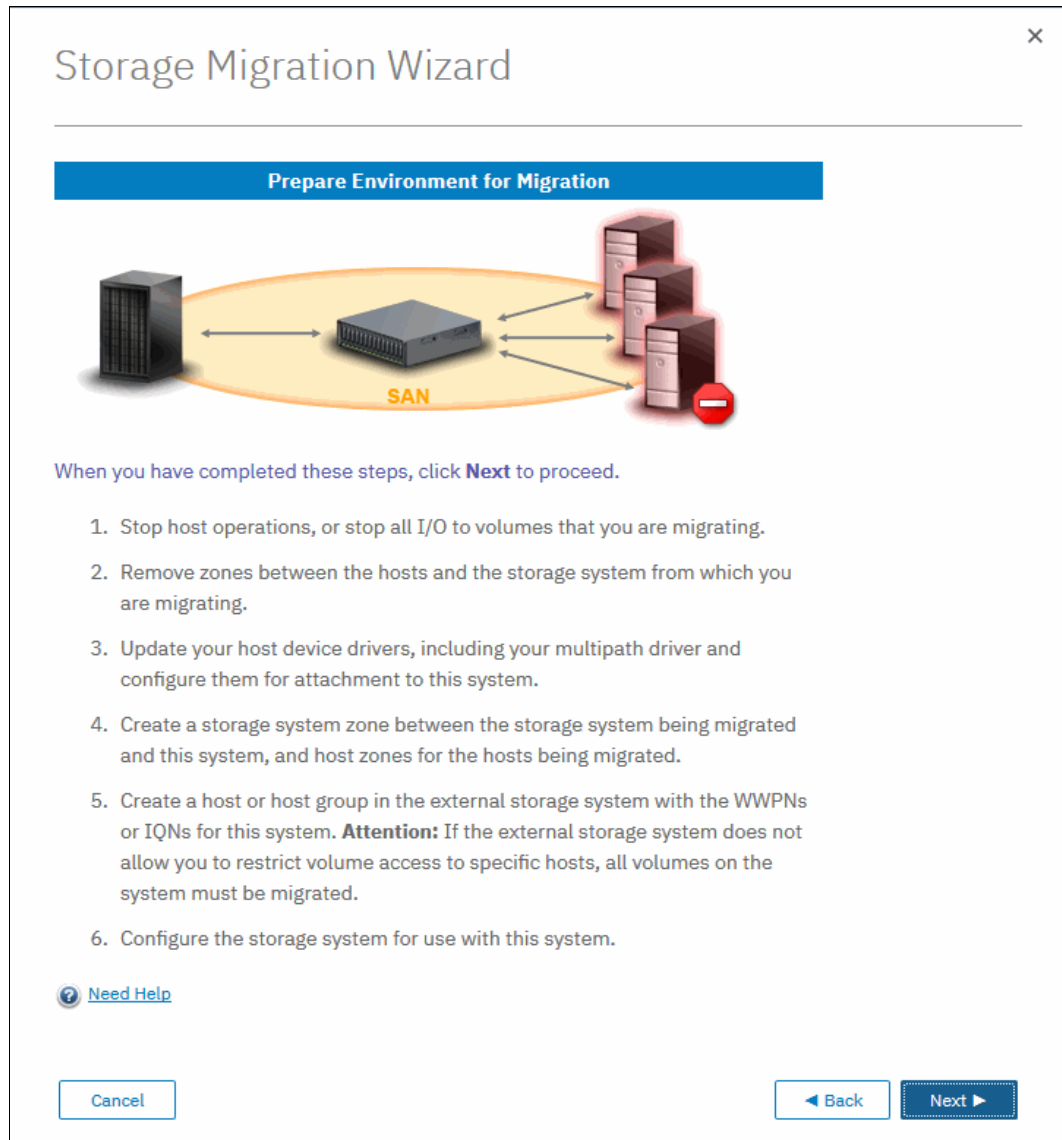


Figure 8-5 Preparing your environment for storage migration

The preparation phase includes the following steps:

- a. Before migrating storage, ensure that all host operations are stopped to prevent applications from generating I/Os to the migrated system.
- b. Remove all existing zones between the hosts and the system you are migrating.
- c. Hosts usually do not support concurrent multipath drivers at the same time. You might need to remove drivers that are not compatible with the IBM SAN Volume Controller, from the hosts and use the recommended device drivers. For more information about supported drivers, see the [IBM SSIC](#).

- d. If you are migrating external storage systems that connect to the system that uses FC or FCoE connections, ensure that you complete suitable zoning changes to simplify migration. Use the following guidelines to ensure that zones are configured correctly for migration:
  - Zoning rules

For every storage system, create one zone that contains this system's ports from every node and all external storage system ports, unless otherwise stated by the zoning guidelines for that storage system.

This system requires single-initiator zoning for all large configurations that contain more than 64 host objects. Each server FC port must be in its own zone, which contains the FC port and this system's ports. In configurations of fewer than 64 hosts, you can have up to 40 FC ports in a host zone if the zone contains similar host bus adapters (HBAs) and operating systems.
  - Storage system zones

In a storage system zone, this system's nodes identify the storage systems. Generally, create one zone for each storage system. Host systems cannot operate on the storage systems directly. All data transfer occurs through this system's nodes.
  - Host zones

In the host zone, the host systems can identify and address this system's nodes. You can have more than one host zone and more than one storage system zone. Create one host zone for each host FC port.
  - Because the IBM SAN Volume Controller should now be seen as a host cluster from the external system to be migrated, you must define the IBM SAN Volume Controller as a host or host group by using the WWPNs or IQNs, on the system to be migrated. Some legacy systems do not permit LUN-to-host mapping and would then present all the LUs to the IBM SAN Volume Controller. In that case, all the LUs should be migrated.
6. If the previous preparation steps are followed, the IBM SAN Volume Controller is now seen as a host from the system to be migrated. LUs can then be mapped to the IBM SAN Volume Controller. Map the external storage system by following the instructions that are shown in Figure 8-6 on page 438.

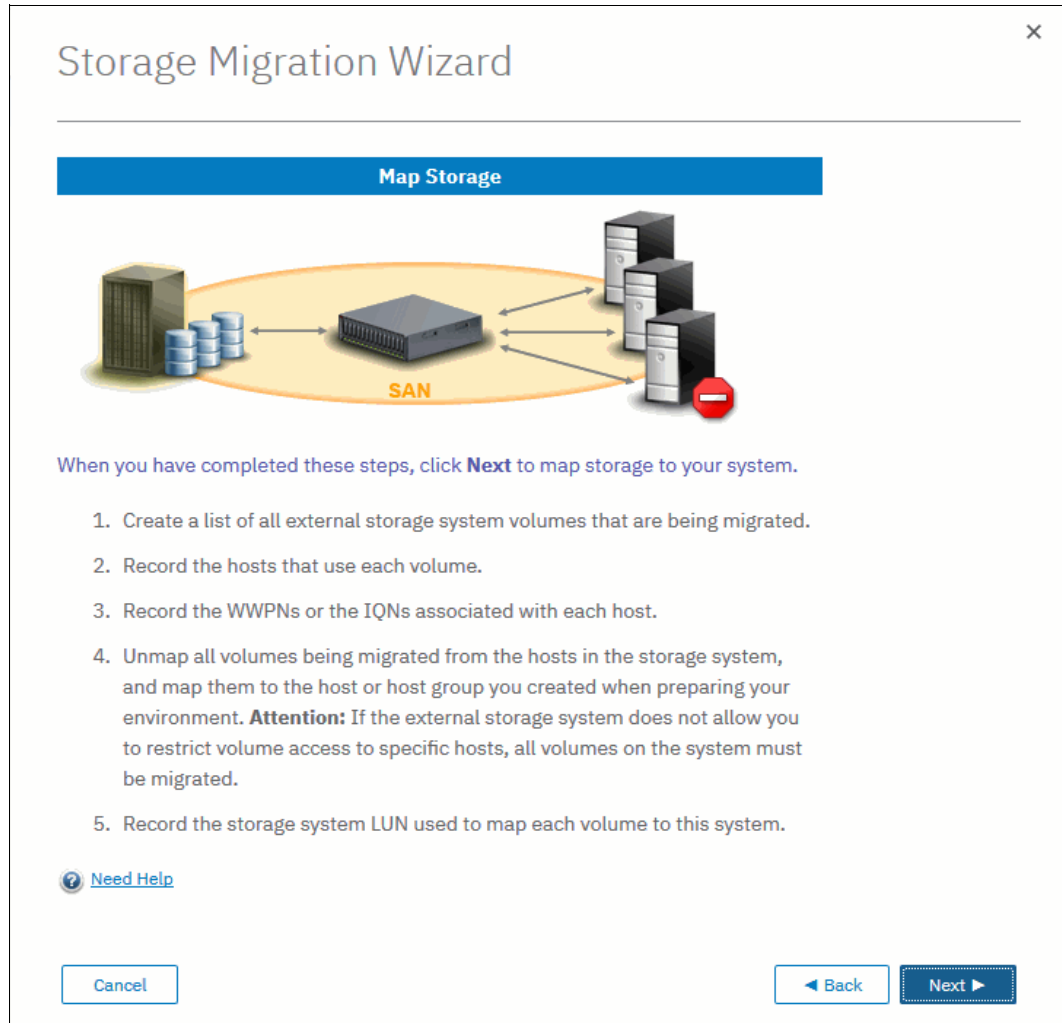


Figure 8-6 Steps to map the LUs to be migrated to the IBM SAN Volume Controller

Before you migrate storage, record the hosts and their WWPNs or IQNs for each volume that is being migrated, and the SCSI LUN when mapped to the IBM SAN Volume Controller.

Table 8-1 shows an example of a table that is used to capture information that relates to the external storage system LUs.

Table 8-1 Example table for capturing external LU information

Volume name or ID	Hosts accessing this LUN	Host WWPNs or IQNs	SCSI LUN when mapped
1 IBM DB2® logs	DB2server	21000024FF2...	0
2 DB2 data	DB2Server	21000024FF2...	1
3 file system	FileServer1	21000024FF2...	2

**Note:** Make sure to record the SCSI ID of the LUs to which the host is originally mapped. Some operating systems do not support changing the SCSI ID during the migration.



Click **Next** and wait for the system to discover external devices.

7. The next window shows all of the MDisks that were found. If the MDisks to be migrated are not in the list, check your zoning or IP configuration, as applicable, and your LUN mappings. Repeat the previous step to trigger the discovery procedure again.

Select the MDisks that you want to migrate, as shown in Figure 8-7. In this example, one MDisk was found and will be migrated: `mdisk8`. Detailed information about an MDisk is available by double-clicking it. To select multiple elements from the table, use the standard **Shift**+left-click or **Ctrl**+left-click actions. Optionally, you can export the discovered MDisks list to a comma-separated value (CSV) file, for further use, by clicking **Export to CSV**.

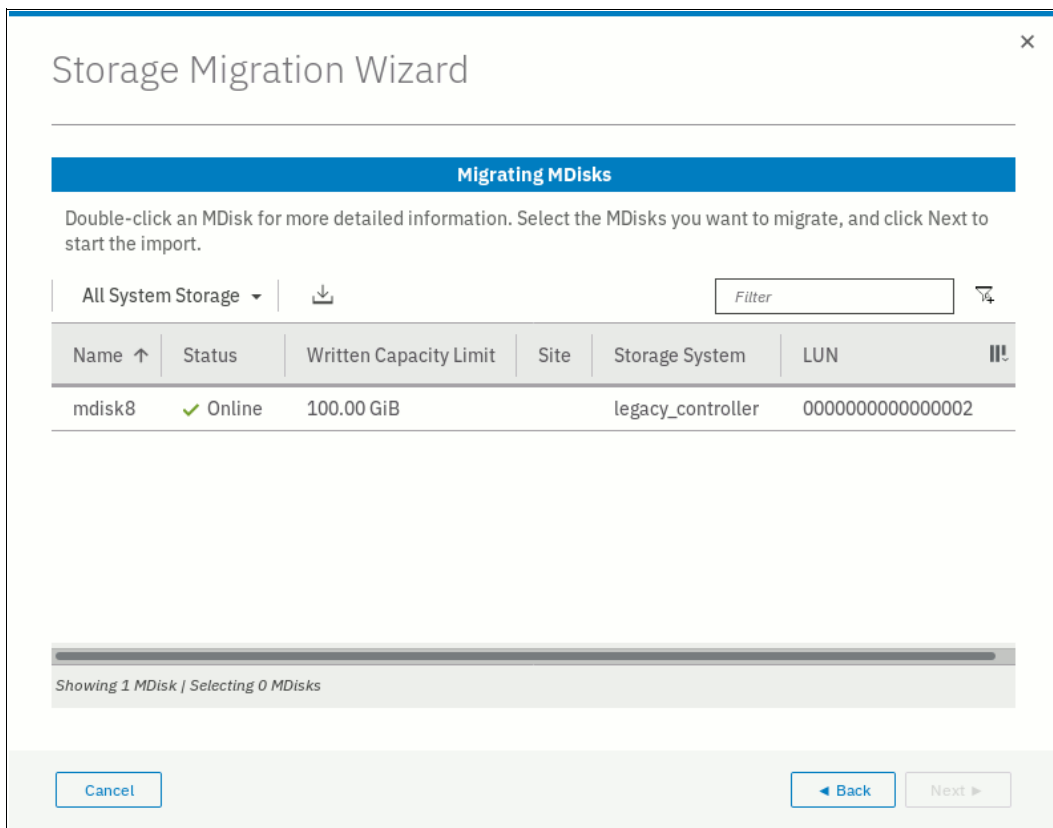


Figure 8-7 Discovering mapped LUs from external storage

**Note:** Select only the MDisks that are applicable to the current migration plan. After step 15 on page 446 of the current migration completes, another migration can be started to migrate any remaining MDisks.

8. Click **Next** and wait for the MDisk to be imported. During this task, the system creates a new storage pool that is called `MigrationPool1_XXXX` and adds the imported MDisk to the storage pool as image-mode volumes.

9. The next window lists all of the hosts that are configured on the system and enables you to configure new hosts. This step is optional and can be bypassed by clicking **Next**. In this example, the host `linuxsrv` is configured, as shown in Figure 8-8. If no host is selected, you can create a host after the migration completes and map the imported volumes to it.

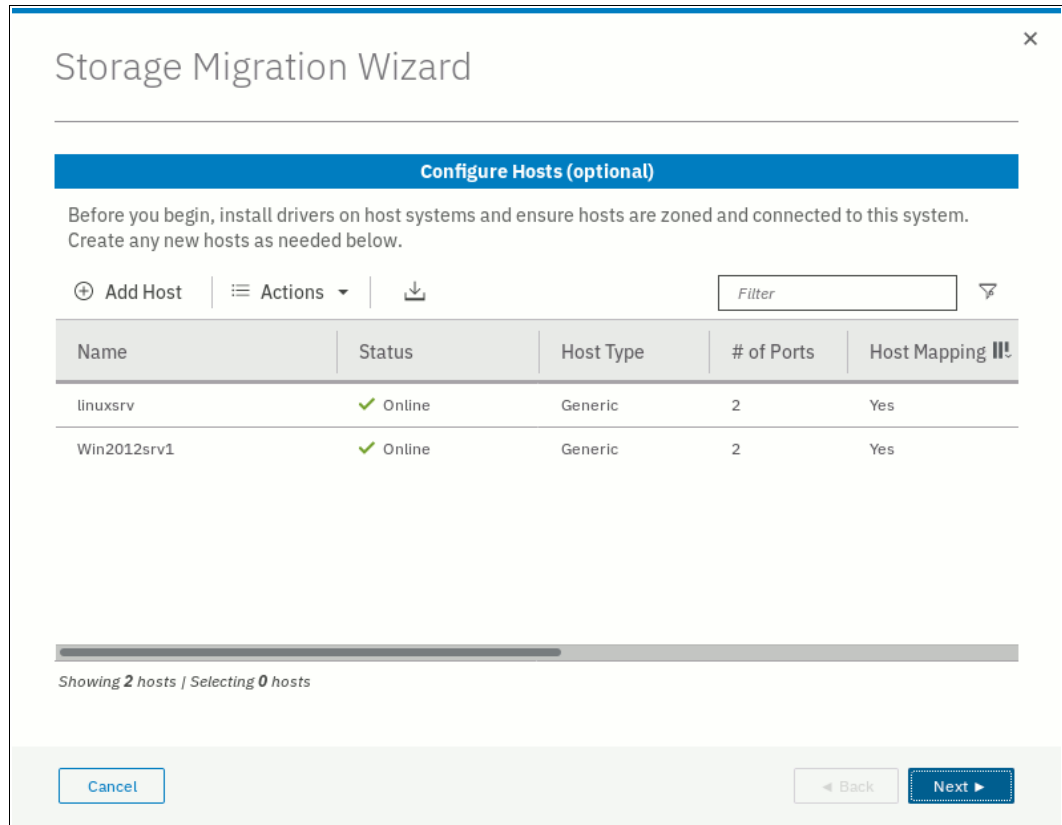


Figure 8-8 Listing of configured hosts to map the imported Volume to

10. If the host that needs access to the migrated data is not configured, select **Add Host** to begin the Add Host wizard. Enter the host connection type, name, and connection details. Optionally, click **Advanced** to modify the host type and I/O group assignment. Figure 8-9 shows the Add Host wizard with the details completed.

For more information about the Add Host wizard, see Chapter 7, “Hosts” on page 351.

**Add Host** [X]

**Required Fields**

Name: ISCSI HOST

Host connections: iSCSI (SCSI)

Host IQN: iqn.1994-05.com.redhat:72dcb719 (+) (-)

**Optional Fields**

CHAP authentication:

CHAP secret: Enter 1 to 79 characters

CHAP username: Enter 1 to 31 characters

Host type: Generic

I/O groups: All

Host cluster: No Host Cluster Selected

Cancel Add

Figure 8-9 If not already defined, you can create a host during the migration process

11. Click **Add**. The host is created and is now listed in the Configure Hosts window, as shown in Figure 8-8 on page 440. Click **Next** to proceed.

12. The next window lists the new volumes and enables you to map them to hosts as shown in Figure 8-10. The volumes are listed with names that were automatically assigned by the system. The names can be changed to reflect something more meaningful to the user by selecting the volume and clicking **Rename** in the **Actions** menu.

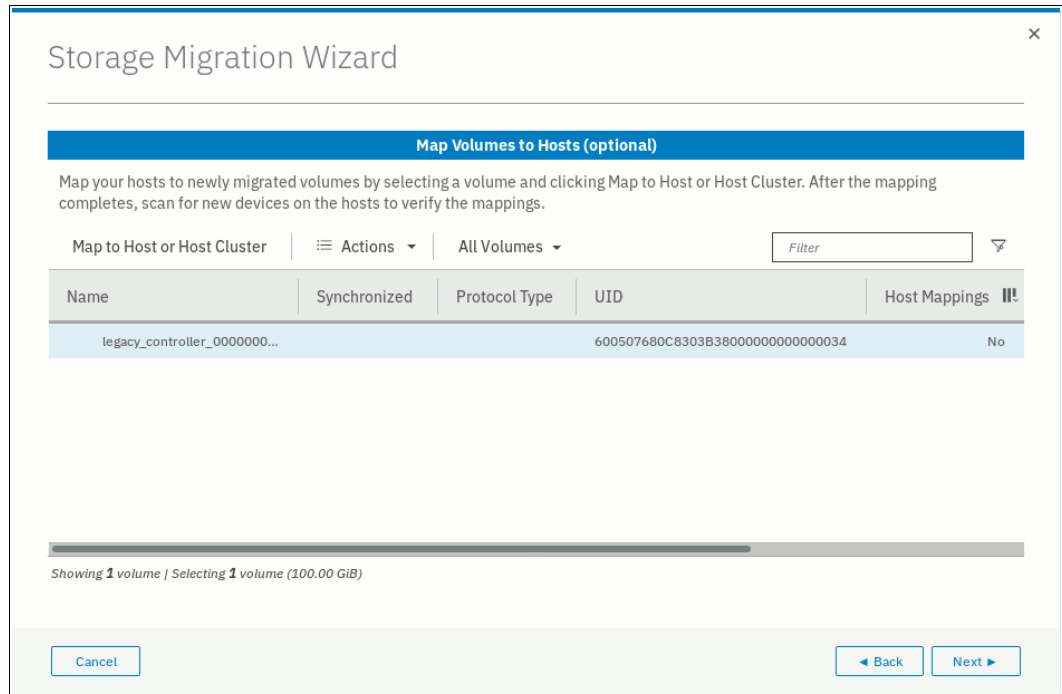


Figure 8-10 Map volumes to hosts

13. Map the volumes to hosts by selecting the volumes and clicking **Map to Host or Host Cluster**, as shown in Figure 8-11. This step is optional and can be bypassed by clicking **Next**.

Create Mapping

Create Mappings to:

Hosts

Host Clusters

Select hosts to map to legacy\_controller\_0000000000000002

Default Contains Filter

Name	Status	Host Type	Host Mappings	Ownership Group
ISCSI HOST	Offline	Generic	No	
Win2012srv1	Online	Generic	Yes	
linuxsrv	Online	Generic	Yes	

Showing 3 hosts | Selecting 1 host

Would you like the system to assign SCSI LUN IDs or manually assign these IDs?

System Assign

Self Assign

Cancel Back Next

Figure 8-11 Select the host to map the new Volume to

You can manually assign a SCSI ID to the LUNs you are mapping. This technique is useful when the host must have the same LUN ID for a LUN before and after it is migrated. To assign the SCSI ID manually, select the **Self Assign** option and follow the instructions as shown in Figure 8-12 on page 444.

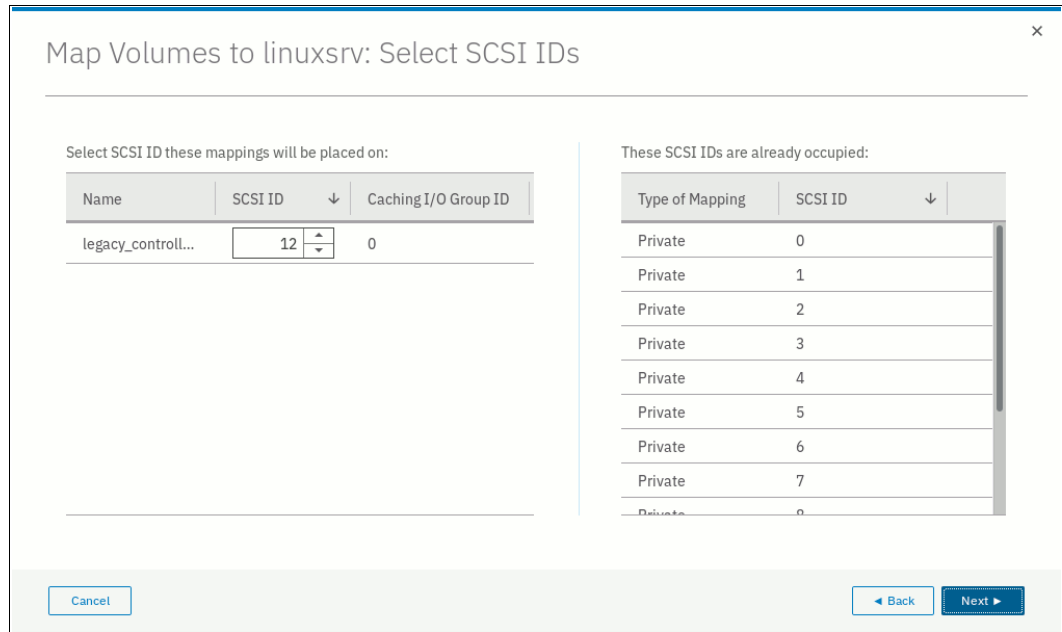


Figure 8-12 Manually assign a LUN SCSI ID to mapped Volume

When your LUN mapping is ready, click **Next**. A new window is displayed with a summary of the new and existing mappings, as shown in Figure 8-13.

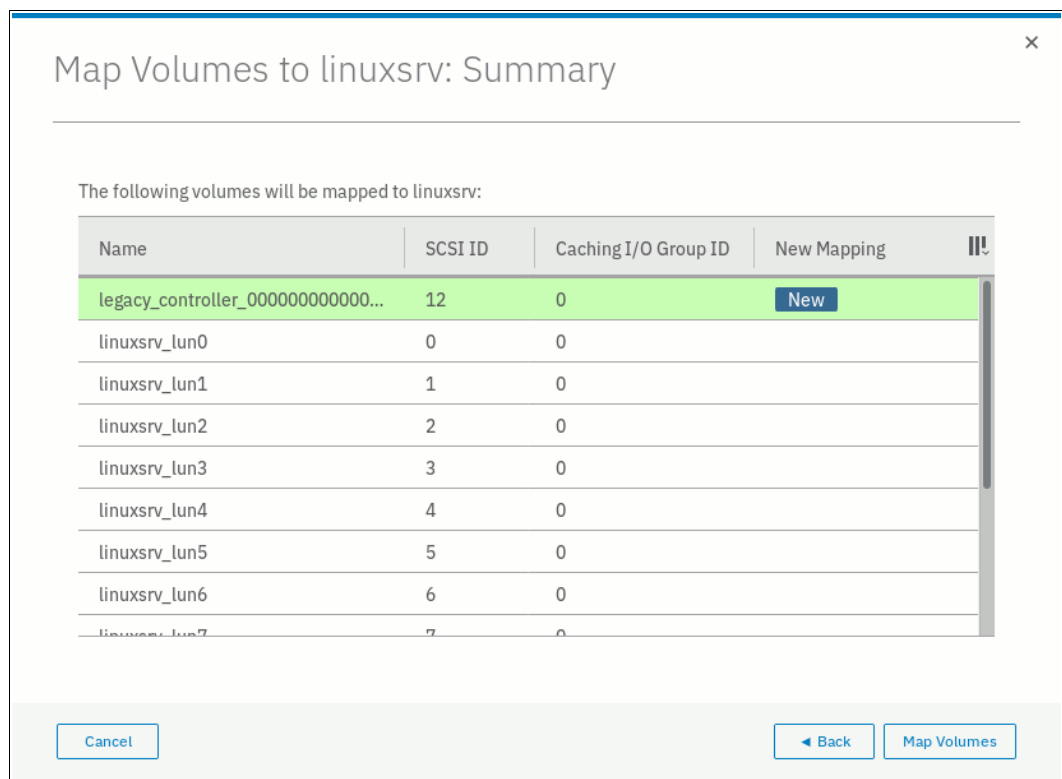


Figure 8-13 Volumes mapping summary before migration

Click **Map Volumes** and wait for the mappings to be created. Continue to map volumes to hosts until all mappings are created. Then, click **Next** to continue with the next migration step.

14. Select the storage pool that you want to migrate the imported volumes into. Ensure that the selected storage pool has enough space to accommodate the migrated volumes before you continue. This is an optional step. You can decide not to migrate to a storage pool and to leave the imported MDisk as an image-mode volume.

However, this technique is not recommended because no volume mirroring is created. Therefore, no protection is available for the imported MDisk, and no data transfer occurs from the storage system to be migrated to the system. Click **Next**, as shown in Figure 8-14.

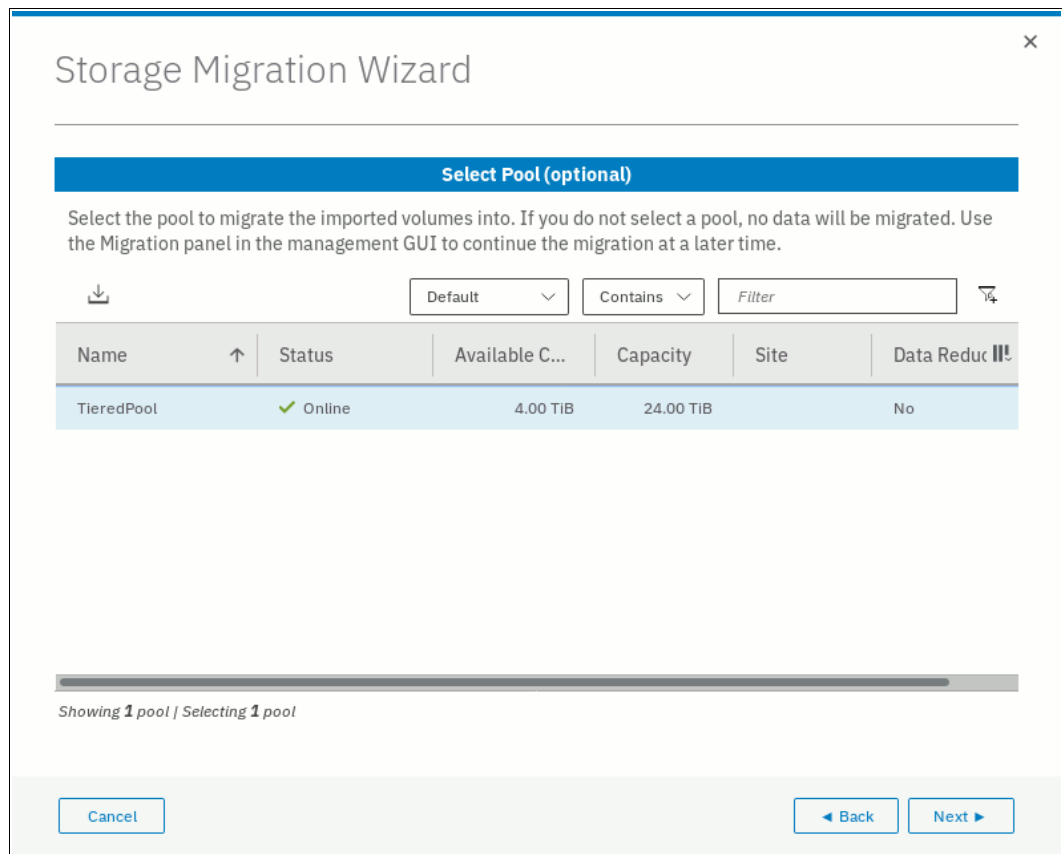


Figure 8-14 Select the pool to migrate the MDisk to

The migration starts. This task continues running in the background and uses the volume mirroring function to place a generic copy of the image-mode volumes in the selected storage pool.

For more information about Volume Mirroring, see Chapter 6, “Volumes” on page 255.

**Note:** With volume mirroring, the system creates two copies (Copy0 and Copy1) of a volume. Typically, Copy0 is in the migration pool, and Copy1 is created in the target pool of the migration. When the host generates a write I/O on the volume, data is written at the same time on both copies. Read I/Os are performed on the primary copy only.

In the background, a mirror synchronization of the two copies is performed and runs until the two copies are synchronized. The speed of this background synchronization can be changed in the volume properties.

For more information about volume mirroring synchronization rate, see Chapter 6, “Volumes” on page 255.

15. Click **Finish** to end the storage migration wizard, as shown in Figure 8-15.

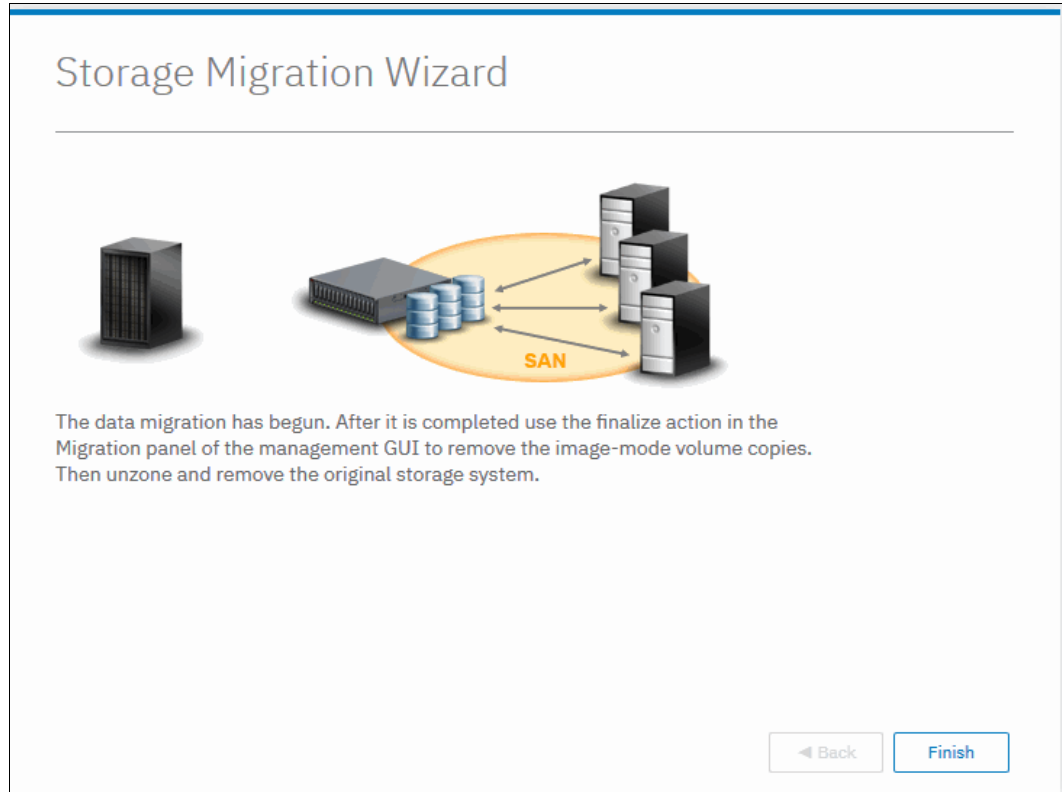


Figure 8-15 Migration is started

The end of the wizard is not the end of the migration task. You can find the progress of the migration in the Storage Migration window, as shown in Figure 8-16. The target storage pool and the progress of the volume copy synchronization is also displayed there.

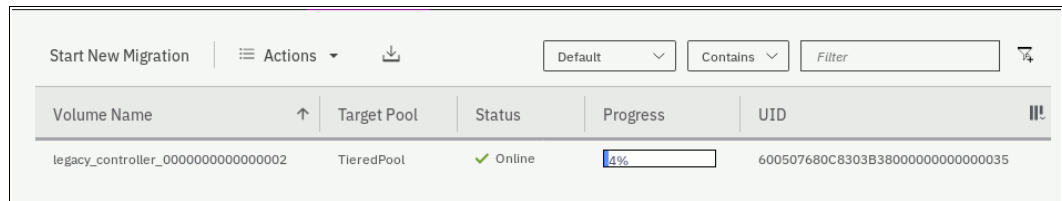


Figure 8-16 The ongoing Migration is listed in the Storage Migration window

16. If you want to check the progress by using the CLI, run the **lsvdisksyncprogress** command because the process is essentially a volume copy, as shown in Example 8-1.

*Example 8-1 Migration progress on the CLI*

```

IBM_2145:ITS0-SV1:superuser>lsvdisksyncprogress
vdisk_id vdisk_name                progress estimated_completion_time
20      legacy_controller_0000000000000002 1      191021123932

```





18. When finalized, the image-mode copies of the volumes are deleted and the associated MDisks are removed from the migration pool. The status of those MDisks returns to unmanaged. You can verify the status of the MDisks by selecting **Pools** → **External Storage**, as shown in Figure 8-19. In the example, `mdisk3` was migrated and finalized, and it displays as unmanaged in the external storage window.

Name	State	Written Capacity Limit	Mode	Site	Pool	Storage System	LUN
> flashsystem	✓ Online	IBM 2145	Serial Number: 2076		Site: Unassigned	WWNN: 500507680B009479	
flashsystem	✓ Online	IBM 2145	Serial Number: 2076		Site: Unassigned	WWNN: 5005076810000F62	
∨ legacy_cont	✓ Online	IBM 2145	Serial Number: 2076		Site: Unassigned	WWNN: 5005076810000F88	
mdisk7	✓ Online	1.00 TiB	Managed		TieredPool	legacy_controller	0000000000000001
mdisk5	✓ Online	1.00 TiB	Managed		TieredPool	legacy_controller	0000000000000000
mdisk8	✓ Online	100.00 GiB	Unmanaged			legacy_controller	0000000000000002
flashsystem	✓ Online	IBM 2145	Serial Number: 2076		Site: Unassigned	WWNN: 500507680B009478	

Figure 8-19 External Storage MDisks window

All the steps that are described in the Storage Migration wizard can be performed manually by using the CLI, but it is highly recommended to use the wizard as a guide.

**Note:** For a real-world demonstration of the storage migration capabilities offered with IBM Spectrum Virtualize, see [this web page](#) (log in required).

The demonstration includes three different step-by-step scenarios showing the integration of an IBM SAN Volume Controller cluster into an environment with one Microsoft Windows Server (image mode), one IBM AIX server (Logical Volume Manager (LVM) mirroring), and one VMware ESXi server (storage vMotion).



## Advanced features for storage efficiency

IBM Spectrum Virtualize running inside the IBM SAN Volume Controller offers several functions for storage optimization and efficiency.

This chapter introduces the basic concepts of those functions. It also provides a short technical overview and implementation recommendations.

For more information about planning and configuration of storage efficiency features, see the following publications:

- ▶ *IBM System Storage SAN Volume Controller, IBM Storwize V7000, and IBM FlashSystem 7200 Best Practices and Performance Guidelines*, SG24-7521
- ▶ *Introduction and Implementation of Data Reduction Pools and Deduplication*, SG24-8430
- ▶ *IBM Real-time Compression in IBM SAN Volume Controller and IBM Storwize V7000*, REDP-4859
- ▶ *Implementing IBM Real-time Compression in SAN Volume Controller and IBM Storwize V7000*, TIPS1083
- ▶ *Implementing IBM Easy Tier with IBM Real-time Compression*, TIPS1072

This chapter includes the following topics:

- ▶ 9.1, “IBM Easy Tier” on page 450
- ▶ 9.2, “Thin-provisioned volumes” on page 466
- ▶ 9.3, “Unmap” on page 468
- ▶ 9.4, “Data Reduction Pools” on page 471
- ▶ 9.5, “Compression with standard pools” on page 482
- ▶ 9.6, “Saving estimation for compression and deduplication” on page 484
- ▶ 9.7, “Overprovisioning and data reduction on external storage” on page 486

## 9.1 IBM Easy Tier

IBM Spectrum Virtualize includes the IBM System Storage Easy Tier function. It enables automated subvolume data placement throughout different storage tiers and automatically moves extents within the same storage tier to intelligently align the system with current workload requirements. It also optimizes the usage of Flash drives or flash arrays.

Many applications exhibit a significant skew in the distribution of I/O workload: a small fraction of the storage is responsible for a disproportionately large fraction of the total I/O workload of an environment.

Easy Tier acts to identify this skew and to automatically place data to take advantage of it. By moving the “hottest” data onto the fastest tier of storage, the workload on the remainder of the storage is significantly reduced. By servicing most of the application workload from the fastest storage, Easy Tier acts to accelerate application performance, and increase overall server utilization. This can reduce costs in servers and application licenses.

Easy Tier also reduces storage cost because the system always places the data with the highest I/O workload on the fastest tier of storage. This means depending on the workload pattern, a large portion of the capacity can be provided by a lower and less expensive tier without impacting application performance.

**Note:** Easy Tier is a licensed function, but it is included in base code. No actions are required to activate the Easy Tier license on IBM SAN Volume Controller.

### 9.1.1 Easy Tier concepts

Easy Tier is a performance optimization function that automatically migrates (or moves) extents that belong to a volume between different storage tiers, based on their I/O load. Movement of the extents is online and unnoticed from the host perspective. As a result of extent movement, the volume no longer has all its data in one tier, but rather in two or more tiers, and each tier provides optimal performance for the extent, as shown in Figure 9-1.

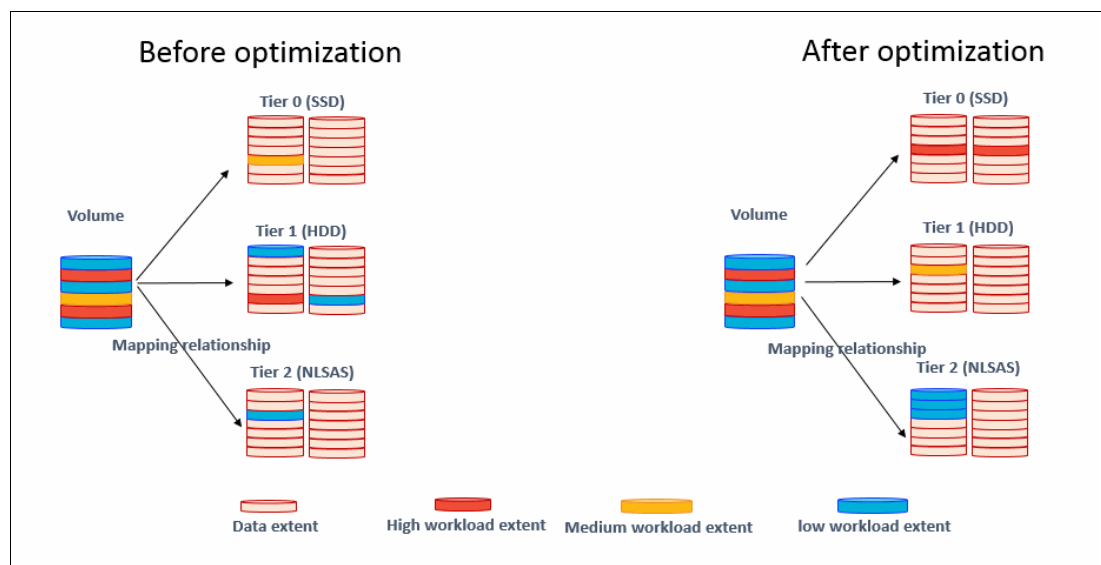


Figure 9-1 Easy Tier

Easy Tier monitors the I/O activity and latency of the extents on all Easy Tier enabled storage pools. Based on the performance log, it creates an extent migration plan and *promotes* (moves) high activity or hot extents to a higher disk tier within the same storage pool. It also *demotes* extents whose activity dropped off (or cooled) by moving them from a higher disk tier managed disk (MDisk) back to a lower tier MDisk.

If a pool only contains MDisks of a single tier, Easy Tier only operates in balancing mode. Extents are moved between MDisks in the same tier to balance I/O workload within that tier.

### Tiers of storage

The MDisks (external logical units [LUs] or Redundant Array of Independent Disks [RAID] arrays) that are presented to the IBM SAN Volume Controller might have different performance attributes because of their technology type, such as Flash or spinning drives and other characteristics.

The system divides available storage into the following tiers:

- ▶ Storage-class memory (SCM)
  - SCM tier exists when the pool contains drives that use persistent memory technologies that improves endurance and speed of current flash storage device technologies.
- ▶ Tier 0 flash
  - Tier 0 flash drives are high-performance flash drives that use enterprise flash technology.
- ▶ Tier 1 flash
  - Tier 1 flash drives represent the read-intensive flash drive technology. Tier 1 flash drives are lower-cost flash drives that typically offer capacities larger than enterprise class flash, but lower performance and write endurance characteristics.
- ▶ Enterprise tier
  - Enterprise tier exists when the pool contains MDisks on enterprise-class hard disk drives (HDDs), which are disk drives that are optimized for performance.
- ▶ Nearline (NL) tier
  - NL tier exists when the pool has MDisks on NL-class disks drives that are optimized for capacity.

The system automatically sets the tier for internal array mode MDisks because it knows the capabilities of array members, physical drives, or modules. External MDisks need manual tier assignment when they are added to a storage pool.

**Note:** The tier of MDisks that are mapped from certain types of IBM System Storage Enterprise Flash is fixed to tier0\_flash, and cannot be changed.

Although IBM SAN Volume Controller can distinguish between five tiers, Easy Tier manages only a three tier storage architecture within each storage pool. MDisk tiers are mapped to Easy Tier tiers, depending on the pool configuration, as listed in Table 9-1.

Table 9-1 Storage Tier to Easy Tier mapping

Configuration	ET Top Tier	ET Middle Tier	ET Bottom Tier
SCM (+ Tier0_Flash)	SCM	(Tier0_Flash)	
SCM + Tier0_Flash (+ Tier1_Flash)	SCM	Tier0_Flash	(Tier1_Flash)

Configuration	ET Top Tier	ET Middle Tier	ET Bottom Tier
SCM + Tier0_Flash (+ Tier1_Flash) + Enterprise + NL ( <b>unsupported</b> )	SCM	Tier0_Flash (+ Tier1_Flash)	Enterprise + NL
SCM + Tier0_Flash + Enterprise/NL	SCM	Tier0_Flash	Enterprise/NL
SCM + Tier0_Flash + Tier1_Flash + Enterprise/NL ( <b>unsupported</b> )	SCM	Tier0_Flash + Tier1_Flash	Enterprise/NL
SCM + Tier1_Flash (+ Enterprise/NL)	SCM	Tier1_Flash	(Enterprise/NL)
SCM + Tier1_Flash + Enterprise + NL	SCM	Tier1_Flash + Enterprise	NL
SCM + Enterprise/NL	SCM	Enterprise/NL	
SCM + Enterprise + NL	SCM	Enterprise	NL
Tier0_Flash (+ Tier1_Flash)	Tier0_Flash	(Tier1_Flash)	
Tier0_Flash + Tier1_Flash + Enterprise/NL	Tier0_Flash	Tier1_Flash	Enterprise/NL
Tier0_Flash + Tier1_Flash + Enterprise + NL	Tier0_Flash	Tier1_Flash + Enterprise	NL
Tier0_Flash + Enterprise (+ NL)	Tier0_Flash	Enterprise	(NL)
Tier0_Flash + NL	Tier0_Flash	NL	
Tier1_Flash (+ Enterprise/NL)		Tier1_Flash	(Enterprise/NL)
Tier1_Flash + Enterprise + NL	Tier1_Flash	Enterprise	NL
Enterprise (+ NL)		Enterprise	(NL)
NL			NL

The table represents all of the possible pool configurations. Some entries in the table contain *optional tiers* (in *italics*), the configurations without the optional tiers are also valid.

Sometimes a single Easy Tier tier contains MDiskS from more than one storage tier. For example, consider a pool with SCM, Tier1\_Flash, Enterprise, and NL. SCM is the top tier and Tier1\_Flash and Enterprise share the middle tier. NL is represented by the bottom tier.

**Note:** Some storage pool configurations with four or more different tiers are not supported. If such a configuration is detected, an error is logged and Easy Tier enters the measure mode, which means no extent migrations are performed.

For more information about planning and configuration considerations or best practices, see *IBM System Storage SAN Volume Controller, IBM Storwize V7000, and IBM FlashSystem 7200 Best Practices and Performance Guidelines, SG24-7521*.

### Easy Tier automatic data placement

Easy Tier continuously monitors volumes for host I/O activity. It collects performance statistics for each extent, and derives exponential moving averages for a rolling 24-hour period of I/O activity. Random and sequential I/O rate, I/O block size and bandwidth for reads and writes, and I/O response time are collected.

A set of algorithms is used to decide where the extents should be located and whether extent relocation is required. Once per day, Easy Tier analyzes the statistics to work out which data should be sent to a higher performing tier or might be sent to a lower tier. Four times per day, it analyzes the statistics to identify if any data needs to be rebalanced between MDisks in the same tier. Once every 5 minutes, Easy Tier checks the statistics to identify if any of the MDisks is overloaded.

Based on this information, Easy Tier generates a migration plan that must be run for optimal data placement. The system then spends as long as needed running the migration plan. The migration rate is limited to make sure host I/O performance is not affected while data is relocated.

The migration plan can consist of the following data movement actions on volume extents, as shown in Figure 9-2. Although each action is only shown once, all movement actions can be performed between any pair of adjacent tiers:

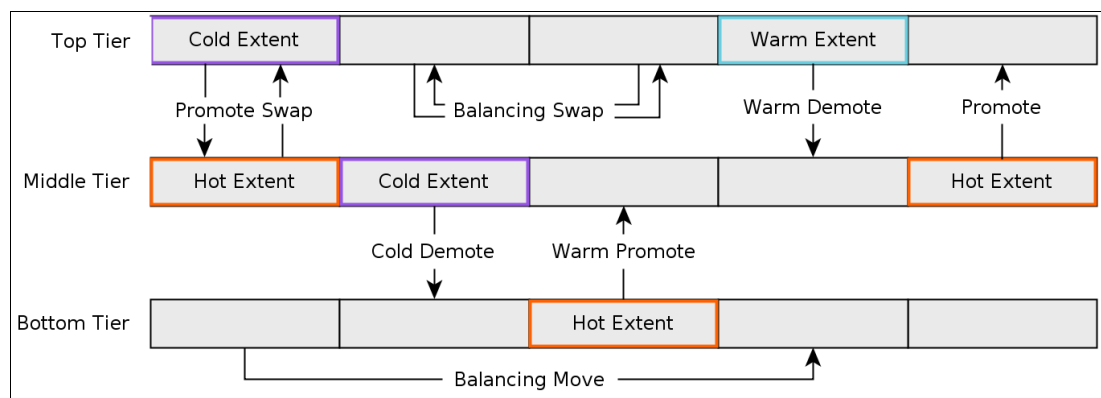


Figure 9-2 Actions on extents

- ▶ **Promote**  
Active data is moved from a lower tier of storage to a higher tier to improve the overall system performance.
- ▶ **Promote Swap**  
Active data is moved from a lower tier of storage to a higher tier to improve overall system performance. Less active data is moved first from the higher tier to the lower tier to make space.
- ▶ **Warm Promote**  
When an MDisk becomes overloaded, active data is moved from a lower tier to a higher tier to reduce the workload on the MDisk. This addresses the situation where a lower tier suddenly becomes active. Instead of waiting for the next migration plan, Easy Tier can react immediately.

Warm promote acts in a similar way to warm demote. If the 5-minute average performance shows that a layer is overloaded, Easy Tier immediately starts to promote extents until the condition is relieved.

► Cold Demote

Inactive or less active data is moved from a higher tier of storage to a lower tier to free up space on the higher tier. In that way, Easy Tier automatically frees extents on the higher storage tier before the extents on the lower tier become hot, which helps the system to be more responsive to new hot data.

► Warm Demote

When an MDisk becomes overloaded, active data is moved from a higher tier to a lower tier to reduce the workload on the MDisk. Easy Tier continuously ensures that the higher performance tier does not suffer from saturation or overload conditions that might affect the overall performance in the pool. This action is triggered when bandwidth or input/output operations per second (IOPS) exceeds a predefined threshold of an MDisk and causes the movement of selected extents from the higher-performance tier to the lower-performance tier to prevent MDisk overload.

► Balancing Move

Data is moved within the same tier from an MDisk with a higher workload to one with a lower workload to balance the workload within the tier. This automatically populates new MDisks that were added to the pool.

► Balancing Swap

Data is moved within the same tier from an MDisk with higher workload to one with a lower workload to balance the workload within the tier. Other less active data is moved first to make space.

Extent migration occurs at a maximum rate of 12 GB every 5 minutes for the entire system. It prioritizes actions as follows:

- Promote and rebalance get equal priority
- Demote is guaranteed 1 GB every 5 minutes, and then gets whatever is left

**Note:** Extent promotion or demotion only occurs between adjacent tiers. In a three-tier storage pool, Easy Tier does not move extents from the top directly to the bottom tier or vice versa without moving to the middle tier first.

The Easy Tier overload protection is designed to avoid overloading any type of MDisk with too much work. To achieve this, Easy Tier must have an indication of the maximum capability of a MDisk.

For an array made of locally attached drives, the system can calculate the performance of the MDisk because it is pre-programmed with performance characteristics for different drives and array configurations. For a storage area network (SAN)-attached MDisk, the system cannot calculate the performance capabilities. Therefore, it is important to follow the best-practice guidelines when configuring external storage, particularly the ratio between physical disks and MDisks presented to the system.

Each MDisk has an Easy Tier load parameter (low, medium, high, very\_high) that can be fine-tuned manually. If you analyze the statistics and find that the system does not appear to be sending enough IOPS to your external MDisk, you can increase the load parameter.



## Easy Tier operating modes

Easy Tier includes the following main operating modes:

- ▶ Off

When off, no statistics are recorded and no cross-tier extent migration occurs. Also, with Easy Tier turned off, no storage pool balancing across MDisks in the same tier is performed, even in single tier pools.

- ▶ Evaluation or measurement only

When in this mode, Easy Tier only collects usage statistics for each extent in a storage pool (if it is enabled on the volume and pool). No extents are moved. This collection is typically done for a single-tier pool that contains only HDDs so that the benefits of adding Flash drives to the pool can be evaluated before any major hardware acquisition.

- ▶ Automatic data placement and storage pool balancing

In this mode, usage statistics are collected and extent migration is performed between tiers (if there is more than one pool in a tier). Also, auto-balance between MDisks in each tier is performed.

**Note:** The auto-balance process automatically balances data when MDisks are added into a pool. However, it does not migrate extents from MDisks to achieve even extent distribution among all old and new MDisks in the storage pool. The Easy Tier migration plan is based on performance, not on the capacity of the underlying MDisks or on the number of extents on them.

## Implementation considerations

Consider the following implementation and operational rules when you use the IBM System Storage Easy Tier function on the SAN Volume Controller:

- ▶ Volumes that are added to storage pools use extents from the “middle” tier of three-tier model, if available. Easy Tier then collects usage statistics to determine which extents to move to “faster” or “slower” tiers. If there are no free extents in the middle tier, extents from the other tiers are used (bottom tier if possible, otherwise top tier).
- ▶ When an MDisk with allocated extents is deleted from a storage pool, extents in use are migrated to MDisks in the same tier as the MDisk that is being removed, if possible. If insufficient extents exist in that tier, extents from another tier are used.
- ▶ Easy Tier monitors extent I/O activity of each copy of a mirrored volume. Easy Tier works with each copy independently of the other copy. This applies to volume mirroring and HyperSwap and Remote Copy (RC).

**Note:** Volume mirroring can have different workload characteristics on each copy of the data because reads are normally directed to the primary copy and writes occur to both copies. Therefore, the number of extents that Easy Tier migrates between the tiers might differ for each copy.

- ▶ For compressed volumes in standard pools, only reads are analyzed by Easy Tier.
- ▶ Easy Tier automatic data placement is not supported on image mode or sequential volumes. However, it supports evaluation mode for such volumes. I/O monitoring is supported and statistics are accumulated.

- ▶ When a volume is migrated out of a storage pool that is managed with Easy Tier, Easy Tier automatic data placement mode is no longer active on that volume. Automatic data placement is also turned off while a volume is being migrated, even when it is between pools that both have Easy Tier automatic data placement enabled. Automatic data placement for the volume is reenabled when the migration is complete.

When the system migrates a volume from one storage pool to another, it attempts to migrate each extent to an extent in the new storage pool from the same tier as the original extent, if possible.

- ▶ When Easy Tier automatic data placement is enabled for a volume, you cannot use the `svctask migrateexts` CLI command on that volume.

## 9.1.2 Implementing and tuning Easy Tier

The Easy Tier function is enabled by default. It starts monitoring I/O activity immediately after storage pools and volumes are created. It also starts extent migration when the necessary I/O statistics are collected.

A few parameters can be adjusted. Also, Easy Tier can be disabled on selected volumes in storage pools.

### MDisk settings

The tier for internal (array) MDisk is detected automatically and depends on the type of drives that are its members. No adjustments are needed.

For an external MDisk, the tier is assigned when it is added to a storage pool. To assign the MDisk, select **Pools** → **External Storage**, select the MDisk (or MDisks) to add and click **Assign**.

**Note:** The tier of MDisks mapped from certain types of IBM System Storage Enterprise Flash is fixed to tier0\_flash and cannot be changed.

You can choose the target storage pool and storage tier that is assigned, as shown in Figure 9-3.

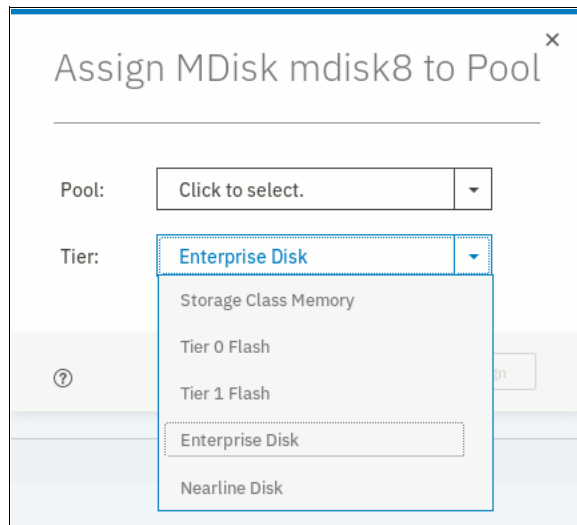


Figure 9-3 Choosing tier when assigning MDisk

To change the storage tier for an MDisk that is assigned, in **Pools** → **External Storage**, right-click one or more selected MDisks and choose **Modify Tier**, as shown in Figure 9-4.

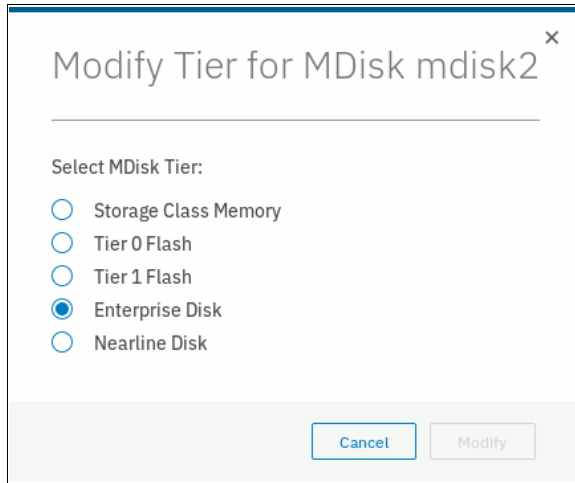


Figure 9-4 Changing MDisk tier

**Note:** Assigning a tier to an external MDisk that does not match the physical back-end storage type is not supported by IBM and can lead to unpredictable consequences.

To determine what tier is assigned to an MDisk, on **Pools** → **External Storage**, select **Actions** → **Customize columns** and select **Tier**. This selection includes the current tier setting into a list of MDisk parameters that are shown in the External Storage pane. You can also find this information in MDisk properties. To view it, right-click the MDisk, select **Properties**, and expand the **View more details** section, as shown in Figure 9-5.

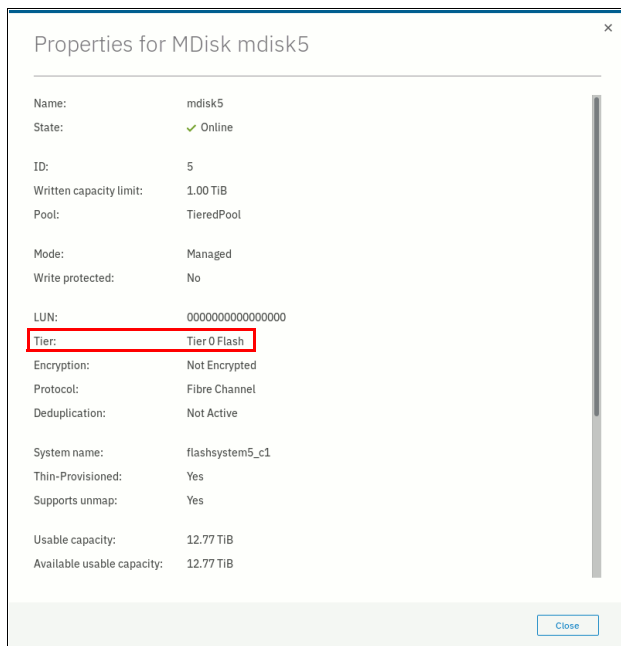


Figure 9-5 MDisk properties

To list MDisk parameters with the CLI, run the `lsmdisk` command. The current tier for each MDisk is shown. To change the external MDisk tier, run `chmdisk` command with the `-tier` parameter, as shown in Example 9-1.

*Example 9-1 Listing and changing tiers for MDisks (partially shown)*

```
IBM_2145:ITS0-SV1:superuser>lsmdisk
id name  status mode      mdisk_grp_id ... tier          encrypt
1 mdisk1 online unmanaged  ... tier0_flash no
2 mdisk2 online managed  0      ... tier_enterprise no
3 mdisk3 online managed  0      ... tier_enterprise no
<...>
IBM_2145:ITS0-SV1:superuser>chmdisk -tier tier1_flash mdisk2
IBM_2145:ITS0-SV1:superuser>
```

For an external MDisk, the system cannot calculate its exact performance capabilities, so it has several predefined levels. In rare cases, statistics analysis might show that Easy Tier is overusing or under-utilizing an MDisk. If so, levels can be adjusted only by using the CLI. Run the `chmdisk` command with `-easytierload` parameter. To reset Easy Tier load to system-default for the chosen MDisk, use `-easytier default`, as shown in Example 9-2.

**Note:** Adjust Easy Tier load settings only if instructed to do so by IBM Technical Support or your solution architect.

*Example 9-2 Changing Easy Tier load*

```
IBM_2145:ITS0-SV1:superuser>chmdisk -easytierload default mdisk2
IBM_2145:ITS0-SV1:superuser>
IBM_2145:ITS0-SV1:superuser>lsmdisk mdisk2 | grep tier
tier tier_enterprise
easy_tier_load high
IBM_2145:ITS0-SV1:superuser>
```

To list the current Easy Tier load setting of an MDisk, run `lsmdisk` with MDisk name or ID as a parameter.

## Storage pool settings

When a storage pool (either standard pool or Data Reduction Pool [DRP]) is created, Easy Tier is enabled by default. The system automatically enables Easy Tier functions when the storage pool contains an MDisk from more than one tier. It also enables automatic rebalancing when the storage pool contains an MDisk from only one tier.

You can disable Easy Tier or switch it to measure-only mode when creating a pool or any moment later. This is not possible with the GUI, only with the system CLI.

To check the current Easy Tier function state on a pool, select **Pools** → **Pools**, right-click the selected pool, choose **Properties**, and expand the **View more details** section, as shown in Figure 9-6 on page 459. This window also displays the amount of data stored on each tier.

Easy Tier can be in one of the following statuses:

- ▶ active

Indicates that a pool is being managed by Easy Tier, and extent migrations between tiers can be performed. Performance-based pool balancing of MDisks in the same tier is also enabled. This is the expected state for a pool with two or more tiers of storage.

- ▶ **balanced**  
Indicates that a pool is being managed by Easy Tier to provide performance-based pool balancing of MDisks in the same tier. This is the expected state for a pool with a single tier of storage.
- ▶ **inactive**  
Indicates that Easy Tier is inactive (disabled).
- ▶ **measured**  
Shows that Easy Tier statistics are being collected but no extent movement can be performed.

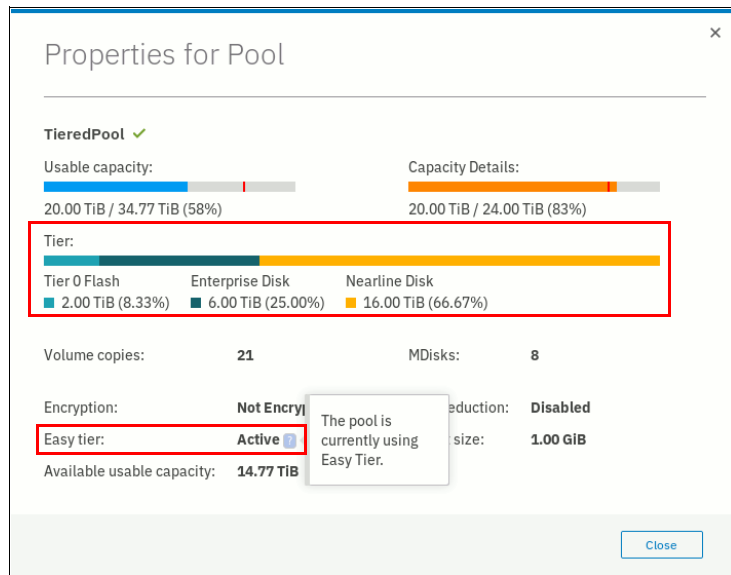


Figure 9-6 Pool properties

To find the status of the Easy Tier function on the pools with the CLI, run the `lsmdiskgrp` command without any parameters. To switch Easy Tier off or back on, run the `chmdiskgrp` command, as shown in Example 9-3. By running `lsmdiskgrp` with pool name/ID as a parameter, you can also determine how much storage of each tier is available within the pool.

Example 9-3 Listing and changing Easy Tier status on pools

```

IBM_2145:ITS0-SV1:superuser>lsmdiskgrp
id name      status mdisk_count ... easy_tier easy_tier_status
0 TieredPool online 1          ... auto      balanced
IBM_2145:ITS0-SV1:superuser>chmdiskgrp -easytier measure TieredPool
IBM_2145:ITS0-SV1:superuser>chmdiskgrp -easytier auto TieredPool
IBM_2145:ITS0-SV1:superuser>

```

## Volume settings

By default, each striped-type volume allows Easy Tier to manage its extents. If you need to fix the volume extent location (for example, to prevent extent demotes and to keep the volume in the higher-performing tier), you can turn off Easy Tier management for a particular volume copy.

**Note:** Thin-provisioned and compressed volumes in a DRP cannot have Easy Tier switched off. It is possible to switch off Easy Tier only at a pool level.

This process can be done by using the CLI only. Run the `lsvdisk` command to check and `chvdisk` command to modify Easy Tier function status on a volume copy, as shown in Example 9-4.

*Example 9-4 Checking and modifying Easy Tier settings on a volume*

---

```
IBM_2145:ITS0-SV1:superuser>lsvdisk vdisk0 |grep easy_tier
easy_tier on
easy_tier_status balanced
IBM_2145:ITS0-SV1:superuser>chvdisk -easytier off vdisk0
IBM_2145:ITS0-SV1:superuser>
```

---

## System-wide settings

There is a system-wide setting called *Easy Tier acceleration* that is disabled by default. Turning it on makes Easy Tier move extents up to four times faster than the default setting. In acceleration mode, Easy Tier can move up to 48 GiB per 5 minutes, whereas in normal mode it moves up to 12 GiB. The following use cases are the most probable use cases for acceleration:

- ▶ When adding capacity to the pool by adding to an existing tier or by adding a tier to the pool, accelerating Easy Tier can quickly spread volumes onto the new MDisks.
- ▶ Migrating the volumes between the storage pools when the target storage pool has more tiers than the source storage pool, so Easy Tier can quickly promote or demote extents in the target pool.

**Note:** Enabling Easy Tier acceleration is advised only during periods of low system activity only after migrations or storage reconfiguration occurred. It is recommended to keep Easy Tier acceleration mode off during normal system operation to avoid performance impacts that are caused by accelerated data migrations.

This setting can be changed non-disruptively, but only by using the CLI. To turn on or off Easy Tier acceleration mode, run the `chsystem` command. Run the `lssystem` command to check its current state, as shown in Example 9-5.

*Example 9-5 The chsystem command*

---

```
IBM_2145:ITS0-SV1:superuser>lssystem |grep easy_tier
easy_tier_acceleration off
IBM_2145:ITS0-SV1:superuser>chsystem -easytieracceleration on
IBM_2145:ITS0-SV1:superuser>
```

---

## 9.1.3 Monitoring Easy Tier activity

When Easy Tier is active, it constantly monitors and records I/O activity, collecting extent heat data. Heat data files are produced approximately once a day and summarize the activity per volume since the last heat data file was produced. Easy Tier activity can be monitored by using the GUI or the external IBM Storage Tier Advisor Tool (STAT) application.

## Monitoring Easy Tier by using the GUI

Select **Monitoring** → **Easy Tier Reports** to view the most recent Easy Tier statistics. Select the storage pool you want to see reports for in the filter section on the left, as shown in Figure 9-7. It takes approximately 24 hours for reports to be available after turning on Easy Tier or after a configuration node failover occurred. If no reports are available, the error message in the figure is shown. In this case, wait until new reports were generated and then revisit the GUI.

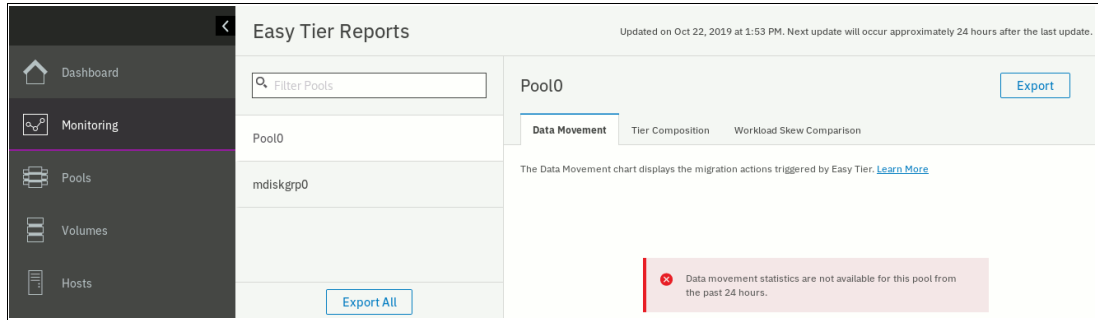


Figure 9-7 Easy Tier reports not available

Three types of reports are available per storage pool: Data Movement, Tier Composition, and Workload Skew Comparison. Select the corresponding tabs in the GUI to view the charts. Alternatively, use the **Export** or **Export All** buttons to download the reports in comma-separated value (CSV) format.

## Data Movement report

The Data Movement report shows the extent migrations Easy Tier performed to relocate data between different tiers of storage and within the same tier for optimal performance, as shown in Figure 9-8. The chart displays the data for the previous 24-hour period, in one-hour increments. You can change the time span and the increments for a more detailed view.

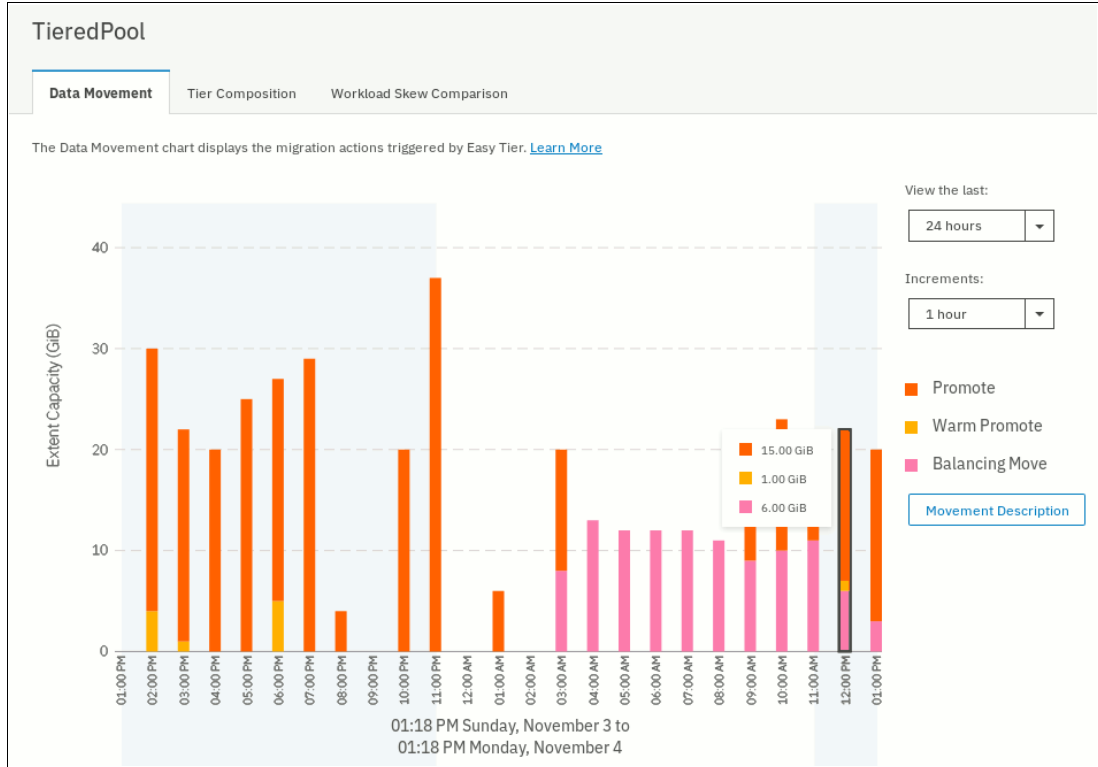


Figure 9-8 Easy Tier Data Movement chart

The X-axis shows a timeline for the selected period by using the selected increments. The Y-axis indicates the amount of extent capacity that was moved. For each time increment, a color-coded bar displays the amount of data moved by each Easy Tier data movement action, such as promote or cold demote. For more information about the different movement actions, see “Easy Tier automatic data placement” on page 453, or click the **Movement Description** button next to the chart to see an explanation in the GUI.



## Tier Composition chart

The Tier Composition chart shows how different types of workloads are distributed between top, middle, and bottom tiers of storage in the selected pool, as shown in Figure 9-9.

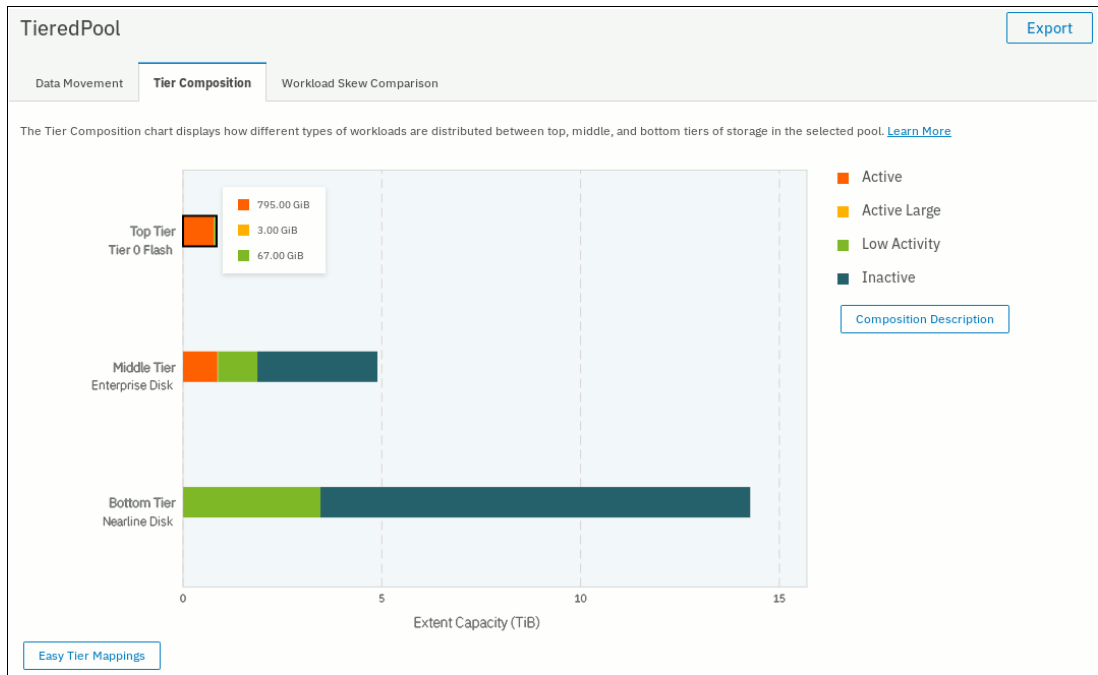


Figure 9-9 Tier Composition chart

A color-coded bar for each tier shows which workload types are present in that tier and to how much of the extent capacity in that tier they can be attributed to. Easy Tier distinguishes between the following workload types. Click the **Composition Description** button to show a short explanation for each workload type in the GUI:

- ▶ **Active**  
Data with more than 0.1 IOPS / Extent access density for small IOPS (< 64 KB block size).
- ▶ **Active Large**  
All data not classified above (> 64 KB block size)
- ▶ **Low Activity**  
Data with less than 0.1 IOPS / Extent access density
- ▶ **Inactive**  
Data with zero IOPS / Extent access density (no recent activity)

Click the **Easy Tier Mappings** button to show which MDisks were assigned to which of the three tiers of Easy Tier, as shown in Figure 9-10 on page 464. How storage tiers in the system are mapped to Easy Tier tiers depends on the available storage tiers in the pool. For a list of all possible mappings, see Table 9-1 on page 451.

Easy Tier Mappings for Pool TieredPool

Filter

ID	MDisk Name	Tier	Easy Tier Group
0	mdisk0	Enterprise Disk	Middle Tier
1	mdisk1	Enterprise Disk	Middle Tier
2	mdisk2	Enterprise Disk	Middle Tier
3	mdisk3	Enterprise Disk	Middle Tier
4	mdisk4	Nearline Disk	Bottom Tier
5	mdisk5	Tier 0 Flash	Top Tier
6	mdisk6	Nearline Disk	Bottom Tier
7	mdisk7	Tier 0 Flash	Top Tier

Showing 8 Easy Tier Mappings | Selecting 0 Easy Tier Mappings

Close

Figure 9-10 Easy Tier Mappings

### Workload Skew Comparison

The Workload Skew Comparison chart displays the percentage of I/O workload that is attributed to a percentage of the total capacity, as shown in Figure 9-11.

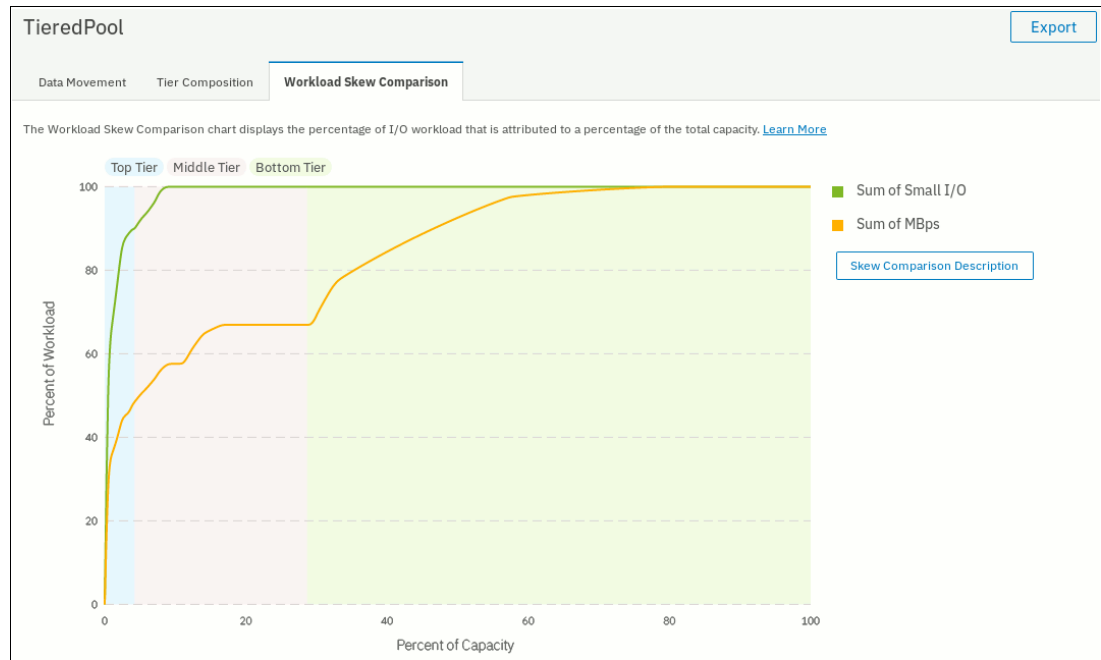


Figure 9-11 Easy Tier Workload Skew Comparison

The X-axis shows the percentage of capacity and the Y-axis shows the corresponding percentage of workload on that capacity. Workload is classified in small I/O (sum of small reads and writes) and MBps (sum of small and large bandwidth). The portion of capacity and workload attributed to a specific tier is color-coded in the chart with a legend above the chart.

Figure 9-11 on page 464 shows that the top tier (Tier1 Flash) contributes only a tiny percentage of capacity to the pool, but is handling around 85% of the IOPS and more than 40% of the bandwidth in that pool. The middle tier (enterprise disk) handles almost of the remaining IOPS and another 20% of the bandwidth. The bottom tier (NL disk) provides the most capacity to the pool, but does almost no small I/O workload.

Use this chart to estimate how much storage capacity in the high tiers must be available to handle most of the workload.

### Monitoring Easy Tier by using the IBM Storage Tier Advisor Tool

The IBM Storage Tier Advisor Tool (STAT) is a Windows console application that can analyze heat data files that are generated by Easy Tier. It produces a graphical display of the amount of “hot” data per volume and predictions of the performance benefits of adding capacity to a tier in a storage pool.

By using this method of monitoring, Easy Tier can provide more insights on top of the information that is available in the GUI.

IBM STAT can be downloaded from this IBM Support [web page](#).

You can download the STAT and install it on your Windows-based computer. The tool is packaged as an ISO file that must be extracted to a temporary location.

The tool installer is at `temporary_location\IMAGES\STAT\Disk1\InstData\NoVM\`. By default, the STAT is installed in `C:\Program Files\IBM\STAT\`.

On the system, the heat data files are found in the `/dumps/easytier` directory on the configuration node, and are named `dpa_heat.node_panel_name.time_stamp.data`. Any heat data file is erased when it exists for longer than seven days.

Heat files must be offloaded and IBM STAT started from a Windows command prompt console with the file specified as a parameter, as shown in Example 9-6.

*Example 9-6 Running STAT in Windows command prompt*

```
C:\Program Files (x86)\IBM\STAT>stat dpa_heat.78DXRY0.191021.075420.data
```

The STAT creates a set of `.html` and `.csv` files that can be used for Easy Tier analysis.

To download a heat data file, select **Settings** → **Support** → **Support Package** → **Download Support Package** → **Download Existing Package**, as shown in Figure 9-12.

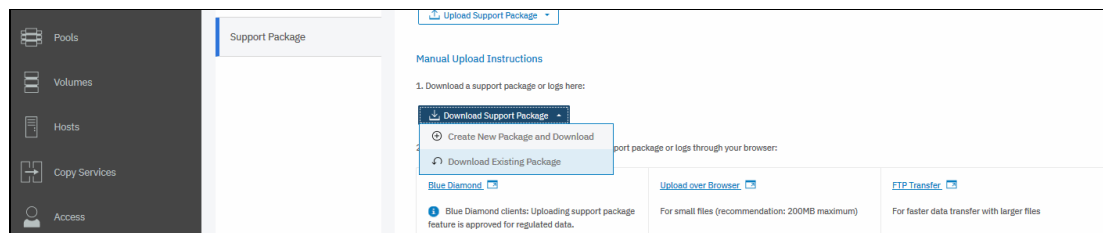


Figure 9-12 Downloading support package

The download window opens and shows all files in the `/dumps` directory and its subfolders on a current configuration node. You can filter the list by using the `easytier` keyword, select the `dpa_heat` file or files that are to be analyzed, and click **Download**, as shown in Figure 9-13 on page 466. Save the files in a convenient location (for example, to a subfolder that holds the IBM STAT executable file).

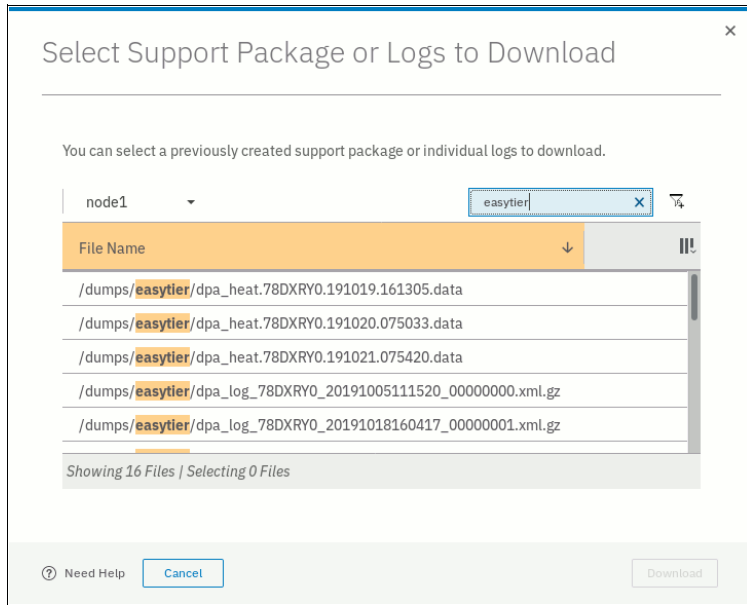


Figure 9-13 Downloading Easy Tier heat data file: *dpa\_heat* files

You can also specify the output directory. IBM STAT creates a set of HTML files, and the user can then open the `index.html` file in a browser to view the results. Also, the following `.csv` files are created and placed in the `Data_files` directory:

- ▶ `<panel_name>_data_movement.csv`
- ▶ `<panel_name>_skew_curve.csv`
- ▶ `<panel_name>_workload_ctg.csv`

These files can be used as input data for other utilities, such as the IBM STAT Charting Utility.

For more information about how to interpret IBM STAT tool output and analyzing CSV files, see *IBM System Storage SAN Volume Controller, IBM Storwize V7000, and IBM FlashSystem 7200 Best Practices and Performance Guidelines, SG24-7521*.

## 9.2 Thin-provisioned volumes

In a shared storage environment, thin provisioning is a method for optimizing the use of available storage. It relies on the allocation of capacity on demand as opposed to the traditional method of allocating all of the capacity at the time of initial provisioning. The use of this principle means that storage environments can achieve significantly higher utilization of physical storage resources by eliminating the unused allocated capacity.

Traditional storage allocation methods often provision large amounts of storage to individual hosts, but some of it remains unused (not written to). This might result in poor usage rates (often as low as 10%) of the underlying physical storage resources. Thin provisioning avoids this issue by presenting more storage capacity to the hosts than it uses from the storage pool. Physical storage resources can be expanded over time to respond to growth.

## 9.2.1 Concepts

IBM SAN Volume Controller supports thin-provisioned volumes in standard pools and in DRPs.

Each volume has a *provisioned capacity* and a *real capacity*. Provisioned capacity is the volume storage capacity that is available to a host. This capacity is detected by host operating systems and applications and can be used when creating a file system. Real capacity is the storage capacity that is reserved to a volume copy from a pool.

In a standard-provisioned volume, the provisioned capacity and real capacity are the same. However, in a thin-provisioned volume, the provisioned capacity can be much larger than the real capacity.

The provisioned capacity of a thin-provisioned volume is typically significantly larger than its real capacity. As more information is written by the host to the volume, more of the real capacity is used. The system identifies read operations to unwritten parts of the provisioned capacity and returns zeros to the server without using any real capacity.

The autoexpand feature prevents a thin-provisioned volume from using up its capacity and going offline. As a thin-provisioned volume uses capacity, the autoexpand feature maintains a fixed amount of unused real capacity, called the *contingency capacity*. For thin-provisioned volumes in standard pools, the autoexpand feature can be turned on and off. For thin-provisioned volumes in DRPs, the autoexpand feature is always enabled.

The capacity of a thin-provisioned volume is split into chunks called *grains*. Write I/O to grains that were not written to cause real capacity to be used to store data and metadata. The grain size of thin-provisioned volumes in standard pools can be 32 KB, 64 KB, 128 KB, or 256 KB.

Generally, smaller grain sizes save space but require more metadata access, which can adversely affect performance. When you use thin-provisioning with FlashCopy, specify the same grain size for the thin-provisioned volume and FlashCopy. The grain size of thin-provisioned volumes in DRPs cannot be changed from the default of 8 KB.

A thin-provisioned volume can be converted non-disruptively to a fully allocated volume, or vice versa, by using the volume mirroring function. For example, you can add a thin-provisioned copy to a fully allocated volume and then remove the fully allocated copy from the volume after they are synchronized.

The fully allocated to thin-provisioned migration procedure uses a zero-detection algorithm so that grains that contain all zeros do not cause any real capacity to be used. Usually, if the system is to detect zeros on the volume, you must use software on the host side to write zeros to all unused space on the disk or file system.

## 9.2.2 Implementation

For more information about creating thin-provisioned volumes, see Chapter 6, “Volumes” on page 255.

### Metadata

In a standard pool, the system uses the real capacity to store data that is written to the volume, and metadata that describes the thin-provisioned configuration of the volume. The metadata that is required for a thin-provisioned volume is less than 0.1% of the provisioned capacity.

Therefore, if the host uses 100% of the provisioned capacity, some extra space is required on your storage pool to store thin provisioning metadata. In a worst case, the real size of a thin provisioned volume can be 100.1% of its virtual capacity.

In a DRP, metadata for a thin-provisioned volume is stored separately from user data, and is not reflected in the volume real capacity. Capacity reporting is handled at the pool level.

## Volume parameters

When creating a thin-provisioned volume in a standard pool, some of its parameters can be modified in Custom mode, as shown in Figure 9-14.

Thin Provisioning	
Real capacity:	<input type="text" value="2"/> <input type="text" value="% of Provisioned capacity"/>
Automatically expand:	<input checked="" type="checkbox"/> Enabled
Warning threshold:	<input checked="" type="checkbox"/> Enabled
	<input type="text" value="80"/> % of Provisioned capacity
Thin-Provisioned Grain Size:	<input type="text" value="256"/> <input type="text" value="KiB"/>

Figure 9-14 Volume parameters for thin-provisioning

Real capacity defines initial volume real capacity and the amount of contingency capacity. When autoexpand is enabled, the system attempts to always maintain the contingency capacity by allocating more real capacity when hosts write to the volume.

The warning threshold can be used to be notified when the volume is about to run out of space.

In a DRP, fine-tuning of these parameters is not required. The real capacity and warning threshold are handled at the pool level. The grain size is always 8 KB and autoexpand is always on.

**Host considerations:** Do not use defragmentation applications on thin-provisioned volumes. The defragmentation process can write data to different areas of a volume, which can cause a thin-provisioned volume to grow up to its provisioned size.

## 9.3 Unmap

IBM Spectrum Virtualize systems running V8.1.0 and later support the SCSI Unmap command. The use of this command enables hosts to notify the storage controller of capacity that is no longer required, which can improve capacity savings and performance of Flash storage.

### 9.3.1 SCSI unmap command

Unmap is a set of Small Computer System Interface (SCSI) primitives that allow hosts to indicate to a storage system that space allocated to a range of blocks on a storage volume is no longer required. This command allows the storage system to take measures and optimize the system so that the space can be reused for other purposes.

When a host writes to a volume, storage is allocated from the storage pool. To free allocated space back to the pool, human intervention is needed on the storage system. The SCSI Unmap feature is used to allow host operating systems to unprovision storage on the storage system, which means that the resources can automatically be freed up in the storage pools and used for other purposes.

One of the most common use cases is a host application, such as VMware, freeing storage within a file system. The storage system can then reorganize the space, such as optimizing the data on the volume or the pool so that space can be reclaimed.

A SCSI unmappable volume is a volume that can have storage unprovision and space reclamation being triggered by the host operating system. The system can pass the SCSI **Unmap** command through to back-end Flash storage and external storage controllers that support the function.

### 9.3.2 Back-end SCSI unmap

The system can generate and send SCSI **Unmap** commands to specific back-end storage controllers and internal Flash storage.

This process occurs when volumes are formatted, deleted, extents are migrated, or an **unmap** command is received from the host. At the time of this writing, SCSI **unmap** commands are sent to only IBM A9000, IBM FlashSystem FS900 AE3, IBM FlashSystem FS9100, IBM Storwize and IBM FlashSystem family systems, and Pure storage systems.

This helps prevent an overprovisioned storage controller from running out of free capacity for write I/O requests. This means when you use supported overprovisioned back-end storage, back-end SCSI **unmap** should normally be left enabled.

Flash storage typically requires empty blocks to serve write I/O requests. This means **unmap** can improve Flash performance by erasing blocks in advance.

This feature is turned on by default. It is recommended to keep back-end **unmap** enabled, especially if a system is virtualizing an overprovisioned storage controller.

To verify that sending **unmap** commands to back-end is enabled, run the CLI command **lssystem**, as shown in Example 9-7.

*Example 9-7 Verifying back-end unmap support status*

---

```
IBM_2145:ITS0-SV1:superuser>lssystem | grep backend_unmap  
backend_unmap on
```

---

### 9.3.3 Host SCSI unmap

The IBM Spectrum Virtualize system can advertise support for SCSI Unmap to hosts. Hosts can then use the set of SCSI **unmap** commands to indicate that formerly used capacity is no longer required on a volume.

When these volumes are in DRPs, that capacity becomes reclaimable capacity and is monitored and collected and eventually redistributed back to the pool for use by the system. Volumes in standard pools do not support automatic space reclamation after data is unmapped and SCSI **unmap** commands are handled as though they were writes with zero data.

The system also sends SCSI **unmap** commands to back-end controllers that support them if host unmaps for corresponding blocks are received (and backend unmap is enabled).

With host SCSI **unmap** enabled, some host types (for example, Windows, Linux, or VMware) change their behavior when creating a file system on a volume, by running SCSI **unmap** commands to the entire capacity of the volume.

The format completes only after all of these **unmap** commands complete. Some host types run a background process (for example, `fstrim` on Linux), which periodically issues SCSI **unmap** commands for regions of a file system that are no longer required. Hosts might also send **unmap** commands when files are deleted in a file system.

Host SCSI **unmap** commands drive more I/O workload to back-end storage. In some circumstances (for example, volumes that are on a heavily loaded NL serial-attached SCSI [SAS] array), this issue can cause an increase in response times on volumes that use the same storage. Also, host formatting time is likely to increase, compared to a system that does not advertise support the SCSI **unmap** command.

If you use DRPs or an overprovisioned back-end that supports **unmap**, it is recommended to turn on SCSI **unmap** support. Host **unmap** support is disabled by default.

If only standard pools are configured and back-end is traditional (fully provisioned), consider keeping host **unmap** support switched off because it does not provide any benefit.

To check and modify current setting for host SCSI **unmap** support, run the `lssystem` and `chsystem` CLI commands, as shown in Example 9-8.

*Example 9-8 Turning host unmap support on*

---

```
IBM_2145:ITS0-SV1:superuser>lssystem | grep host_unmap
host_unmap off
IBM_2145:ITS0-SV1:superuser>chsystem -hostunmap on
IBM_2145:ITS0-SV1:superuser>
```

---

**Note:** You can switch host **unmap** support on and off non-disruptively on the system side. However, hosts must rediscover storage, or (in the worst case) be restarted for them to stop sending **unmap** commands.

### 9.3.4 Offload IO throttle

Throttles are a mechanism to control the amount of resources that are used when the system is processing I/Os on supported objects. If a throttle limit is defined, the system processes the I/O for that object or delays the processing of the I/O to free resources for more critical I/O operations.

Offload commands, such as **unmap** and **XCOPY** free up hosts and speed up the copy process by offloading the operations of specific types of hosts to a storage system. These commands are used by hosts to format new file systems, or copy volumes without the host needing to read and then write data.

Offload commands can sometimes create I/O intensive workloads, potentially taking bandwidth from production volumes and affecting performance, especially if the underlying storage cannot handle the amount of I/O that is generated.



Throttles can be used to delay processing for offloads to free bandwidth for other more critical operations. This can improve performance, but limits the rate at which host features, such as VMware VMotion, can copy data. It can also increase the time it takes to format file systems on a host.

**Note:** For systems that are managing any NL storage, it might be recommended to set the offload throttle to 100 MBps.

To implement offload throttle, you can run the `mkthrottle` command with `-type offload` parameter. In the GUI, select **Monitoring** → **System** and click **System Actions** → **Edit System Offload Throttle**, as shown in Figure 9-15.

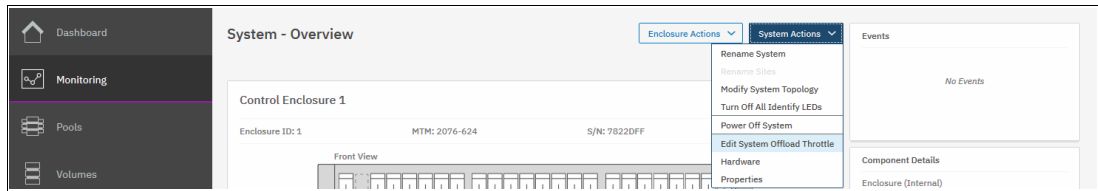


Figure 9-15 Setting offload throttle

## 9.4 Data Reduction Pools

DRPs provide a set of techniques that can be used to reduce the amount of usable capacity that is required to store data. This helps increase storage efficiency and reduce storage costs. Available techniques include thin-provisioning, compression, and deduplication.

DRPs automatically reclaim used capacity that is no longer needed by host systems and return it back to the pool as available capacity for future reuse.

**Note:** This book provides only an overview of DRP aspects. For more information, see *Introduction and Implementation of Data Reduction Pools and Deduplication*, SG24-8430.

### 9.4.1 Introduction to DRP

The system can use different data reduction methods simultaneously, which increases the capacity savings across the entire storage pool.

DRPs support five types of volumes:

- ▶ Fully allocated  
This type provides no data reduction.
- ▶ Thin-provisioned  
This type provides data reduction by allocating storage on demand when writing to the volume.
- ▶ Thin and Compressed  
In addition to being thin-provisioned, data is compressed before being written to storage.
- ▶ Thin and Deduplicated  
In addition being thin-provisioned, duplicates of data blocks are detected and are replaced with references to the first copy.

- ▶ Thin and Compressed and Deduplicated

This type achieves maximum data reduction by combining thin-provisioning, compression, and deduplication.

Volumes in a DRP track when capacity is freed from hosts and possible unused capacity that can be collected and reused within the storage pool. When a host no longer needs the data that is stored on a volume, the host system uses SCSI **unmap** commands to release that capacity from the volume. When these volumes are in DRPs, that capacity becomes reclaimable capacity and is monitored and collected and eventually redistributed back to the pool for use by the system.

**Note:** If the usable capacity usage of a DRP exceeds more than 85%, I/O performance can be affected. The system needs 15% of usable capacity available in DRPs to ensure that capacity reclamation can be performed efficiently.

At its core, a DRP uses a log structured array (LSA) to allocate capacity. An LSA allows a tree-like directory to be used to define the physical placement of data blocks independent of size and logical location.

Each volume has a range of logical block addresses (LBAs), starting from 0 and ending with the block address that fills the capacity. The LSA allows the system to allocate data sequentially when written to volumes (in any order) and provides a directory that provides a lookup to match volume LBA with physical address within the array. A volume in a DRP contains directory metadata to store the mapping from logical address on the volume to physical location on the backend storage.

This directory often is too large to store in memory; therefore, it must be read from storage as required. The lookup and maintenance of this metadata results in I/O amplification. I/O amplification occurs when a single host-generated read or write I/O results in more than one back-end storage I/O request. For example, a read I/O request might have to read some directory metadata in addition to the data. A write I/O request might have to read directory metadata write updated directory metadata, journal metadata, and the data.

Conversely, data reduction by design reduces the size of data by using compression and deduplication, so less data is written to the backend storage.

## 9.4.2 DRP benefits

DRPs are a new type of storage pool that implement techniques, such as thin-provisioning, compression, and deduplication to reduce the amount of physical capacity that is required to store data. Savings in storage capacity requirements translate into reduction in the cost of storing the data.

The cost reductions that are achieved through software can facilitate the transition to all Flash storage. Flash storage has lower operating costs, lower power consumption, higher density, and is cheaper to cool than disk storage. However, the cost of Flash storage is still higher. Data reduction can reduce the total cost of ownership of an All-Flash system to be competitive with HDD.

One benefit of DRP is in the form of capacity savings that are achieved by deduplication and compression. Real-time deduplication identifies duplicate data blocks during write I/O operations and stores a reference to the first copy of the data instead of writing the data to the storage pool a second time. It does this by maintaining a fingerprint database that contains hashes of data blocks that are written to the pool. If new data written by hosts matches an entry in this database a reference is generated in the directory metadata instead of writing the new data.

Compression reduces the size of the host data that is written to the storage pool. DRP uses the same Lempel-Ziv (LZ) based IBM Real-time Compression (RtC) and decompression algorithm that is used by Random Access Compression Engine (RACE) in standard pools. However, DRP offers a new implementation of data compression that is fully integrated into the IBM Spectrum Virtualize I/O stack. It makes optimal use of node resources, such as memory and CPU cores and uses hardware acceleration on supported platforms more efficiently. DRP compression also operates on smaller block sizes, which results in more consistent and predictable performance.

Deduplication and compression can be combined, in which case data is first deduplicated and then compressed. Therefore, deduplication references are created on the compressed data stored on the physical domain.

DRPs support end-to-end SCSI **unmap** functions. Hosts use the set of SCSI **unmap** commands to indicate the formerly used capacity is no longer required on a target volume. Reclaimable capacity is unused capacity that is created when data is overwritten, volumes are deleted, or when data is marked as unneeded by a host by using the SCSI **unmap** command. That capacity can be collected and reused on the system.

DRPs, the directory, and the reduction techniques were designed around optimizing for Flash and future solid-state storage technologies. All metadata operations are 4 KB, which is ideal for Flash storage to maintain low and consistent latency. All data read operations are 8 KB (before reduction) and again, designed to minimize latency because Flash is well-suited for small block workload with high IOPS. All write operations are coalesced into 256 KB sequential writes to simplify the garbage collection on Flash devices and to gain full stride writes from RAID arrays.

DRP is designed to work well alongside Easy Tier. The directory metadata of DRPs does not fit in memory and is therefore stored on disk that uses dedicated metadata volumes, which are separate from the data. The metadata volumes are small but frequently accessed by small block I/O requests. Performance gains are expected because they are optimal candidates for promotion to the fastest tier of storage through Easy Tier. In contrast, data volumes are large but frequently rewritten data is grouped to consolidate “heat”. Easy Tier can work as usual to accurately identify active data.

### 9.4.3 Planning for DRPs

Before configuring and using DRPs in production environments, it is important to plan for capacity and performance. Because DRPs have different performance characteristics than standard pools, sizing models cannot be used directly without modifications.

For more information about how to estimate capacity savings that are achieved by compression and deduplication, see 9.6, “Saving estimation for compression and deduplication” on page 484.

The following software and hardware requirements must be met for DRP compression and deduplication:

- ▶ Enabled Compression license
- ▶ SV1 nodes include hardware compression accelerator cards to use compression
- ▶ The system runs V8.1.3.2 or higher

RACE compression in standard pools and DRP compressed volumes can coexist in the same I/O group. However, deduplication is not supported in the same I/O group as RACE compressed volumes.

In most cases, it is recommended to enable compression for all thin-provisioned and deduplicated volumes. Overhead in DRPs is caused by metadata handling, which is the same for compressed volumes and thin-provisioned volumes without compression. An exception are systems without support for hardware compression.

In systems with compressing backend storage controllers, certain system configurations make determining accurate physical capacity usage difficult. If the system contains compressing backend storage controllers and DRPs with thin-provisioned volumes without compression, it is difficult to monitor capacity usage.

In this case, overcommitting and losing access to write operations is possible. To prevent this issue from occurring, use compressed volumes (with or without deduplication) or fully allocated volumes. Separate compressed volumes and fully allocated volumes by using separate pools. For more information, see Chapter 9.7, “Overprovisioning and data reduction on external storage” on page 486.

A maximum of four DRPs can be in a system. When this limit is reached, only standard pools can be created.

A DRP uses a customer data volume per I/O group to store volume data. The maximum size of a customer data volume is limited to 128,000 extents per I/O group. This places a limit on the maximum physical capacity in a pool after data reduction that depends on the extent size, number of DRPs, and number of I/O groups, as listed in Table 9-2. DRPs have a minimum extent size of 1024 MB.

*Table 9-2 Maximum physical capacity after data reduction*

<b>Extent Size</b>	<b>1 DRP - 1 I/O group</b>	<b>1 DRP - 4 I/O groups</b>	<b>4 DRP - 4 I/O groups</b>
1024 MB	128 TiB	512 TiB	2 PiB
2048 MB	256 TiB	1 PiB	4 PiB
4096 MB	512 TiB	2 PiB	8 PiB
8192 MB	1 PiB	4 PiB	16 PiB

Overwriting data, unmapping data, and deleting volumes causes reclaimable capacity in the pool to increase. Garbage collection is performed in the background to convert reclaimable capacity to available capacity. This operation requires free capacity in the pool to operate efficiently and without affecting I/O performance. A general guideline is to ensure that the provisioned capacity with the DRP does not exceed 85% of the total usable capacity of the DRP.

To ensure garbage collection is working properly, a DRP has a minimum capacity limit, depending on extent size and number of IO groups, as listed in Table 9-3. Even when there are no volumes in the pool, some of the space is used to store metadata. The required metadata capacity depends on the total capacity of the storage pool and on the extent size. This issue should be considered when planning capacity.

Table 9-3 Minimum capacity in a single Data Reduction Pool

Extent size	1 I/O group	4 I/O groups
1024 MB	255 GiB	1 TiB
2048 MB	0.5 TiB	2 TiB
4096 MB	1 TiB	4 TiB
8192 MB	2 TiB	8 TiB

**Note:** The default extent size in a DRP is 4 GB. If the estimated total capacity in the pool exceeds the documented limits, choose a larger extent size. If the estimated total capacity is relatively small, consider the use of a smaller extent size for a smaller metadata overhead and a lower minimum capacity limit.

For more information about the use of data reduction on the system and the backend storage, see 9.7, “Overprovisioning and data reduction on external storage” on page 486.

## 9.4.4 Implementing DRP with compression and deduplication

To use all data reduction technologies on the system, you must create a DRP and volumes within the DRP, and then map these volumes to hosts that support SCSI **unmap** commands. Also consider enabling support for host SCSI **unmap** commands as described in 9.3.3, “Host SCSI unmap” on page 469. The implementation process for DRP is similar to standard pools, but has its own specifics.

### Creating pools and volumes

To create a DRP, enable **Data Reduction** in the Create Pool window, which is available by selecting **Pools** → **Pools**. For more information about how to create a storage pool and populate it with MDisks, see Chapter 5, “Storage pools” on page 199.

To create a volume within a DRP, select **Volumes** → **Volumes**, and click **Create Volumes**.

Figure 9-16 shows the Create Volumes window. In the Capacity Savings menu, the following selections are available: None, Thin-Provisioned, and Compressed. If Compressed or Thin-Provisioned is selected, the Deduplicated option also becomes available and can be selected.

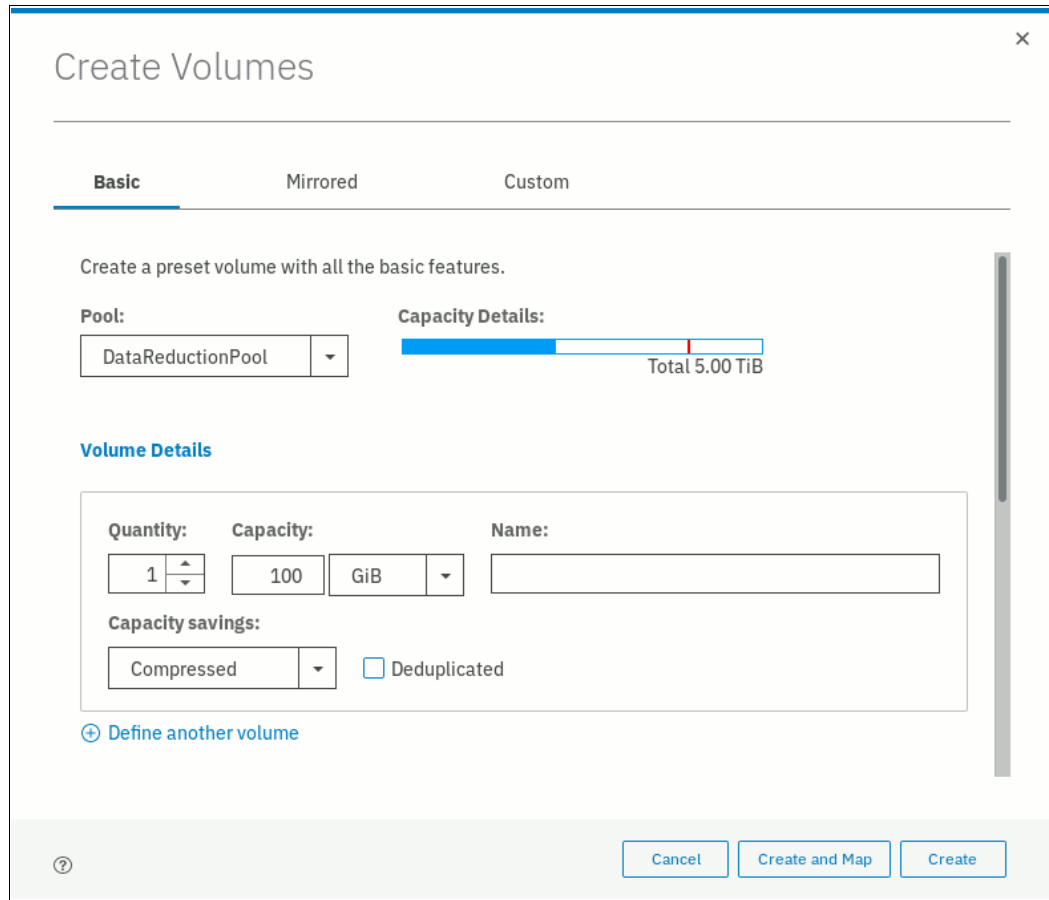


Figure 9-16 Create compressed volume

## Capacity monitoring

Capacity monitoring in DRPs is done mostly on the system level and on the storage pool level. Use the **Dashboard** in the GUI to view a summary of the capacity usage and capacity savings of the entire system.

The Pools page in the management GUI is used for reporting on the storage pool level and displays the Usable Capacity and Capacity Details. Usable capacity indicates the amount of capacity that is available for storing data on a pool after formatting and RAID techniques are applied. Capacity details is the capacity that is available for volumes before any capacity savings methods are applied. Select **Pools** → **Pools** to monitor this capacity, as shown in Figure 9-17.

Name	State	Usable Capacity	Capacity Details	Data Reduction
DataReductionPool	✓ Online	1.15 TiB / 5.00 TiB (23%)	2.13 TiB / 5.00 TiB (43%)	Yes
TieredPool	✓ Online	20.00 TiB / 24.00 TiB (83%)	20.00 TiB / 24.00 TiB (83%)	No

Figure 9-17 Data Reduction Pool capacity overview

To see more detailed capacity reporting including the warning threshold and capacity savings, open the pool properties dialog by right-clicking a pool and selecting **Properties**. This dialog shows the savings achieved by thin-provisioning, compression, deduplication, and the total data reduction savings in the pool, as shown in Figure 9-18. The Reclaimable capacity also is shown. This unused capacity is created when data is overwritten, volumes are deleted, or when data is marked as unneeded by a host by using the SCSI **unmap** command. This capacity is converted to available capacity by the garbage collection background process.

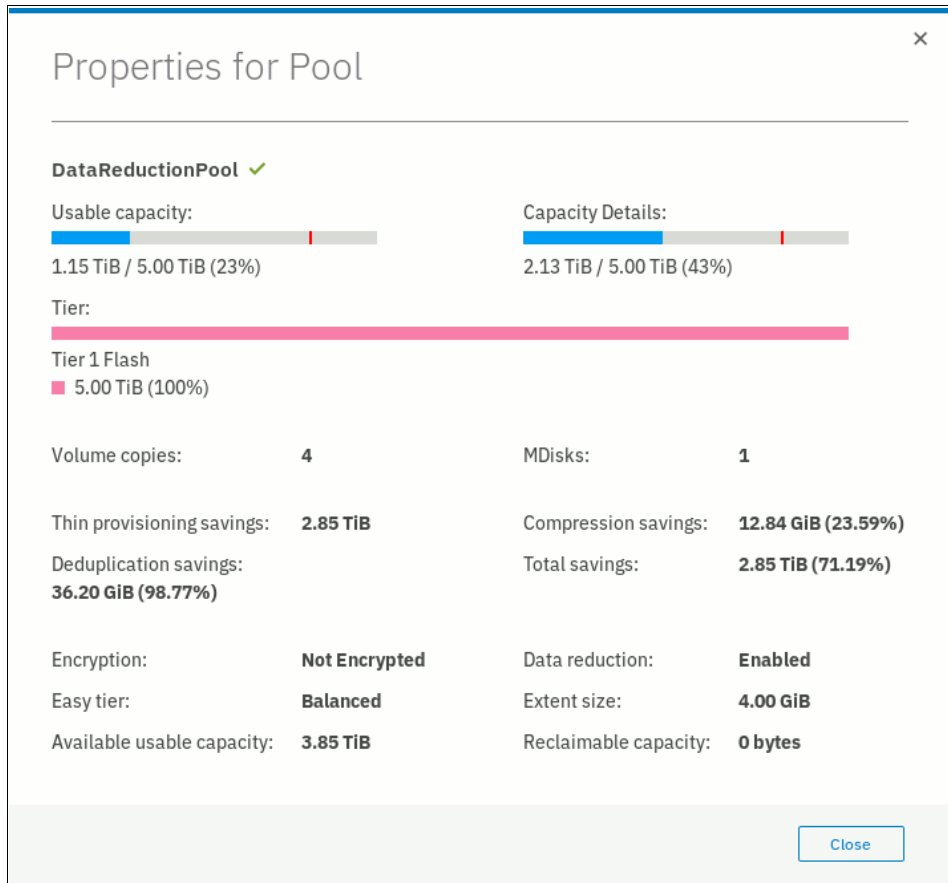


Figure 9-18 Capacity reporting in a Data Reduction Pool

Thin-provisioned, compressed, and deduplicated volumes do not provide a detailed per-volume capacity reporting, as shown in the Volumes by Pool pane in Figure 9-19. Only the Capacity (which is the provisioned capacity that is available to hosts) is shown. Real capacity, Used capacity, and Compression savings are not applicable for volumes with capacity savings. Only fully allocated volumes display those parameters.

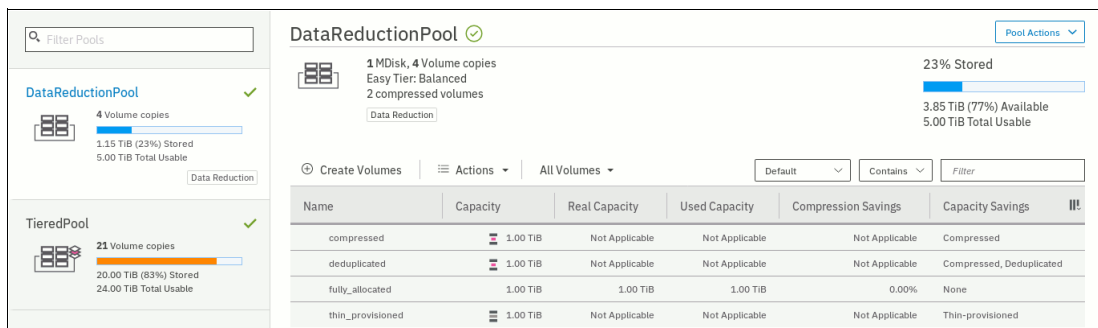


Figure 9-19 Volumes in DRP pool

Per-volume compression savings are not visible directly, but they can be accurately estimated by using the IBM Comprestimator. For more information, see 9.6.1, “Evaluate compression savings by using IBM Comprestimator” on page 484. This tool can be used for compressed volumes to analyze the volume level compression savings.

The CLI can be used for limited capacity reporting on the volume level. The `used_capacity_before_reduction` indicates the total amount of data that is written to a thin-provisioned or compressed volume copy in a data reduction storage pool before data reduction occurs. This field is empty for fully allocated volume copies and volume copies not in a DRP.

To find this value, run the `lsvdisk` command with volume name or ID as a parameter, as shown in Example 9-9. It shows a thin-provisioned volume without compression and deduplication with a virtual size of 1 TiB provisioned to the host. A 53 GB file was written from the host.

*Example 9-9 Data Reduction Pool volume capacity reporting on the CLI*

---

```
IBM_2145:ITS0-SV1:superuser>lsvdisk thin_provisioned
id 34
name vdisk1

capacity 1.00TB

used_capacity
real_capacity
free_capacity

tier tier_scm
tier_capacity 0.00MB
tier tier0_flash
tier_capacity 0.00MB
tier tier1_flash
tier_capacity 0.00MB
tier tier_enterprise
tier_capacity 0.00MB
tier tier_nearline
tier_capacity 0.00MB

compressed_copy no
uncompressed_used_capacity
deduplicated_copy no
used_capacity_before_reduction 53.04GB
```

---

The used, real, free capacity, and the capacity that is stored on each storage tier is not shown for volumes (except fully allocated volumes) in DRPs.

Capacity reporting on the pool level is available by running the CLI command `lsmdiskgrp` with the pool ID or name as a parameter, as shown in Example 9-10.

*Example 9-10 DRP capacity reporting on the CLI*

---

```
IBM_2145:ITS0-SV1:superuser>lsmdiskgrp 1 | grep -E "capacity|compression|tier
tier"
capacity 5.00TB
free_capacity 2.87TB
virtual_capacity 4.00TB
```



```
used_capacity 1.14TB
real_capacity 1.14TB
tier tier_scm
tier_capacity 0.00MB
tier_free_capacity 0.00MB
tier tier0_flash
tier_capacity 0.00MB
tier_free_capacity 0.00MB
tier tier1_flash
tier_capacity 5.00TB
tier_free_capacity 3.85TB
tier tier_enterprise
tier_capacity 0.00MB
tier_free_capacity 0.00MB
tier tier_nearline
tier_capacity 0.00MB
tier_free_capacity 0.00MB
compression_active no
compression_virtual_capacity 0.00MB
compression_compressed_capacity 0.00MB
compression_uncompressed_capacity 0.00MB
child_mdisk_grp_capacity 0.00MB
used_capacity_before_reduction 143.68GB
used_capacity_after_reduction 94.64GB
overhead_capacity 52.00GB
deduplication_capacity_saving 36.20GB
reclaimable_capacity 0.00MB
physical_capacity 5.00TB
physical_free_capacity 3.85TB
```

---

**Note:** Compression-related properties are not valid for DRPs. They are used exclusively for RACE compression in standard pools.

For more information about every reported value, see [IBM Knowledge Center](#).

## Migrating to and from DRP

Data migration from or to a DRP is done by using volume mirroring. A second copy in the target pool is added to the source volume and the original copy is optionally removed after the synchronization process completes. If the volume has two copies, one of the copies must be removed or a more complex migration scheme that uses FlashCopy, RC, host mirroring, or similar must be used.

Depending on the system configuration and the type of migration, a one-step migration or a two-step migration is necessary. The reason is that compressed volumes in standard pools cannot coexist with deduplicated volumes in DRPs. Therefore, a two-step migration is required in the following scenarios:

- ▶ Migrating to a deduplicated volume in a DRP if compressed volumes in standard pools exist in the same I/O group
- ▶ Migrating to a compressed volume in a standard pool if deduplicated volumes exist in the same I/O group

**Note:** The two-step migration requires all volumes that cannot coexist with the wanted volume type to be migrated in one step.

The migration processes work as follows:

- One-step migration process

To create a second copy, right-click the source volume and choose **Add Volume Copy**, as shown in Figure 9-20. Choose the target pool of the migration for the second copy and select the wanted capacity savings.

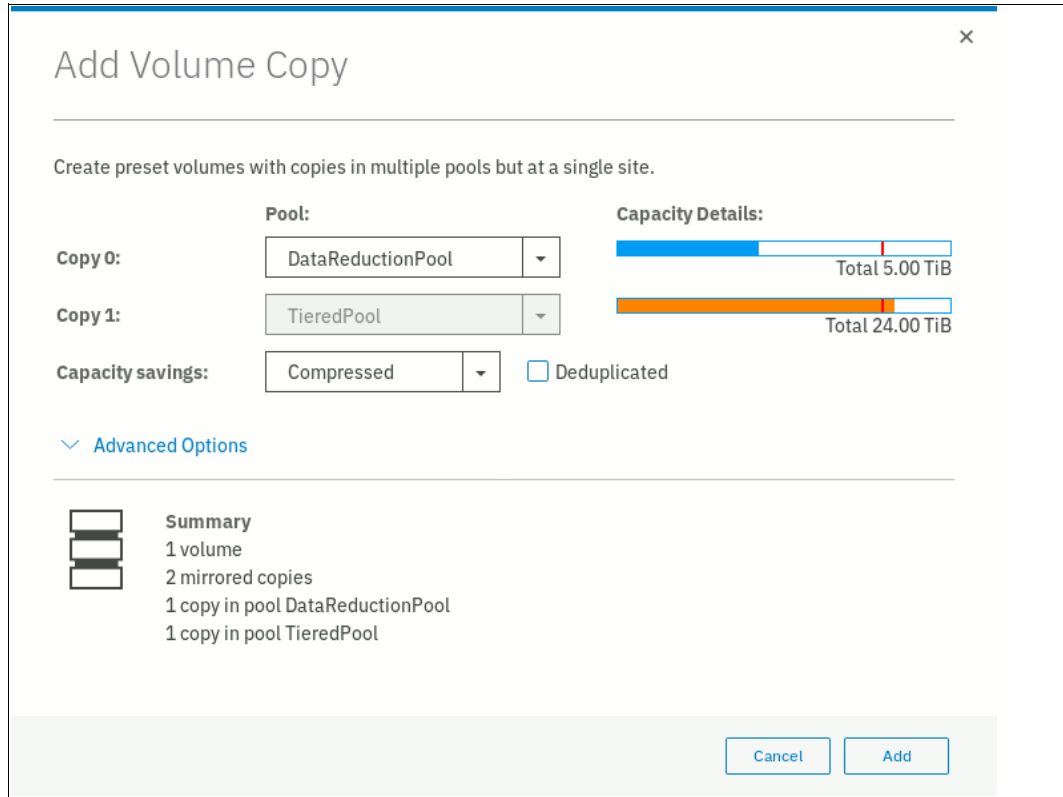


Figure 9-20 Add volume copy dialog

After you click **Add**, synchronization starts. The time synchronization takes to complete depends on the size of the volume, system performance, and the configured migration rate. You can increase the synchronization rate by right-clicking the volume and selecting **Modify Mirror Sync Rate**.

When both copies are synchronized, Yes is displayed for both copies in the Synchronized column in the Volumes pane. You can track the synchronization process by using the Running tasks pane, as shown in Figure 9-21. After it reaches 100% and the copies are in-sync, you can complete migration by deleting the source copy.

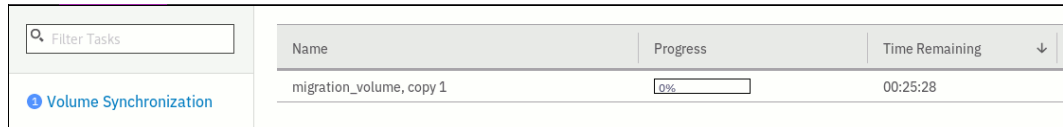


Figure 9-21 Synchronization progress

► Two-step migration process

Some volume types in DRPs cannot coexist in the same I/O group as compressed volumes in standard pools. The following two-step migration process is used to work around these limitations:

- a. Add a second copy to each source volume in the target pool as described in the one-step migration process but select capacity savings options that can coexist with the volumes in the current system configuration. For example, do not enable deduplication for the new volume copy if compressed volumes in standard pools exist in the same I/O group.

Wait for synchronization to complete and then delete the source copy.

Complete this process for all volumes that must be migrated and then verify that no volumes that cannot coexist with the wanted volume type are left in the same I/O group. For example, verify that no more compressed volumes in standard pools are left in the same I/O group when migrating to deduplicated volumes in a DRP.

- b. For all of the volumes handled in the previous step, add a second copy in the target pool, but select the wanted capacity savings options. Wait for synchronization to complete and then delete the source copy to complete migration.

Alternatively, right-click one of the volumes, click **Modify Capacity Savings**, and select the wanted options. A second copy is created and the source copy is deleted automatically after synchronization completes.

## Garbage collection and volume deletion

DRP includes built-in capabilities to enable garbage collection of unused storage capacity. Garbage Collection is a DRP process that reduces the amount of data that is stored on external storage systems and internal drives by reclaiming previously used storage resources that are no longer needed by host systems.

When a DRP is created, the system monitors the pool for reclaimable capacity from host unmap operations. When space is freed from a host operating system, it is a process called *unmapping*. Hosts indicate that the allocated capacity is no longer required on a target volume. The freed space is collected and reused by the system automatically, without having to reallocate the capacity manually.

Removing thin-provisioned or compressed volume copies in a DRP is an asynchronous operation. Volume copies that are removed from the system enter the *deleting* state, during which the used capacity of the copies is converted to reclaimable capacity in the pool by using a background deletion process.

The removal process of deduplicated volume copies also searches and moves deduplication references that other volumes might have to the deleting volume copies. This process is done to ensure that deduplicated data that is on deleted copies continues to be available for other volumes in the system.

After this process is complete, the volume copies are deleted and they disappear from the system configuration. In a second step, garbage collection can give the reclaimable capacity that is generated in the first step back to the pool as available capacity. This means that used capacity of a removed volume is not available for reuse immediately after the removal.

The time it takes to delete a thin-provisioned or compressed volume copy depends on the size of the volume, system configuration, and workload. For deduplicated copies, the duration also depends on the amount and size of other deduplicated copies in the pool.

Therefore, it might take a much time to delete a small deduplicated copy if many other deduplicated volumes are in the same pool. The deletion process is a background process that does not affect system I/O performance.

The deleting state of a volume or volume copy can be seen by running the `lsvdisk` CLI command. The GUI hides volumes in this state, but it shows deleting volume copies if the volume contains another copy.

**Note:** Removing thin-provisioned or compressed volume copies in a DRP might take much time to complete. Used capacity is not immediately returned to the pool as available capacity.

When one copy of a mirrored volume is in the deleting state, it is not possible to add a new copy to the volume before the deletion finishes. If a new copy must be added without waiting for the deletion to complete, first split the copy that must be deleted into a new volume and then delete the new volume and add a new second copy to the original volume. To split a copy into a new volume, right-click the copy and select **Split into New Volume** in the GUI or run the `splitvdiskcopy` command on the CLI.

## 9.5 Compression with standard pools

RACE technology was first introduced in the RtC Appliances. It is integrated into the system software stack as the RtC solution.

RACE or RtC is used for compression of the volume copies in standard pools. DRPs use a new compression implementation that is described in 9.4.2, “DRP benefits” on page 472.

Consider the following points:

- ▶ SA2 and SV2 nodes support only compression in DRPs
- ▶ SV1 nodes need hardware compression accelerator cards to use compression

For more information about RtC compression, see *IBM Real-time Compression in IBM SAN Volume Controller and IBM Storwize V7000*, REDP-4859.

### 9.5.1 Real-time Compression concepts

At a high level, the IBM RACE component compresses data that is written into the storage system dynamically. This compression occurs transparently, so Fibre Channel (FC) and internet Small Computer Systems Interface (iSCSI) connected hosts are not aware of the compression. RACE is an inline compression technology, which means that each host write is compressed as it passes through IBM Spectrum Virtualize to the disks.

This real-time technology has a clear benefit over other compression technologies that are post-processing based. These other technologies do not provide immediate capacity savings. Therefore, they are not a good fit for primary storage workloads, such as databases and active data set applications.

RACE is based on the LZ lossless data compression algorithm and operates by using a real-time method. When a host sends a write request, the request is acknowledged by the write cache of the system, and then destaged to the storage pool. As part of its destaging, the write request passes through the compression engine and is then stored in compressed format onto the storage pool. Therefore, writes are acknowledged immediately after they are received by the write cache with compression occurring as part of the destaging to internal or external physical storage.

Capacity is saved when the data is written by the host because the host writes are smaller when they are written to the storage pool. IBM Rtc is a self-tuning solution that adapts to the workload that runs on the system at any particular moment.

## 9.5.2 Implementing Rtc

To create a compressed volume, choose **Capacity Savings** → **Compressed** in the Create Volumes window, as shown in Figure 9-22. For more information about creating volumes, see Chapter 6, “Volumes” on page 255.



Figure 9-22 Creating compressed virtual disk

It is possible to add compression to an existing thin-provisioned or fully allocated volume. To do so, you must change the capacity savings settings of the volume by right-clicking it and selecting **Modify Capacity Settings**. In the menu, select **Compression** as the Capacity Savings option, as shown in Figure 9-23.

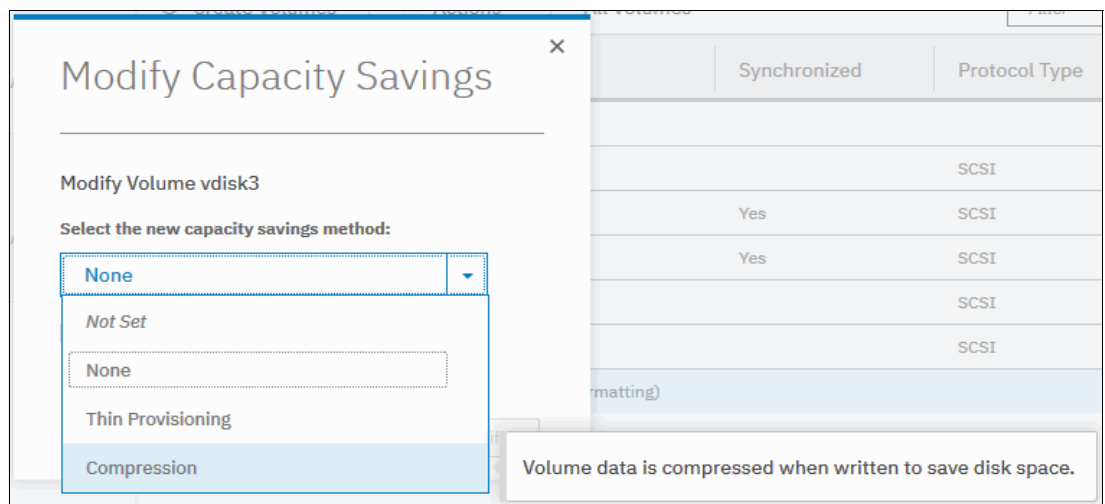


Figure 9-23 Selecting capacity setting

The system automatically compresses data on the volume and simultaneously starts compressing new data that is written by hosts. Because this process is nondisruptive, the volume stays online and remains accessible by hosts and applications.

This capability enables clients to regain space from the storage pool, which can then be reused for other applications.

With the virtualization of external storage systems, the ability to compress stored data significantly enhances and accelerates the benefit to users. This capability enables them to see a tremendous return on their investment.

On the initial purchase of a system with RtC, clients can defer their purchase of new storage. When storage is needed, IT purchases a lower amount of the required storage before compression.

For more information about volume migration for compressed volumes on standard pools to DRPs, see “Migrating to and from DRPs” on page 479.

## 9.6 Saving estimation for compression and deduplication

This section describes the tools that are used for sizing the environment for compression and deduplication.

### 9.6.1 Evaluate compression savings by using IBM Comprestimator

IBM Comprestimator is a utility that estimates the capacity savings that can be achieved when compression is used for storage volumes. The utility is integrated into the system by using the GUI and the CLI. It can also be installed and used on host systems.

Comprestimator provides a quick and accurate estimation of compression and thin-provisioning benefits. The utility performs read-only operations and therefore does not affect the data that is stored on the volume.

If the compression savings prove to be beneficial in your environment, volume mirroring can be used to convert volumes to compressed volumes.

In the **Volumes** pane, right-click any volume, and select **Space Savings** → **Estimate Compression Savings**. This action shows you the results and the date of the latest estimation cycle, as shown in Figure 9-24. If no analysis was done, the system suggests running it. A new estimation of all volumes can be started from this dialog. To run or rerun analysis on a single volume, select **Analyze** in **Space Savings** submenu.

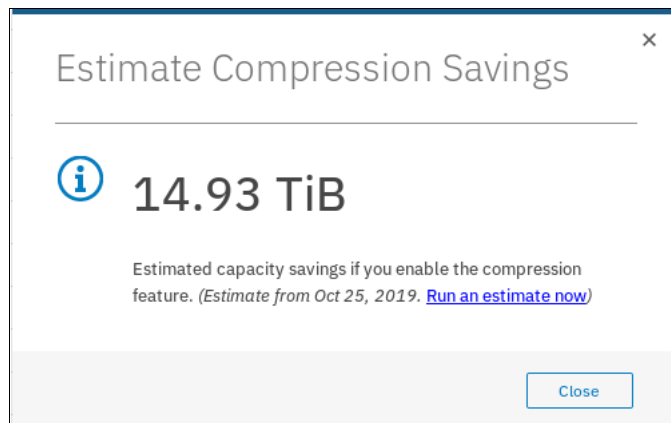


Figure 9-24 Estimate Compression Savings

To analyze all the volumes on the system from the CLI, run the **analyzevdiskbysystem** command.

The use of the command analyzes all the volumes that are created on the system. Volumes that are created during or after the analysis are not included and can be analyzed individually. The time it takes to analyze all the volumes on system depends on the number of volumes that are being analyzed and results can be expected at about a minute per volume. For example, if a system has 50 volumes, compression savings analysis takes approximately 50 minutes.

You can analyze a single volume by specifying its name or ID as a parameter in the **analyzevdisk** CLI command.

To check the progress of the analysis, run the **lsvdiskanalysisprogress** command. This command displays the total number of volumes on the system, total number of volumes that are remaining to be analyzed, and estimated time of completion.

The **lsvdiskanalysis** command is run to display information for thin provisioning and compression estimation analysis report for all volumes.

If you use a version of IBM Spectrum Virtualize that is older than V7.6 or if you want to estimate the compression savings of another IBM or non-IBM storage system, the separate IBM Comprestimator Utility can be installed on a host that is connected to the device that must be analyzed. For more information and the latest version of this utility, see this IBM Support [web page](#).

Consider the following recommended best practices for using Comprestimator:

- ▶ Run the Comprestimator utility before implementing an IBM Spectrum Virtualize solution and DRPs.
- ▶ Download the latest version of the Comprestimator utility if you are not using one that is included in your IBM Spectrum Virtualize solution.
- ▶ Use Comprestimator to analyze volumes that contain as much active data as possible rather than volumes that are nearly empty or newly created. This use ensures more accuracy when sizing your environment for compression and DRPs.

**Note:** Comprestimator can run for a long period (a few hours) when it is scanning a relatively empty device. The utility randomly selects and reads 256 KB samples from the device. If the sample is empty (that is, full of null values), it is skipped. A minimum number of samples with data is required to provide an accurate estimation. When a device is mostly empty, many random samples are empty. As a result, the utility runs for a longer time as it tries to gather enough non-empty samples that are required for an accurate estimate. The scan is stopped if the number of empty samples is over 95%.

## 9.6.2 Evaluating compression and deduplication

To help with the profiling and analysis of user workloads that are to be migrated to the new system, IBM provides a highly accurate Data Reduction Estimation Tool (DRET) that supports deduplication and compression. The tool operates by scanning target workloads on any legacy array (from IBM or third party) and then merging all scan results to provide an integrated system level data reduction estimate. It provides a report of what it expects the deduplication and compression savings to be from data written to a disk.

The DRET utility uses advanced mathematical and statistical algorithms to perform an analysis with low memory footprint. The utility runs on a host that can access the devices to be analyzed. It performs only read operations so it has no effect on the data that is stored on the device.

The following sections provide information about installing DRET on a host and using it to analyze devices on it. Depending on the environment configuration, DRET in many cases is used on more than one host to analyze more data types.

When DRET is used to analyze a block device that is used by a file system, all underlying data in the device is analyzed, regardless of whether this data belongs to files that were already deleted from the file system. For example, you can fill a 100 GB file system and make it 100% used, then delete all the files in the file system to make it 0% used. When scanning the block device that is used for storing the file system in this example, DRET accesses the data that belongs to the files that are deleted.

**Important:** The preferred method of using DRET is to analyze volumes that contain as much active data as possible rather than volumes that are mostly empty of data. This increases the accuracy level and reduces the risk of analyzing old data that is deleted, but might still have traces on the device.

For more information and the latest version of this utility, see this IBM Support [web page](#).

## 9.7 Overprovisioning and data reduction on external storage

Starting from IBM Spectrum Virtualize V8.1.x, overprovisioning on selected back-end controllers is supported. This means if back-end storage performs data deduplication or data compression on LUs provisioned from it, they still can be used as external MDisks on the system. However, other configuration and monitoring considerations must be taken into account.

Overprovisioned MDisks from controllers that are supported by this feature can be used as managed mode MDisks in the system and can be added to storage pools (including DRPs).

Implementation steps for overprovisioned MDisks are same as for fully allocated storage controllers. The system detects if the MDisk is overprovisioned, its total physical capacity, used, and remaining physical capacity. It detects if SCSI **unmap** commands are supported by the back-end. By sending SCSI **unmap** commands to overprovisioned MDisks, the system marks data that is no longer in use. Then, the garbage collection processes on the back-end can free unused capacity and convert it to free space.

At the time of this writing, the following back-end controllers are supported by overprovisioned MDisks:

- ▶ IBM A9000 V12.1.0 and above
- ▶ FlashSystem 900 V1.4
- ▶ FlashSystem 9000 AE2 expansions
- ▶ IBM Storwize or FlashSystem family system with code V8.1.0 and above
- ▶ Pure Storage



Extra caution is required when planning and monitoring capacity for such configurations. See Table 9-4 for an overview of configuration guidelines when overprovisioned external storage controllers are used.

Table 9-4 Using data reduction at two levels

System	Backend	Comments
DRP	Fully allocated	<i>Recommended.</i> Use DRP on the system to plan for compression and deduplication. DRP at the top level provides the best application capacity reporting.
Fully allocated	Overprovisioned, single tier of storage	<i>Recommended with appropriate precautions.</i> Track physical capacity use carefully to avoid out-of-space. The system can report physical use, but does not manage to avoid out-of-space. No visibility of each application's use at the system layer. If the backend runs out-of-space, there is limited ability to recover. Consider creating a sacrificial emergency space volume.
DRP with compression	Overprovisioned	<i>Recommended with appropriate precautions.</i> Assume 1:1 compression in backend storage and do not overcommit capacity in the backend. Small extra savings are achieved from compressing DRP metadata.
Fully allocated	Overprovisioned, multiple tiers of storage	<i>Use with great care.</i> Easy Tier is unaware of physical capacity in tiers of a hybrid pool and it tends to fill the top tier with hottest data. Changes in compressibility of data in the top tier can overcommit the storage leading to out-of-space.
DRP with thin-provisioned or fully allocated volumes	Overprovisioned	<i>Avoid.</i> Very difficult to understand physical capacity use of the uncompressed volumes. High risk of overcommitting the backend. If a mix of DRP and fully allocated volumes is required, use separate pools.
DRP	DRP	<i>Avoid.</i> Creates two levels of IO amplification and capacity overhead. DRP at the bottom layer provides no benefit.

When DRPs are used with a compressing backend controller, use compression in DRP and avoid overcommitting the backend by assuming a 1:1 compression ratio in backend storage. Small, extra savings are realized from compressing metadata.

**Note:** Fully allocated volumes that are above overprovisioned MDisk configurations must be avoided or used with extreme caution because it can lead to overcommitting back-end storage.

The concept of provisioning groups is used for capacity reporting and monitoring of overprovisioned external storage controllers. A provisioning group is an object that represents a set of MDisks that share physical resources. Each overprovisioned MDisk is part of a provisioning group that defines the physical storage resources available to a set of MDisks.

Storage controllers report the usable capacity of an overprovisioned MDisk based on its provisioning group. If multiple MDisks are part of the same provisioning group, all of these MDisks share the physical storage resources and report the same usable capacity. However, this usable capacity is not available to each MDisk individually because it is shared between all of these MDisks.

Provisioning groups are used differently depending on the back-end storage, as shown in the following examples:

- ▶ A9000 and FlashSystem 900: The entire subsystem forms one provisioning group.
- ▶ Storwize family systems: The storage pool forms a provisioning group, which allows more than one independent provisioning group in a system.
- ▶ RAID with compressing drives: An array is a provisioning group that presents the physical storage that is in use much like an external array.

Capacity usage should primarily be monitored on the overprovisioned backend storage controller itself.

From the system, capacity usage can be monitored on over-provisioned MDisks by using one of the following methods:

- ▶ GUI Dashboard, as shown in Figure 9-25.

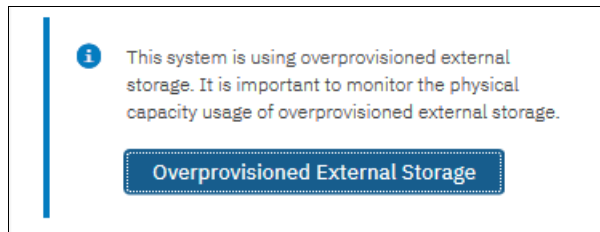


Figure 9-25 Dashboard button for overprovisioned storage monitoring

Click **Overprovisioned External Storage** to show an overview of overprovisioned MDisks and provisioning groups used by the system, as shown in Figure 9-26.

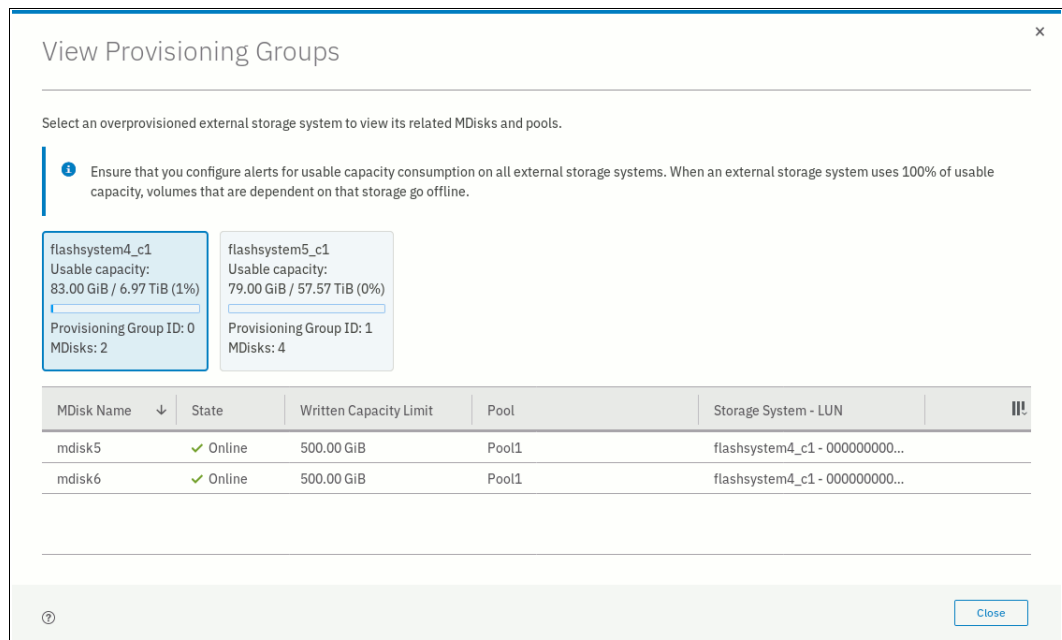


Figure 9-26 View Provisioning Groups

- MDisk properties window that opens by right-clicking an MDisk in **Pools** → **MDisks by Pools** and **Properties** menu option, as shown in Figure 9-27.

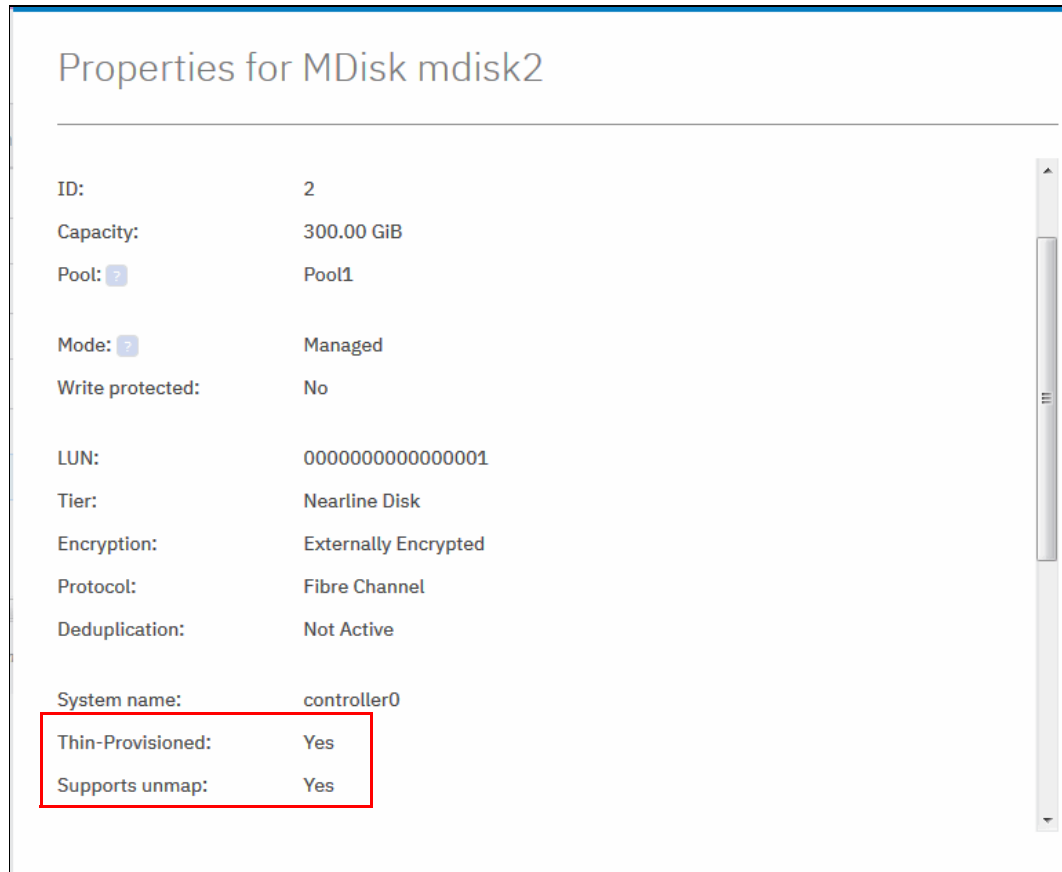


Figure 9-27 Thin-provisioned MDisk properties

- CLI command `lsmdisk` with MDisk name or ID as a parameter, as shown in Example 9-11.

*Example 9-11 Ismdisk parameters for thin provisioned MDisks*

```
IBM_2145:ITS0-SV1:superuser>lsmdisk mdisk2
id 2
name mdisk2
status online
mode managed
<...>
dedupe no
<...>
over_provisioned yes
supports_unmap yes
provisioning_group_id
physical_capacity 299.00GB
physical_free_capacity 288.00GB
write_protected no
allocated_capacity 11.00GB
effective_used_capacity 300.00GB
```

The overprovisioning status and SCSI Unmap support for the selected MDisk are displayed.

The `physical_capacity` and `physical_free_capacity` parameters belong to the MDisk's provisioning group. They indicate the total physical storage capacity and formatted available physical space in the provisioning group that contains this MDisk.

**Note:** It is not recommended to create multiple storage pools from MDisks in a single provisioning group.



# Advanced Copy Services

This chapter describes the Advanced Copy Services that are a group of functions that provide different methods of data copy. It also describes the storage software capabilities to support the interaction with hybrid clouds. These functions are enabled by IBM Spectrum Virtualize software that runs inside IBM SAN Volume Controller.

This chapter includes the following topics:

- ▶ 10.1, “IBM FlashCopy” on page 492
- ▶ 10.2, “Managing FlashCopy by using the GUI” on page 523
- ▶ 10.3, “Transparent Cloud Tiering” on page 562
- ▶ 10.4, “Implementing Transparent Cloud Tiering” on page 565
- ▶ 10.5, “Volume mirroring and migration options” on page 576
- ▶ 10.6, “Remote Copy” on page 578
- ▶ 10.7, “Remote Copy commands” on page 606
- ▶ 10.8, “Native IP replication” on page 613
- ▶ 10.9, “Managing Remote Copy by using the GUI” on page 634
- ▶ 10.10, “Remote Copy memory allocation” on page 669
- ▶ 10.11, “Troubleshooting Remote Copy” on page 670

## 10.1 IBM FlashCopy

Through the IBM FlashCopy function of the IBM Spectrum Virtualize, you can perform a *point-in-time (PiT) copy* of one or more volumes. This section describes the inner workings of FlashCopy and provides more information about its configuration and use.

You can use FlashCopy to help you solve critical and challenging business needs that require duplication of data of your source volume. Volumes can remain online and active while you create consistent copies of the data sets. Because the copy is performed at the block level, it operates below the host operating system and its cache. Therefore, the copy is not apparent to the host unless it is mapped.

While the FlashCopy operation is performed, the source volume is frozen briefly to initialize the FlashCopy bitmap after which I/O can resume. Although several FlashCopy options require the data to be copied from the source to the target in the background (which can take time to complete), the resulting data on the target volume is presented so that the copy appears to complete immediately. This feature means that the copy can immediately be mapped to a host and is directly accessible for read *and* write operations.

### 10.1.1 Business requirements for FlashCopy

When you are deciding whether FlashCopy addresses your needs, you must adopt a combined business and technical view of the problems that you want to solve. First, determine the needs from a business perspective. Then, determine whether FlashCopy can address the technical needs of those business requirements.

The business applications for FlashCopy are wide-ranging. Common use cases for FlashCopy include, but are not limited to, the following examples of rapidly creating:

- ▶ Consistent backups of dynamically changing data
- ▶ Consistent copies of production data to facilitate data movement or migration between hosts
- ▶ Copies of production data sets for application development and testing, auditing purposes and data mining, and for quality assurance

Regardless of your business needs, FlashCopy within the IBM Spectrum Virtualize is flexible and offers a broad feature set, which makes it applicable to several scenarios.

#### **Back up improvements with FlashCopy**

FlashCopy does not reduce the time that it takes to perform a backup to traditional backup infrastructure. However, it can be used to minimize and under certain conditions, eliminate application downtime that is associated with performing backups. FlashCopy can also transfer the resource usage of performing intensive backups from production systems.

After the FlashCopy is performed, the resulting image of the data can be backed up to tape, as though it were the source system. After the copy to tape is completed, the image data is redundant and the target volumes can be discarded. For time-limited applications, such as these examples, “no copy” or incremental FlashCopy is used most often. The use of these methods puts less load on your servers infrastructure.

When FlashCopy is used for backup purposes, the target data often is managed as read-only *at the operating system level*. This approach provides extra security by ensuring that your target data was not modified and remains true to the source.

## Restore with FlashCopy

FlashCopy can perform a restore from any FlashCopy mapping. Therefore, you can restore (or copy) from the target to the source of your regular FlashCopy relationships. When restoring data from FlashCopy, this method can be qualified as reversing the direction of the FlashCopy mappings.

This capability has the following benefits:

- ▶ Pairing mistakes are not a concern. You trigger a restore.
- ▶ The process appears instantaneous.
- ▶ You can maintain a pristine image of your data while you are restoring what was the primary data.

This approach can be used for various applications, such as recovering your production database application after an errant batch process that caused extensive damage.

**Preferred practices:** Although restoring from a FlashCopy is quicker than a traditional tape media restore, you must not use restoring from a FlashCopy as a substitute for good backup and archiving practices. Instead, keep one to several iterations of your FlashCopies so that you can near-instantly recover your data from the most recent history, and keep your long-term backup and archive as suitable for your business.

In addition to the restore option that copies the original blocks from the target volume to modified blocks on the source volume, the target can be used to perform a restore of individual files. To do that, you make the target available on a host. It is suggested to not make the target available to the source host because seeing duplicates of disks causes problems for most host operating systems. Copy the files to the source by using normal host data copy methods for your environment.

For more information about how to use reverse FlashCopy, see 10.1.12, “Reverse FlashCopy” on page 514.

## Moving and migrating data with FlashCopy

FlashCopy can be used to facilitate the movement or migration of data between hosts while minimizing downtime for applications. By using FlashCopy, application data can be copied from source volumes to new target volumes while applications remain online. After the volumes are fully copied and synchronized, the application can be brought down and then immediately brought back up on the new server that is accessing the new FlashCopy target volumes.

This method differs from the other migration methods, which are described later in this chapter. Common uses for this capability are host and back-end storage hardware refreshes.

## Application testing with FlashCopy

It is often important to test a new version of an application or operating system that is using actual production data. This testing ensures the highest quality possible for your environment. FlashCopy makes this type of testing easy to accomplish without putting the production data at risk or requiring downtime to create a constant copy.

You can create a FlashCopy of your source and use that for your testing. This copy is a duplicate of your production data down to the block level so that even physical disk identifiers are copied. Therefore, it is impossible for your applications to tell the difference.

You can also use the FlashCopy feature to create restart points for long running batch jobs. This option means that if a batch job fails several days into its run, it might be possible to restart the job from a saved copy of its data rather than rerunning the entire multiday job.

### 10.1.2 FlashCopy principles and terminology

The FlashCopy function creates a PiT or time-zero (T0) copy of data that is stored on a source volume to a target volume by using a Copy on Write (CoW) and copy on-demand mechanism.

When a FlashCopy operation starts, a checkpoint creates a *bitmap table* that indicates that no part of the source volume was copied. Each bit in the bitmap table represents one region of the source volume and its corresponding region on the target volume. Each region is called a *grain*.

The relationship between two volumes defines the way data are copied and is called a *FlashCopy mapping*.

FlashCopy mappings between multiple volumes can be grouped in a Consistency group to ensure their PiT (or T0) is identical for all of them. A simple one-to-one FlashCopy mapping does not need to belong to a consistency group.

Figure 10-1 shows the basic terms that are used with FlashCopy. All elements are explained later in this chapter.

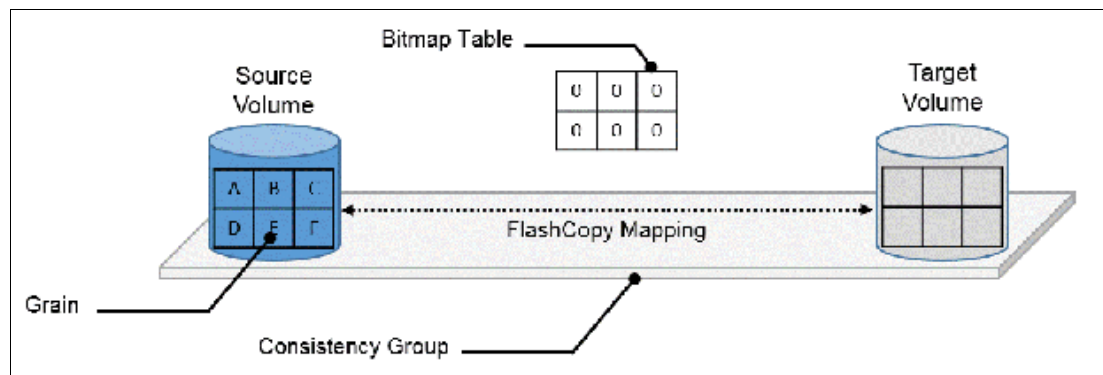


Figure 10-1 FlashCopy terminology

### 10.1.3 FlashCopy mapping

The relationship between the source volume and the target volume is defined by a FlashCopy mapping. The FlashCopy mapping can have three different types, four attributes, and seven different states.

The FlashCopy mapping can be one of the following types:

- ▶ Snapshot: Sometimes referred to as *nocopy*, a snapshot is a PiT copy of a volume without a background copy of the data from the source volume to the target. Only the changed blocks on the source volume are copied. The target copy cannot be used without an active link to the source.
- ▶ Clone: Sometimes referred to as *full copy*, a clone is a PiT copy of a volume with background copy of the data from the source volume to the target. All blocks from the source volume are copied to the target volume. The target copy becomes a usable independent volume.



- ▶ **Backup:** Sometimes referred to as *incremental*, a backup FlashCopy mapping consists of a PiT full copy of a source volume, plus periodic increments or “deltas” of data that changed between two points in time.

The FlashCopy mapping has four property attributes (clean rate, copy rate, autodelete, incremental) and seven different states that are described later in this chapter. Users can perform the following actions on a FlashCopy mapping:

- ▶ **Create:** Define a source and target, and set the properties of the mapping.
- ▶ **Prepare:** The system must be prepared before a FlashCopy copy starts. It flushes the cache and makes it “transparent” for a short time, so no data is lost.
- ▶ **Start:** The FlashCopy mapping is started and the copy begins immediately. The target volume is immediately accessible.
- ▶ **Stop:** The FlashCopy mapping is stopped (by the system or by the user). Depending on the state of the mapping, the target volume is usable or not usable.
- ▶ **Modify:** Some properties of the FlashCopy mapping can be modified after creation.
- ▶ **Delete:** Delete the FlashCopy mapping. This action does not delete volumes (source or target) from the mapping.

The source and target volumes must be the same size. The minimum granularity that IBM Spectrum Virtualize supports for FlashCopy is an entire volume. It is not possible to use FlashCopy to copy only part of a volume.

**Important:** As with any PiT copy technology, you are bound by operating system and application requirements for interdependent data and the restriction to an entire volume.

The source and target volumes must belong to the same SAN Volume Controller system, but they do not have to be in the same I/O group or storage pool.

Volumes that are members of a FlashCopy mapping cannot have their size increased or decreased while they are members of the FlashCopy mapping.

All FlashCopy operations occur on FlashCopy mappings. FlashCopy does not alter the volumes. However, multiple operations can occur at the same time on multiple FlashCopy mappings because of the use of Consistency Groups.

## 10.1.4 Consistency groups

To overcome the issue of dependent writes across volumes and to create a consistent image of the client data, a FlashCopy operation must be performed on multiple volumes as an atomic operation. To accomplish this method, the IBM Spectrum Virtualize supports the concept of *consistency groups*.

Consistency groups address the requirement to preserve PiT data consistency across multiple volumes for applications that include related data that spans multiple volumes. For these volumes, consistency groups maintain the integrity of the FlashCopy by ensuring that “dependent writes” are run in the application’s intended sequence. Also, consistency groups provide an easy way to manage several mappings.

FlashCopy mappings can be part of a consistency group, even if only one mapping exists in the consistency group. If a FlashCopy mapping is not part of any consistency group, it is referred to as *stand-alone*.

## Dependent writes

It is crucial to use consistency groups when a data set spans multiple volumes. Consider the following typical sequence of writes for a database update transaction:

1. A write is run to update the database log, which indicates that a database update is about to be performed.
2. A second write is run to perform the update to the database.
3. A third write is run to update the database log, which indicates that the database update completed successfully.

The database ensures the correct ordering of these writes by waiting for each step to complete before the next step is started. However, if the database log (updates 1 and 3) and the database (update 2) are on separate volumes, it is possible for the FlashCopy of the database volume to occur before the FlashCopy of the database log. This sequence can result in the target volumes seeing writes 1 and 3 but not 2 because the FlashCopy of the database volume occurred before the write was completed.

In this case, if the database was restarted by using the backup that was made from the FlashCopy target volumes, the database log indicates that the transaction completed successfully. In fact, it did not complete successfully because the FlashCopy of the volume with the database file was started (the bitmap was created) before the write completed to the volume. Therefore, the transaction is lost and the integrity of the database is in question.

Most of the actions that the user can perform on a FlashCopy mapping are the same for consistency groups.

## 10.1.5 Crash consistent copy and hosts considerations

FlashCopy consistency groups do not provide application consistency. It ensures only that volume points-in-time are consistent between them.

Because FlashCopy is at the block level, it is necessary to understand the interaction between your application and the host operating system. From a logical standpoint, it is easiest to think of these objects as “layers” that sit on top of one another. The application is the topmost layer, and beneath it is the operating system layer.

Both of these layers have various levels and methods of caching data to provide better speed. Therefore, because the IBM SAN Volume Controller and FlashCopy sit below these layers, they are unaware of the cache at the application or operating system layers.

To ensure the integrity of the copy that is made, it is necessary to flush the host operating system and application cache for any outstanding reads or writes before the FlashCopy operation is performed. Failing to flush the host operating system and application cache produces what is referred to as a *crash consistent* copy.

The resulting copy requires the same type of recovery procedure, such as log replay and file system checks, that is required following a host crash. FlashCopies that are crash consistent often can be used after file system and application recovery procedures.

Various operating systems and applications provide facilities to stop I/O operations and ensure that all data is flushed from host cache. If these facilities are available, they can be used to prepare for a FlashCopy operation. When this type of facility is unavailable, the host cache must be flushed manually by quiescing the application and unmounting the file system or drives.

The target volumes are overwritten with a complete image of the source volumes. Before the FlashCopy mappings are started, any data that is held on the host operating system (or application) caches for the target volumes must be discarded. The easiest way to ensure that no data is held in these caches is to unmount the target volumes before the FlashCopy operation starts.

**Preferred practice:** From a practical standpoint, when you have an application that is backed by a database and you want to make a FlashCopy of that application's data, it is sufficient in most cases to use the write-suspend method that is available in most modern databases. This is possible because the database maintains strict control over I/O.

This method is as opposed to flushing data from the application and backing database, which is always the suggested method because it is safer. However, this method can be used when facilities do not exist or your environment includes time sensitivity.

### IBM Spectrum Protect Snapshot

IBM FlashCopy is not application aware and a third-party tool is needed to link the application to the FlashCopy operations.

IBM Spectrum Protect Snapshot protects data with integrated, application-aware snapshot backup and restore capabilities that use FlashCopy technologies in the IBM Spectrum Virtualize.

You can protect data that is stored by IBM DB2 SAP, Oracle, Microsoft Exchange, and Microsoft SQL Server applications. You can create and manage volume-level snapshots for file systems and custom applications.

In addition, it enables you to manage frequent, near-instant, nondisruptive, application-aware backups and restores that use integrated application and VMware snapshot technologies. IBM Spectrum Protect Snapshot can be widely used in IBM and non-IBM storage systems.

**Note:** To see how IBM Spectrum Protect Snapshot can help your business, see [IBM Knowledge Center](#).

### 10.1.6 Grains and bitmap: I/O indirection

When a FlashCopy operation starts, a checkpoint is made of the source volume. No data is copied at the time that a start operation occurs. Instead, the checkpoint creates a bitmap that indicates that no part of the source volume was copied. Each bit in the bitmap represents one region of the source volume. Each region is called a *grain*.

You can think of the bitmap as a simple table of ones or zeros. The table tracks the difference between a source volume grains and a target volume grains. At the creation of the FlashCopy mapping, the table is filled with zeros, which indicates that no grain is copied yet.

When a grain is copied from source to target, the region of the bitmap that refers to that grain is updated (for example, from “0” to “1”), as shown in Figure 10-2.

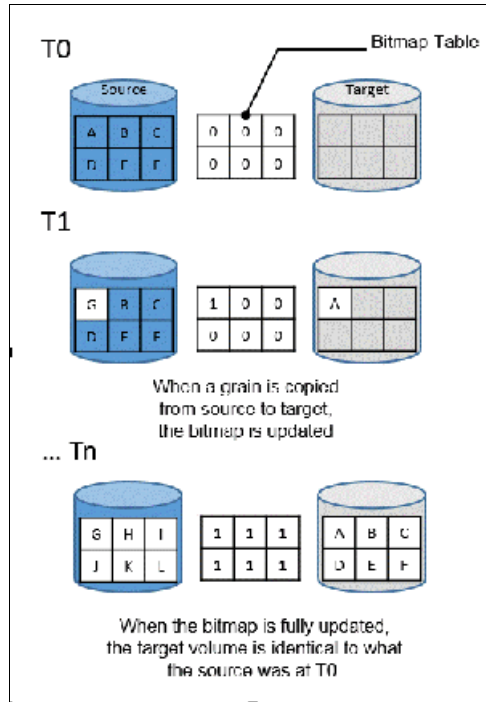


Figure 10-2 A simplified representation of grains and bitmap

The grain size can be 64 KB or 256 KB (the default is 256 KB). The grain size cannot be selected by the user when a FlashCopy mapping is created from the GUI. The FlashCopy bitmap contains 1 bit for each grain. The bit records whether the associated grain is split by copying the grain from the source to the target.

After a FlashCopy mapping is created, the grain size for that FlashCopy mapping cannot be changed. When a FlashCopy mapping is created, the grain size of that mapping is used if the grain size parameter is not specified and one of the volumes in the mapping is part of a FlashCopy mapping.

If neither volume in the new mapping is part of another FlashCopy mapping and at least one of the volumes in the mapping is a compressed volume, the default grain size is 64 KB for performance considerations. Other than in this situation, the default grain size is 256 KB.

### Copy on Write and Copy on Demand

IBM Spectrum Virtualize FlashCopy uses CoW mechanism to copy data from a source volume to a target volume.

As shown in Figure 10-3, when data is written on a source volume, the grain where the to-be-changed blocks is stored is first copied to the target volume and then modified on the source volume. The bitmap is updated to track the copy.

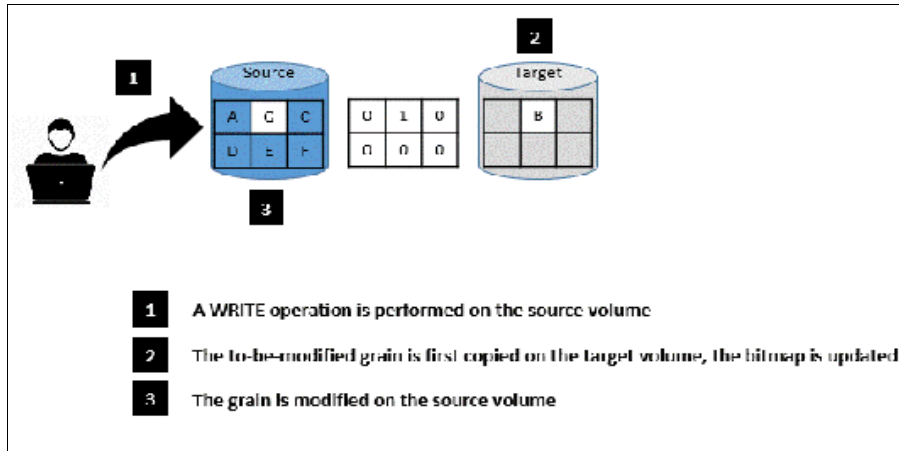


Figure 10-3 Copy on Write steps

With IBM FlashCopy, the target volume is immediately accessible for read *and* write operations. Therefore, a target volume can be modified, even if it is part of a FlashCopy mapping. As shown in Figure 10-4, when a Write operation is performed on the *target* volume, the grain that contains the blocks to be changed is first copied from the source (*Copy on-Demand*). It is then modified with the new value. The bitmap is modified so the grain from the source is *not* copied again, even if it is changed or if a background copy is enabled.

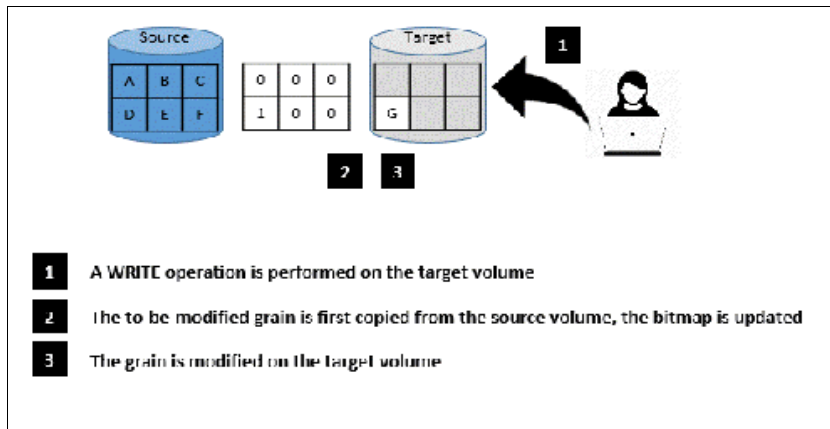


Figure 10-4 Copy on-Demand steps

**Note:** If all the blocks of the grain to be modified are changed, there is no need to copy the source grain first. There is no copy on demand and it is directly modified.

## FlashCopy indirection layer

The FlashCopy indirection layer governs the I/O to the source and target volumes when a FlashCopy mapping is started, which is done by using the FlashCopy bitmap. The purpose of the FlashCopy indirection layer is to enable the source and target volumes for read and write I/O immediately after the FlashCopy is started.

The indirection Layer intercepts any I/O coming from a host (read or write operation) and addressed to a FlashCopy volume (source or target). It determines whether the addressed volume is a source or a target, its direction (read or write), and the state of the bitmap table for the FlashCopy mapping that the addressed volume is in. It then decides what operation to perform. The different I/O indirections are described next.

### Read from the source Volume

When a user performs a read operation on the source volume, there is no redirection. The operation is similar to what is done with a volume that is not part of a FlashCopy mapping.

### Write on the source Volume

Performing a write operation on the source volume modifies a block or a set of blocks, which modifies a grain on the source. It generates one of the following actions, depending on the state of the grain to be modified.

Consider the following points:

- ▶ If the bitmap indicates that the grain was copied, the source grain is changed and the target volume and the bitmap table remain unchanged, as shown in Figure 10-5.

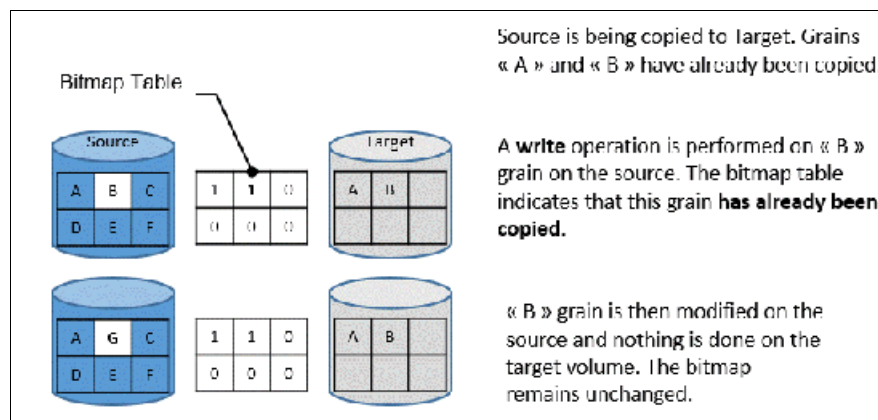


Figure 10-5 Modifying an already copied grain on the Source

- ▶ If the bitmap indicates that the grain is not yet copied, the grain is first copied on the target (CoW), the bitmap table is updated, and the grain is modified on the source, as shown in Figure 10-6.

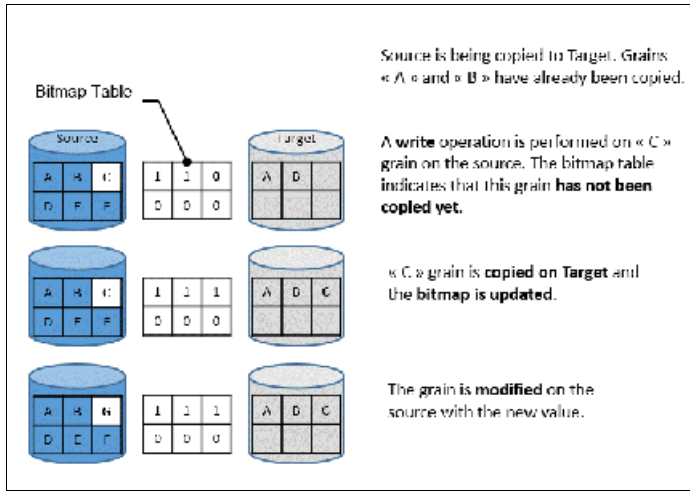


Figure 10-6 Modifying a non-copied grain on the Source

### Write on a Target Volume

Because FlashCopy target volumes are immediately accessible in Read and Write mode, it is possible to perform write operations on the target volume when the FlashCopy mapping is started. Performing a write operation on the target generates one of the following actions, depending on the bitmap:

- ▶ If the bitmap indicates the grain to be modified on the target was not yet copied, it is first copied from the source (copy on demand). The bitmap is updated, and the grain is modified on the target with the new value, as shown in Figure 10-7. The source volume remains unchanged.

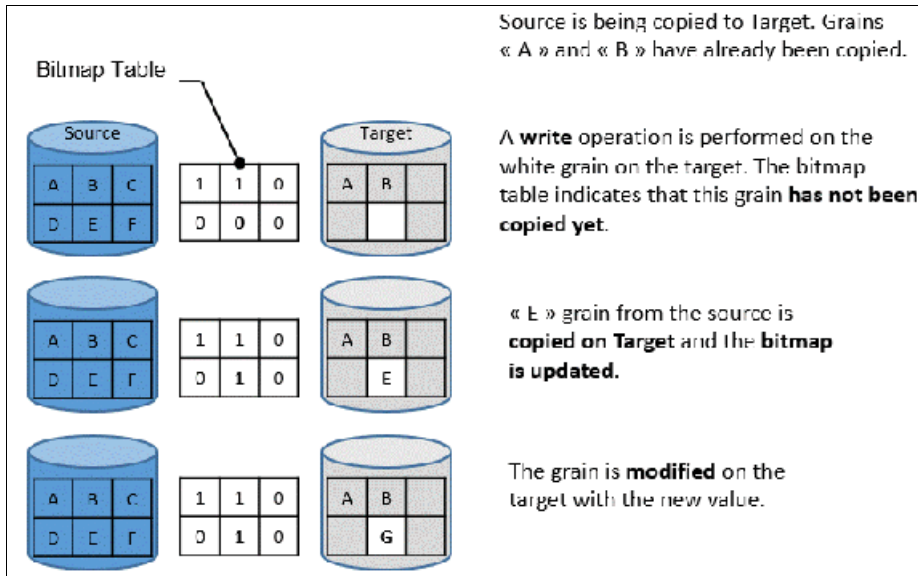


Figure 10-7 Modifying a non-copied grain on the target

**Note:** If the entire grain is to be modified and not only part of it (some blocks only), the copy on demand is bypassed. The bitmap is updated, and the grain on the target is modified but not copied first.

- If the bitmap indicates the grain to be modified on the target was copied, it is directly changed. The bitmap is *not* updated, and the grain is modified on the target with the new value, as shown in Figure 10-8.

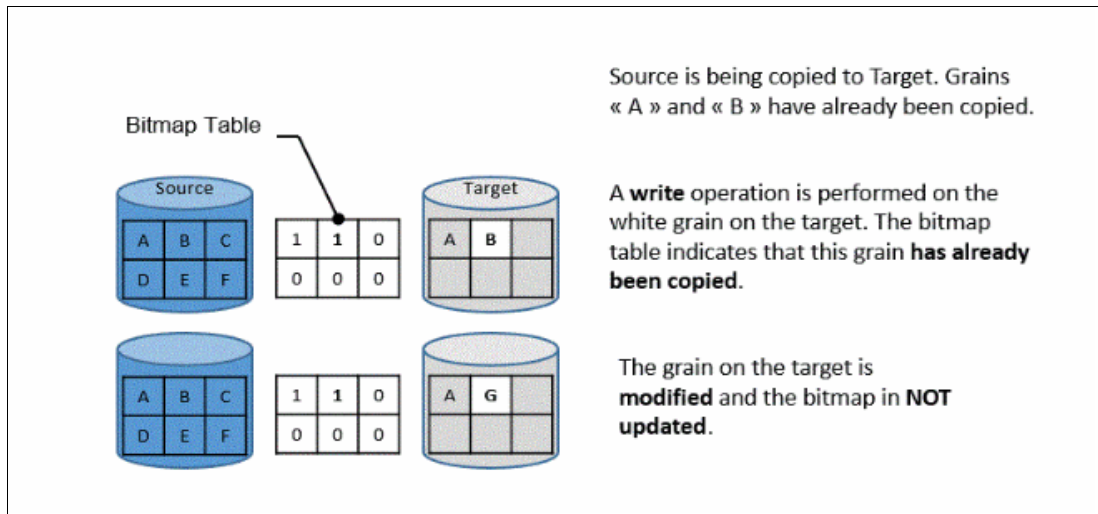


Figure 10-8 Modifying an already copied grain on the Target

**Note:** The bitmap is not updated in that case. Otherwise, it might be copied from the source late if a background copy is ongoing or if write operations are made on the source. That process over-writes the changed grain on the target.



## Read from a target volume

Performing a read operation on the target volume returns the value in the grain on the source or on the target, depending on the bitmap. Consider the following points:

- ▶ If the bitmap indicates that the grain was copied from the source or that the grain was modified on the target, the grain on the target is read, as shown in Figure 10-9.
- ▶ If the bitmap indicates that the grain was not yet copied from the source or was not modified on the target, the grain on the source is read, as shown in Figure 10-9.

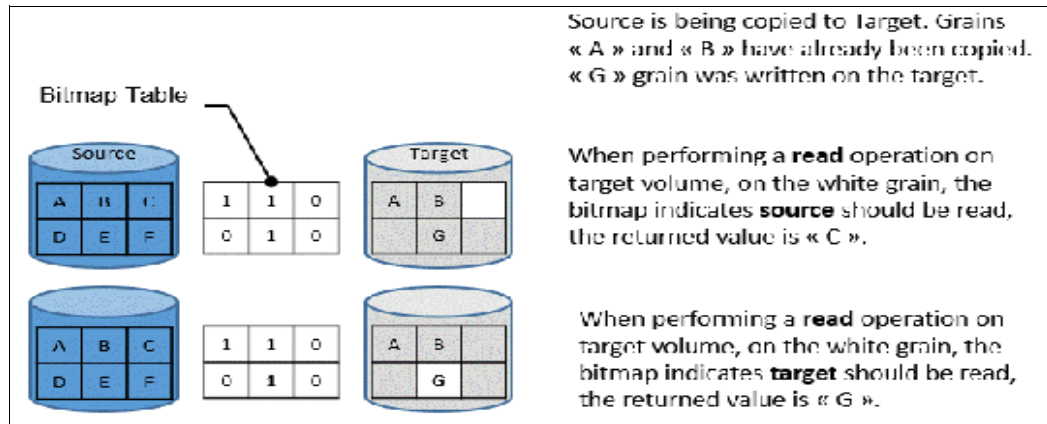


Figure 10-9 Reading a grain on target

If source has multiple targets, the Indirection layer algorithm behaves differently on Target I/Os. For more information about multi-target operations, see 10.1.11, “Multiple target FlashCopy” on page 509.

## 10.1.7 Interaction with cache

IBM Spectrum Virtualize based systems have their cache divided into upper and lower cache. Upper cache serves mostly as write cache and hides the write latency from the hosts and application. Lower cache is a read/write cache and optimizes I/O to and from disks. Figure 10-10 shows the IBM Spectrum Virtualize cache architecture.

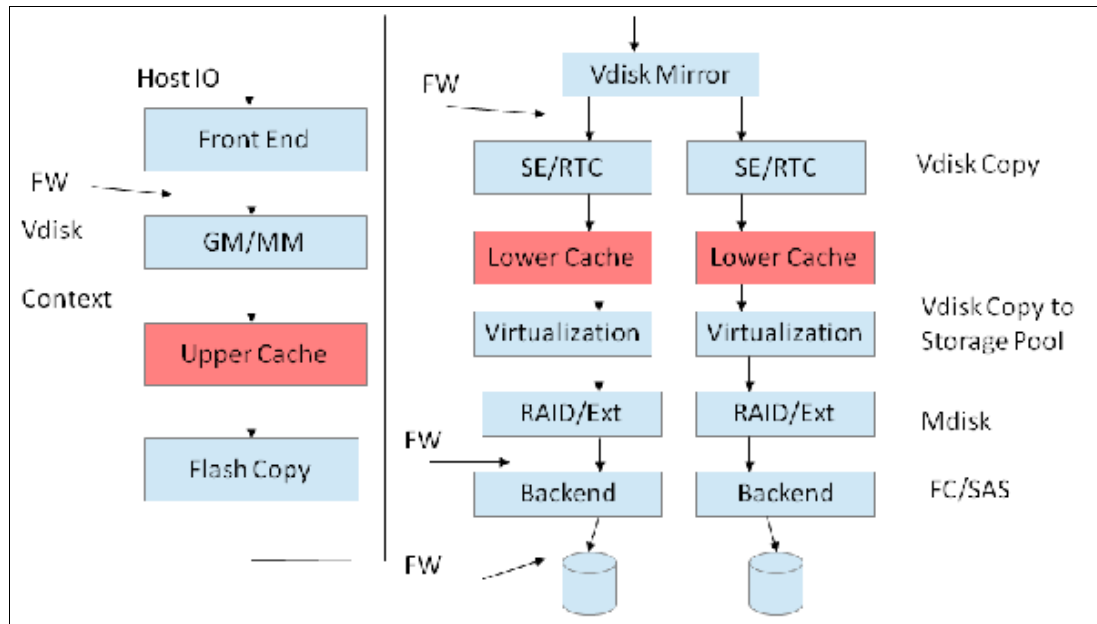


Figure 10-10 New cache architecture

This copy-on-write process introduces significant latency into write operations. To isolate the active application from this extra latency, the FlashCopy indirection layer is placed logically between upper and lower cache. Therefore, the extra latency that is introduced by the copy-on-write process is encountered only by the internal cache operations and not by the application.

The two-level cache provides more performance improvements to the FlashCopy mechanism. Because the FlashCopy layer is above lower cache in the IBM Spectrum Virtualize software stack, it can benefit from read prefetching and coalescing writes to backend storage. Preparing FlashCopy benefits from the two-level cache because upper cache write data does not have to go directly to backend storage, but to lower cache layer instead.

## 10.1.8 Background Copy Rate

The Background Copy Rate is a property of a FlashCopy mapping. A grain copy from the source to the target can occur when triggered by a write operation on the source or target volume, or when background copy is enabled. With background copy enabled, the target volume eventually becomes a clone of the source volume at the time the mapping was started (T0). When the copy is completed, the mapping can be removed between the two volumes and you can end up with two independent volumes.

The background copy rate property determines the speed at which grains are copied as a background operation, immediately after the FlashCopy mapping is started. That speed is defined by the user when the FlashCopy mapping is created, and can be changed dynamically for each individual mapping, whatever its state. Mapping copy rate values can be 0 - 150, with the corresponding speeds that are listed in Table 10-1.

Table 10-1 Copy rate values

User-specified copy rate attribute value	Data copied/sec	256 KB grains/sec	64 KB grains/sec
1 - 10	128 kibibytes (KiB)	0.5	2
11 - 20	256 KiB	1	4
21 - 30	512 KiB	2	8
31 - 40	1 MiB	4	16
41 - 50	2 MiB	8	32
51 - 60	4 MiB	16	64
61 - 70	8 MiB	32	128
71 - 80	16 MiB	64	256
81 - 90	32 MiB	128	512
91 - 100	64 MiB	256	1024
101 - 110	128 MiB	512	2048
111 - 120	256 MiB	1024	4096
121 - 130	512 MiB	2048	8192
131 - 140	1 GiB	4096	16384
141 - 150	2 GiB	8192	32768

When the background copy function is not performed (copy rate = 0), the target volume remains a valid copy of the source data only while the FlashCopy mapping remains in place.

The *grains per second* numbers represent the maximum number of grains that the SAN Volume Controller copies per second. This amount assumes that the bandwidth to the managed disks (MDisks) can accommodate this rate.

If the IBM SAN Volume Controller cannot achieve these copy rates because of insufficient width from the nodes to the MDisks, the background copy I/O contends for resources on an equal basis with the I/O that is arriving from the hosts. Background copy I/O and I/O that is arriving from the hosts tend to see an increase in latency and a consequential reduction in throughput.

Background copy and foreground I/O continue to progress, and do not stop, hang, or cause the node to fail.

The background copy is performed by one of the nodes that belong to the I/O group in which the source volume is stored. This responsibility is moved to the other node in the I/O group if the node that performs the background and stopping copy fails.

## 10.1.9 Incremental FlashCopy

When a FlashCopy mapping is stopped (because the entire source volume was copied onto the target volume or a user manually stopped it), the bitmap table is reset. Therefore, when the same FlashCopy is started again, the copy process is restarted from the beginning.

Running the `-incremental` option when creating the FlashCopy mapping allows the system to keep the bitmap as it is when the mapping is stopped. Therefore, when the mapping is started again (at another PiT), the bitmap is reused and only changes between the two copies are applied to the target.

A system that provides Incremental FlashCopy capability allows the system administrator to refresh a target volume without having to wait for a full copy of the source volume to be complete. At the point of refreshing the target volume, if the data changed on the source or target volumes for a particular grain, the grain from the source volume is copied to the target.

The advantages of Incremental FlashCopy are useful only if a previous full copy of the source volume was obtained. Incremental FlashCopy helps with only further recovery time objectives (RTOs, which are time needed to recover data from a previous state), it does not help with the initial RTO.

For example, as shown in Figure 10-11 on page 507, a FlashCopy mapping was defined between a source volume and a target volume by using the `-incremental` option.

Consider the following points:

- ▶ The mapping is started on Copy1 date. A *full copy* of the source volume is made, and the bitmap is updated every time that a grain is copied. At the end of Copy1, all grains are copied and the target volume is an exact replica of the source volume at the beginning of Copy1. Although the mapping is stopped because of the `-incremental` option, the bitmap is maintained.
- ▶ Changes are made on the source volume and the bitmap is updated, although the FlashCopy mapping is not active. For example, grains E and C on the source are changed in G and H, their corresponding bits are changed in the bitmap. The target volume is untouched.
- ▶ The mapping is started again on Copy2 date. The bitmap indicates that only grains E and C were changed; therefore, only G and H are copied on the target volume. The other grains do not need to be copied because they were copied the first time. The copy time is much quicker than for the first copy as only a fraction of the source volume is copied.

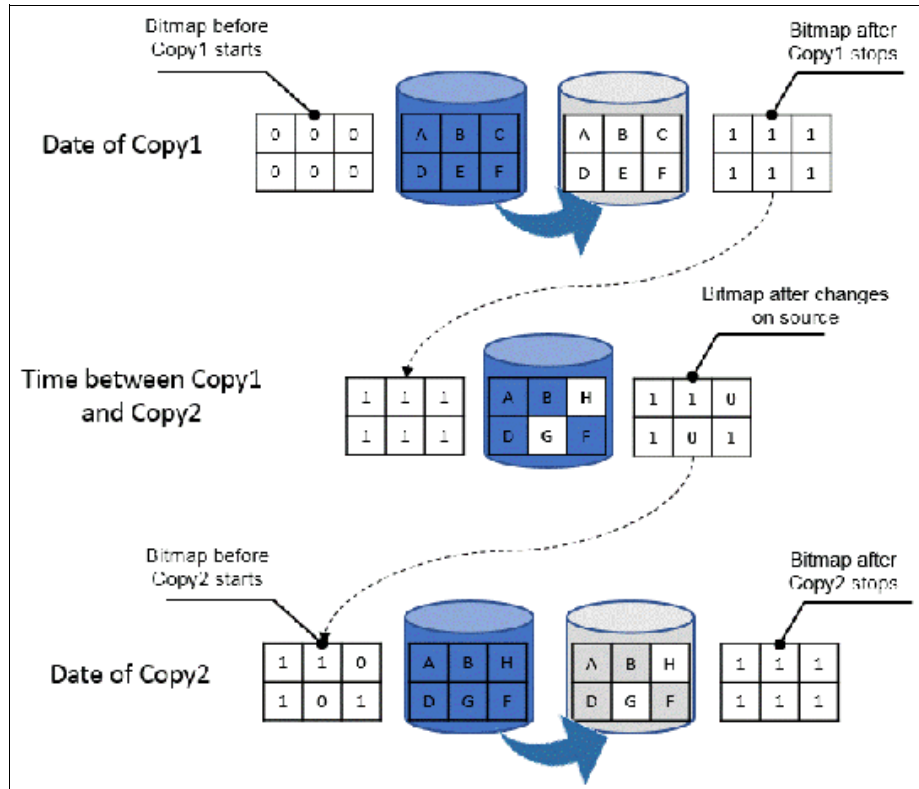


Figure 10-11 Incremental FlashCopy example

### 10.1.10 Starting FlashCopy mappings and consistency groups

You can prepare, start, or stop FlashCopy on a stand-alone mapping or a consistency group.

When the CLI is used to perform FlashCopy on volumes, run a **prestartfcmap** or **prestartfcconsistgrp** command *before* you start a FlashCopy (regardless of the type and options specified). These commands put the cache into write-through mode and provides a flushing of the I/O that is bound for your volume. After FlashCopy is started, an effective copy of a source volume to a target volume is created.

The content of the source volume is presented immediately on the target volume and the original content of the target volume is lost.

FlashCopy commands can then be run to the FlashCopy consistency group and therefore, simultaneously for all of the FlashCopy mappings that are defined in the consistency group. For example, when a FlashCopy **start** command is run to the consistency group, all of the FlashCopy mappings in the consistency group are started at the same time. This simultaneous start results in a PiT copy that is consistent across all of the FlashCopy mappings that are contained in the consistency group.

Rather than running **prestartfcmap** or **prestartfcconsistgrp**, you can also use the **-prep** parameter in the **startfcmap** or **startfcconsistgrp** command to prepare and start FlashCopy in one step.

**Important:** After an individual FlashCopy mapping is added to a consistency group, it can be managed as part of the group only. Operations, such as prepare, start, and stop, are no longer allowed on the individual mapping.

## FlashCopy mapping states

At any point, a mapping is in one of the following states:

- ▶ Idle or copied

The source and target volumes act as independent volumes, even if a mapping exists between the two. Read and write caching is enabled for the source and the target volumes. If the mapping is incremental and the background copy is complete, the mapping records only the differences between the source and target volumes. If the connection to both nodes in the I/O group that the mapping is assigned to is lost, the source and target volumes are offline.

- ▶ Copying

The copy is in progress. Read and write caching is enabled on the source and the target volumes.

- ▶ Prepared

The mapping is ready to start. The target volume is online, but is not accessible. The target volume cannot perform read or write caching. Read and write caching is failed by the Small Computer System Interface (SCSI) front end as a hardware error. If the mapping is incremental and a previous mapping completed, the mapping records only the differences between the source and target volumes. If the connection to both nodes in the I/O group that the mapping is assigned to is lost, the source and target volumes go offline.

- ▶ Preparing

The target volume is online, but not accessible. The target volume cannot perform read or write caching. Read and write caching is failed by the SCSI front end as a hardware error. Any changed write data for the source volume is flushed from the cache. Any read or write data for the target volume is discarded from the cache. If the mapping is incremental and a previous mapping completed, the mapping records only the differences between the source and target volumes. If the connection to both nodes in the I/O group that the mapping is assigned to is lost, the source and target volumes go offline.

- ▶ Stopped

The mapping is stopped because you issued a stop command or an I/O error occurred. The target volume is offline and its data is lost. To access the target volume, you must restart or delete the mapping. The source volume is accessible and the read and write cache is enabled. If the mapping is incremental, the mapping is recording write operations to the source volume. If the connection to both nodes in the I/O group that the mapping is assigned to is lost, the source and target volumes go offline.

- ▶ Stopping

The mapping is copying data to another mapping. If the background copy process is complete, the target volume is online while the stopping copy process completes. If the background copy process is incomplete, data is discarded from the target volume cache. The target volume is offline while the stopping copy process runs. The source volume is accessible for I/O operations.

- ▶ Suspended

The mapping started, but it did not complete. Access to the metadata is lost, which causes the source and target volume to go offline. When access to the metadata is restored, the mapping returns to the copying or stopping state and the source and target volumes return online. The background copy process resumes. If the data was not flushed and was written to the source or target volume before the suspension, it is in the cache until the mapping leaves the suspended state.

### Summary of FlashCopy mapping states

Table 10-2 lists the various FlashCopy mapping states and the corresponding states of the source and target volumes.

Table 10-2 FlashCopy mapping state summary

State	Source		Target	
	Online/Offline	Cache state	Online/Offline	Cache state
Idling/Copied	Online	Write-back	Online	Write-back
Copying	Online	Write-back	Online	Write-back
Stopped	Online	Write-back	Offline	N/A
Stopping	Online	Write-back	► Online if copy complete ► Offline if copy incomplete	N/A
Suspended	Offline	Write-back	Offline	N/A
Preparing	Online	Write-through	Online but not accessible	N/A
Prepared	Online	Write-through	Online but not accessible	N/A

### 10.1.11 Multiple target FlashCopy

A volume can be the source of multiple target volumes. A target volume can also be the source of another target volume. However, a target volume can have only one source volume. A source volume can have multiple target volumes in one or multiple consistency groups. A consistency group can contain multiple FlashCopy mappings (source-target relations). A source volume can belong to multiple consistency groups. Figure 10-12 on page 510 shows these different possibilities.

Every source-target relation is a FlashCopy mapping and is maintained with its own bitmap table. No consistency group bitmap table exists.

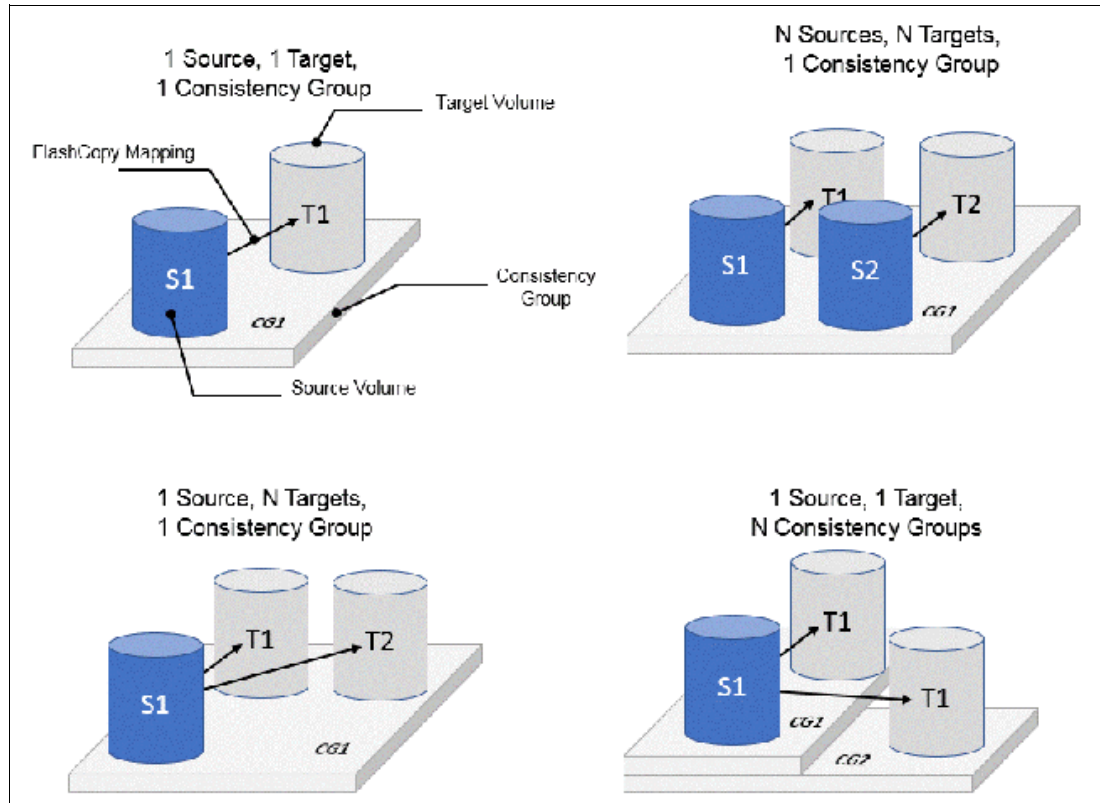


Figure 10-12 Consistency groups and mappings combinations

When a source volume is in a FlashCopy mapping with multiple targets, in multiple consistency groups, it allows the copy of a single source at multiple points in time and therefore, keeps multiple versions of a single volume.

### Consistency group with multiple target FlashCopy

A consistency group aggregates FlashCopy mappings, not volumes. Therefore, where a source volume has multiple FlashCopy mappings, they can be in the same or separate consistency groups.

If a particular volume is the source volume for multiple FlashCopy mappings, you might want to create separate consistency groups to separate each mapping of the same source volume. Regardless of whether the source volume with multiple target volumes is in the same consistency group or in separate consistency groups, the resulting FlashCopy produces multiple identical copies of the source data.

### Dependencies

When a source volume has multiple target volumes, a mapping is created for each source-target relationship. When data is changed on the source volume, it is first copied to the target volume because of the copy-on-write mechanism that is used by FlashCopy.

You can create up to 256 targets for a single source volume. Therefore, a single write operation on the source volume might result in 256 write operations (one per target volume). This configuration generates a large workload that the system cannot handle, which leads to a heavy performance impact on front-end operations.



To avoid any significant effect on performance because of multiple targets, FlashCopy creates dependencies between the targets. Dependencies can be considered as “hidden” FlashCopy mappings that are not visible to and cannot be managed by the user. A dependency is created between the most recent target and the previous one (in order of start time). Figure 10-13 shows an example of a source volume with three targets.

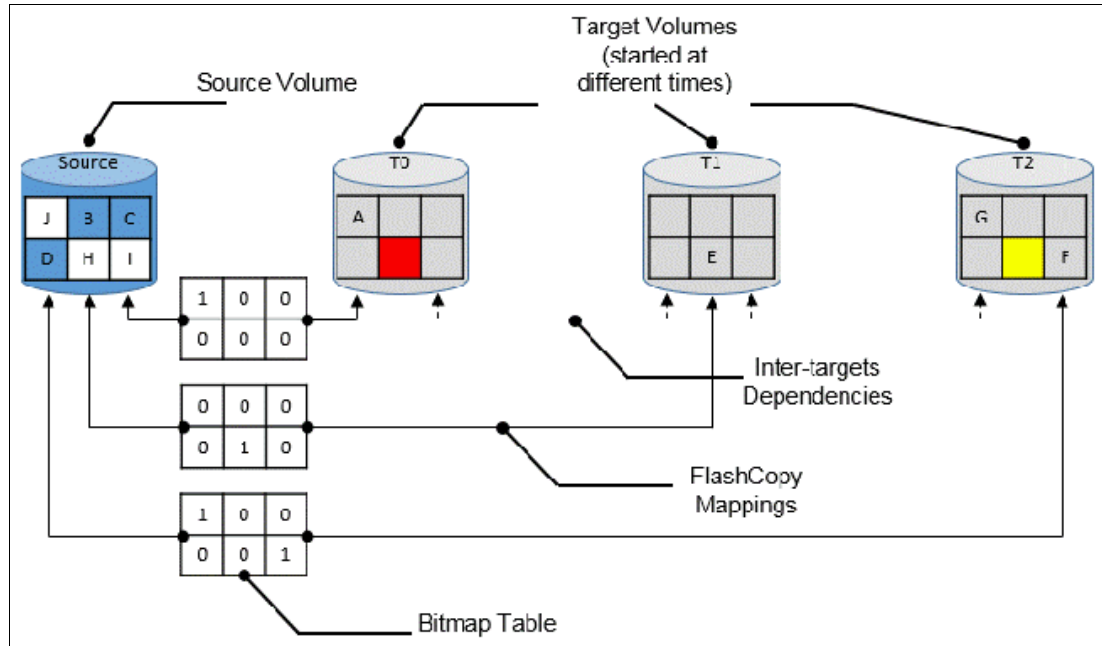


Figure 10-13 FlashCopy dependencies example

When the three targets are started, Target T0 was started first and considered the “oldest.” Target T1 was started next and is considered “next oldest,” and finally, Target T2 was started last and considered the “most recent” or “newest.” The “next oldest” target for T2 is T1. The “next oldest” target for T1 is T0. T1 is newer than T2, and T0 is newer than T1.

### Source read with multiple target FlashCopy

No specific behavior is shown for read operations on source volumes when multiple targets exist for that volume. The data is always read from the source.

### Source write with multiple target FlashCopy (Copy on Write)

A write to the source volume does not cause its data to be copied to all of the targets. Instead, it is copied to the most recent target volume only. For example, consider the sequence of events that are listed in Table 10-3 for a source volume and three targets started at different times. In this example, no background copy exists. The “most recent” target is indicated with an asterisk.

Table 10-3 Sequence example of write IOs on a source with multiple targets

	Source volume	Target T0	Target T1	Target T2
Time 0: mapping with T0 is started	A B C D E F	___* ---	Not started	Not started
Time 1: change of “A” is made on source (->“G”)	G B C D E F	A __* ---	Not started	Not started
Time 2: mapping with T1 is started	G B C D E F	A __ ---	___* ---	Not started

	Source volume	Target T0	Target T1	Target T2
Time 3: change of "E" is made on source (->"H")	G B C D H F	A __ __	__ __* _ E_	Not started
Time 4: mapping with T2 is started	G B C D H F	A __ __	__ __ _ E_	__ __* __
Time 5: change of "F" is made on source (->"I")	G B C D H I	A __ __	__ __ _ E_	__ __* __ F
Time 6: change of "G" is made on source (->"J")	J B C D H I	A __ __	__ __ _ E_	G __* __ F
Time 7: stop of Source-T2 mapping	J B C D H I	A __ __	G __* _ E F	Stopped
Time 8: stop of Source-T1 mapping	J B C D H I	A __* _ E F	Stopped	Stopped
* "most recent" target				

An intermediate target disk (not the oldest or the newest) treats the set of newer target volumes and the true source volume as a type of composite source. It treats all older volumes as a kind of target (and behaves like a source to them).

### **Target read with multiple target FlashCopy**

Target reading with multiple targets depends on whether the grain was copied. Consider the following points:

- ▶ If the grain that is read is copied from the source to the target, the read returns data from the target that is read.
- ▶ If the grain is not yet copied, each of the newer mappings is examined in turn. The read is performed from the first copy (the oldest) that is found. If none is found, the read is performed from the source.

For example, in Figure 10-13 on page 511, if the yellow grain on T2 is read, it returns "H" because no newer target than T2 exists. Therefore, the source is read.

As another example, in Figure 10-13 on page 511, if the red grain on T0 is read, it returns "E" because two newer targets exist for T0, and T1 is the oldest of those targets.

### **Target write with multiple target FlashCopy (Copy on Demand)**

A write to an intermediate or the newest target volume must consider the state of the grain within its own mapping and the state of the grain of the next oldest mapping. Consider the following points:

- ▶ If the grain in the target that is written is copied and if the grain of the next oldest mapping is not yet copied, the grain must be copied before the write can proceed to preserve the contents of the next oldest mapping.

For example, in Figure 10-13 on page 511, if the grain "G" is changed on T2, it must be copied to T1 (next oldest not yet copied) first and then changed on T2.

- ▶ If the grain in the target that is being written is not yet copied, the grain is copied from the oldest copied grain in the mappings that are newer than the target, or from the source if none is copied. For example, in Figure 10-13 on page 511, if the red grain on T0 is written, it is first copied from T1 (data "E"). After this copy is done, the write can be applied to the target.

Table 10-4 lists the indirection layer algorithm in a multi-target FlashCopy.

Table 10-4 Summary table of the FlashCopy indirection layer algorithm

Accessed volume	Was the grain copied?	Host I/O operation	
		Read	Write
Source	No	Read from the source volume.	Copy grain to most recently started target for this source, then write to the source.
	Yes	Read from the source volume.	Write to the source volume.
Target	No	If any newer targets exist for this source in which this grain was copied, read from the oldest of these targets. Otherwise, read from the source.	Hold the write. Check the dependency target volumes to see whether the grain was copied. If the grain is not copied to the next oldest target for this source, copy the grain to the next oldest target. Then, write to the target.
	Yes	Read from the target volume.	Write to the target volume.

### Stopping process in a multiple target FlashCopy: Cleaning Mode

When a mapping that contains a target that includes dependent mappings is stopped, the mapping enters the stopping state. It then begins copying all grains that are uniquely held on the target volume of the mapping that is being stopped to the next oldest mapping that is in the copying state. The mapping remains in the stopping state until all grains are copied, and then enters the stopped state. This mode is referred to as the *Cleaning Mode*.

For example, if the mapping Source-T2 was stopped, the mapping enters the stopping state while the cleaning process copies the data of T2 to T1 (next oldest). After all of the data is copied, Source-T2 mapping enters the stopped state, and T1 is no longer dependent upon T2. However, T0 remains dependent upon T1.

For example, as shown in Table 10-3 on page 511, if you stop the Source-T2 mapping on “Time 7,” then the grains that are not yet copied on T1 are copied from T2 to T1. Reading T1 is then like reading the source at the time T1 was started (“Time 2”).

As another example, with Table 10-3 on page 511, if you stop the Source-T1 mapping on “Time 8,” the grains that are not yet copied on T0 are copied from T1 to T0. Reading T0 is then similar to reading the source at the time T0 was started (“Time 0”).

If you stop the Source-T1 mapping while Source-T0 mapping and Source-T2 are still in copying mode, the grains that are not yet copied on T0 are copied from T1 to T0 (next oldest). T0 now depends upon T2.

Your target volume is still accessible while the cleaning process is running. When the system is operating in this mode, it is possible that host I/O operations can prevent the cleaning process from reaching 100% if the I/O operations continue to copy new data to the target volumes.

### ***Cleaning rate***

The data rate at which data is copied from the target of the mapping being stopped to the next oldest target is determined by the *cleaning rate*. This property of FlashCopy mapping can be changed dynamically. It is measured as is the copyrate property, but both properties are independent. Table 10-5 lists the relationship of the cleaning rate values to the attempted number of grains to be split per second.

Table 10-5 *Cleaning rate values*

<b>User-specified copy rate attribute value</b>	<b>Data copied/sec</b>	<b>256 KB grains/sec</b>	<b>64 KB grains/sec</b>
1 - 10	128 KiB	0.5	2
11 - 20	256 KiB	1	4
21 - 30	512 KiB	2	8
31 - 40	1 MiB	4	16
41 - 50	2 MiB	8	32
51 - 60	4 MiB	16	64
61 - 70	8 MiB	32	128
71 - 80	16 MiB	64	256
81 - 90	32 MiB	128	512
91 - 100	64 MiB	256	1024
101 - 110	128 MiB	512	2048
111 - 120	256 MiB	1024	4096
121 - 130	512 MiB	2048	8192
131 - 140	1 GiB	4096	16384
141 - 150	2 GiB	8192	32768

### **10.1.12 Reverse FlashCopy**

Reverse FlashCopy enables FlashCopy targets to become restore points for the source without breaking the FlashCopy mapping, and without having to wait for the original copy operation to complete. A FlashCopy source supports multiple targets (up to 256), and therefore, multiple rollback points.

A key advantage of the IBM Spectrum Virtualize Multiple Target Reverse FlashCopy function is that the reverse FlashCopy does not destroy the original target. This feature enables processes that use the target, such as a tape backup or tests, to continue uninterrupted.

IBM Spectrum Virtualize also can create an optional copy of the source volume to be made before the reverse copy operation starts. This ability to restore back to the original source data can be useful for diagnostic purposes.

The production disk is instantly available with the backup data. Figure 10-14 shows an example of Reverse FlashCopy with a simple FlashCopy mapping (single target).

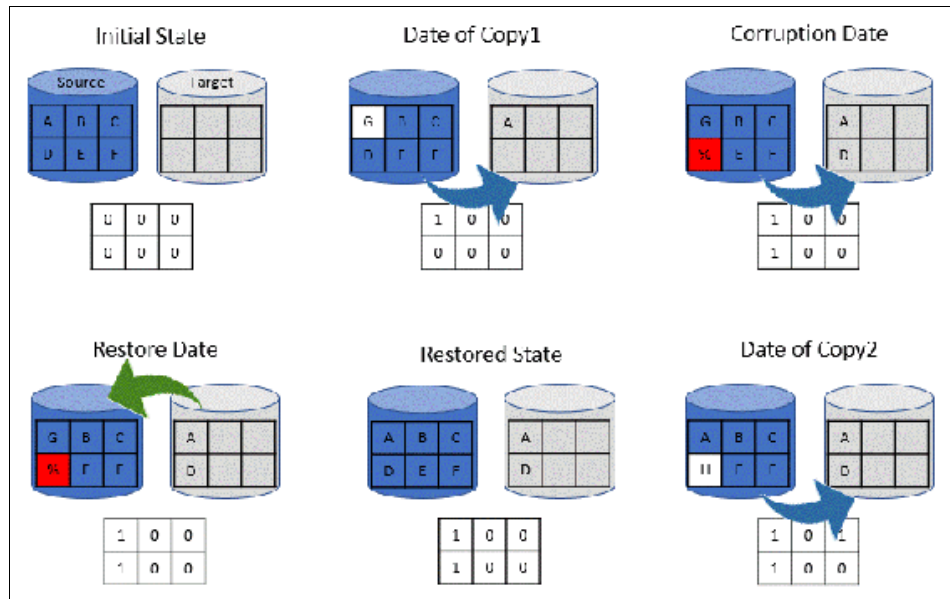


Figure 10-14 A reverse FlashCopy example for data restoration

This example assumes that a simple FlashCopy mapping was created between the “source” volume and “target” volume, and no background copy is set.

When the FlashCopy mapping starts (Date of Copy1), if source volume is changed (write operations on grain “A”), the modified grains are first copied to target, the bitmap table is updated, and the source grain is modified (from “A” to “G”).

At a specific time (“Corruption Date”), data is modified on another grain (grain “D” below), so it is first written on the target volume and the bitmap table is updated. Unfortunately, the new data is corrupted on source volume.

The storage administrator can then use the Reverse FlashCopy feature by completing the following steps:

1. Create a mapping from target to source (if not already created). Because FlashCopy recognizes that the target volume of this new mapping is a source in another mapping, it does not create another bitmap table. It uses the existing bitmap table instead, with its updated bits.
2. Start the new mapping. Because of the existing bitmap table, only the *modified* grains are copied.

After the restoration is complete, at the “Restored State” time, source volume data is similar to what it was before the Corruption Date. The copy can resume with the restored data (Date of Copy2) and for example, data on the source volume can be modified (“D” grain is changed in “H” grain in the example below). In this last case, because “D” grain was copied, it is not copied again on target volume.

Consistency groups are reversed by creating a set of reverse FlashCopy mappings and adding them to a new reverse consistency group. Consistency groups cannot contain more than one FlashCopy mapping with the same target volume.

### 10.1.13 FlashCopy and image mode Volumes

FlashCopy can be used with image mode volumes. Because the source and target volumes must be the same size, you must create a target volume with the same size as the image mode volume when you are creating a FlashCopy mapping. To accomplish this task by using the CLI, run the `svcinfo lsvdisk -bytes volumename` command. The size in bytes is then used to create the volume that is used in the FlashCopy mapping.

This method provides an exact number of bytes because image mode volumes might not line up one-to-one on other measurement unit boundaries. Example 10-1 shows the size of the ITS0-RS-TST volume. The ITS0-TST01 volume is then created, which specifies the same size.

*Example 10-1 Listing the size of a volume in bytes and creating a volume of equal size*

---

```
IBM_2145:ITS0-SV1:superuser>lsvdisk -bytes ITS0-RS-TST
id 42
name ITS0-RS-TST
IO_group_id 0
IO_group_name io_grp0
status online
mdisk_grp_id 0
mdisk_grp_name Pool0
capacity 21474836480
type striped
formatted no
formatting yes
mdisk_id
mdisk_name
FC_id
.....

IBM_2145:ITS0-SV1:superuser>mkvdisk -mdiskgrp Pool0 -iogrp 0 -size 21474836480
-unit b -name ITS0-TST01
Virtual Disk, id [43], successfully created
IBM_2145:ITS0-SV1:superuser>

IBM_2145:ITS0-SV1:superuser>lsvdisk -delim " "
42 ITS0-RS-TST 0 io_grp0 online 0 Pool0 20.00GB striped
600507680C9B8000480000000000002C 0 1 not_empty 0 no 0 0 Pool0 yes no 42
ITS0-RS-TST
43 ITS0-TST01 0 io_grp0 online 0 Pool0 20.00GB image
600507680C9B8000480000000000002D 0 1 not_empty 0 no 0 0 Pool0 yes no 43 ITS0-TST01
IBM_2145:ITS0-SV1:superuser>
```

---

**Tip:** Alternatively, you can run the `expandvdisksize` and `shrinkvdisksize` volume commands to modify the size of the volume.

These actions must be performed before a mapping is created.

## 10.1.14 FlashCopy mapping events

This section describes the events that modify the states of a FlashCopy. It also describes the mapping events that are listed in Table 10-6.

**Overview of a FlashCopy sequence of events:** The FlashCopy sequence includes the following tasks:

1. Associate the source data set with a target location (one or more source and target volumes).
2. Create a FlashCopy mapping for each source volume to the corresponding target volume. The target volume must be equal in size to the source volume.
3. Discontinue access to the target (application dependent).
4. Prepare (pre-trigger) the FlashCopy:
  - a. Flush the cache for the source.
  - b. Discard the cache for the target.
5. Start (trigger) the FlashCopy:
  - a. Pause I/O (briefly) on the source.
  - b. Resume I/O on the source.
  - c. Start I/O on the target.

Table 10-6 Mapping events

Mapping event	Description
Create	<p>A FlashCopy mapping is created between the specified source volume and the specified target volume. The operation fails if any one of the following conditions is true:</p> <ul style="list-style-type: none"> <li>▶ The source volume is a member of 256 FlashCopy mappings.</li> <li>▶ The node has insufficient bitmap memory.</li> <li>▶ The source and target volumes are different sizes.</li> </ul>
Prepare	<p>The <b>prestartfcmap</b> or <b>prestartfcconsistgrp</b> command is directed to a consistency group for FlashCopy mappings that are members of a normal consistency group or to the mapping name for FlashCopy mappings that are stand-alone mappings. The <b>prestartfcmap</b> or <b>prestartfcconsistgrp</b> command places the FlashCopy mapping into the Preparing state.</p> <p>The <b>prestartfcmap</b> or <b>prestartfcconsistgrp</b> command can corrupt any data that was on the target volume because cached writes are discarded. Even if the FlashCopy mapping is never started, the data from the target might be changed logically during the act of preparing to start the FlashCopy mapping.</p>
Flush done	<p>The FlashCopy mapping automatically moves from the preparing state to the prepared state after all cached data for the source is flushed and all cached data for the target is no longer valid.</p>

Mapping event	Description
Start	<p>When all of the FlashCopy mappings in a consistency group are in the prepared state, the FlashCopy mappings can be started. To preserve the cross-volume consistency group, the start of all of the FlashCopy mappings in the consistency group must be synchronized correctly concerning I/Os that are directed at the volumes by running the <b>startfcmap</b> or <b>startfcconsistgrp</b> command.</p> <p>The following actions occur during the running of the <b>startfcmap</b> command or the <b>startfcconsistgrp</b> command:</p> <ul style="list-style-type: none"> <li>▶ New reads and writes to all source volumes in the consistency group are paused in the cache layer until all ongoing reads and writes beneath the cache layer are completed.</li> <li>▶ After all FlashCopy mappings in the consistency group are paused, the internal cluster state is set to enable FlashCopy operations.</li> <li>▶ After the cluster state is set for all FlashCopy mappings in the consistency group, read and write operations continue on the source volumes.</li> <li>▶ The target volumes are brought online.</li> </ul> <p>As part of the <b>startfcmap</b> or <b>startfcconsistgrp</b> command, read and write caching is enabled for the source and target volumes.</p>
Modify	<p>The following FlashCopy mapping properties can be modified:</p> <ul style="list-style-type: none"> <li>▶ FlashCopy mapping name</li> <li>▶ Clean rate</li> <li>▶ Consistency group</li> <li>▶ Copy rate (for background copy or stopping copy priority)</li> <li>▶ Automatic deletion of the mapping when the background copy is complete</li> </ul>
Stop	<p>The following separate mechanisms can be used to stop a FlashCopy mapping:</p> <ul style="list-style-type: none"> <li>▶ Issue a command</li> <li>▶ An I/O error occurred</li> </ul>
Delete	<p>This command requests that the specified FlashCopy mapping is deleted. If the FlashCopy mapping is in the copying state, the <b>force</b> flag must be used.</p>
Flush failed	<p>If the flush of data from the cache cannot be completed, the FlashCopy mapping enters the stopped state.</p>
Copy complete	<p>After all of the source data is copied to the target and there are no dependent mappings, the state is set to copied. If the option to automatically delete the mapping after the background copy completes is specified, the FlashCopy mapping is deleted automatically. If this option is not specified, the FlashCopy mapping is not deleted automatically and can be reactivated by preparing and starting again.</p>
Bitmap online/offline	<p>The node failed.</p>



## 10.1.15 Thin-provisioned FlashCopy

FlashCopy source and target volumes can be thin-provisioned.

### Source or target thin-provisioned

The most common configuration is a fully allocated source and a thin-provisioned target. By using this configuration, the target uses a smaller amount of real storage than the source.

With this configuration, use a copyrate equal to 0 only. In this state, the virtual capacity of the target volume is identical to the capacity of the source volume, but the real capacity (the one used on the storage system) is lower, as shown on Figure 10-15. The real size of the target volume increases with writes that are performed on the source volume, on not already copied grains. Eventually, if the entire source volume is written (unlikely), the real capacity of the target volume is identical to the source's volume.

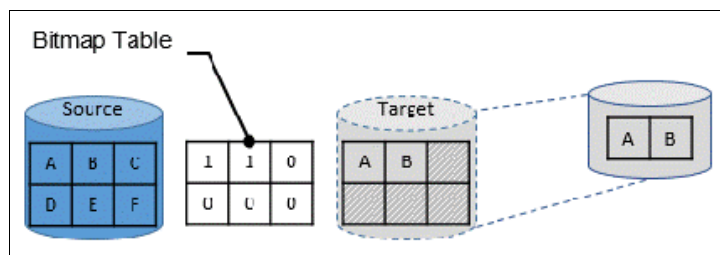


Figure 10-15 Thin-provisioned target volume

### Source and target thin-provisioned

When the source and target volumes are thin-provisioned, only the data that is allocated to the source is copied to the target. In this configuration, the background copy option has no effect.

**Performance:** The best performance is obtained when the grain size of the thin-provisioned volume is the same as the grain size of the FlashCopy mapping.

### Thin-provisioned incremental FlashCopy

The implementation of thin-provisioned volumes does not preclude the use of incremental FlashCopy on the same volumes. It does not make sense to have a fully allocated source volume and then use incremental FlashCopy (which is always a full copy the first time) to copy this fully allocated source volume to a thin-provisioned target volume. However, this action is not prohibited.

Consider the following optional configurations:

- ▶ A thin-provisioned source volume can be copied incrementally by using FlashCopy to a thin-provisioned target volume. Whenever the FlashCopy is performed, only data that was modified is recopied to the target. If space is allocated on the target because of I/O to the target volume, this space is not reclaimed with subsequent FlashCopy operations.
- ▶ A fully allocated source volume can be copied incrementally by using FlashCopy to another fully allocated volume at the same time as it is being copied to multiple thin-provisioned targets (taken at separate points in time). By using this combination, a single full backup can be kept for recovery purposes, and the backup workload is separated from the production workload. At the same time, older thin-provisioned backups can be retained.

## 10.1.16 Serialization of I/O by FlashCopy

In general, the FlashCopy function in the IBM Spectrum Virtualize introduces no explicit serialization into the I/O path. Therefore, many concurrent I/Os are allowed to the source and target volumes.

However, a lock exists for each grain and this lock can be in shared or exclusive mode. For multiple targets, a common lock is shared, and the mappings are derived from a particular source volume. The lock is used in the following modes under the following conditions:

- ▶ The lock is held in shared mode during a read from the target volume, which touches a grain that was not copied from the source.
- ▶ The lock is held in exclusive mode while a grain is being copied from the source to the target.

If the lock is held in shared mode and another process wants to use the lock in shared mode, this request is granted unless a process is waiting to use the lock in exclusive mode.

If the lock is held in shared mode and it is requested to be exclusive, the requesting process must wait until all holders of the shared lock free it.

Similarly, if the lock is held in exclusive mode, a process wanting to use the lock in shared or exclusive mode must wait for it to be freed.

## 10.1.17 Event handling

When a FlashCopy mapping is not copying or stopping, the FlashCopy function does not affect the handling or reporting of events for error conditions that are encountered in the I/O path. Event handling and reporting are affected only by FlashCopy when a FlashCopy mapping is copying or stopping; that is, actively moving data.

These scenarios are described next,

### Node failure

Normally, two copies of the FlashCopy bitmap are maintained. One copy of the FlashCopy bitmap is on each of the two nodes that make up the I/O group of the source volume. When a node fails, one copy of the bitmap for all FlashCopy mappings whose source volume is a member of the failing node's I/O group becomes inaccessible.

FlashCopy continues with a single copy of the FlashCopy bitmap that is stored as non-volatile in the remaining node in the source I/O group. The system metadata is updated to indicate that the missing node no longer holds a current bitmap. When the failing node recovers or a replacement node is added to the I/O group, the bitmap redundancy is restored.

### Path failure (Path Offline state)

In a fully functioning system, all of the nodes have a software representation of every volume in the system within their application hierarchy.

Because the storage area network (SAN) that links IBM SAN Volume Controller nodes to each other and to the MDisks is made up of many independent links, it is possible for a subset of the nodes to be temporarily isolated from several of the MDisks. When this situation occurs, the MDisks are said to be *Path Offline* on certain nodes.

**Other nodes:** Other nodes might see the MDisks as Online because their connection to the MDisks still exists.

### ***Path Offline for the source Volume***

If a FlashCopy mapping is in the copying state and the source volume goes path offline, this path offline state is propagated to all target volumes up to, but not including, the target volume for the newest mapping that is 100% copied but remains in the copying state. If no mappings are 100% copied, all of the target volumes are taken offline. Path offline is a state that exists on a per-node basis. Other nodes might not be affected. If the source volume comes online, the target and source volumes are brought back online.

### ***Path Offline for the target Volume***

If a target volume goes path offline but the source volume is still online and if any dependent mappings exist, those target volumes also go path offline. The source volume remains online.

## **10.1.18 Asynchronous notifications**

FlashCopy raises informational event log entries for certain mapping and consistency group state transitions. These state transitions occur as a result of configuration events that complete asynchronously. The informational events can be used to generate Simple Network Management Protocol (SNMP) traps to notify the user.

Other configuration events complete synchronously, and no informational events are logged as a result of the following events:

▶ **PREPARE\_COMPLETED**

This state transition is logged when the FlashCopy mapping or consistency group enters the prepared state as a result of a user request to prepare. The user can now start (or stop) the mapping or consistency group.

▶ **COPY\_COMPLETED**

This state transition is logged when the FlashCopy mapping or consistency group enters the idle\_or\_copied state when it was in the copying or stopping state. This state transition indicates that the target disk now contains a complete copy and no longer depends on the source.

▶ **STOP\_COMPLETED**

This state transition is logged when the FlashCopy mapping or consistency group enters the stopped state as a result of a user request to stop. It is logged after the automatic copy process completes. This state transition includes mappings where no copying needed to be performed. This state transition differs from the event that is logged when a mapping or group enters the stopped state as a result of an I/O error.

## **10.1.19 Interoperation with Metro Mirror and Global Mirror**

A volume can be part of any copy relationship; that is, FlashCopy, Metro Mirror (MM)], or Remote Mirror. Therefore, FlashCopy can work with MM/Global Mirror (GM) to provide better protection of the data.

For example, you can perform an MM copy to duplicate data from Site\_A to Site\_B, and then perform a daily FlashCopy to back up the data to another location.

**Note:** A volume cannot be part of FlashCopy, MM, or Remote Mirror, if it is set to Transparent Cloud Tiering (TCT) function.

Table 10-7 lists the supported combinations of FlashCopy and Remote Copy (RC). In the table, *RC* refers to MM and GM.

Table 10-7 *FlashCopy and remote copy interaction*

Component	RC primary site	RC secondary site
FlashCopy Source	Supported	Supported latency: When the FlashCopy relationship is in the preparing and prepared states, the cache at the RC secondary site operates in write-through mode. This process adds latency to the latent RC relationship.
FlashCopy Target	This is a supported combination and has the following restrictions: <ul style="list-style-type: none"> <li>▶ Running a <b>stop -force</b> might cause the RC relationship to be fully resynchronized.</li> <li>▶ Code level must be 6.2.x or later.</li> <li>▶ I/O group must be the same.</li> </ul>	This is a supported combination with the major restriction that the FlashCopy mapping cannot be copying, stopping, or suspended. Otherwise, the restrictions are the same as at the RC primary site.

## 10.1.20 FlashCopy attributes and limitations

The FlashCopy function in IBM Spectrum Virtualize features the following attributes:

- ▶ The target is the T0 copy of the source, which is known as *FlashCopy mapping target*.
- ▶ FlashCopy produces an exact copy of the source volume, including any metadata that was written by the host operating system, Logical Volume Manager (LVM), and applications.
- ▶ The source volume and target volume are available (almost) immediately following the FlashCopy operation.
- ▶ The source and target volumes:
  - Must be the same “virtual” size
  - Must be on the same SAN Volume Controller system
  - Do not need to be in the same I/O group or storage pool
- ▶ The storage pool extent sizes can differ between the source and target.
- ▶ The target volumes can be the source volumes for other FlashCopy mappings (*cascaded FlashCopy*). However, a target volume can have only one source copy.
- ▶ Consistency groups are supported to enable FlashCopy across multiple volumes at the same time.
- ▶ The target volume can be updated independently of the source volume.
- ▶ Bitmaps that are governing I/O redirection (I/O indirection layer) are maintained in both nodes of the IBM SAN Volume Controller I/O group to prevent a single point of failure (SPOF).

- ▶ FlashCopy mapping and consistency groups can be automatically withdrawn after the completion of the background copy.
- ▶ Thin-provisioned FlashCopy (or Snapshot in the GUI) use disk space only when updates are made to the source or target data, and not for the entire capacity of a volume copy.
- ▶ FlashCopy licensing is based on the virtual capacity of the source volumes.
- ▶ Incremental FlashCopy copies all of the data when you first start FlashCopy, and then only the changes when you stop and start FlashCopy mapping again. Incremental FlashCopy can substantially reduce the time that is required to re-create an independent image.
- ▶ Reverse FlashCopy enables FlashCopy targets to become restore points for the source without breaking the FlashCopy relationship, and without having to wait for the original copy operation to complete.
- ▶ The size of the source and target volumes cannot be altered (increased or decreased) while a FlashCopy mapping is defined.

IBM FlashCopy limitations for IBM Spectrum Virtualize V8.3 are listed in Table 10-8.

*Table 10-8 FlashCopy limitations in V8.3*

Property	Maximum number
FlashCopy mappings per system	5000
FlashCopy targets per source	256
FlashCopy mappings per consistency group	512
FlashCopy consistency groups per system	500
Total FlashCopy volume capacity per I/O group	4096 TiB

## 10.2 Managing FlashCopy by using the GUI

It is often easier to work with the FlashCopy function from the GUI if you have a reasonable number of host mappings. However, in enterprise data centers with many host mappings, use the CLI to run your FlashCopy commands.

### 10.2.1 FlashCopy presets

The IBM Spectrum Virtualize GUI interface provides three FlashCopy presets (Snapshot, Clone, and Backup) to simplify the more common FlashCopy operations.

Although these presets meet most FlashCopy requirements, they do not support all possible FlashCopy options. If more specialized options are required that are not supported by the presets, the options must be performed by using CLI commands.

This section describes the preset options and their use cases.

## Snapshot

This preset creates a copy-on-write PiT copy. The snapshot is not intended to be an independent copy. Instead, the copy is used to maintain a view of the production data at the time that the snapshot is created. Therefore, the snapshot holds only the data from regions of the production volume that changed since the snapshot was created. Because the snapshot preset uses thin provisioning, only the capacity that is required for the changes is used.

Snapshot uses the following preset parameters:

- ▶ Background copy: None
- ▶ Incremental: No
- ▶ Delete after completion: No
- ▶ Cleaning rate: No
- ▶ Primary copy source pool: Target pool

### Use case

The user wants to produce a copy of a volume without affecting the availability of the volume. The user does not anticipate many changes to be made to the source or target volume; a significant proportion of the volumes remains unchanged.

By ensuring that only changes require a copy of data to be made, the total amount of disk space that is required for the copy is reduced. Therefore, many Snapshot copies can be used in the environment.

Snapshots are useful for providing protection against corruption or similar issues with the validity of the data, but they do not provide protection from physical controller failures. Snapshots can also provide a vehicle for performing repeatable testing (including “what-if” modeling that is based on production data) without requiring a full copy of the data to be provisioned.

For example, in Figure 10-16, the source volume user can still work on the original data volume (as with a production volume) and the target volumes can be accessed instantly. Users of target volumes can modify the content and perform “what-if” tests; for example, versioning. Storage administrators do not need to perform full copies of a volume for temporary tests. However, the target volumes must remain linked to the source. When the link is broken (FlashCopy mapping stopped or deleted), the target volumes become unusable.

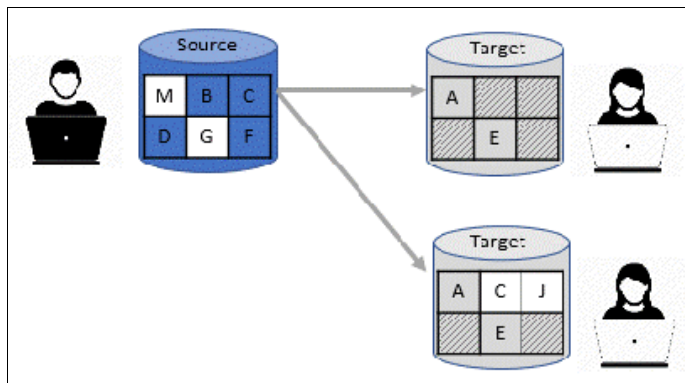


Figure 10-16 FlashCopy snapshot preset example

## Clone

The clone preset creates a replica of the volume, which can be changed without affecting the original volume. After the copy completes, the mapping that was created by the preset is automatically deleted.

Clone uses the following preset parameters:

- ▶ Background copy rate: 50
- ▶ Incremental: No
- ▶ Delete after completion: Yes
- ▶ Cleaning rate: 50
- ▶ Primary copy source pool: Target pool

#### ***Use case***

Users want a copy of the volume that they can modify without affecting the original volume. After the clone is established, it is not expected that it is refreshed or that the original production data must be referenced again. If the source is thin-provisioned, the target is thin-provisioned for the auto-create target.

## **Backup**

The backup preset creates an incremental PiT replica of the production data. After the copy completes, the backup view can be refreshed from the production data, with minimal copying of data from the production volume to the backup volume.

Backup uses the following preset parameters:

- ▶ Background Copy rate: 50
- ▶ Incremental: Yes
- ▶ Delete after completion: No
- ▶ Cleaning rate: 50
- ▶ Primary copy source pool: Target pool

#### ***Use case***

The user wants to create a copy of the volume that can be used as a backup if the source becomes unavailable, such as because of loss of the underlying physical controller. The user plans to periodically update the secondary copy, and does not want to suffer from the resource demands of creating a copy each time.

Incremental FlashCopy times are faster than full copy, which helps to reduce the window where the new backup is not yet fully effective. If the source is thin-provisioned, the target is also thin-provisioned in this option for the auto-create target.

Another use case, which is not supported by the name, is to create and maintain (periodically refresh) an independent image that can be subjected to intensive I/O (for example, data mining) without affecting the source volume's performance.

**Note:** IBM Spectrum Virtualize in general and FlashCopy in particular are not backup solutions on their own. For example, FlashCopy backup preset does not schedule a regular copy of your volumes. Instead, it over-writes the mapping target and does not make a copy of it before starting a new “backup” operation. It is the user's responsibility to handle the target volumes (for example, saving them to tapes) and the scheduling of the FlashCopy operations.

## 10.2.2 FlashCopy window

This section describes the tasks that you can perform at a FlashCopy level by using the IBM Spectrum Virtualize GUI.

When the IBM Spectrum Virtualize GUI is used, FlashCopy components can be seen in different windows. Three windows are related to FlashCopy and are available by using the **Copy Services** menu, as shown in Figure 10-17.

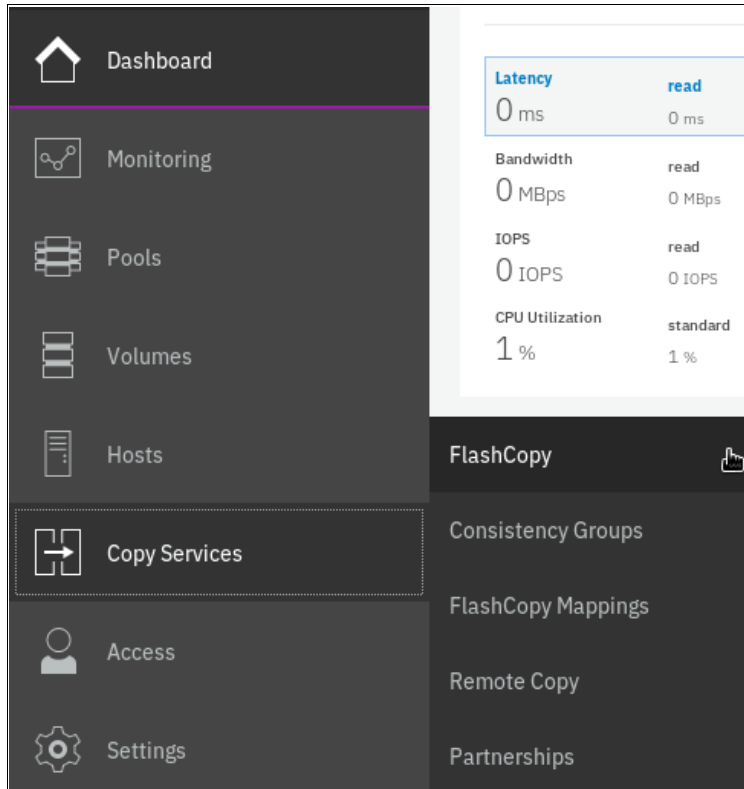


Figure 10-17 Copy Services menu

The FlashCopy window is accessible by clicking **Copy Services** → **FlashCopy**. It displays all of the volumes that are defined in the system. Volumes that are part of a FlashCopy mapping appear, as shown in Figure 10-18. By clicking a source volume, you can display the list of its target volumes.

Volume Name	Status	Progress	Capacity	Group	Flash Time
ITSO-FC-VOL-01			10.00 GiB		
ITSO-FC-VOL-01_03	✓ Copied	100%			Oct 22, 2019, 3:21:06 PM
ITSO-FC-VOL-01_05	⌛ Copying	0%			Oct 22, 2019, 3:21:21 PM
ITSO-FC-VOL-01_04	⌛ Copying	0%			Oct 22, 2019, 3:21:16 PM
ITSO-FC-VOL-01_01	✓ Copied	100%		fccstgrp1	Oct 18, 2019, 2:20:11 PM

Figure 10-18 Source and target volumes displayed in the FlashCopy window

All volumes are listed in this window, and target volumes appear twice (as a regular volume and as a target volume in a FlashCopy mapping).



Consider the following points:

- ▶ The Consistency Group window is accessible by clicking **Copy Services** → **Consistency Groups**. Use the Consistency Groups window (as shown in Figure 10-19) to list the FlashCopy mappings that are part of consistency groups and part of no consistency groups.

Mapping Name	↑	Status	Source Volume	Target Volume	Progress	Flash Tim
▼ Not in a Group						
fcmap6		✓ Copied	ITSO-FC-VOL-01	ITSO-FC-VOL-01_03	100%	Oct 22,
fcmap7		⌂ Copying	ITSO-FC-VOL-01	ITSO-FC-VOL-01_04	0%	Oct 22,
fcmap8		⌂ Copying	ITSO-FC-VOL-01	ITSO-FC-VOL-01_05	0%	Oct 22,
▼ fccstgrp1 Idle or Copied						
fcmap0		✓ Copied	ITSO-FC-VOL-01	ITSO-FC-VOL-01_01	100%	Oct 18,
fcmap1		✓ Copied	ITSO-FC-VOL-04	ITSO-FC-VOL-04_01	100%	Oct 18,
fcmap2		✓ Copied	ITSO-FC-VOL-03	ITSO-FC-VOL-03_01	100%	Oct 18,
fcmap3		✓ Copied	ITSO-FC-VOL-05	ITSO-FC-VOL-05_01	100%	Oct 18,
fcmap4		✓ Copied	ITSO-FC-VOL-02	ITSO-FC-VOL-02_01	100%	Oct 18,

Figure 10-19 Consistency Groups window

- ▶ The FlashCopy Mappings window is accessible by clicking **Copy Services** → **FlashCopy Mappings**. Use the FlashCopy Mappings window (as shown in Figure 10-20) to display the list of mappings between source volumes and target volumes.

Mapping Name	↑	Status	Source Volume	Target Volume	Progress	Group	Flash Tim
fcmap0		✓ Copied	ITSO-FC-VOL-01	ITSO-FC-VOL-01_01	100%	fccstgrp1	Oct 18, 2019, 2:
fcmap1		✓ Copied	ITSO-FC-VOL-04	ITSO-FC-VOL-04_01	100%	fccstgrp1	Oct 18, 2019, 2:
fcmap2		✓ Copied	ITSO-FC-VOL-03	ITSO-FC-VOL-03_01	100%	fccstgrp1	Oct 18, 2019, 2:
fcmap3		✓ Copied	ITSO-FC-VOL-05	ITSO-FC-VOL-05_01	100%	fccstgrp1	Oct 18, 2019, 2:
fcmap4		✓ Copied	ITSO-FC-VOL-02	ITSO-FC-VOL-02_01	100%	fccstgrp1	Oct 18, 2019, 2:
fcmap6		✓ Copied	ITSO-FC-VOL-01	ITSO-FC-VOL-01_03	100%		Oct 22, 2019, 3:
fcmap7		⌂ Copying	ITSO-FC-VOL-01	ITSO-FC-VOL-01_04	0%		Oct 22, 2019, 3:
fcmap8		⌂ Copying	ITSO-FC-VOL-01	ITSO-FC-VOL-01_05	0%		Oct 22, 2019, 3:

Showing 8 FC mappings | Selecting 0 FC mappings

Figure 10-20 FlashCopy mapping window

## 10.2.3 Creating a FlashCopy mapping

This section describes creating FlashCopy mappings for volumes and their targets.

Open the FlashCopy window from the **Copy Services** menu, as shown in Figure 10-21. Select the volume for which you want to create the FlashCopy mapping. Right-click the volume or click the **Actions** menu.

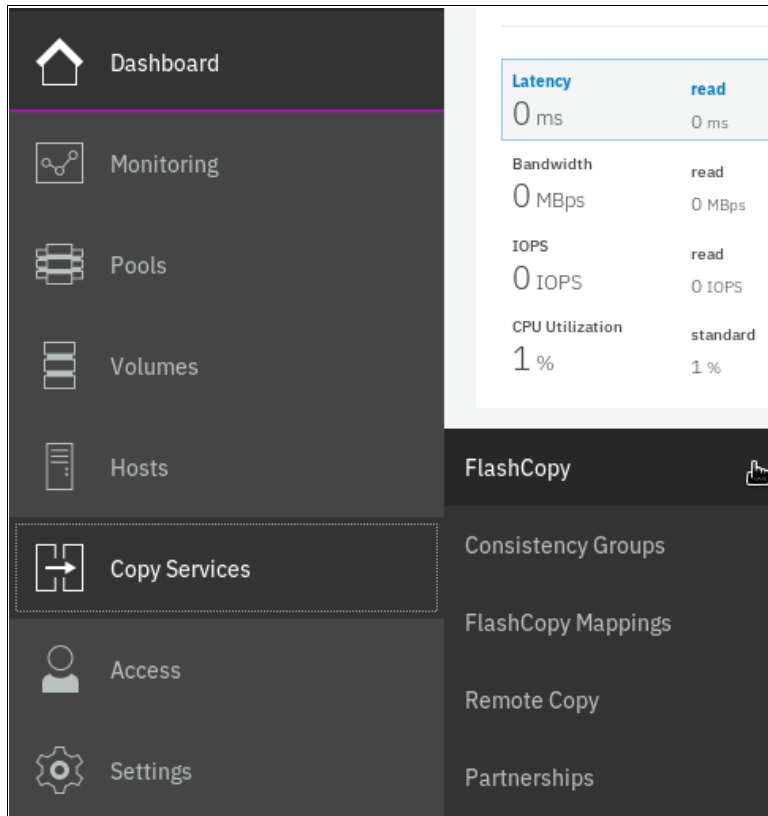


Figure 10-21 FlashCopy window

**Multiple FlashCopy mappings:** To create multiple FlashCopy mappings at the same time, select multiple volumes by pressing and holding **Ctrl** and clicking the entries that you want.

Depending on whether you created the target volumes for your FlashCopy mappings or you want the system to create the target volumes for you, the following options are available:

- ▶ If you created the target volumes, see “Creating a FlashCopy mapping with existing target Volumes” on page 529.
- ▶ If you want the system to create the target volumes for you, see “Creating a FlashCopy mapping and target volumes” on page 534.

## Creating a FlashCopy mapping with existing target Volumes

Complete the following steps to use existing target volumes for the FlashCopy mappings:

**Attention:** When starting a FlashCopy mapping from a source volume to a target volume, data that is on the target is over-written. The system does not prevent you from selecting a target volume that is mapped to a host and contains data.

1. Right-click the volume that you want to create a FlashCopy mapping for, and select **Advanced FlashCopy** → **Use Existing Target Volumes**, as shown in Figure 10-22.

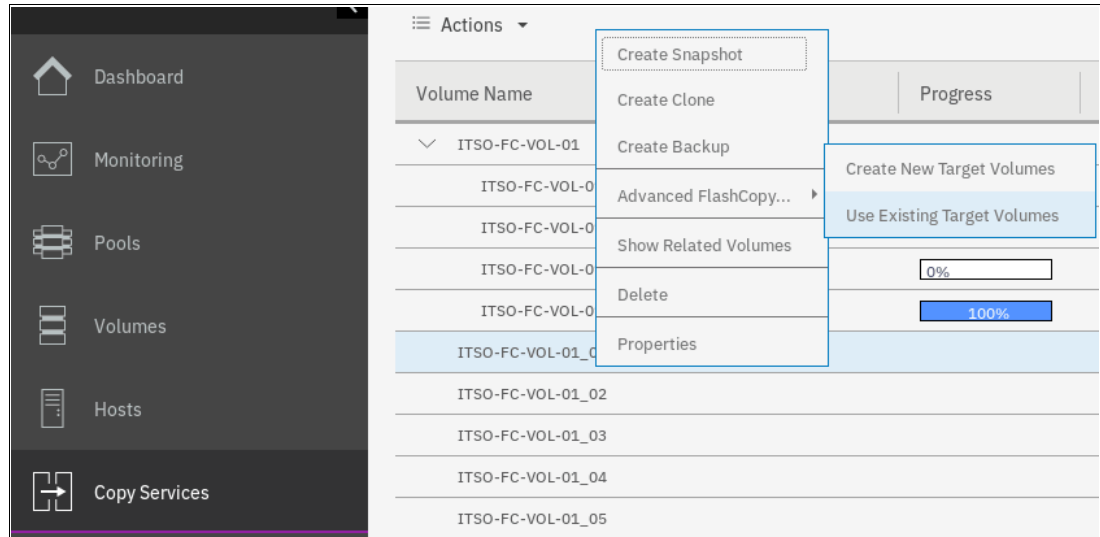


Figure 10-22 Creating a FlashCopy mapping with an existing target

The Create FlashCopy Mapping window opens, as shown in Figure 10-23 on page 530. In this window, you create the mapping between the selected source volume and the target volume you want to create a mapping with. Then, click **Add**.

**Important:** The source volume and the target volume must be of equal size. Therefore, only targets of the same size are shown in the list for a source volume.

Volumes that are a target in a FlashCopy mapping cannot be a target in a new mapping. Therefore, only volumes that are not targets can be selected.

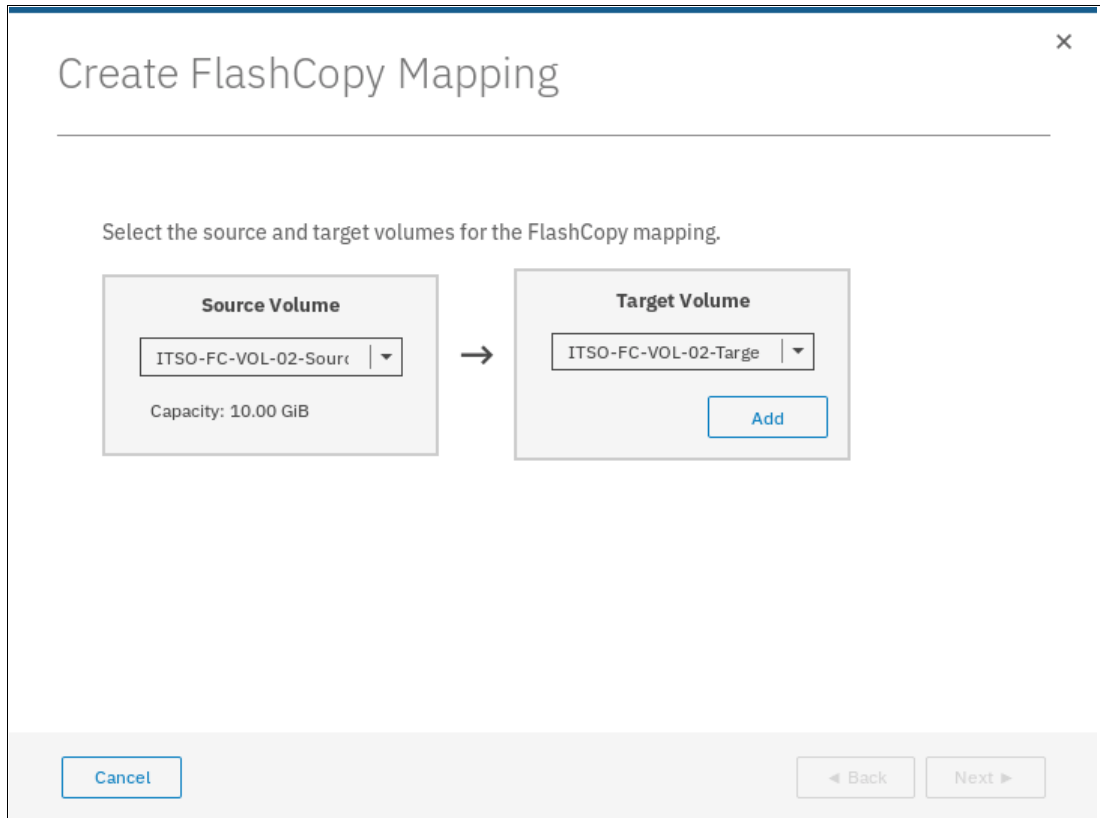


Figure 10-23 Selecting source and target for a FlashCopy mapping

To remove a mapping that was created, click **X** (see Figure 10-24 on page 531).

2. Click **Next** after you create all of the mappings that you need, as shown in Figure 10-24.

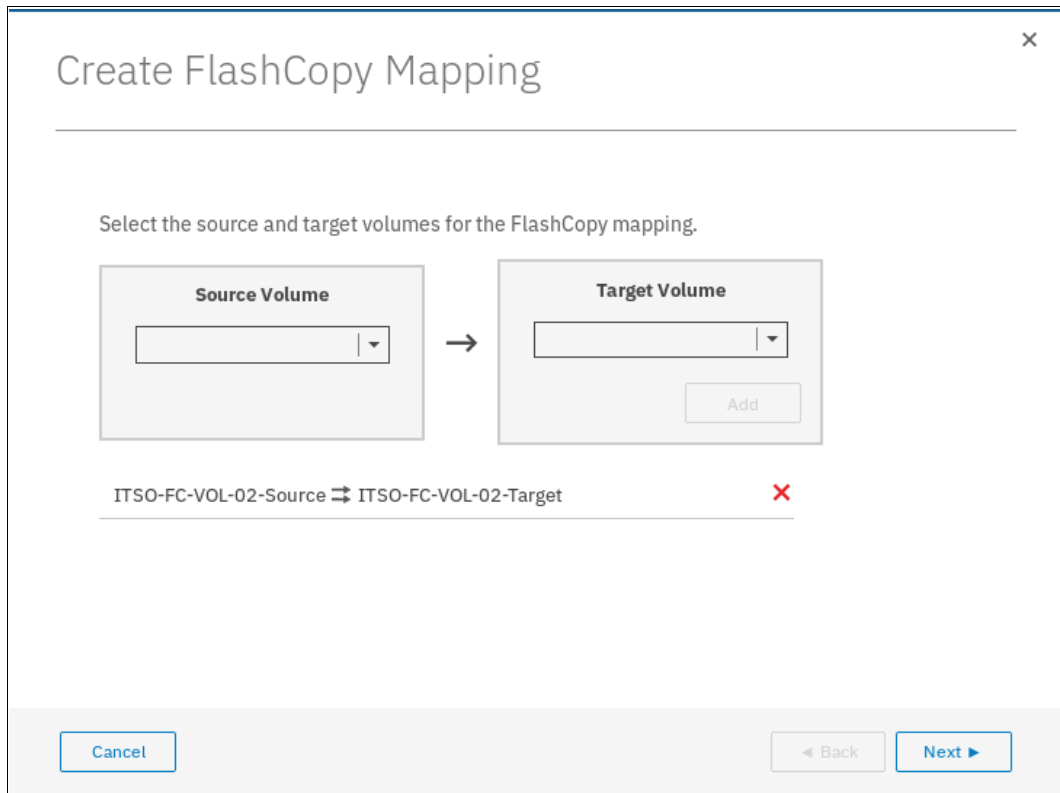


Figure 10-24 Viewing source and target at creation time

3. In the next window, select one FlashCopy preset. The GUI provides the following presets to simplify common FlashCopy operations, as shown in Figure 10-25 on page 532. For more information about the presets, see 10.2.1, “FlashCopy presets” on page 523:

- Snapshot: Creates a PiT snapshot copy of the source volume.
- Clone: Creates a PiT replica of the source volume.
- Backup: Creates an incremental FlashCopy mapping that can be used to recover data or objects if the system experiences data loss. These backups can be copied multiple times from source and target volumes.

**Note:** If you want to create a simple Snapshot of a volume, you likely want the target volume to be defined as thin-provisioned to save space on your system. If you use an existing target, ensure it is thin-provisioned first. The use of the Snapshot preset does not make the system check whether the target volume is thin-provisioned.

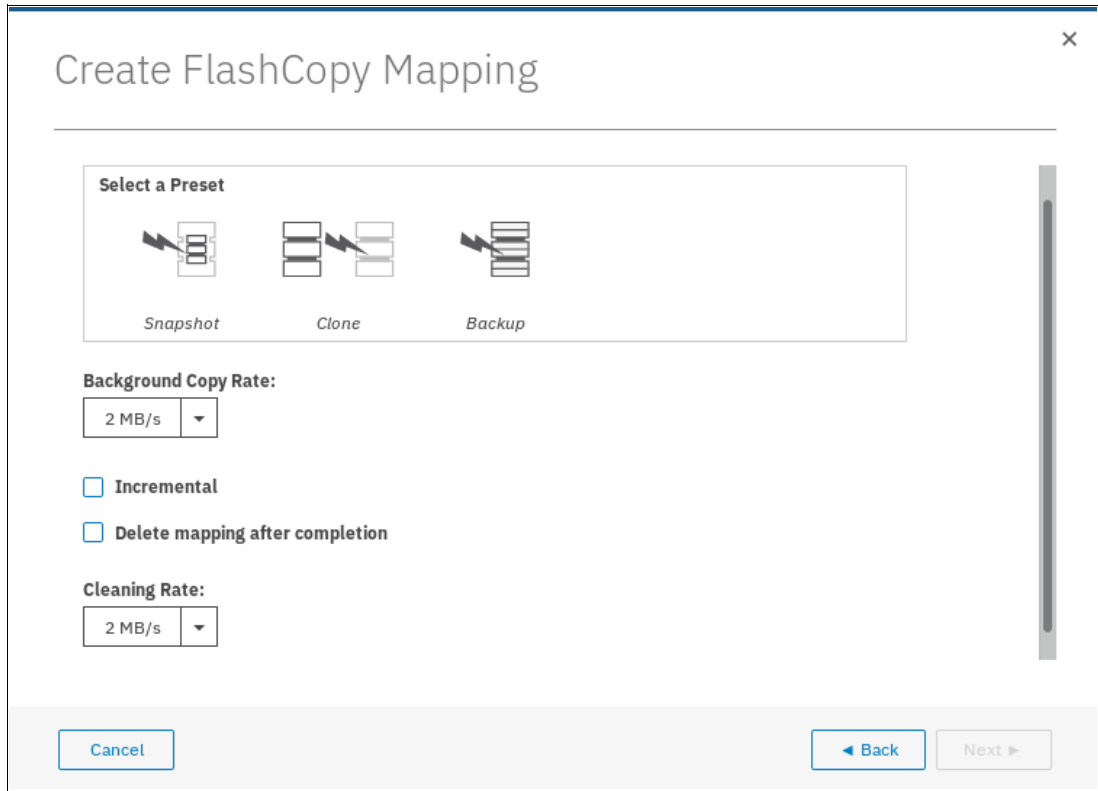


Figure 10-25 FlashCopy mapping preset selection

When selecting a preset, some options, such as Background Copy Rate, Incremental, and Delete mapping after completion, are automatically changed or selected. You can still change the automatic settings, but this is not recommended for the following reasons:

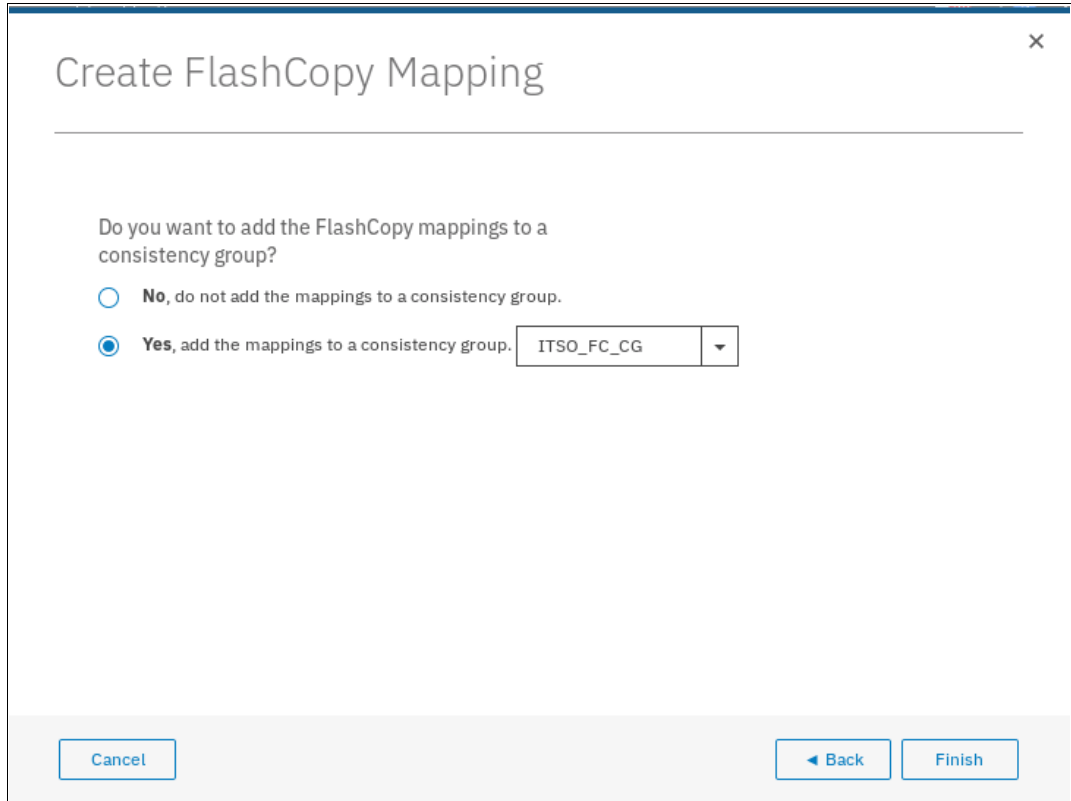
- If you select the **Backup** preset but then clear **Incremental** or select **Delete mapping after completion**, you lose the benefits of the incremental FlashCopy and must copy the entire source volume each time you start the mapping.
- If you select the **Snapshot** preset but then change the **Background Copy Rate**, you will have a full copy of your source volume.

For more information about the Background Copy Rate and the Cleaning Rate, see Table 10-1 on page 505, or Table 10-5 on page 514.

When your FlashCopy mapping setup is ready, click **Next**.

4. You can choose whether to add the mappings to a consistency group, as shown in Figure 10-26.

If you want to include this FlashCopy mapping in a consistency group, select **Yes, add the mappings to a consistency group** and select the consistency group from the drop-down menu.



Do you want to add the FlashCopy mappings to a consistency group?

No, do not add the mappings to a consistency group.

Yes, add the mappings to a consistency group. ITSO\_FC\_CG

Cancel      < Back      Finish

Figure 10-26 Select or not a consistency group for the FlashCopy mapping

5. It is possible to add a FlashCopy mapping to a consistency group or to remove a FlashCopy mapping from a consistency group after they are created. If you do not know at this stage what to do, you can change it later. Click **Finish**.

The FlashCopy mapping is now ready for use. It is visible in the three different windows: FlashCopy, FlashCopy mappings, and Consistency Groups.

**Note:** Creating a FlashCopy mapping does *not* automatically start any copy. You must manually start the mapping.

## Creating a FlashCopy mapping and target volumes

Complete the following steps to create target volumes for FlashCopy mapping:

1. Right-click the volume that you want to create a FlashCopy mapping for and select **Advanced FlashCopy** → **Create New Target Volumes**, as shown in Figure 10-27.

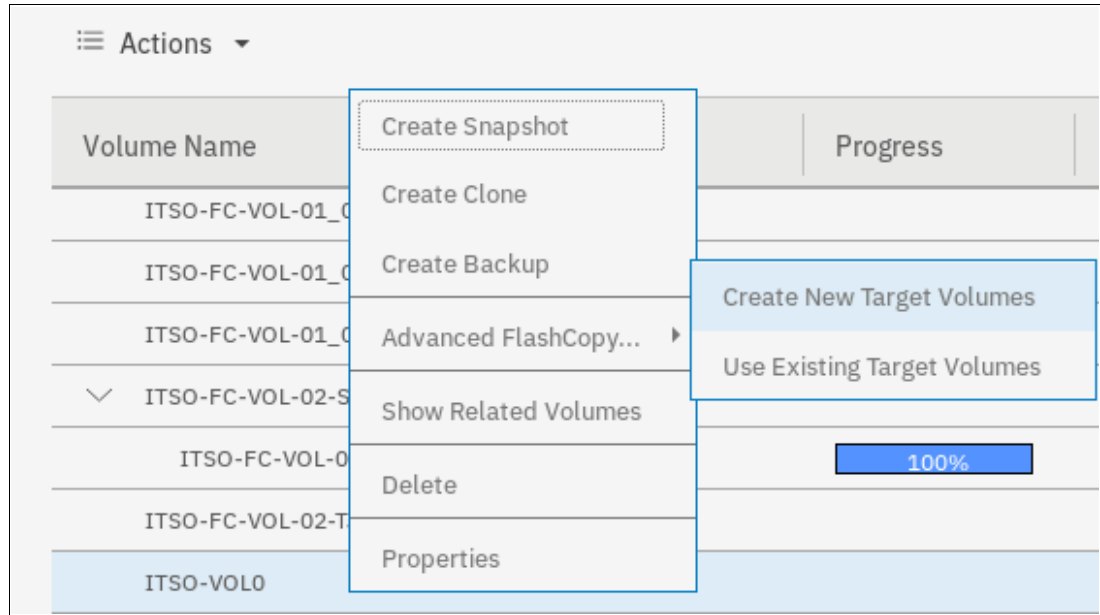


Figure 10-27 Creating a FlashCopy mapping and creating targets

2. In the next window, select one FlashCopy preset. The GUI provides the following presets to simplify common FlashCopy operations, as shown in Figure 10-28 on page 535. For more information about the presets, see 10.2.1, “FlashCopy presets” on page 523:
  - Snapshot: Creates a PiT snapshot copy of the source volume.
  - Clone: Creates a PiT replica of the source volume.
  - Backup: Creates an incremental FlashCopy mapping that can be used to recover data or objects if the system experiences data loss. These backups can be copied multiple times from source and target volumes.

**Note:** If you want to create a simple Snapshot of a volume, you likely want the target volume to be defined as thin-provisioned to save space on your system. If you use an existing target, ensure it is thin-provisioned first. The use of the Snapshot preset does not make the system check whether the target volume is thin-provisioned.



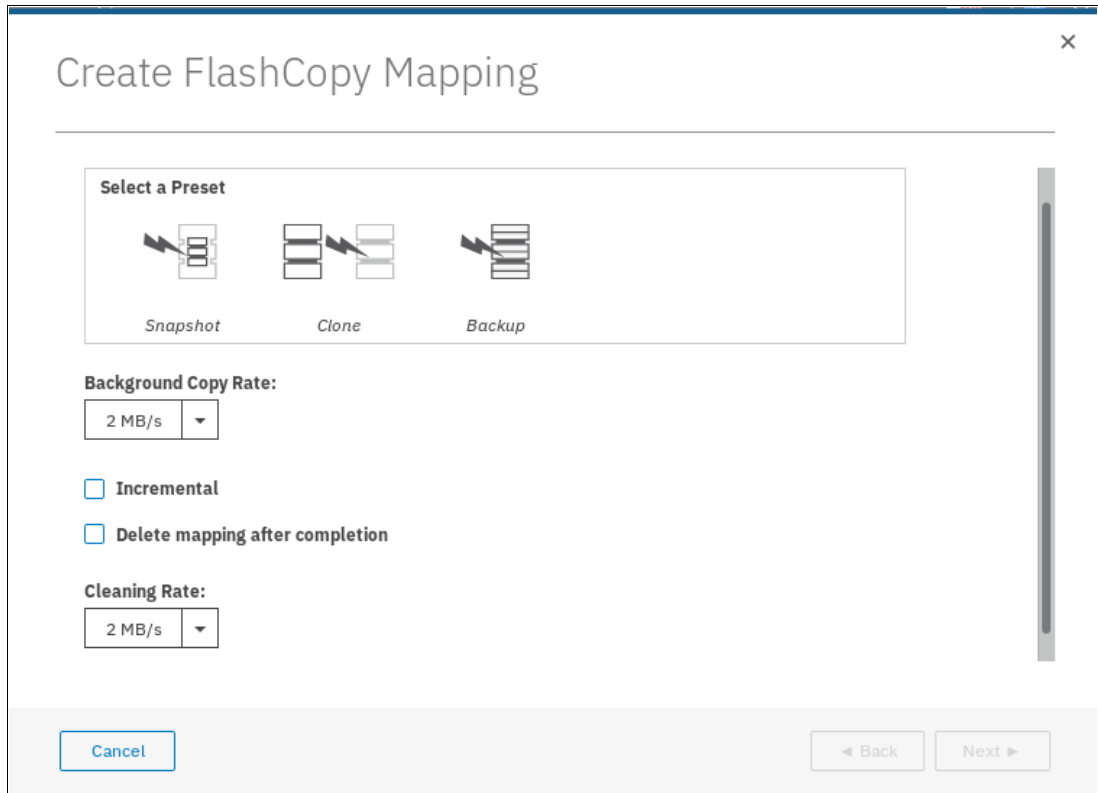


Figure 10-28 FlashCopy mapping preset selection

When selecting a preset, some options, such as Background Copy Rate, Incremental, and Delete mapping, after completion are automatically changed or selected. You can still change the automatic settings, but this is not recommended for the following reasons:

- If you select the **Backup** preset but then clear **Incremental** or select **Delete mapping after completion**, you lose the benefits of the incremental FlashCopy. You must copy the entire source volume each time you start the mapping.
- If you select the **Snapshot** preset but then change the **Background Copy Rate**, you have a full copy of your source volume.

For more information about the Background Copy Rate and the Cleaning Rate, see Table 10-1 on page 505, or Table 10-5 on page 514.

When your FlashCopy mapping setup is ready, click **Next**.

3. You can choose whether to add the mappings to a consistency group, as shown in Figure 10-29.

If you want to include this FlashCopy mapping in a consistency group, select **Yes, add the mappings to a consistency group**, and select the consistency group from the drop-down menu.

Create FlashCopy Mapping

Do you want to add the FlashCopy mappings to a consistency group?

No, do not add the mappings to a consistency group.

Yes, add the mappings to a consistency group. ITSO\_FC\_CG

Cancel      < Back      Next >

Figure 10-29 Select a consistency group for the FlashCopy mapping

4. It is possible to add a FlashCopy mapping to a consistency group or to remove a FlashCopy mapping from a consistency group after they are created. If you do not know at this stage what to do, you can change it later. Click **Next**.

5. The system prompts the user to select the pool that is used to automatically create targets, as shown in Figure 10-30. Click **Next**.

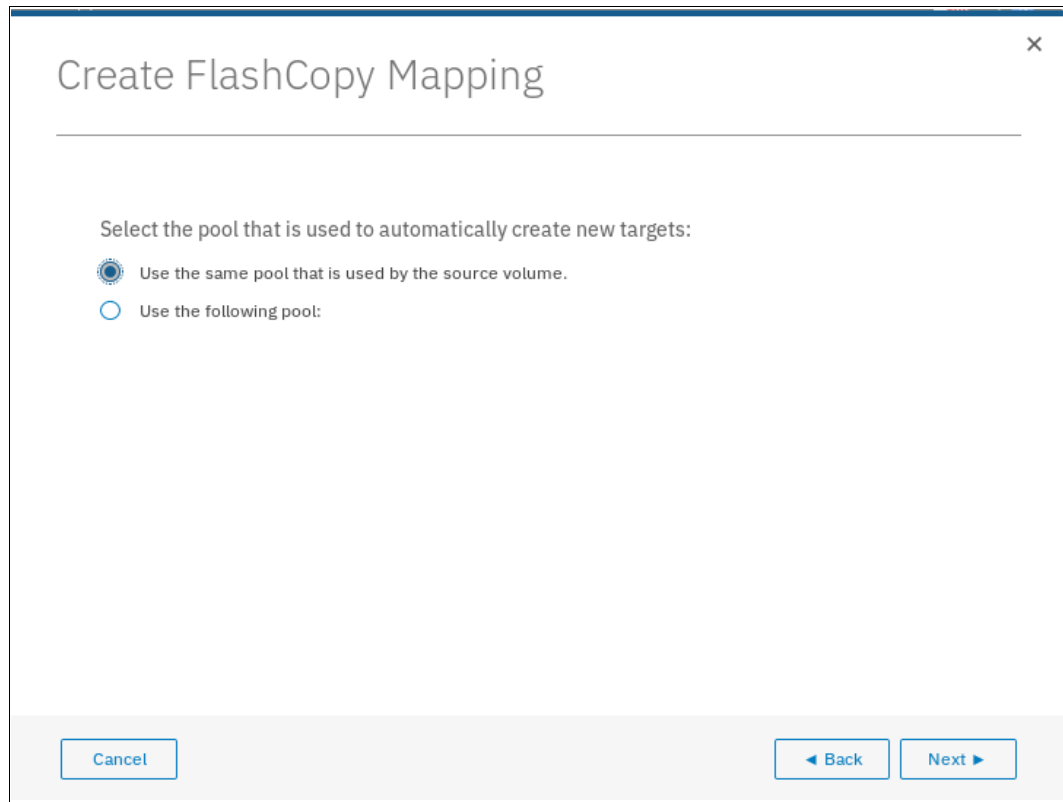


Figure 10-30 Select the pool

6. The system prompts the user how to define the new volumes that are created, as shown in Figure 10-31 on page 538. It can be None, Thin-provisioned, or Inherit from source volume. If Inherit from source volume is selected, the system checks the type of the source volume and then creates a target of the same type. Click **Finish**.

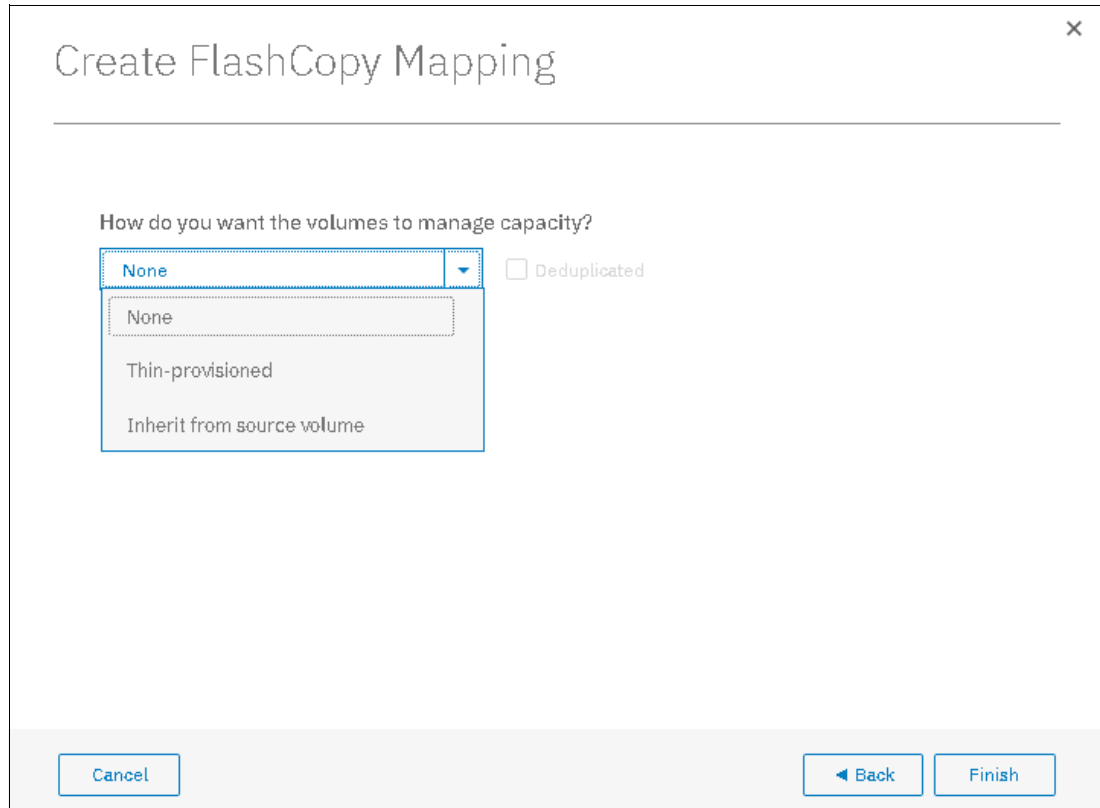


Figure 10-31 Select the type of volumes for the created targets

**Note:** If you selected multiple source volumes to create FlashCopy mappings, selecting **Inherit properties from source Volume** applies to each newly created target volume. For example, if you selected a compressed volume and a generic volume as sources for the new FlashCopy mappings, the system creates a compressed target and a generic target.

The FlashCopy mapping is now ready for use. It is visible in the three different windows: FlashCopy, FlashCopy mappings, and consistency groups.

## 10.2.4 Single-click snapshot

The *snapshot* creates a PiT backup of production data. The snapshot is not intended to be an independent copy. Instead, it is used to maintain a view of the production data at the time that the snapshot is created. Therefore, the snapshot holds only the data from regions of the production volume that changed since the snapshot was created. Because the snapshot preset uses thin provisioning, only the capacity that is required for the changes is used.

Snapshot uses the following preset parameters:

- ▶ Background copy: No
- ▶ Incremental: No
- ▶ Delete after completion: No
- ▶ Cleaning rate: No
- ▶ Primary copy source pool: Target pool

To create and start a snapshot, complete the following steps:

1. Open the FlashCopy window from the **Copy Services** → **FlashCopy** menu.
2. Select the volume that you want to create a snapshot of, and right-click it or click **Actions** → **Create Snapshot**, as shown in Figure 10-32.

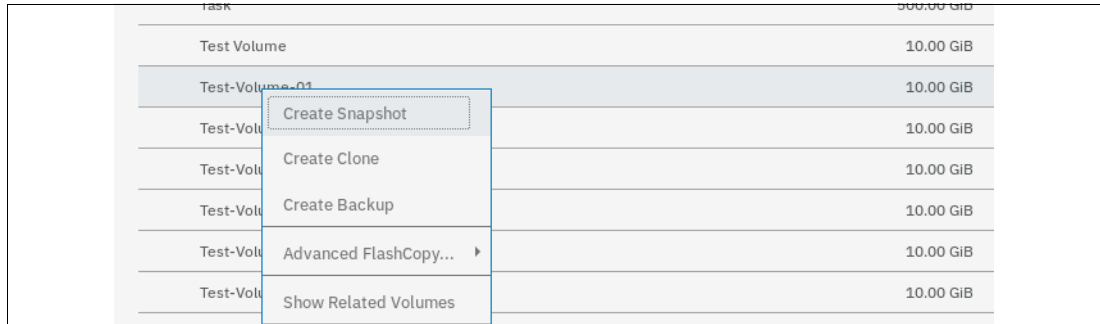


Figure 10-32 Single-click snapshot creation and start

3. You can select multiple volumes at a time, which creates as many snapshots automatically. The system then automatically groups the FlashCopy mappings in a new consistency group, as shown in Figure 10-33.

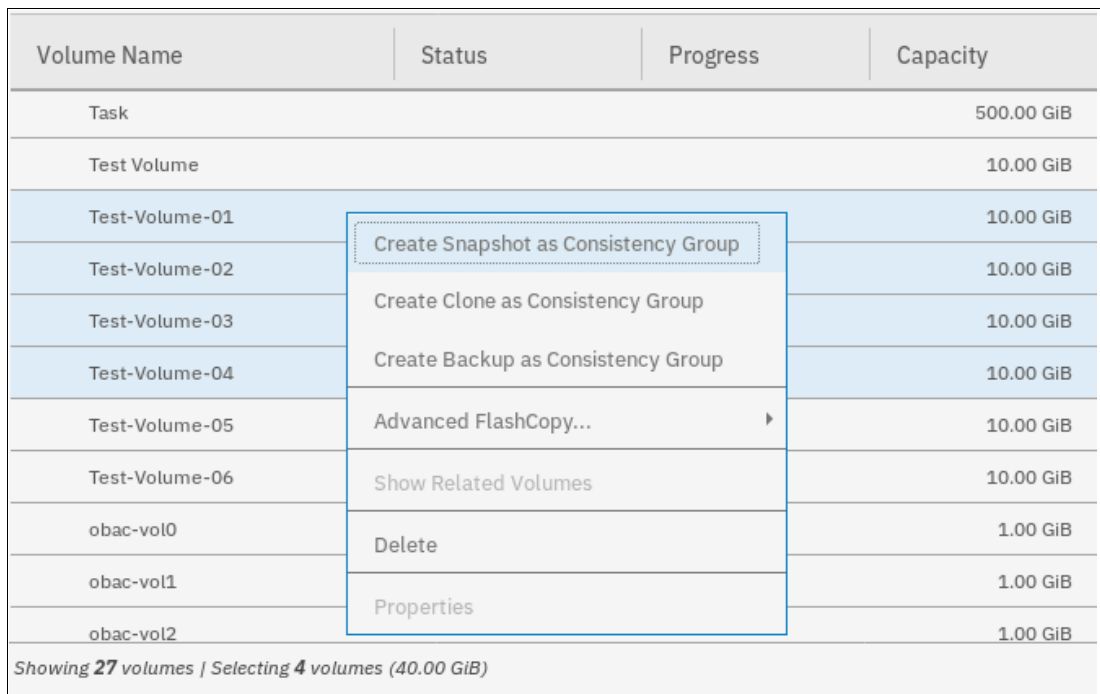


Figure 10-33 Selection single-click snapshot creation and start

For each selected source volume, the following actions occur:

- A FlashCopy mapping is automatically created. It is named by default fcmappXX.
- A target volume is created. By default the source name is appended with a \_XX suffix.
- A consistency group is created for each mapping, unless multiple volumes were selected. Consistency groups are named by default fccstgrpX.

The newly created consistency group is automatically started.

## 10.2.5 Single-click clone

The *clone preset* creates a replica of the volume, which can be changed without affecting the original volume. After the copy completes, the mapping that was created by the preset is automatically deleted.

The clone preset uses the following parameters:

- ▶ Background copy rate: 50
- ▶ Incremental: No
- ▶ Delete after completion: Yes
- ▶ Cleaning rate: 50
- ▶ Primary copy source pool: Target pool

To create and start a snapshot, complete the following steps:

1. Open the FlashCopy window from the **Copy Services** → **FlashCopy** menu.
2. Select the volume that you want to create a snapshot of, and right-click it or click **Actions** → **Create Clone**, as shown in Figure 10-34.

Volume Name	Status	Progress	Capacity
Task			500.00 GiB
Test Volume			10.00 GiB
Test-Volume-01			10.00 GiB
Test-Volume-02			10.00 GiB
Test-Volume-03			10.00 GiB
Test-Volume-04			10.00 GiB
Test-Volume-05			10.00 GiB
Test-Volume-06			10.00 GiB
obac-vol0			1.00 GiB
obac-vol1			1.00 GiB
obac-vol2			1.00 GiB

Showing 27 volumes | Selecting 1 volume (10.00 GiB)

Figure 10-34 Single-click clone creation and start

3. You can select multiple volumes at a time, which creates as many snapshots automatically. The system then automatically groups the FlashCopy mappings in a new consistency group, as shown in Figure 10-35.

Volume Name	Status	Progress	Capacity
Task			500.00 GiB
Test Volume			10.00 GiB
Test-Volume-01			10.00 GiB
Test-Volume-02			10.00 GiB
Test-Volume-03		Create Snapshot as Consistency Group	10.00 GiB
Test-Volume-04		Create Clone as Consistency Group	10.00 GiB
Test-Volume-05		Create Backup as Consistency Group	10.00 GiB
Test-Volume-06		Advanced FlashCopy...	10.00 GiB
obac-vol0		Show Related Volumes	1.00 GiB
obac-vol1		Delete	1.00 GiB
obac-vol2		Properties	1.00 GiB
Showing 27 volumes   Selecting 4 v			

Figure 10-35 Selection single-click clone creation and start

For each selected source volume, the following actions occur:

- A FlashCopy mapping is automatically created. It is named by default fcmappXX.
- A target volume is created. The source name is appended with an \_XX suffix.
- A consistency group is created for each mapping, unless multiple volumes were selected. Consistency groups are named by default fccstgrpX.
- The newly created consistency group is automatically started.

## 10.2.6 Single-click backup

The backup creates a PiT replica of the production data. After the copy completes, the backup view can be refreshed from the production data, with minimal copying of data from the production volume to the backup volume. The backup preset uses the following parameters:

- ▶ Background Copy rate: 50
- ▶ Incremental: Yes
- ▶ Delete after completion: No
- ▶ Cleaning rate: 50
- ▶ Primary copy source pool: Target pool

To create and start a backup, complete the following steps:

1. Open the FlashCopy window from the **Copy Services** → **FlashCopy** menu.
2. Select the volume that you want to create a backup of, and right-click it or click **Actions** → **Create Backup**, as shown in Figure 10-36.

Volume Name	Status	Progress	Capacity
Task			500.00 GiB
Test Volume			10.00 GiB
Test-Volume-01			10.00 GiB
Test-Volume-02			10.00 GiB
Test-Volume-03			10.00 GiB
Test-Volume-04			10.00 GiB
Test-Volume-05			10.00 GiB
Test-Volume-06			10.00 GiB
obac-vol0			1.00 GiB
obac-vol1			1.00 GiB
obac-vol2			1.00 GiB

Showing 27 volumes | Selecting 1 volume (10.00 GiB)

Figure 10-36 Single-click backup creation and start

3. You can select multiple volumes at a time, which creates as many snapshots automatically. The system then automatically groups the FlashCopy mappings in a new consistency group, as shown Figure 10-37 on page 543.



Volume Name	Status	Progress	Capacity
Task			500.00 GiB
Test Volume			10.00 GiB
Test-Volume-01			10.00 GiB
Test-Volume-02			10.00 GiB
Test-Volume-03			10.00 GiB
Test-Volume-04			10.00 GiB
Test-Volume-05			10.00 GiB
Test-Volume-06			10.00 GiB
obac-vol0			1.00 GiB
obac-vol1			1.00 GiB
obac-vol2			1.00 GiB
Showing 27 volumes / Selecting			

Create Snapshot as Consistency Group

Create Clone as Consistency Group

Create Backup as Consistency Group

Advanced FlashCopy...

Show Related Volumes

Delete

Figure 10-37 Selection single-click backup creation and start

For each selected source volume, the following actions occur:

- A FlashCopy mapping is automatically created. It is named by default fcmappXX.
- A target volume is created. It is named after the source name with a \_XX suffix.
- A consistency group is created for each mapping, unless multiple volumes were selected. Consistency groups are named by default fccstgrpX.
- The newly created consistency group is automatically started.

## 10.2.7 Creating a FlashCopy consistency group

To create a FlashCopy consistency group in the GUI, complete the following steps:

1. Open the Consistency Groups window by clicking **Copy Services** → **Consistency Groups**. Click **Create Consistency Group**, as shown in Figure 10-38.

Mapping Name	Status	Source Volume	Target Volume	Progress
⊕ Create Consistency Group   ⋮ Actions ▾   Default ▾   Contains ▾   Filter				
⌵ Not in a Group				
fcmapp0	✓ Copied	ITSO-FC-VOL-01	ITSO-FC-VOL-01_01	100%
fcmapp1	✓ Idle	ITSO-VOL0	ITSO-VOL0_02	0%
fcmapp7	⌛ Copying	ITSO-FC-VOL-01	ITSO-FC-VOL-01_04	0%
fcmapp8	⌛ Copying	ITSO-FC-VOL-01	ITSO-FC-VOL-01_05	0%

Figure 10-38 Creating a consistency group

2. Enter the FlashCopy Consistency Group name that you want to use and the ownership group; then, click **Create**, as shown in Figure 10-39.

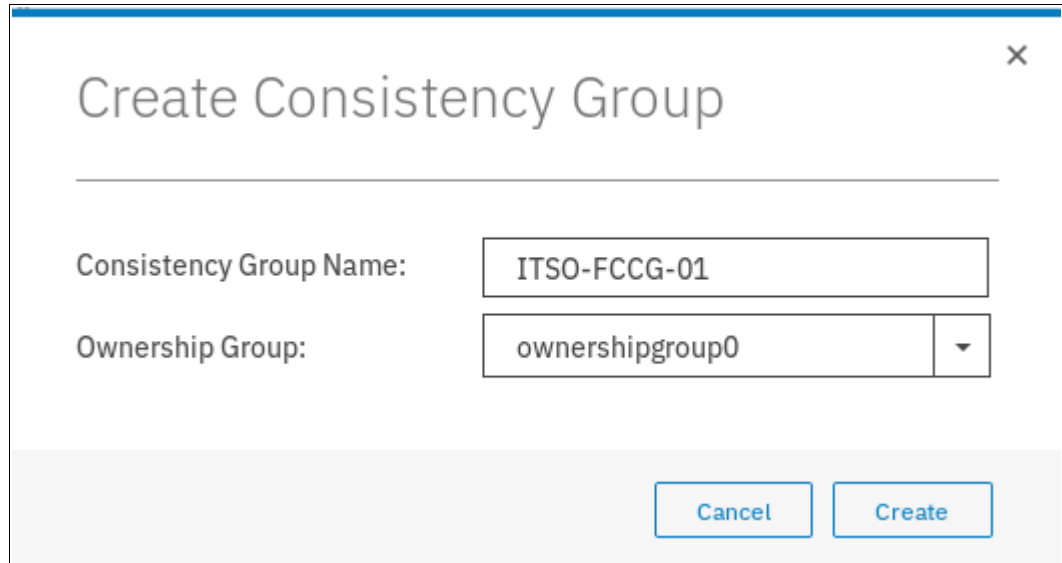


Figure 10-39 Enter the name and ownership group of new consistency group

**Consistency Group name:** You can use the letters A - Z and a - z, the numbers 0 - 9, and the underscore ( \_ ) character. The volume name can be 1 - 63 characters.

## 10.2.8 Creating FlashCopy mappings in a Consistency Group

To create a FlashCopy Consistency Group in the GUI, complete the following steps:

1. Open the Consistency Groups window by clicking **Copy Services** → **Consistency Groups**. This example assumes that source and target volumes were previously created.
2. Select the Consistency Group in which you want to create the FlashCopy mapping. If you prefer not to create a FlashCopy mapping in a Consistency Group, select **Not in a Group**, and right-click the selected consistency group or click **Actions** → **Create FlashCopy Mapping**, as shown in Figure 10-40 on page 545.

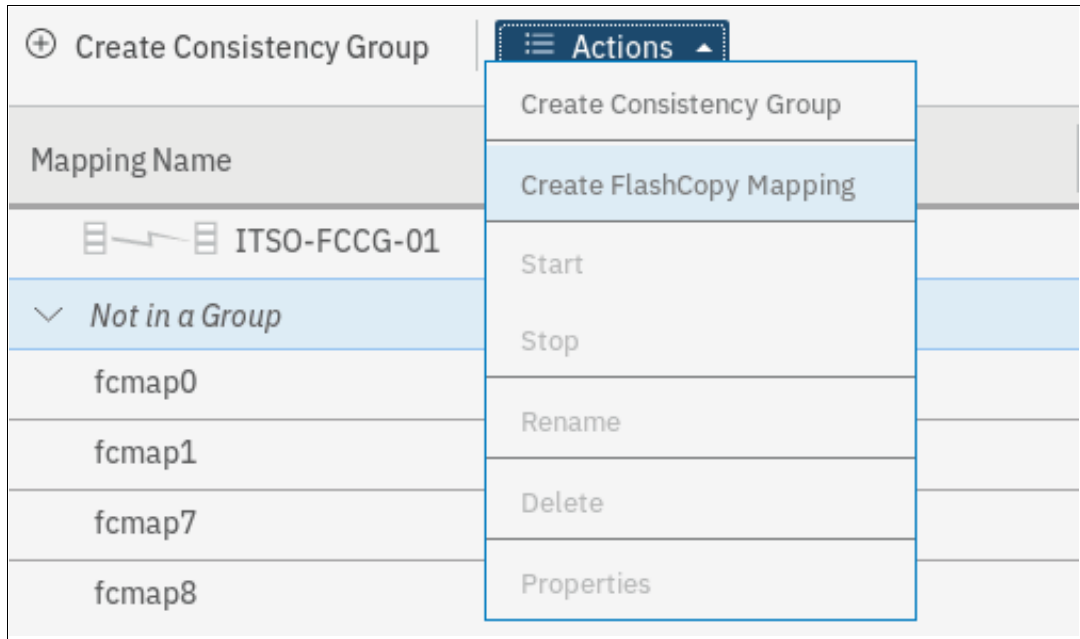


Figure 10-40 Creating a FlashCopy mapping

3. Select a volume in the source volume column by using the drop-down menu. Then, select a volume in the target volume column by using the drop-down menu. Click **Add**, as shown in Figure 10-41.

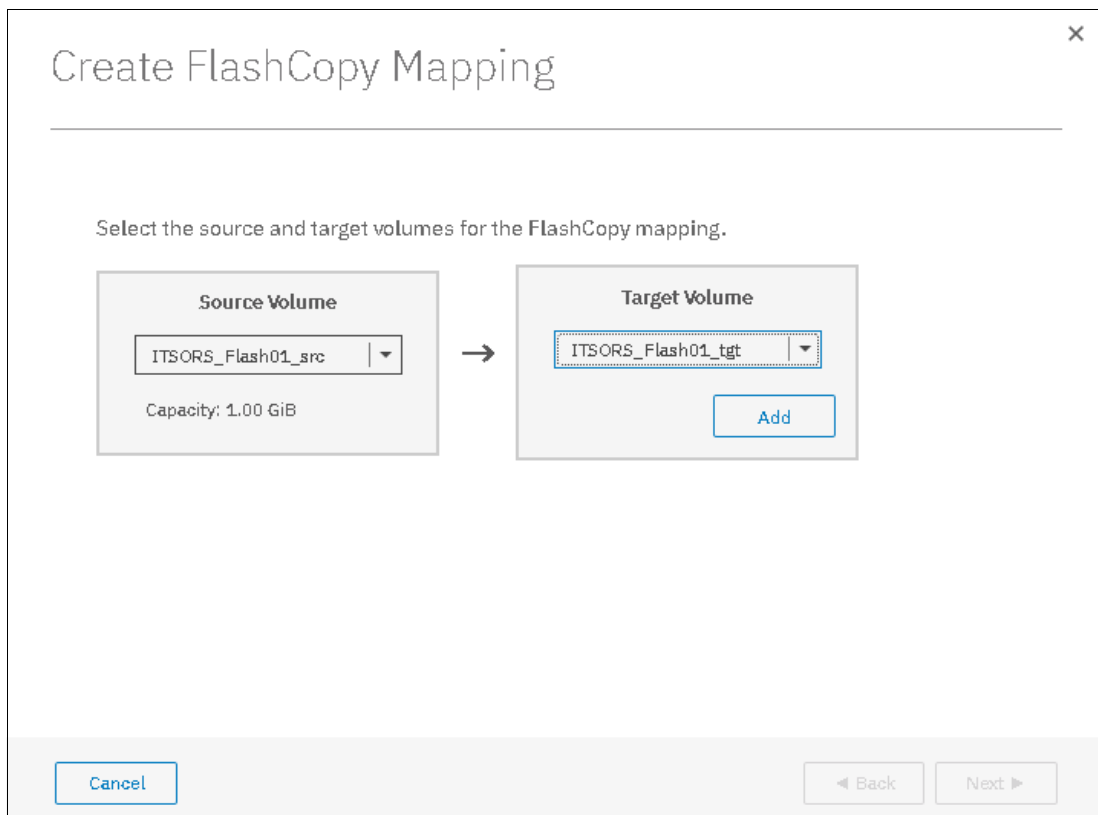


Figure 10-41 Select source and target volumes for the FlashCopy mapping

Repeat this step to create other mappings. To remove a mapping that was created, click **X**. Click **Next**.

**Important:** The source and target volumes must be of equal size. Therefore, only the targets with the suitable size are shown for a source volume.

Volumes that are target volumes in another FlashCopy mapping cannot be target of a new FlashCopy mapping. Therefore, they do not appear in the list.

- In the next window, select one FlashCopy preset. The GUI provides the following presets to simplify common FlashCopy operations, as shown in Figure 10-42. For more information about the presets, see 10.2.1, “FlashCopy presets” on page 523:
  - **Snapshot:** Creates a PiT snapshot copy of the source volume.
  - **Clone:** Creates a PiT replica of the source volume.
  - **Backup:** Creates an incremental FlashCopy mapping that can be used to recover data or objects if the system experiences data loss. These backups can be copied multiple times from source and target volumes.

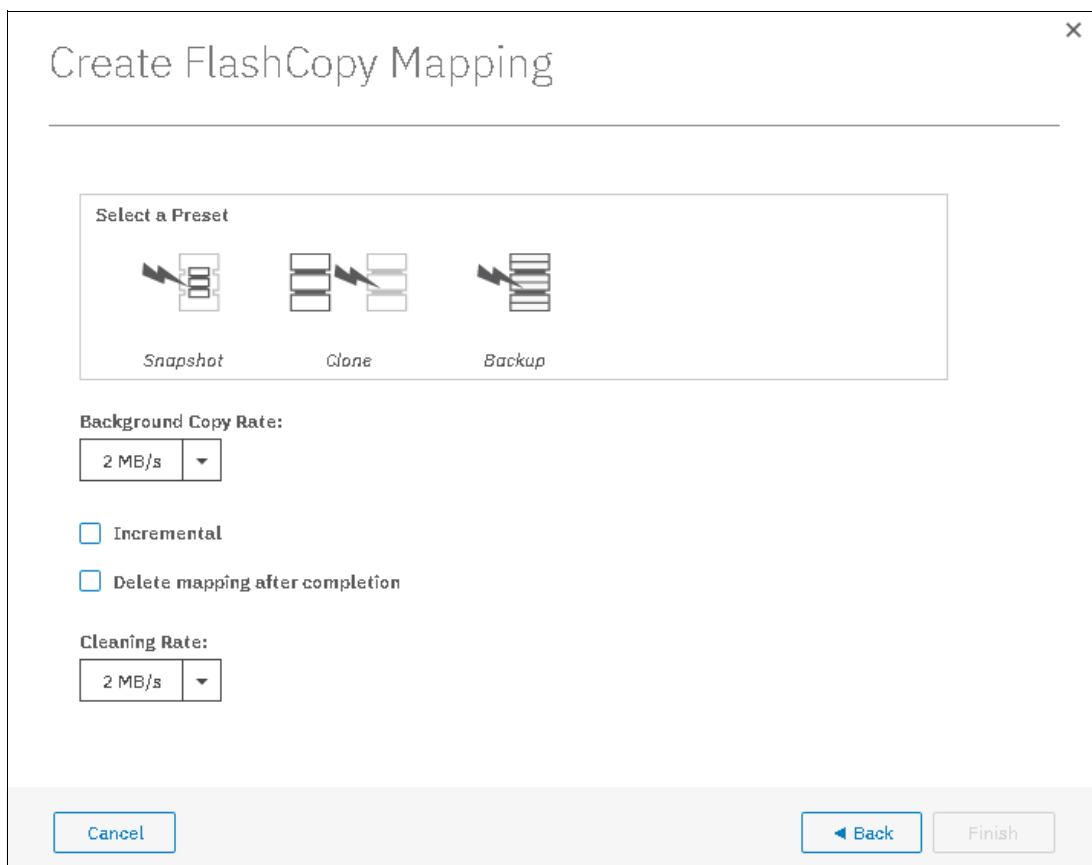


Figure 10-42 FlashCopy mapping preset selection

When selecting a preset, some options, such as Background Copy Rate, Incremental, and Delete mapping after completion, are automatically changed or selected. You can still change the automatic settings, but this is not recommended for the following reasons:

- If you select the **Backup** preset but then clear **Incremental** or select **Delete mapping after completion**, you lose the benefits of the incremental FlashCopy. You must copy the entire source volume each time you start the mapping.

- If you select the **Snapshot** preset but then change the **Background Copy Rate**, you have a full copy of your source volume.

For more information about the Background Copy Rate and the Cleaning Rate, see Table 10-1 on page 505, or Table 10-5 on page 514.

5. When your FlashCopy mapping setup is ready, click **Finish**.

## 10.2.9 Showing related Volumes

To show related volumes for a specific FlashCopy mapping, complete the following steps:

1. Open the Copy Services FlashCopy Mappings window.
2. Right-click a FlashCopy mapping and select **Show Related Volumes**, as shown in Figure 10-43. Also, depending on which window you are inside Copy Services, you can right-click at mappings and select **Show Related Volumes**.

Mapping Name	Status	Source Volume
ITSO-FCCG-01	Idle or Copied	
fcmap2	Copied	ITSO-FC-VOL-02-S
Not in a Group		
fcmap0	ed	ITSO-FC-VOL-01
fcmap1		ITSO-VOL0
fcmap7	ying	ITSO-FC-VOL-01
fcmap8	ying	ITSO-FC-VOL-01
Selected 1 FlashCopy mappin		

Create Consistency Group

Create FlashCopy Mapping

Move to Consistency Group

Remove from Consistency Group

Start

Stop

Rename Mapping

Delete Mapping

Show Related Volumes

Figure 10-43 Showing related volumes for a mapping, a consistency group or another volume

3. In the related volumes window, you can see the related mapping for a volume, as shown in Figure 10-44 on page 548. If you click one of these volumes, you can see its properties.

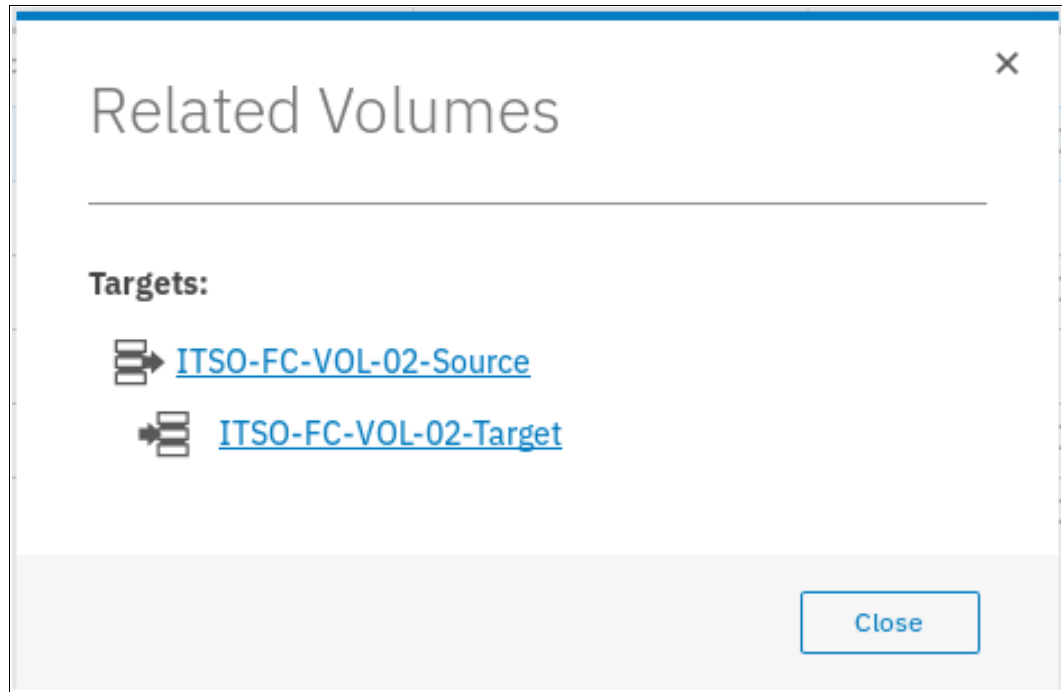


Figure 10-44 Showing related volumes list

### 10.2.10 Moving FlashCopy mappings across Consistency Groups

To move one or multiple FlashCopy mappings to a Consistency Group, complete the following steps:

1. Open the FlashCopy, Consistency Groups, or FlashCopy Mappings window.
2. Right-click the FlashCopy mappings that you want to move and select **Move to Consistency Group**, as shown in Figure 10-45.

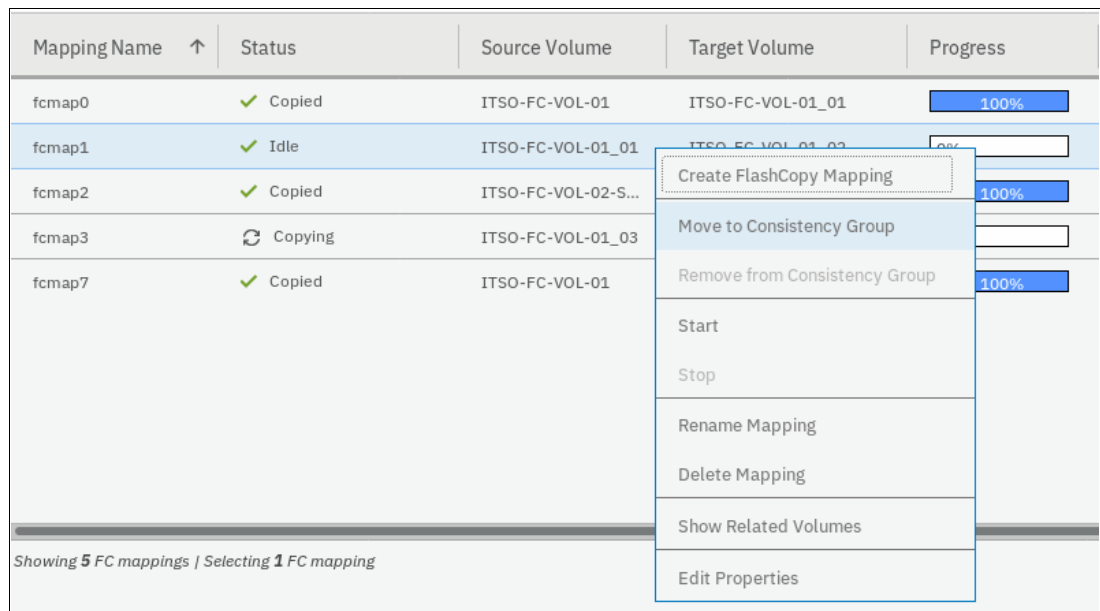


Figure 10-45 Moving a FlashCopy mapping to a consistency group

**Note:** You cannot move a FlashCopy mapping that is in a copying, stopping, or suspended state. The mapping should be idle-or-copied or stopped to be moved.

3. In the Move FlashCopy Mapping to Consistency Group window, select the Consistency Group for the FlashCopy mappings selection by using the drop-down menu, as shown in Figure 10-46.

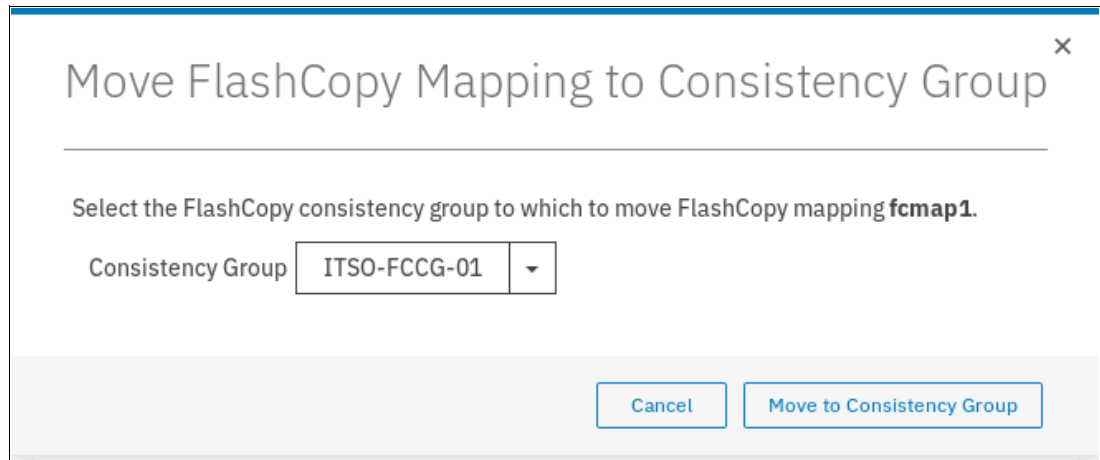


Figure 10-46 Selecting the consistency group where to move the FlashCopy mapping

4. Click **Move to Consistency Group** to confirm your changes.

### 10.2.11 Removing FlashCopy mappings from Consistency Groups

To remove one or multiple FlashCopy mappings from a Consistency Group, complete the following steps:

1. Open the FlashCopy Consistency Groups, or FlashCopy Mappings window.
2. Right-click the FlashCopy mappings that you want to remove and select **Remove from Consistency Group**, as shown in Figure 10-47 on page 550.

**Note:** Only FlashCopy mappings that belong to a consistency group can be removed.

Mapping Name	↑	Status	Source Volume	Target Volume
ITSO-FCCG-01		Idle or Copied		
fcmap1		✓ Copied	ITSO-FC-VOL-01	ITSO-FC-VOL-01_02
fcmap2		✓ Copied		02-Target
Not in a Group				
fcmap0		✓ Copied		01_01
fcmap3		⌛ Copying		01_03_0
fcmap7		✓ Copied		01_04
Selected 1 FlashCopy mapping				

Figure 10-47 Removing FlashCopy mappings from a consistency group

- In the Remove FlashCopy Mapping from Consistency Group window, click **Remove**, as shown in Figure 10-48.

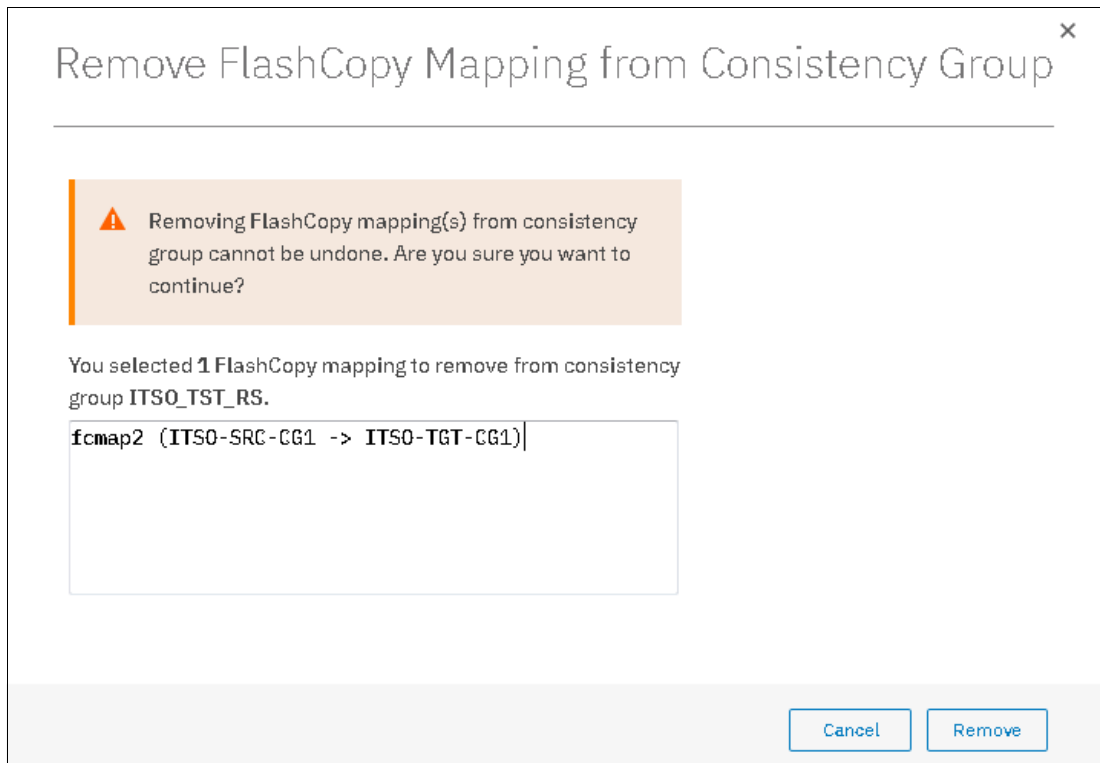


Figure 10-48 Confirm the selection of mappings to be removed



## 10.2.12 Modifying a FlashCopy mapping

To modify a FlashCopy mapping, complete the following steps:

1. Open the FlashCopy Consistency Groups, or FlashCopy Mappings window.
2. Right-click the FlashCopy mapping that you want to edit and select **Edit Properties**, as shown in Figure 10-49.

Mapping Name ↑	Status	Source Volume	Target Volume	Progress
fcmap0	✓ Copied	ITSO-FC-VOL-01	ITSO-FC-VOL-01_01	100%
fcmap1	✓ Copied	ITSO-FC-VOL-01_01	ITSO-FC-VOL-01_02	100%
fcmap2	✓ Copied	ITSO-F	et	100%
fcmap3	🔄 Copying	ITSO-F	01	0%
fcmap7	✓ Copied	ITSO-F		100%

Create FlashCopy Mapping

Move to Consistency Group

Remove from Consistency Group

Start

Stop

Rename Mapping

Delete Mapping

Show Related Volumes

Edit Properties

Showing 5 FC mappings | Selecting 1 FC mapping

Figure 10-49 Editing a FlashCopy mapping properties

**Note:** It is not possible to select multiple FlashCopy mappings to edit their properties all at the same time.

3. In the Edit FlashCopy Mapping window, you can modify the background copy rate and the cleaning rate for a selected FlashCopy mapping, as shown in Figure 10-50.

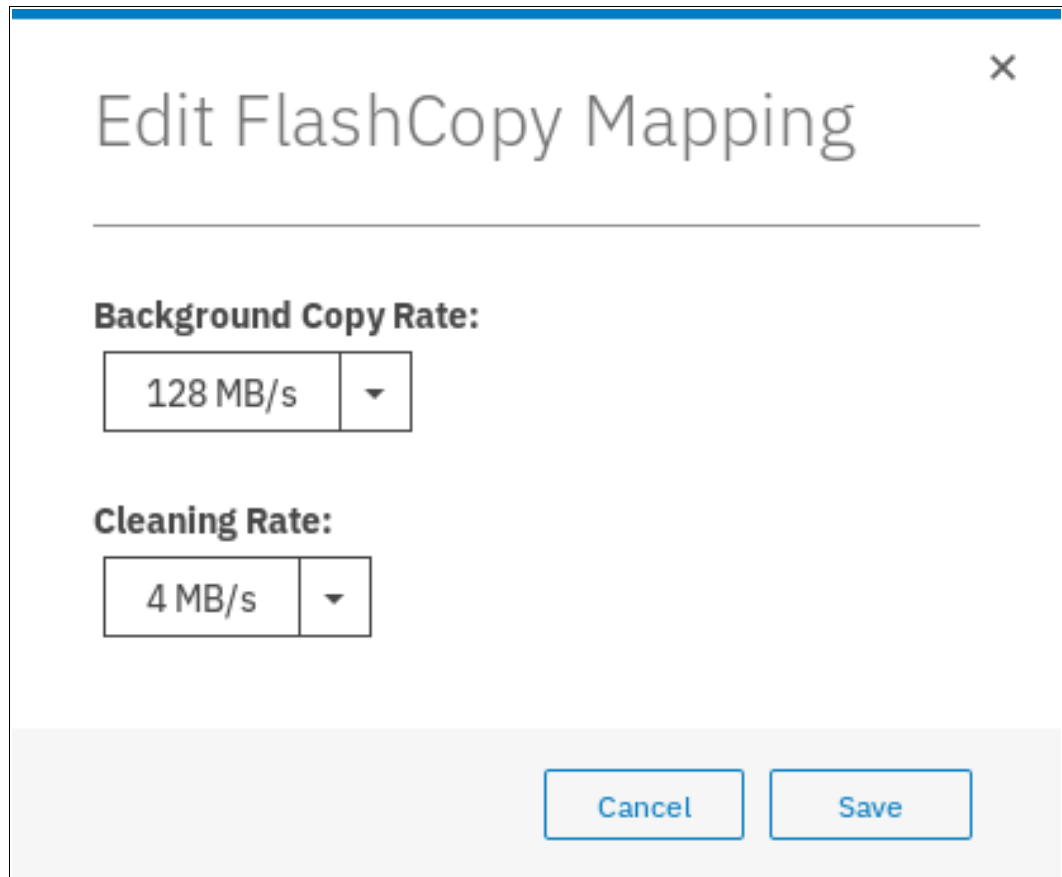


Figure 10-50 Editing copy and cleaning rates of a FlashCopy mapping

For more information about the Background Copy Rate and the Cleaning Rate, see Table 10-1 on page 505, or Table 10-5 on page 514.

4. Click **Save** to confirm your changes.

### 10.2.13 Renaming FlashCopy mappings

To rename one or multiple FlashCopy mappings, complete the following steps:

1. Open the FlashCopy Consistency Groups, or FlashCopy Mappings window.
2. Right-click the FlashCopy mappings that you want to rename and select **Rename Mapping**, as shown in Figure 10-51 on page 553.

Mapping Name ↑	Status	Source Volume	Target Volume	Progress
fcmap0	✓ Copied	ITSO-FC-VOL-01	ITSO-FC-VOL-01_01	100%
fcmap1	✓ Copied	ITSO-FC-VOL-01_01	ITSO-FC-VOL-01_02	100%
fcmap2	✓ Copied	ITSO-	get	100%
fcmap3	🔄 Copying	ITSO-	_01	0%
fcmap7	✓ Copied	ITSO-		100%

Showing 5 FC mappings | Selecting 1 FC mapping

- Create FlashCopy Mapping
- Move to Consistency Group
- Remove from Consistency Group
- Start
- Stop
- Rename Mapping
- Delete Mapping
- Show Related Volumes
- Edit Properties

Figure 10-51 Renaming FlashCopy mappings

- In the Rename FlashCopy Mapping window, enter the new name that you want to assign to each FlashCopy mapping and click **Rename**, as shown in Figure 10-52.

**FlashCopy mapping name:** You can use the letters A - Z and a - z, the numbers 0 - 9, and the underscore (\_) character. The FlashCopy mapping name can be 1 - 63 characters.

## Rename FlashCopy Mapping

---

\*New Name

fcmap1

Reset
Cancel
Rename

Figure 10-52 Renaming the selected FlashCopy mappings

### Renaming a Consistency Group

To rename a Consistency Group, complete the following steps:

- Open the Consistency Groups window.

- Right-click the consistency group you want to rename and select **Rename**, as shown in Figure 10-53.

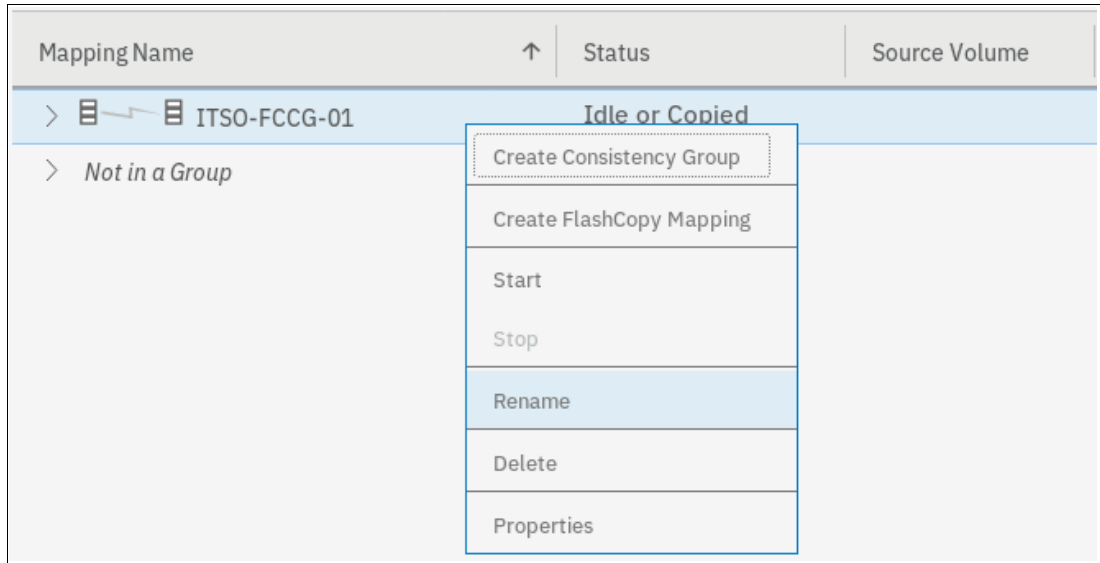


Figure 10-53 Renaming a consistency group

- Enter the new name that you want to assign to the Consistency Group and click **Rename**, as shown in Figure 10-54.

**Note:** It is not possible to select multiple consistency groups to edit their names all at the same time.



Figure 10-54 Renaming the selected consistency group

**Consistency Group name:** The name can consist of the letters A - Z and a - z, the numbers 0 - 9, the dash (-), and the underscore (\_) character. The name can be 1 - 63 characters. However, the name cannot start with a number, a dash, or an underscore.

## 10.2.14 Deleting FlashCopy mappings

To delete one or multiple FlashCopy mappings, complete the following steps:

1. Open the FlashCopy Consistency Groups, or FlashCopy Mappings window.
2. Right-click the FlashCopy mappings that you want to delete and select **Delete Mapping**, as shown in Figure 10-55.

Mapping Name ↑	Status	Source Volume	Target Volume	Progress
fcmap0	✓ Copied	ITSO-FC-VOL-01	ITSO-FC-VOL-01_01	100%
fcmap1	✓ Copied	ITSO-FC-VOL-01_01		100%
fcmap2	✓ Copied	ITSO-FC-VOL-02-S...		100%
fcmap3	🔄 Copying	ITSO-FC-VOL-01_03		
fcmap7	✓ Copied	ITSO-FC-VOL-01		100%

Create FlashCopy Mapping

Move to Consistency Group

Remove from Consistency Group

Start

Stop

Rename Mapping

**Delete Mapping**

Show Related Volumes

Edit Properties

Showing 5 FC mappings | Selecting 1 FC mapping

Figure 10-55 Deleting FlashCopy mappings

3. The Delete FlashCopy Mapping window opens, as shown in Figure 10-56. In the **Verify the number of FlashCopy mappings that you are deleting** field, enter the number of volumes that you want to remove. This verification was added to help avoid deleting the wrong mappings.

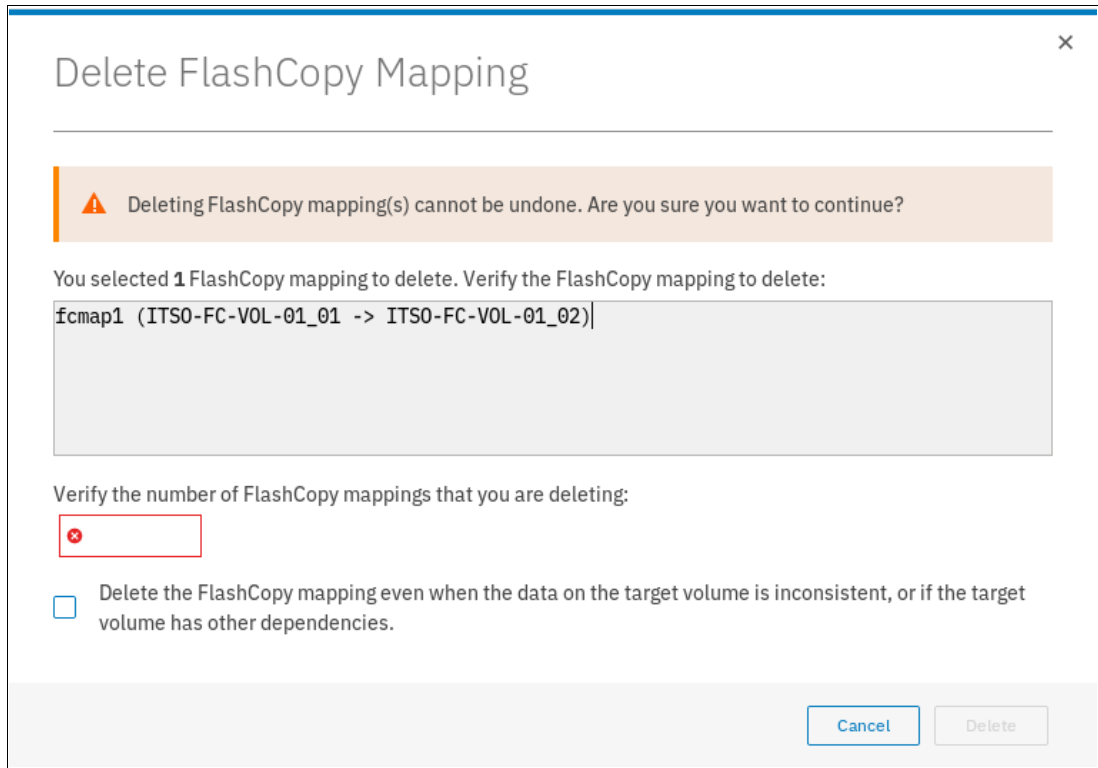


Figure 10-56 Confirming the selection of FlashCopy mappings to be deleted

4. If you still have target volumes that are inconsistent with the source volumes and you want to delete these FlashCopy mappings, select the **Delete the FlashCopy mapping even when the data on the target volume is inconsistent, or if the target volume has other dependencies** option. Click **Delete**.

### 10.2.15 Deleting a FlashCopy consistency group

**Important:** Deleting a consistency group does not delete the FlashCopy mappings that it contains.

To delete a consistency group, complete the following steps:

1. Open the Consistency Groups window.

2. Right-click the consistency group that you want to delete and select **Delete**, as shown in Figure 10-57.

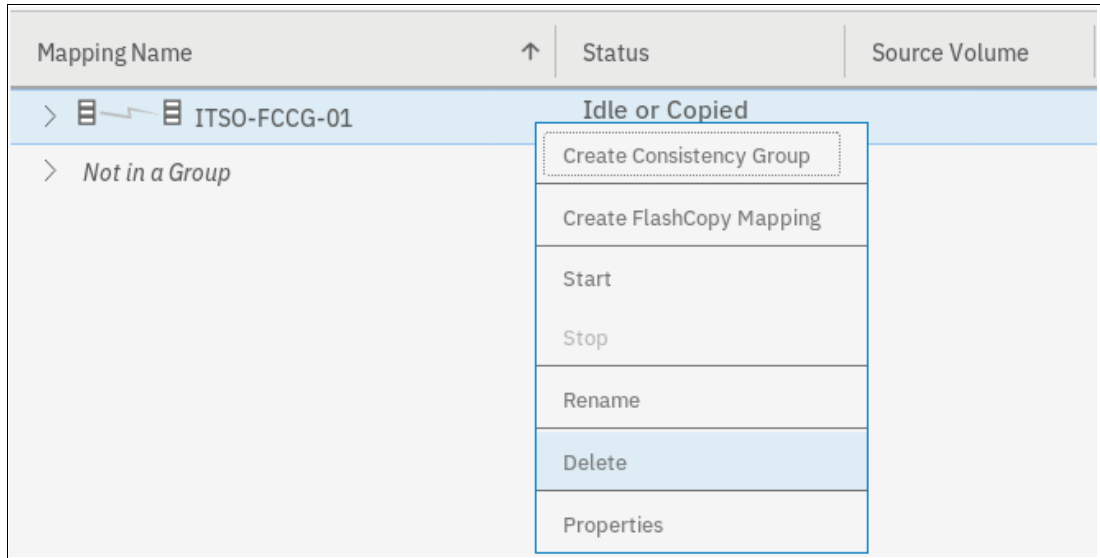


Figure 10-57 Deleting a consistency group

3. A warning message is displayed, as shown in Figure 10-58. Click **Yes**.

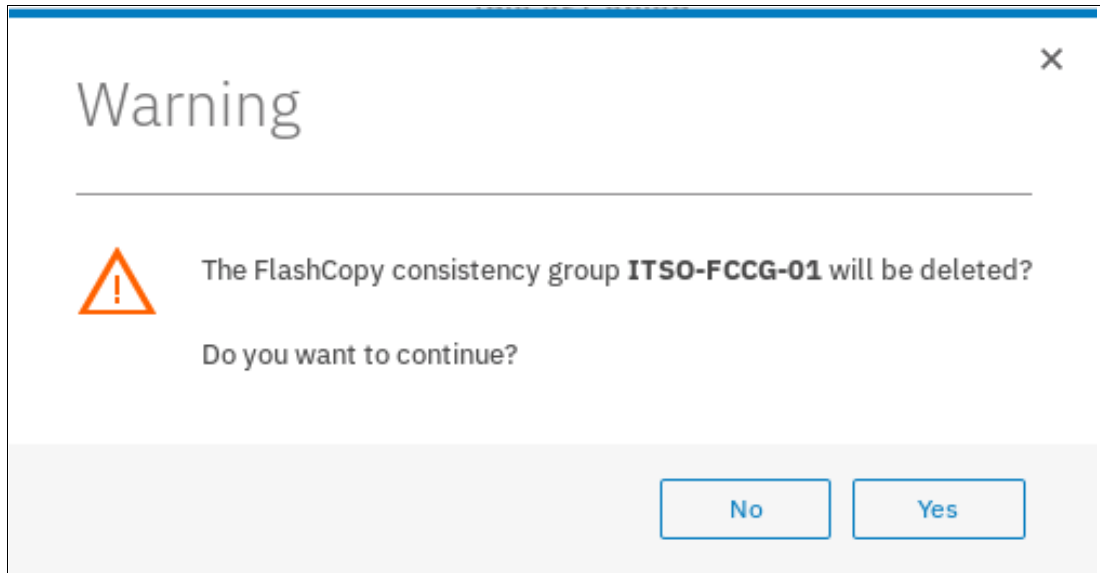


Figure 10-58 Confirming the consistency group deletion

## 10.2.16 Starting FlashCopy mappings

**Important:** Only FlashCopy mappings that do not belong to a consistency group can be started individually. If FlashCopy mappings are part of a consistency group, they can be started only all together by using the consistency group **start** command.

It is the **start** command that defines the “PiT”. It is the moment that is used as a reference (T0) for all subsequent operations on the source and the target volumes. To start one or multiple FlashCopy mappings that do not belong to a consistency group, complete the following steps:

1. Open the FlashCopy Consistency Groups, or FlashCopy Mappings window.
2. Right-click the FlashCopy mappings that you want to start and select **Start**, as shown in Figure 10-59.

Mapping Name	Status	Source Volume	Target Volume	Progress
fcmap0	✓ Copied	ITSO-FC-VOL-01	ITSO-FC-VOL-01_01	100%
fcmap1	✓ Copied	ITSO-FC-VOL-01_01	ITSO-FC-VOL-01_02	100%
fcmap2	✓ Copied		D-FC-VOL-02-Target	100%
fcmap3	🔄 Copying		D-FC-VOL-01_03_01	0%
fcmap7	✓ Copied		D-FC-VOL-01_04	100%

Create FlashCopy Mapping

Move to Consistency Group

Remove from Consistency Group

**Start**

Stop

Rename Mapping

Delete Mapping

Show Related Volumes

Edit Properties

Showing 5 FC mappings | Selecting 1 FC mapping

Figure 10-59 Starting FlashCopy mappings

You can check the FlashCopy state and the progress of the mappings in the Status and Progress columns of the table, as shown in Figure 10-60.

Mapping Name	Status	Source Volume	Target Volume	Progress	Group
fcmap0	✓ Copied	ITSO-FC-VOL-01	ITSO-FC-VOL-01_01	100%	
fcmap1	🔄 Copying	ITSO-FC-VOL-01_01	ITSO-FC-VOL-01_02	3%	
fcmap2	✓ Copied	ITSO-FC-VOL-02-S...	ITSO-FC-VOL-02-Target	100%	ITSO-FCCG-01
fcmap3	🔄 Copying	ITSO-FC-VOL-01_03	ITSO-FC-VOL-01_03_01	0%	
fcmap7	✓ Copied	ITSO-FC-VOL-01	ITSO-FC-VOL-01_04	100%	

Figure 10-60 FlashCopy mappings status and progress examples

FlashCopy Snapshots depend on the source volume and should be in a “copying” state if the mapping is started.



FlashCopy clones and the first occurrence of FlashCopy backup can take some time to complete, depending on the copyrate value and the size of the source volume. The next occurrences of FlashCopy backups are faster because only the changes that were made during two occurrences are copied.

For more information about FlashCopy starting operations and states, see 10.1.10, “Starting FlashCopy mappings and consistency groups” on page 507.

## 10.2.17 Stopping FlashCopy mappings

**Important:** Only FlashCopy mappings that do not belong to a consistency group can be stopped individually. If FlashCopy mappings are part of a consistency group, they can be stopped all together only by using the consistency group `stop` command.

The only reason to stop a FlashCopy mapping is for incremental FlashCopy. When the first occurrence of an incremental FlashCopy is started, a full copy of the source volume is made. When 100% of the source volume is copied, the FlashCopy mapping does not stop automatically and a manual stop can be performed. The target volume is available for read and write operations, during the copy, and after the mapping is stopped.

In any other case, stopping a FlashCopy mapping interrupts the copy and resets the bitmap table. Because only part of the data from the source volume was copied, the copied grains might be meaningless without the remaining grains. Therefore, the target volumes are placed offline and are unusable, as shown in Figure 10-61.

Mapping Name ↑	Status	Source Volume	Target Volume	Progress	Group
fcmap0	✓ Copied	ITSO-FC-VOL-01	ITSO-FC-VOL-01_01	100%	
fcmap1	✓ Copied	ITSO-FC-VOL-01_01	ITSO-FC-VOL-01_02	100%	
fcmap2	✓ Copied	ITSO-FC-VOL-02-S...	ITSO-FC-VOL-02-Target	100%	ITSO-FCCG-01
fcmap3	▲ Stopped	ITSO-FC-VOL-01_03	ITSO-FC-VOL-01_03_01	0%	
fcmap7	✓ Copied	ITSO-FC-VOL-01	ITSO-FC-VOL-01_04	100%	

Figure 10-61 Showing target volumes state and FlashCopy mappings status

To stop one or multiple FlashCopy mappings that do not belong to a consistency group, complete the following steps:

1. Open the FlashCopy Consistency Groups, or FlashCopy Mappings window.
2. Right-click the FlashCopy mappings that you want to stop and select **Stop**, as shown in Figure 10-62.

Mapping Name ↑	Status	Source Volume	Target Volume	Progress
fcmap0	✓ Copied	ITSO-FC-VOL-01	ITSO-FC-VOL-01_01	100%
fcmap1	✓ Copied	ITSO-FC-VOL-01_01	ITSO-FC-VOL-01_02	100%
fcmap2	✓ Copied	ITSO-FC-VOL-02-S...	ITSO-FC-VOL-02-Target	100%
fcmap3	🔄 Copying	ITSO-FC-VOL-01_03	ITSO-FC-VOL-01_03_01	0%
fcmap7	✓ Copied		-01_04	100%

Create FlashCopy Mapping

Move to Consistency Group

Remove from Consistency Group

Start

**Stop**

Rename Mapping

Delete Mapping

Showing 5 FC mappings | Selecting 1 FC mapping

Figure 10-62 Stopping FlashCopy mappings

**Note:** FlashCopy mappings can be in a stopping state for some time if you created dependencies between several targets. It is in a cleaning mode. For more information about dependencies and stopping process, see “Stopping process in a multiple target FlashCopy: Cleaning Mode” on page 513.

## 10.2.18 Memory allocation for FlashCopy

Copy Services features require that small amounts of volume cache be converted from cache memory into bitmap memory to allow the functions to operate at an I/O group level. If not enough bitmap space is allocated when you try to use one of the functions, you cannot complete the configuration. The total memory that can be dedicated to these functions is not defined by the physical memory in the system. The memory is constrained by the software functions that use the memory.

For every FlashCopy mapping that is created on an IBM Spectrum Virtualize system, a bitmap table is created to track the copied grains. By default, the system allocates 20 MiB of memory for a minimum of 10 TiB of FlashCopy source volume capacity and 5 TiB of incremental FlashCopy source volume capacity.

Depending on the grain size of the FlashCopy mapping, the memory capacity usage is different. One MiB of memory provides the following volume capacity for the specified I/O group:

- ▶ For clones and snapshots FlashCopy with 256 KiB grains size, 2 TiB of total FlashCopy source volume capacity

- ▶ For clones and snapshots FlashCopy with 64 KiB grains size, 512 GiB of total FlashCopy source volume capacity
- ▶ For incremental FlashCopy, with 256 KiB grains size, 1 TiB of total incremental FlashCopy source volume capacity
- ▶ For incremental FlashCopy, with 64 KiB grains size, 256 GiB of total incremental FlashCopy source volume capacity

Review Table 10-9 to calculate the memory requirements and confirm that your system can accommodate the total installation size.

Table 10-9 Memory allocation for FlashCopy services

Minimum allocated bitmap space	Default allocated bitmap space	Maximum allocated bitmap space	Minimum <sup>1</sup> functionality when using the default values
0	20 MiB	2 GiB	10 TiB of FlashCopy source volume capacity  5 TiB of incremental FlashCopy source volume capacity
<sup>1</sup> The actual amount of functionality might increase based on settings, such as grain size and strip size.			

FlashCopy includes the FlashCopy function, Global Mirror with Change Volumes (GMCV), and active-active (HyperSwap) relationships.

For multiple FlashCopy targets, you must consider the number of mappings. For example, for a mapping with a grain size of 256 KiB, 8 KiB of memory allows one mapping between a 16 GiB source volume and a 16 GiB target volume. Alternatively, for a mapping with a 256 KiB grain size, 8 KiB of memory allows two mappings between one 8 GiB source volume and two 8 GiB target volumes.

When creating a FlashCopy mapping, if you specify an I/O group other than the I/O group of the source volume, the memory accounting goes toward the specified I/O group, not toward the I/O group of the source volume.

When creating FlashCopy relationships or mirrored volumes, more bitmap space is allocated automatically by the system, if required.

For FlashCopy mappings, only one I/O group uses bitmap space. By default, the I/O group of the source volume is used.

When you create a reverse mapping, such as when you run a restore operation from a snapshot to its source volume, a bitmap is created.

When you configure change volumes for use with GM, two internal FlashCopy mappings are created for each change volume.

You can modify the resource allocation for each I/O group of an IBM SAN Volume Controller system by selecting **Settings** → **System** and clicking the **Resources** menu, as shown in Figure 10-63.

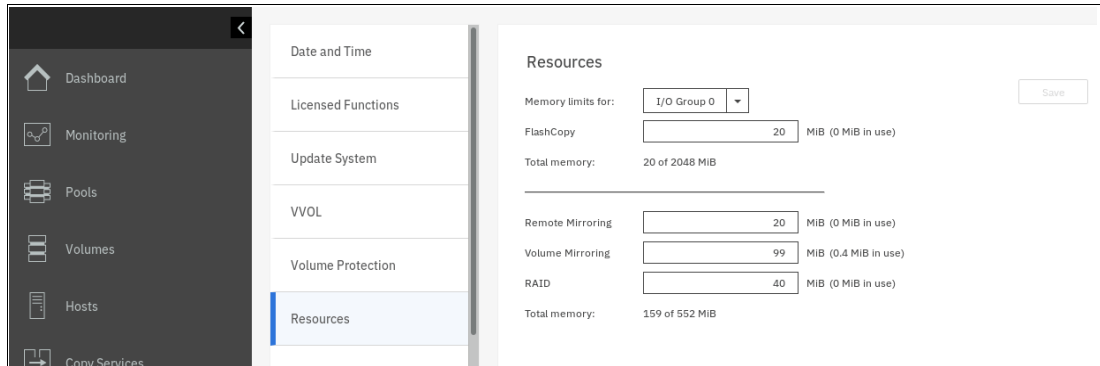


Figure 10-63 Modifying resources allocation per I/O group

## 10.3 Transparent Cloud Tiering

Introduced in V7.8, TCT is a function of IBM Spectrum Virtualize that uses IBM FlashCopy mechanisms to produce a PiT snapshot of the data. TCT helps to increase the flexibility to protect and transport data to public or private cloud infrastructure. This technology is built on top of IBM Spectrum Virtualize software capabilities. TCT uses the cloud to store snapshot targets and provides alternatives to restore snapshots from the private and public cloud of an entire volume or set of volumes.

TCT can help to solve business needs that require duplication of data of your source volume. Volumes can remain online and active while you create snapshot copies of the data sets. TCT operates below the host operating system and its cache. Therefore, the copy is not apparent to the host.

IBM Spectrum Virtualize features built-in software algorithms that allow the TCT function to securely interact; for example, with Information Dispersal Algorithms (IDA), which is essentially the interface to IBM Cloud Object Storage.

*Object Storage* is a general term that refers to the entity in which IBM Cloud Object Storage organizes, manages, and stores units of data. To transform these snapshots of traditional data into Object Storage, the storage nodes and the IDA import the data and transform it into several metadata and slices. The object can be read by using a subset of those slices. When an Object Storage entity is stored as IBM Cloud Object Storage, the objects must be manipulated or managed as a whole unit. Therefore, objects cannot be accessed or updated partially.

IBM Spectrum Virtualize uses internal software components to support HTTP-based REST application programming interface (API) to interact with an external cloud service provider (CSP) or private cloud.

For more information about the IBM Cloud Object Storage portfolio, see this [web page](#).

**Demonstration:** The IBM Client Demonstration Center has a demonstration available at this [web page](#) (log in required).

### 10.3.1 Considerations for using Transparent Cloud Tiering

TCT can help to address certain business needs. When considering whether to use TCT, adopt a combination of business and technical views of the challenges and determine whether TCT can solve both of those needs.

The use of TCT can help businesses to manipulate data as shown in the following examples:

- ▶ Creating a consistent snapshot of dynamically changing data
- ▶ Creating a consistent snapshot of production data to facilitate data movement or migration between systems that are running at different locations
- ▶ Creating a snapshot of production data sets for application development and testing
- ▶ Creating a snapshot of production data sets for quality assurance
- ▶ Using secure data tiering to off-premises cloud providers

From a technical standpoint, ensure that you evaluate the network capacity and bandwidth requirements to support your data migration to off-premises infrastructure. To maximize productivity, you must match your amount of data that must be transmitted off cloud plus your network capacity.

From a security standpoint, ensure that your on-premises or off-premises cloud infrastructure supports your requirements in terms of methods and level of encryption.

Regardless of your business needs, TCT within the IBM Spectrum Virtualize can provide opportunities to manage the exponential data growth and to manipulate data at low cost.

Today, many CSPs offers several *storage-as-services* solutions, such as content repository, backup, and archive. Combining all of these services, your IBM Spectrum Virtualize can help you solve many challenges that are related to rapid data growth, scalability, and manageability at attractive costs.

### 10.3.2 Transparent Cloud Tiering as backup solution and data migration

TCT can also be used as backup and data migration solution. In certain conditions, can be easily applied to eliminate the downtime that is associated with the needs to import and export data.

When TCT is applied as your backup strategy, IBM Spectrum Virtualize uses the same FlashCopy functions to produce *PiT* snapshot of an entire volume or set of volumes.

To ensure the integrity of the snapshot, it might be necessary to flush the host operating system and application cache of any outstanding reads or writes before the snapshot is performed. Failing to flush the host operating system and application cache can produce inconsistent and useless data.

Many operating systems and applications provide mechanism to stop I/O operations and ensure that all data is flushed from host cache. If these mechanisms are available, they can be used in combination with snapshot operations. When these mechanisms are not available, it might be necessary to flush the cache manually by quiescing the application and unmounting the file system or logical drives.

When choosing cloud Object Storage as a backup solution, be aware that the Object Storage must be managed as a whole. Backup and restore of individual files, folders, and partitions, are not possible.

To interact with CSPs or a private cloud, the IBM Spectrum Virtualize requires interaction with the correct architecture and specific properties. Conversely, CSPs offered attractive prices per Object Storage in cloud and deliver an easy-to-use interface. Normally, cloud providers offer low-cost prices for Object Storage space, and charges are applied for the cloud outbound traffic only.

### 10.3.3 Restoring data by using Transparent Cloud Tiering

TCT can also be used to restore data from any snapshot that is stored in cloud providers. When the cloud accounts' technical and security requirements are met, the storage objects in the cloud can be used as a data recovery solution. The recovery method is similar to back up, except that the reverse direction is applied.

TCT running on IBM Spectrum Virtualize queries for Object Storage stored in a cloud infrastructure. It enables users to restore the objects into a new volume or set of volumes.

This approach can be used for various applications, such as recovering your production database application after an errant batch process that caused extensive damage.

**Note:** Always consider the bandwidth characteristics and network capabilities when choosing to use TCT.

Restoring individual files by using TCT is not possible. Object Storage is unlike a file or a block; therefore, Object Storage must be managed as a whole unit piece of storage, and not partially. Cloud Object Storage is accessible by using an HTTP-based REST API.

### 10.3.4 Transparent Cloud Tiering restrictions

The following restrictions must be considered before TCT is used:

- ▶ Because the Object Storage is normally accessed by using the HTTP protocol on top of a TCP/IP stack, all traffic that is associated with cloud service flows through the node management ports.
- ▶ The size of cloud-enabled volumes cannot change. If the size of the volume changes, a snapshot must be created, so new Object Storage is constructed.
- ▶ TCT cannot be applied to volumes that are part of traditional copy services, such as FlashCopy, MM, GM, and HyperSwap.
- ▶ Volume containing two physical copies in two different storage pools cannot be part of TCT.
- ▶ Cloud Tiering snapshots cannot be taken from a volume that is part of migration activity across storage pools.
- ▶ Because VMware vSphere Virtual Volumes (VVOLs) are managed by a specific VMware application, these volumes are not candidates for TCT.
- ▶ File system volumes, such as volumes that are provisioned by the IBM Storwize V7000 Unified platform, are not qualified for TCT.

## 10.4 Implementing Transparent Cloud Tiering

This section describes the steps and requirements to implement TCT by using your IBM Spectrum Virtualize.

### 10.4.1 Domain Name System configuration

Because most of IBM Cloud Object Storage is managed and accessible by using the HTTP protocol, the Domain Name System (DNS) setting is an important requirement to ensure consistent resolution of domain names to internet resources.

Using your IBM Spectrum Virtualize management GUI, click **Settings** → **System** → **DNS** and insert your DNS Internet Protocol Version 4 (IPv4) or Internet Protocol Version 6 (IPv6). The DNS name can be anything that you want, and is used as a reference. Click **Save** after you complete the choices, as shown in Figure 10-64.

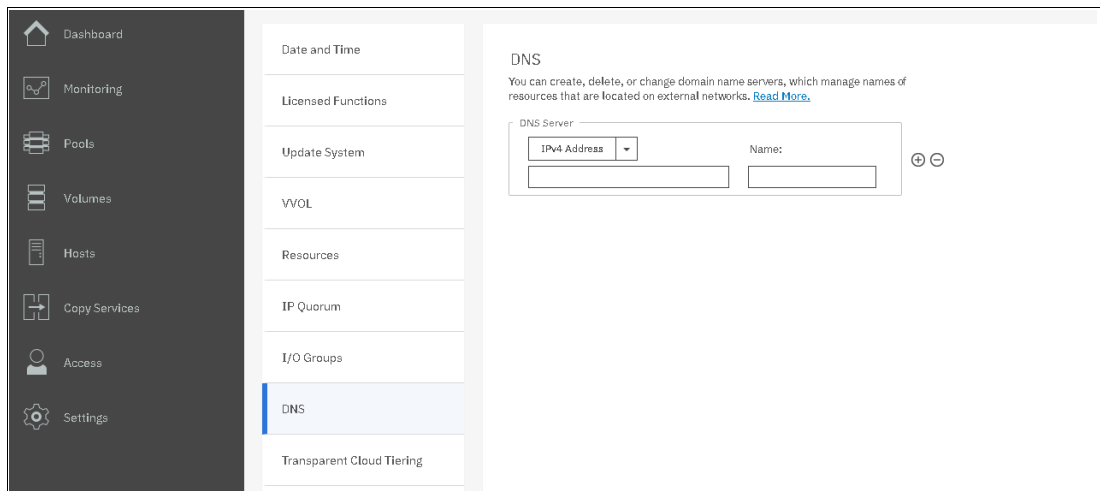


Figure 10-64 DNS settings

## 10.4.2 Enabling Transparent Cloud Tiering

After you complete the DNS settings, you can enable the TCT function in your IBM Spectrum Virtualize system by completing the following steps:

1. Using the IBM Spectrum Virtualize GUI, click **Settings** → **System** → **Transparent Cloud Tiering** and then, click **Enable Cloud Connection**, as shown in Figure 10-65. The TCT wizard starts and shows the welcome warning.

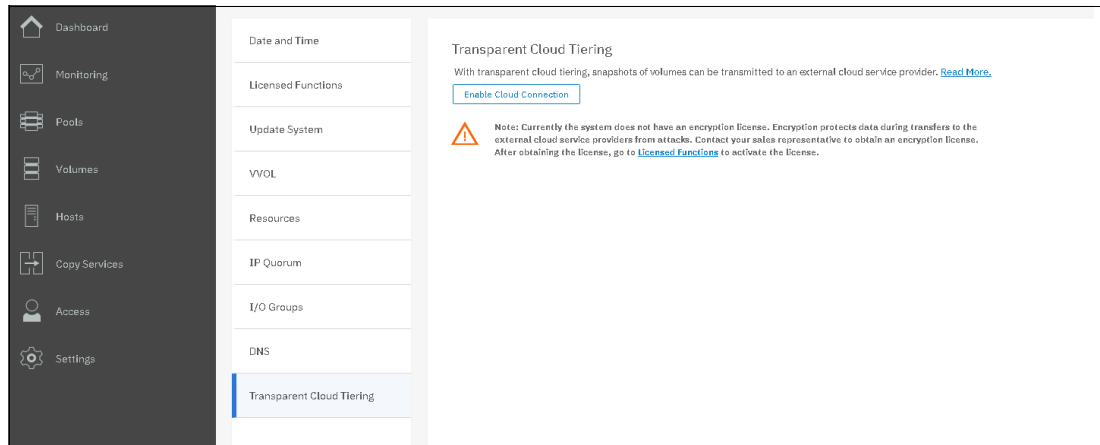


Figure 10-65 Enabling Cloud Tiering

**Note:** It is important to implement encryption before enabling cloud connecting. Encryption protects your data from attacks during the transfer to the external cloud service. Because the HTTP protocol is used to connect to cloud infrastructure, it is likely to start transactions by using the internet. For purposes of this writing, our system does not have encryption enabled.

2. Click **Next** to continue. You must select one of three CSPs:
  - IBM Cloud
  - OpenStack Swift
  - Amazon S3



Figure 10-66 shows the available options.

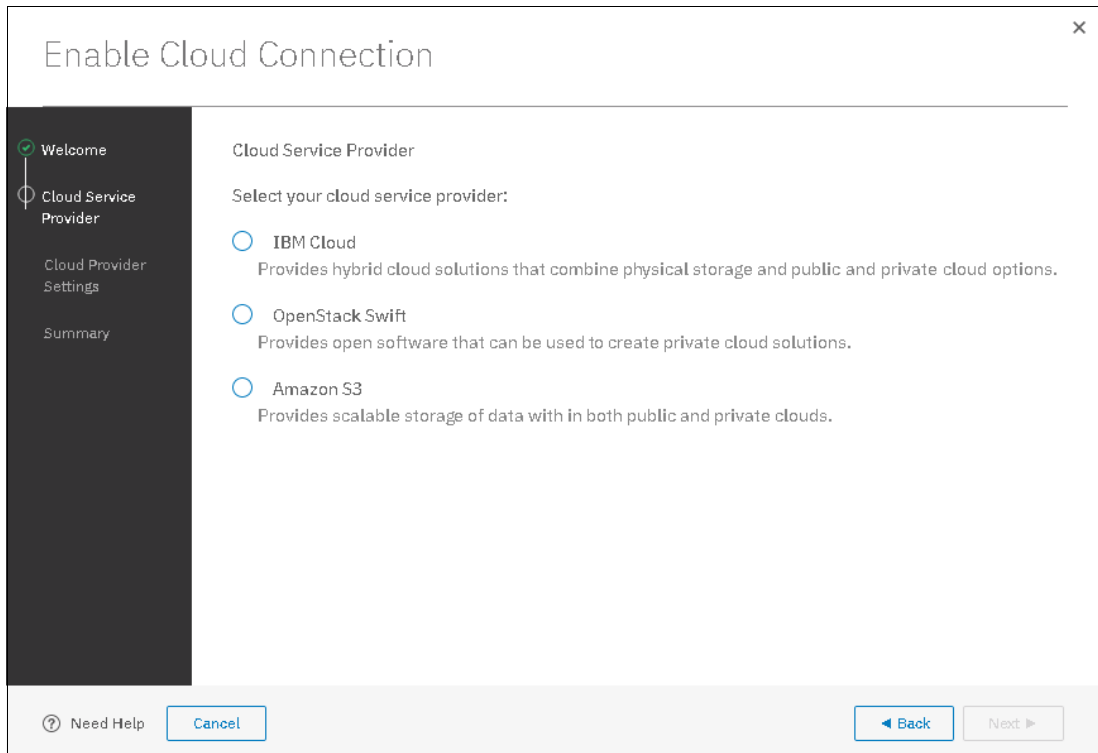


Figure 10-66 Selecting cloud service provider

3. In the next window, you must complete the settings of the Cloud Provider, credentials, and security access keys. The required settings can change depending on your CSP. An example of an empty form for an IBM Cloud connection is shown in Figure 10-67 on page 568.

Enable Cloud Connection

Cloud Provider Settings

IBM Cloud account

Object Storage URL:

Tenant:

User name:

API key:

Show characters

Container prefix:

Encryption  Enable

Bandwidth:

Upload:

No limit  Limit to:  Mbps

Download:

No limit  Limit to:  Mbps

Need Help

Figure 10-67 Entering cloud service provider information

4. Review your settings and click **Finish**, as shown in Figure 10-68.

Enable Cloud Connection

Summary

Provider: OpenStack Swift

Endpoint: http://9.71.48.122:8080/auth/v1.0

Keystone: Disabled

Encryption: Disabled

Max Upload bandwidth: No limit

Max Download bandwidth: No limit

Figure 10-68 Cloud Connection summary

- The cloud credentials can be viewed and updated at any time by using the function icons in left side of the GUI and clicking **Settings** → **Systems** → **Transparent Cloud Tiering**. From this window, you can also verify the status, the data usage statistics, and the upload and download bandwidth limits set to support this functionality.

In the account information window, you can visualize your cloud account information. This window also enables you to remove the account.

An example of visualizing your cloud account information is shown in Figure 10-69.

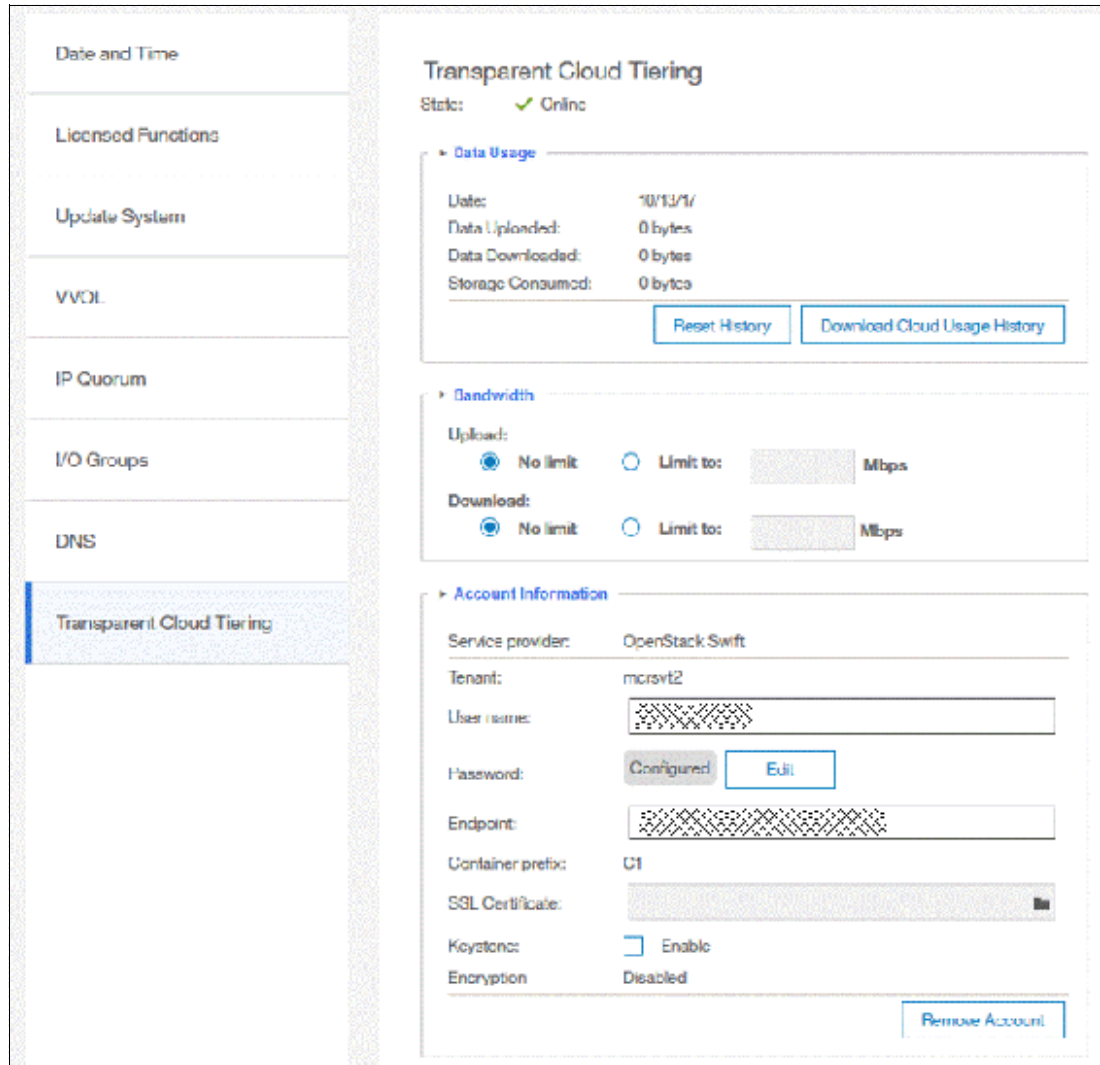


Figure 10-69 Enabled Transparent Cloud Tiering window

### 10.4.3 Creating cloud snapshots

To manage the cloud snapshots, the IBM Spectrum Virtualize provides a section in the GUI named Cloud Volumes. This section shows you how to add the volumes that are going to be part of the TCT. As described in 10.3.4, “Transparent Cloud Tiering restrictions” on page 564, cloud snapshot is available only for volumes that do not have a relationship to the list of restrictions previously mentioned.

Any volume can be added to the cloud volumes. However, snapshots work only for volumes that are not related to any other copy service.

To create and cloud snapshots, complete the following steps:

1. Click **Volumes** → **Cloud Volumes**, as shown in Figure 10-70.

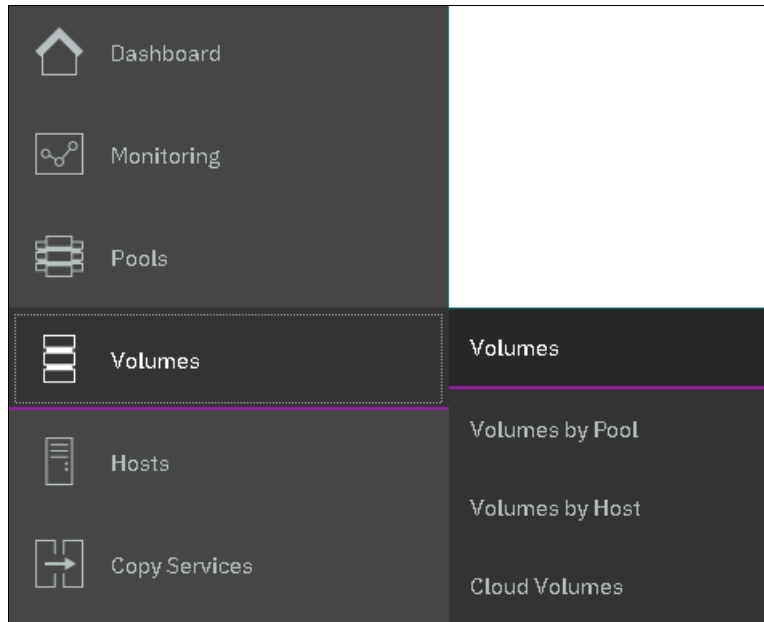


Figure 10-70 Cloud volumes menu

2. A new window opens, and you can use the GUI to select one or more volumes that you need to enable a cloud snapshot or you can add volumes to the list, as shown in Figure 10-71.

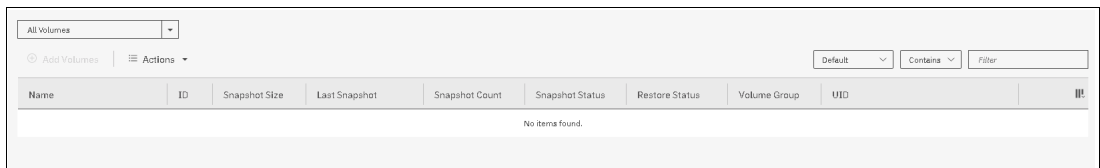


Figure 10-71 Cloud volumes window

3. Click **Add Volumes** to enable cloud-snapshot on volumes. A new window opens, as shown in Figure 10-72. Select the volumes that you want to enable Cloud Tiering for and click **Next**.

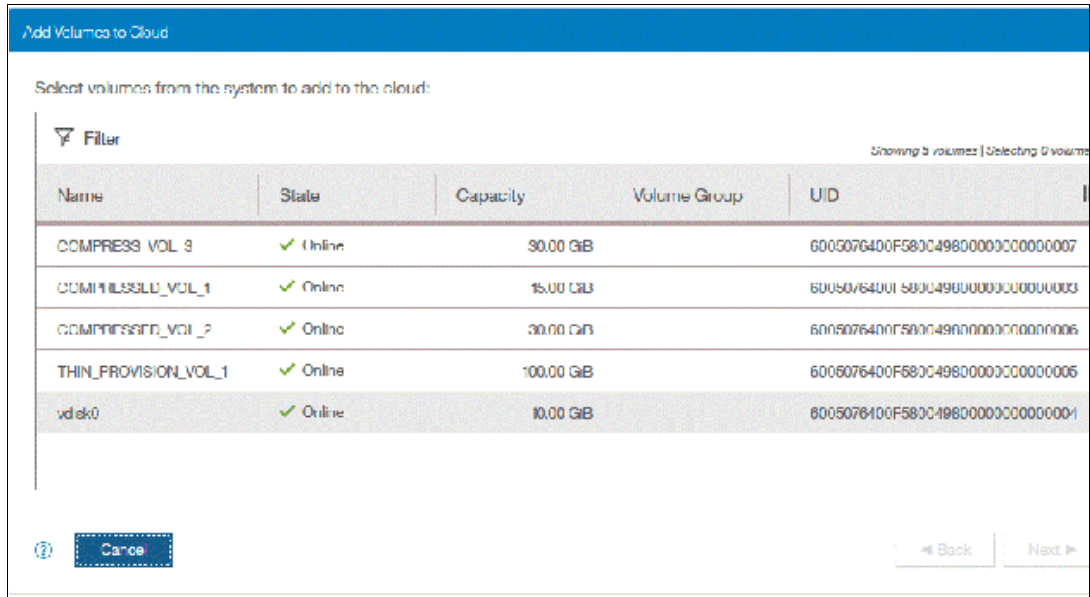


Figure 10-72 Adding volumes to Cloud Tiering

4. IBM Spectrum Virtualize GUI provides two options for you to select. If the first option is selected, the system decides what type of snapshot is created based on previous objects for each selected volume. If a full copy (full snapshot) of a volume was created, the system makes an incremental copy of the volume.

The second option creates a full snapshot of one or more selected volumes. You can select the second option for a first occurrence of a snapshot and click **Finish**, as shown in Figure 10-73. You can also select the second option, even if another full copy of the volume exists.

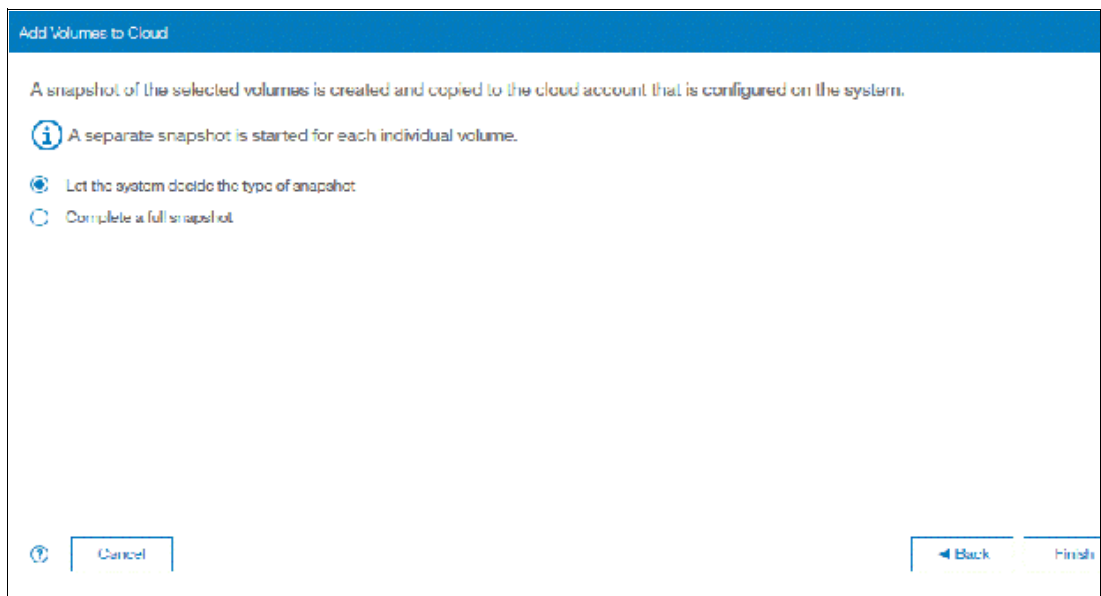


Figure 10-73 Selecting if a full copy is made or if the system decides

The **Cloud Volumes** window shows complete information about the volumes and their snapshots. The GUI shows the following information:

- Name of the volume
- ID of the volume assigned by the IBM Spectrum Virtualize
- Snapshot size
- Date and time that the last snapshot was created
- Number of snapshots that are taken for every volume
- Snapshot status
- Restore status
- Volume group for a set of volumes
- Volume unique identifier (UID)

Figure 10-74 shows an example of a Cloud Volumes list.

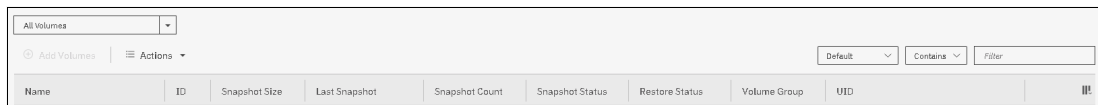


Figure 10-74 Cloud Volumes list example

5. Click the **Actions** menu in the Cloud Volumes window to create and manage snapshots. Also, you can use the menu to cancel, disable, and restore snapshots to volumes, as shown in Figure 10-75.

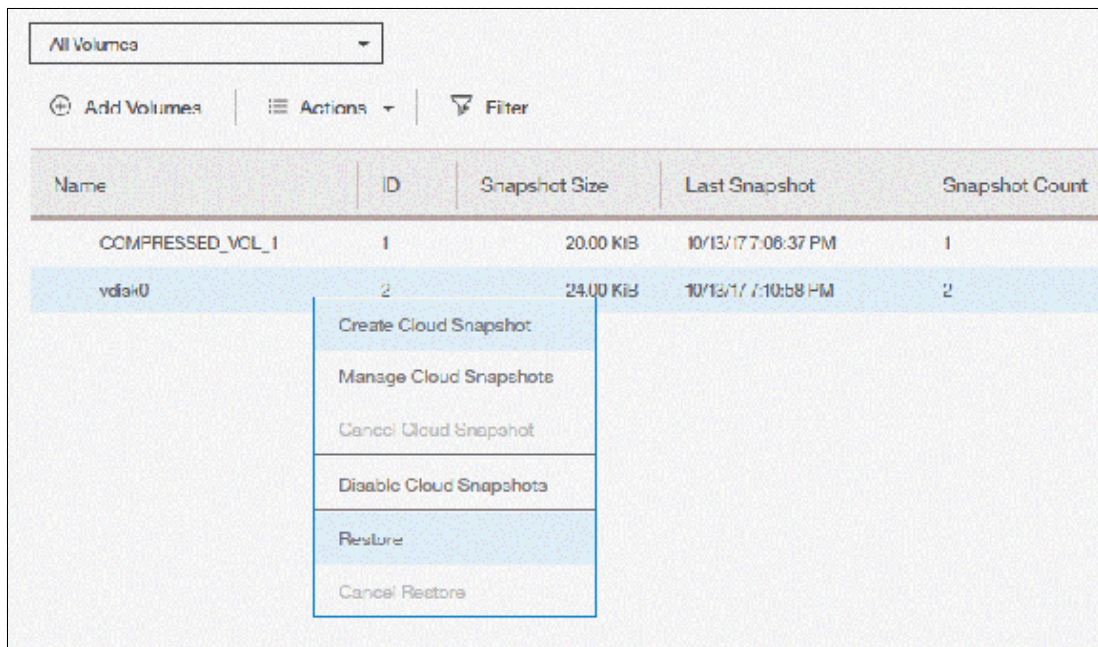


Figure 10-75 Available actions in Cloud Volumes window

#### 10.4.4 Managing cloud snapshots

To manage volume cloud snapshots, open the Cloud Volumes window, right-click the volume that you want to manage the snapshots from, and select **Manage Cloud Snapshot**.

“Managing” a snapshot is deleting one or multiple versions. The list of PiT copies appear and provide details about their status, type, and snapshot date, as shown in Figure 10-76 on page 573.



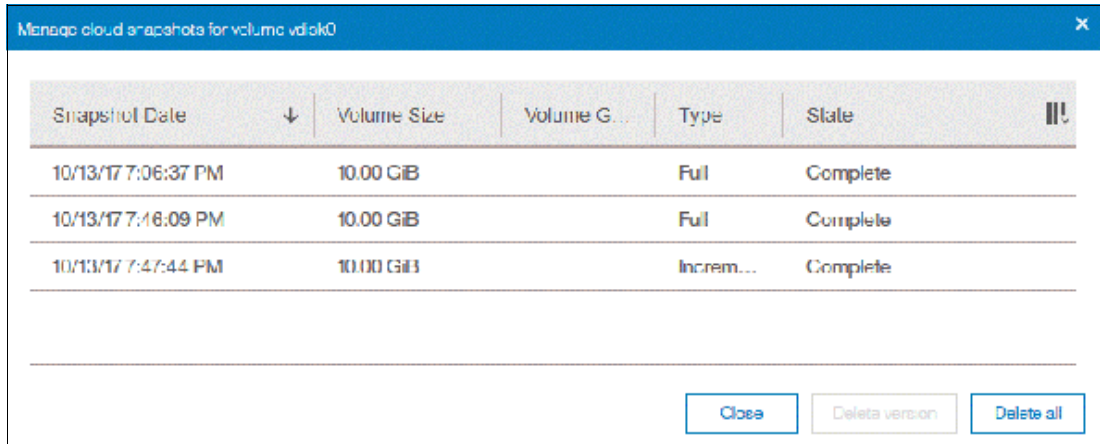


Figure 10-76 Deleting versions of a volume's snapshots

From this window, an administrator can delete old snapshots (old PiT copies) if they are no longer needed. The most recent copy cannot be deleted. If you want to delete the most recent copy, you must first disable Cloud Tiering for the specified volume.

### 10.4.5 Restoring cloud snapshots

This option allows IBM Spectrum Virtualize to restore snapshots from the cloud to the selected volumes or to create volumes with the restored data.

If the cloud account is shared among systems, IBM Spectrum Virtualize queries the snapshots that are stored in the cloud, and enables you to restore to a new volume. To restore a volume's snapshot, complete the following steps:

1. Open the Cloud Volumes window.
2. Right-click a volume and select **Restore**, as shown in Figure 10-77.

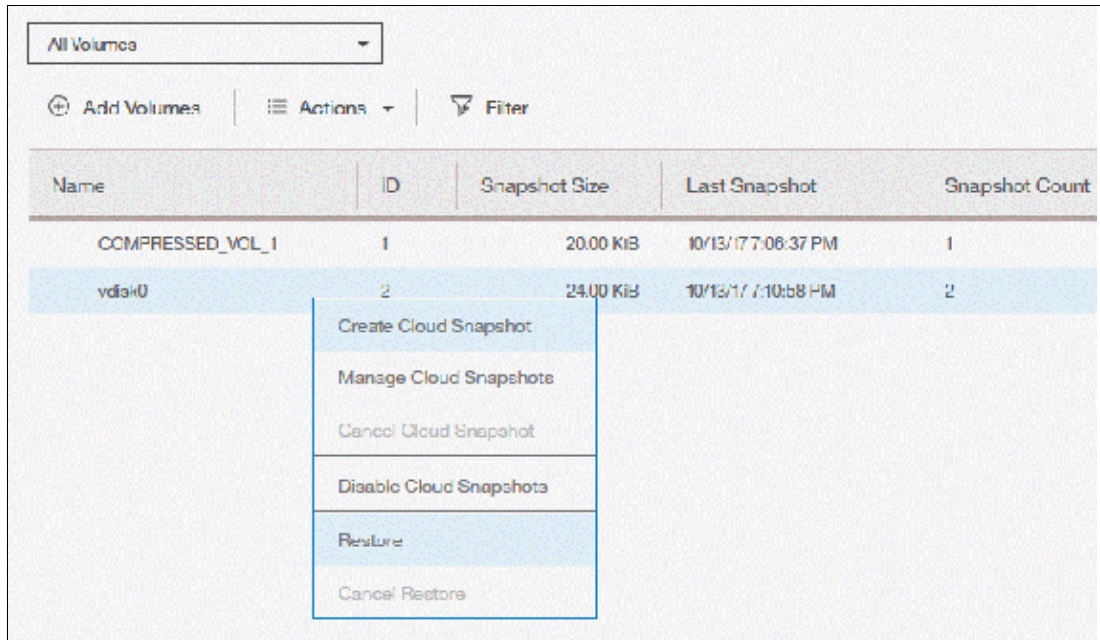


Figure 10-77 Selecting a volume to restore a snapshot from

3. A list of available snapshots is displayed. The snapshots date (PiT), their type (full or incremental), their state, and their size are shown (see Figure 10-78). Select the version that you want to restore and click **Next**.

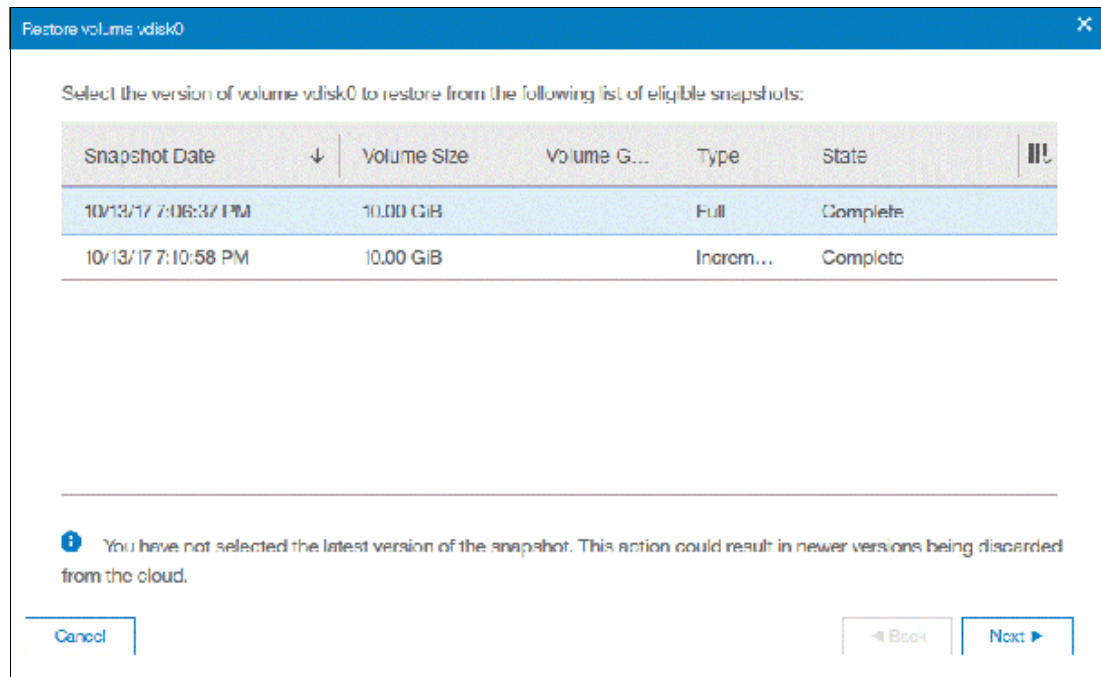


Figure 10-78 Selecting a snapshot version to restore

If the snapshot version that you selected has later generations (more recent Snapshot dates), the newer copies are removed from the cloud.

4. The IBM Spectrum Virtualize GUI provides two options to restore the snapshot from cloud. You can restore the snapshot from cloud directly to the selected volume, or create a volume to restore the data on, as shown in Figure 10-79. Make a selection and click **Next**.

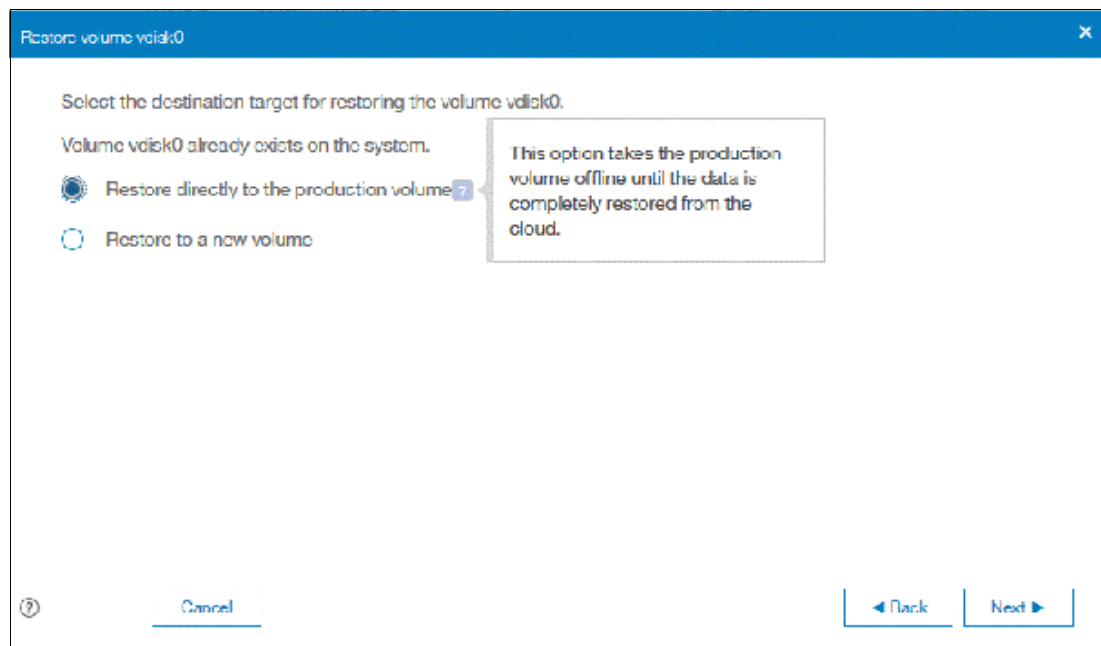


Figure 10-79 Restoring a snapshot on an existing volume or on a new volume

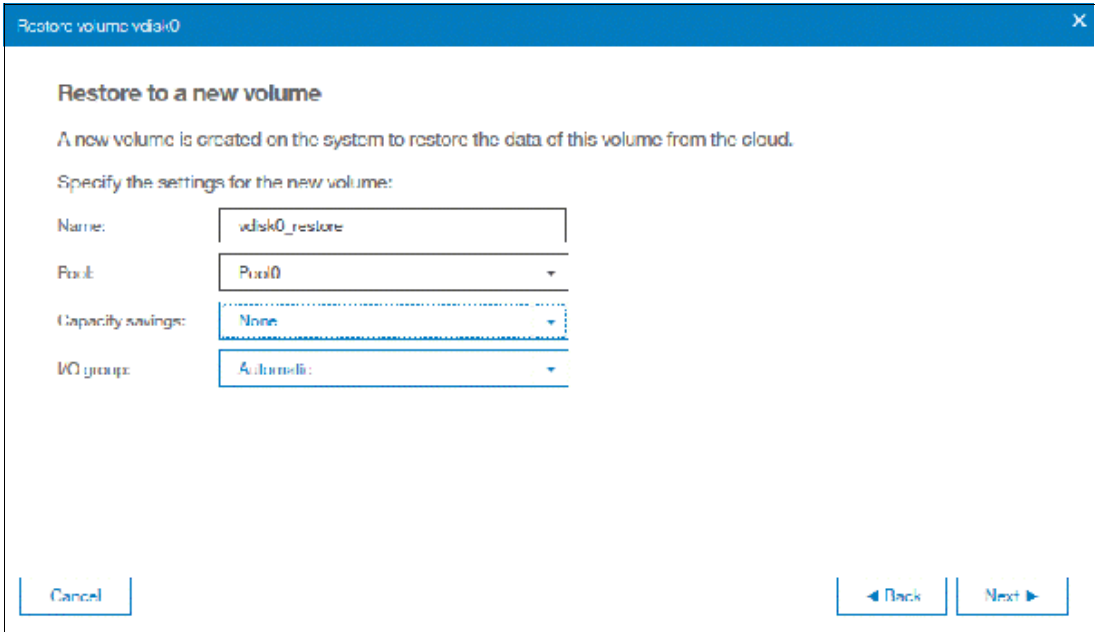


**Note:** Restoring a snapshot on the volume overwrites the data on the volume. The volume is taken offline (no read or write access) and the data from the PiT copy of the volume are written. The volume returns back online when all data is restored from the cloud.

5. If you selected the **Restore to a new Volume** option, you must enter the following information for the volume to be created with the snapshot data, as shown in Figure 10-80:
  - Name
  - Storage Pool
  - Capacity Savings (None, Compressed or Thin-provisioned)
  - I/O group

You are not asked to enter the volume size because the new volume's size is identical to the snapshot copy size

Enter the settings for the new volume and click **Next**.



Restore to a new volume

A new volume is created on the system to restore the data of this volume from the cloud.

Specify the settings for the new volume:

Name:

Pool:

Capacity savings:

I/O group:

Figure 10-80 Restoring a snapshot to a new volume

6. A Summary window is displayed so you can review the restoration settings, as shown in Figure 10-81 on page 576. Click **Finish**. The system creates a volume or overwrites the selected volume. The more recent snapshots (later versions) of the volume are deleted from the cloud.

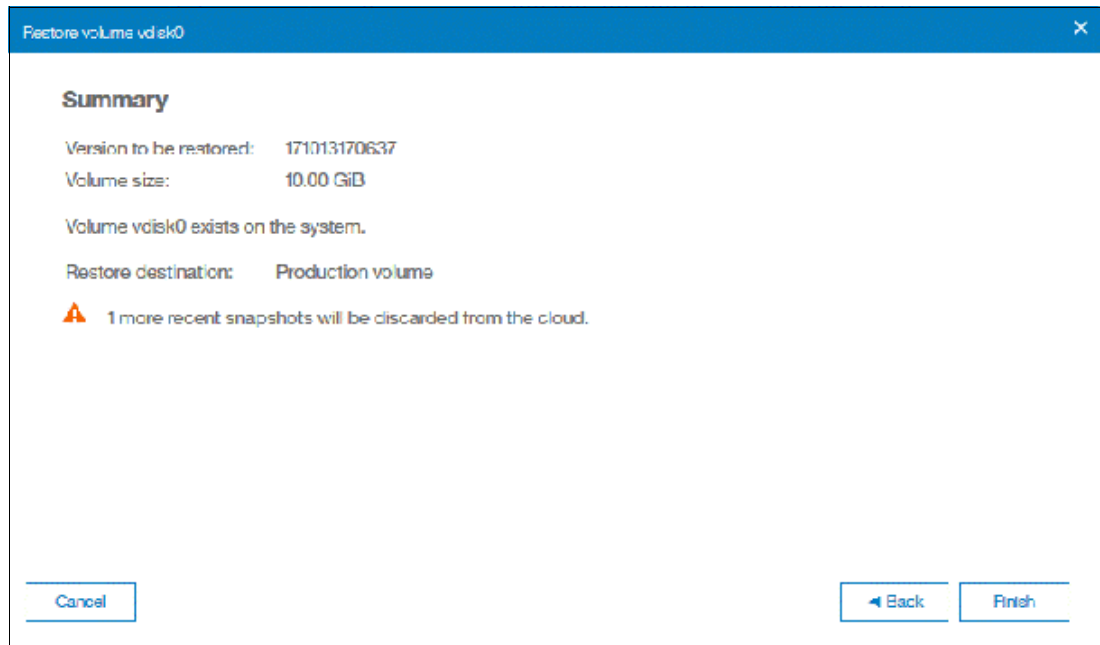


Figure 10-81 Restoring a snapshot summary

If you chose to restore the data from the cloud to a new volume, the new volume appears immediately in the volumes window. However, it is taken offline until all the data from the snapshot is written. The new volume is independent. It is not defined as a target in a FlashCopy mapping with the selected volume, for example. It also is not mapped to a host.

## 10.5 Volume mirroring and migration options

*Volume mirroring* is a simple Redundant Array of Independent Disks (RAID) 1-type function that enables a volume to remain online, even when the storage pool that is backing it becomes inaccessible. Volume mirroring is designed to protect the volume from storage infrastructure failures by seamless mirroring between storage pools.

Volume mirroring is provided by a specific volume mirroring function in the I/O stack. It cannot be manipulated like a FlashCopy or other types of copy volumes. However, this feature provides migration functionality, which can be obtained by splitting the mirrored copy from the source or by using the *migrate to* function. Volume mirroring cannot control backend storage mirroring or replication.

With volume mirroring, host I/O completes when both copies are written. This feature is enhanced with a tunable latency tolerance. This tolerance provides an option to give preference to losing the redundancy between the two copies. This tunable timeout value is Latency or Redundancy.

The Latency tuning option, which is set by running the **chvdisk -mirrorwritepriority Latency** command, is the default. It prioritizes host I/O latency, which yields a preference to host I/O over availability. However, you might need to give preference to redundancy in your environment when availability is more important than I/O response time. Run the **chvdisk -mirrorwritepriority redundancy** command to set the redundancy option.

Regardless of which option you choose, volume mirroring can provide extra protection for your environment.

Migration offers the following options:

► **Export to Image mode**

By using this option, you can move storage from managed mode to image mode, which is useful if you use the IBM SAN Volume Controller as a migration device. For example, vendor A's product cannot communicate with vendor B's product, but you must migrate data from vendor A to vendor B. By using Export to Image mode, you can migrate data by using Copy Services functions and then return control to the native array while maintaining access to the hosts.

► **Import to Image mode**

By using this option, you can import a storage MDisk or logical unit number (LUN) with its data from an external storage system without putting metadata on it so that the data remains intact. After you import it, all copy services functions can be used to migrate the storage to other locations while the data remains accessible to your hosts.

► **Volume migration by using volume mirroring and then by using Split into New Volume**

By using this option, you can use the available RAID 1 functions. You create two copies of data that initially has a set relationship (one volume with two copies, one primary and one secondary) but then break the relationship (two volumes, both primary and no relationship between them) to make them independent copies of data.

You can use this option to migrate data between storage pools and devices. You might use this option if you want to move volumes to multiple storage pools. Each volume can have two copies at a time, which means that you can add only one copy to the original volume, and then you must split those copies to create another copy of the volume.

► **Volume migration by using move to another pool**

By using this option, you can move any volume between storage pools without any interruption to the host access. This option is a quicker version of the Volume Mirroring and Split into New Volume option. You might use this option if you want to move volumes in a single step, or you do not have a volume mirror copy.

**Migration:** Although these migration methods do not disrupt access, a brief outage does occur to install the host drivers for your IBM SAN Volume Controller if they are not yet installed.

With volume mirroring, you can move data to different MDisks within the same storage pool or move data between different storage pools. The use of volume mirroring over volume migration is beneficial because with volume mirroring, storage pools do not need to have the same extent size as is the case with volume migration.

**Note:** Volume mirroring does not create a second volume before you split copies. Volume mirroring adds a second copy of the data under the same volume. Therefore, you have one volume that is presented to the host with two copies of data that are connected to this volume. Only splitting copies creates another volume, and then both volumes have only one copy of the data.

Starting with V7.3 and the introduction of the new cache architecture, mirrored volume performance was significantly improved. Now, lower cache is beneath the volume mirroring layer, which means that both copies have their own cache.

This approach helps when you have copies of different types; for example, generic and compressed, because now both copies use its independent cache and performs its own read prefetch. Destaging of the cache can now be done independently for each copy, so one copy does not affect performance of a second copy.

Also, because the IBM Spectrum Virtualize destage algorithm is MDisk aware, it can tune or adapt the destaging process, depending on MDisk type and usage, for each copy independently.

For more information about Volume Mirroring, see Chapter 6, “Volumes” on page 255.

## 10.6 Remote Copy

This section describes the RC services, which are a synchronous RC called MM, asynchronous RC that is called GM, and GMCV. RC in an IBM Spectrum Virtualize system is similar to RC in the IBM System Storage DS8000 family at a functional level, but the implementation differs.

IBM Spectrum Virtualize provides a single point of control when RC is enabled in your network (regardless of the disk subsystems that are used) if those disk subsystems are supported by the system.

The general application of RC services is to maintain two real-time synchronized copies of a volume. Often, the two copies are geographically dispersed between two IBM Spectrum Virtualize systems. However, it is possible to use MM or GM within a single system (within an I/O group). If the master copy fails, you can enable an auxiliary copy for I/O operations.

**Tips:** Intracluster MM/GM uses more resources within the system when compared to an intercluster MM/GM relationship, where resource allocation is shared between the systems. Use intercluster MM/GM when possible. For mirroring volumes in the same system, it is better to use volume mirroring or the FlashCopy feature.

A typical application of this function is to set up a dual-site solution that uses two IBM SAN Volume Controller or IBM Storwize systems. The first site is considered the *primary site* or *production site*, and the second site is considered the *backup site* or *failover site*. The failover site is activated when a failure at the first site is detected.

When MM or GM are used, a certain amount of bandwidth is required for the system intercluster heartbeat traffic. The amount of traffic depends on how many nodes are in each of the two clustered systems.

Table 10-10 lists the amount of heartbeat traffic (in megabits per second) that is generated by various sizes of clustered systems.

Table 10-10 Intersystem heartbeat traffic in Mbps

IBM Spectrum Virtualize system 1	IBM Spectrum Virtualize system 2			
	2 nodes	4 nodes	6 nodes	8 nodes
2 nodes	5	6	6	6
4 nodes	6	10	11	12

IBM Spectrum Virtualize system 1	IBM Spectrum Virtualize system 2			
6 nodes	6	11	16	17
8 nodes	6	12	17	21

### 10.6.1 IBM SAN Volume Controller and IBM Storwize system layers

An IBM Spectrum Virtualize based system can be in one of the two layers: the *replication* layer or the *storage* layer. The system layer affects how the system interacts with IBM SAN Volume Controller systems and other external IBM Spectrum Virtualize based systems. The IBM SAN Volume Controller is always set to replication layer. This parameter *cannot* be changed.

In the storage layer, an IBM Storwize/FlashSystem 9100 family system has the following characteristics and requirements:

- ▶ The system can perform MM and GM replication with other storage-layer systems.
- ▶ The system can provide external storage for replication-layer systems or IBM SAN Volume Controller.
- ▶ The system cannot use a storage-layer system as external storage.

In the replication layer, an IBM SAN Volume Controller or an IBM Storwize system has the following characteristics and requirements:

- ▶ Can perform MM and GM replication with other replication-layer systems
- ▶ Cannot provide external storage for a replication-layer system
- ▶ Can use a storage-layer system as external storage

An IBM Storwize family system is in the storage layer by default, but the layer can be changed. For example, you might want to change an IBM Storwize V7000 to a replication layer if you want to virtualize other IBM Storwize systems or to replicate to an IBM SAN Volume Controller system.

**Note:** Before you change the system layer, the following conditions must be met on the system at the time of layer change:

- ▶ No other IBM Storwize or FlashSystem 9100 can exist as a backend or host entity.
- ▶ No system partnerships can exist.
- ▶ No IBM Storwize or FlashSystem 9100 system can be visible on the SAN fabric.

The layer can be changed during normal host I/O.

In your IBM Storwize system, run the `lssystem` command to check the current system layer, as shown in Example 10-2.

*Example 10-2 Output from the lssystem command showing the system layer*

```
IBM_Storwize:mcr-fab2-cluster-07:superuser>lssystem
id 00000100298056C2
name mcr-fab2-cluster-07
...
8.3.1.0 (build 150.15.1910111241200)
...
```

layer storage

...

**Note:** Consider the following rules for creating remote partnerships between the IBM SAN Volume Controller and IBM Storwize Family systems:

- ▶ An IBM SAN Volume Controller is always in the replication layer.
- ▶ By default, the IBM Storwize systems are in the storage layer, but can be changed to the replication layer.
- ▶ A system can form partnerships only with systems in the same layer.
- ▶ Starting in V6.4, an IBM SAN Volume Controller or Storwize system in the replication layer can virtualize an IBM Storwize system in the storage layer.

## 10.6.2 Multiple IBM Spectrum Virtualize systems replication

Each IBM Spectrum Virtualize system can maintain up to three partner system relationships, which enables as many as four systems to be directly associated with each other. This system partnership capability enables the implementation of disaster recovery (DR) solutions.

**Note:** For more information about restrictions and limitations of native IP replication, see 10.8.2, “IP partnership limitations” on page 615.

Figure 10-82 shows an example of a multiple system mirroring configuration.

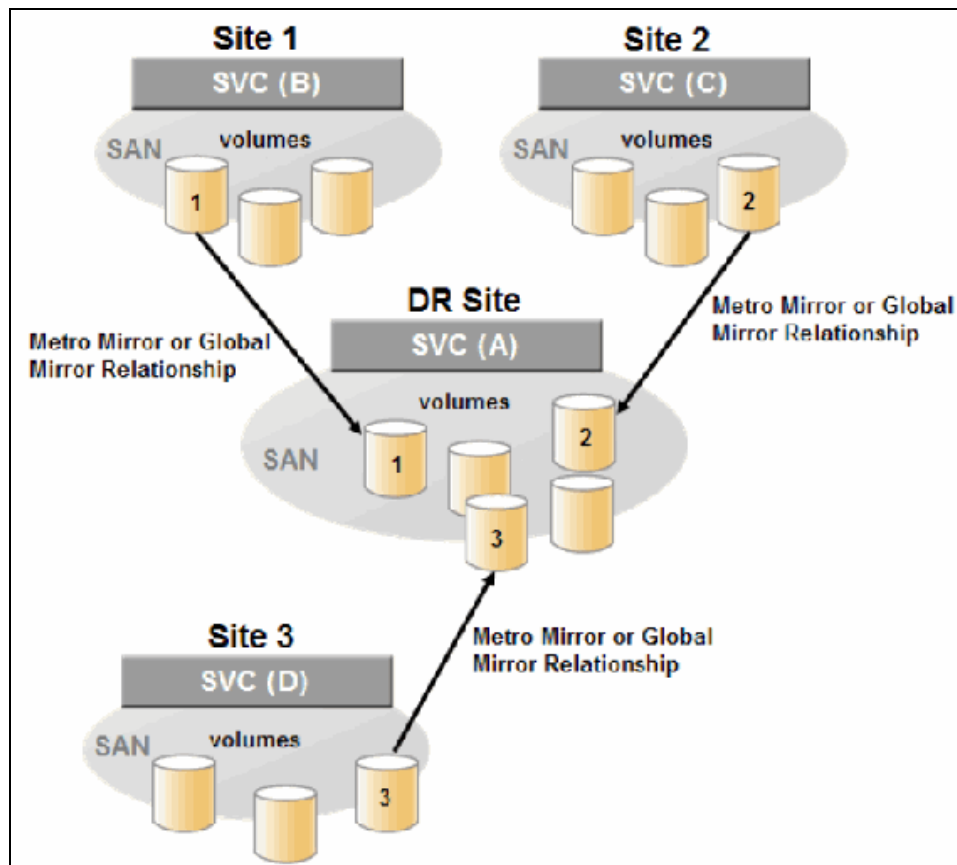


Figure 10-82 Multiple system mirroring configuration example

## Supported multiple system mirroring topologies

Multiple system mirroring supports various partnership topologies, as shown in the example in Figure 10-83. This example is a star topology ( $A \rightarrow B$ ,  $A \rightarrow C$ , and  $A \rightarrow D$ ).

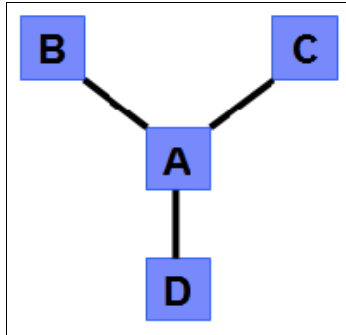


Figure 10-83 Star topology

Figure 10-83 shows four systems in a star topology, with System A at the center. System A can be a central DR site for the three other locations.

By using a star topology, you can migrate applications by using a process, such as the one described in the following example:

1. Suspend application at A.
2. Remove the  $A \rightarrow B$  relationship.
3. Create the  $A \rightarrow C$  relationship (or the  $B \rightarrow C$  relationship).
4. Synchronize to system C, and ensure that  $A \rightarrow C$  is established:
  - $A \rightarrow B$ ,  $A \rightarrow C$ ,  $A \rightarrow D$ ,  $B \rightarrow C$ ,  $B \rightarrow D$ , and  $C \rightarrow D$
  - $A \rightarrow B$ ,  $A \rightarrow C$ , and  $B \rightarrow C$

Figure 10-84 shows an example of a triangle topology ( $A \rightarrow B$ ,  $A \rightarrow C$ , and  $B \rightarrow C$ ).

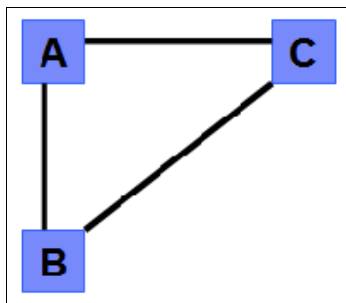


Figure 10-84 Triangle topology

Figure 10-85 shows an example of an IBM SAN Volume Controller fully connected topology (A → B, A → C, A → D, B → D, and C → D).

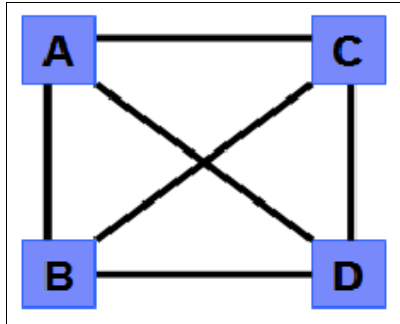


Figure 10-85 Fully connected topology

Figure 10-85 shows a fully connected mesh in which every system has a partnership to each of the three other systems. This topology enables volumes to be replicated between any pair of systems; for example, A → B, A → C, and B → C.

Figure 10-86 shows a daisy-chain topology.

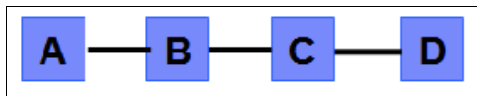


Figure 10-86 Daisy-chain topology

Although systems can have up to three partnerships, volumes can be part of only one RC relationship; for example, A → B.

**System partnership intermix:** All of these topologies are valid for the intermix of the IBM SAN Volume Controller with the IBM Storwize V7000/FlashSystem 9100 if the IBM Storwize V7000/FlashSystem 9100 is set to the replication layer.

### 10.6.3 Importance of write ordering

Many applications that use block storage are required to survive failures, such as loss of power or a software crash, and to not lose data that existed before the failure. Because many applications must perform many update operations in parallel, maintaining write ordering is key to ensure the correct operation of applications after a disruption.

An application that performs many database updates is designed with the concept of dependent writes. With dependent writes, it is important to ensure that an earlier write completed before a later write is started. Reversing or performing the order of writes differently than the application intended can undermine the application's algorithms and can lead to problems, such as detected or undetected data corruption.

The IBM Spectrum Virtualize MM and GM implementation operates in a manner that is designed to always keep a consistent image at the secondary site. The GM implementation uses complex algorithms that identify sets of data and number those sets of data in sequence. The data is then applied at the secondary site in the defined sequence.



Operating in this manner ensures that if the relationship is in a `Consistent_Synchronized` state, the GM target data is at least crash consistent and supports quick recovery through application crash recovery facilities.

For more information about dependent writes, see 10.1.13, “FlashCopy and image mode Volumes” on page 516.

## Remote Copy consistency groups

An RC consistency group can contain an arbitrary number of relationships up to the maximum number of MM/GM relationships that is supported by the IBM Spectrum Virtualize system. MM/GM commands can be issued to an RC consistency group.

Therefore, these commands can be issued simultaneously for all MM/GM relationships that are defined within that consistency group, or to a single MM/GM relationship that is not part of an RC consistency group. For example, when a `startrcconsistgrp` command is issued to the consistency group, all of the MM/GM relationships in the consistency group are started at the same time.

Figure 10-87 shows the concept of RC consistency groups.

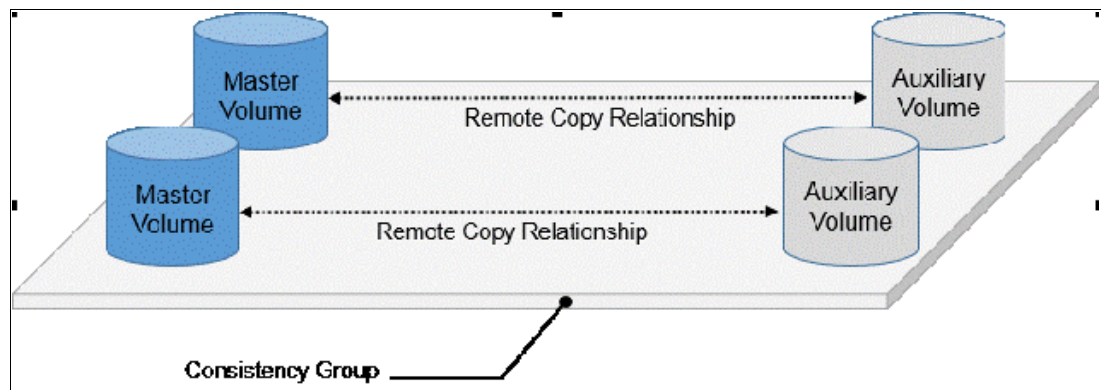


Figure 10-87 Remote Copy consistency group

Certain uses of MM/GM require the manipulation of more than one relationship. RC consistency groups can group relationships so that they are manipulated in unison.

Consider the following points:

- ▶ MM/GM relationships can be part of a consistency group, or they can be stand-alone and, therefore, are handled as single instances.
- ▶ A consistency group can contain zero or more relationships. An empty consistency group with zero relationships in it has little purpose until it is assigned its first relationship, except that it has a name.
- ▶ All relationships in a consistency group must have corresponding master and auxiliary volumes.
- ▶ All relationships in one consistency group must be the same type; for example, only MM or only GM.

Although consistency groups can be used to manipulate sets of relationships that do not need to satisfy these strict rules, this manipulation can lead to unwanted side effects. The rules behind a consistency group mean that certain configuration commands are prohibited. These configuration commands are not prohibited if the relationship is not part of a consistency group.

For example, consider the case of two applications that are independent, yet they are placed into a single consistency group. If an error occurs, synchronization is lost and a background copy process is required to recover synchronization. While this process is progressing, MM/GM rejects attempts to enable access to the auxiliary volumes of either application.

If one application finishes its background copy more quickly than the other application, MM/GM still refuses to grant access to its auxiliary volumes, even though it is safe in this case. The MM/GM policy is to refuse access to the entire consistency group if any part of it is inconsistent. Stand-alone relationships and consistency groups share a common configuration and state models. All of the relationships in a non-empty consistency group feature the same state as the consistency group.

## 10.6.4 Remote Copy intercluster communication

In the traditional Fibre Channel (FC), the intercluster communication between systems in a MM/GM partnership is performed over the SAN. This section describes this communication path.

For more information about intercluster communication between systems in an IP partnership, see 10.8.6, “States of IP partnership” on page 619.

### Zoning

At least two FC ports of every node of each system must communicate with each other to create the partnership. Switch zoning is critical to facilitate intercluster communication.

### Intercluster communication channels

When an IBM Spectrum Virtualize system partnership is defined on a pair of systems, the following intercluster communication channels are established:

- ▶ A single control channel, which is used to exchange and coordinate configuration information
- ▶ I/O channels between each of these nodes in the systems

These channels are maintained and updated as nodes and links appear and disappear from the fabric, and are repaired to maintain operation where possible. If communication between the systems is interrupted or lost, an event is logged (and the MM/GM relationships stop).

**Alerts:** You can configure the system to raise SNMP traps to the enterprise monitoring system to alert on events that indicate an interruption in internode communication occurred.

### Intercluster links

All IBM SAN Volume Controller nodes maintain a database of other devices that are visible on the fabric. This database is updated as devices appear and disappear.

Devices that advertise themselves as IBM SAN Volume Controller or IBM Storwize V7000 nodes are categorized according to the system to which they belong. Nodes that belong to the same system establish communication channels between themselves and exchange messages to implement clustering and the functional protocols of IBM Spectrum Virtualize.

Nodes that are in separate systems do not exchange messages after initial discovery is complete, unless they are configured together to perform an RC relationship.

The intercluster link carries control traffic to coordinate activity between two systems. The link is formed between one node in each system. The traffic between the designated nodes is distributed among logins that exist between those nodes.

If the designated node fails (or all of its logins to the remote system fail), a new node is chosen to carry control traffic. This node change causes the I/O to pause, but it does not put the relationships in a `ConsistentStopped` state.

**Note:** Run the `chsystem` command with `-partnerfcportmask` to dedicate several FC ports only to system-to-system traffic to ensure that RC is not affected by other traffic, such as host-to-node traffic or node-to-node traffic within the same system.

## 10.6.5 Metro Mirror overview

MM establishes a synchronous relationship between two volumes of equal size. The volumes in an MM relationship are referred to as the *master* (primary) volume and the *auxiliary* (secondary) volume. Traditional FC MM is primarily used in a metropolitan area or geographical area, up to a maximum distance of 300 km (186.4 miles) to provide synchronous replication of data.

With synchronous copies, host applications write to the master volume, but they do not receive confirmation that the write operation completed until the data is written to the auxiliary volume. This action ensures that both the volumes have identical data when the copy completes. After the initial copy completes, the MM function always maintains a fully synchronized copy of the source data at the target site.

MM has the following characteristics:

- ▶ Zero recovery point objective (RPO)
- ▶ Synchronous
- ▶ Production application performance that is affected by round-trip latency

Increased distance directly affects host I/O performance because the writes are synchronous. Use the requirements for application performance when you are selecting your MM auxiliary location.

Consistency groups can be used to maintain data integrity for dependent writes, which is similar to FlashCopy consistency groups.

IBM Spectrum Virtualize provides intracluster and intercluster MM, which are described next.

### Intracluster Metro Mirror

Intracluster MM performs the intracluster copying of a volume, in which both volumes belong to the same system and I/O group within the system. Because it is within the same I/O group, sufficient bitmap space must exist within the I/O group for both sets of volumes and licensing on the system.

**Important:** Performing MM across I/O groups within a system is not supported.

### Intercluster Metro Mirror

Intercluster MM performs intercluster copying of a volume, in which one volume belongs to a system and the other volume belongs to a separate system.

Two IBM Spectrum Virtualize systems must be defined in a partnership, which must be performed on both systems to establish a fully functional MM partnership.

By using standard single-mode connections, the supported distance between two systems in an MM partnership is 10 km (6.2 miles), although greater distances can be achieved by using extenders. For extended distance solutions, contact your IBM representative.

**Limit:** When a local fabric and a remote fabric are connected for MM purposes, the inter-switch link (ISL) hop count between a local node and a remote node cannot exceed seven.

### 10.6.6 Synchronous Remote Copy

MM is a fully synchronous RC technique that ensures that writes are committed at the master and auxiliary volumes before write completion is acknowledged to the host, but only if writes to the auxiliary volumes are possible.

Events, such as a loss of connectivity between systems, can cause mirrored writes from the master volume and the auxiliary volume to fail. In that case, MM suspends writes to the auxiliary volume and enables I/O to the master volume to continue to avoid affecting the operation of the master volumes.

Figure 10-88 shows how a write to the master volume is mirrored to the cache of the auxiliary volume before an acknowledgment of the write is sent back to the host that issued the write. This process ensures that the auxiliary is synchronized in real time if it is needed in a failover situation.

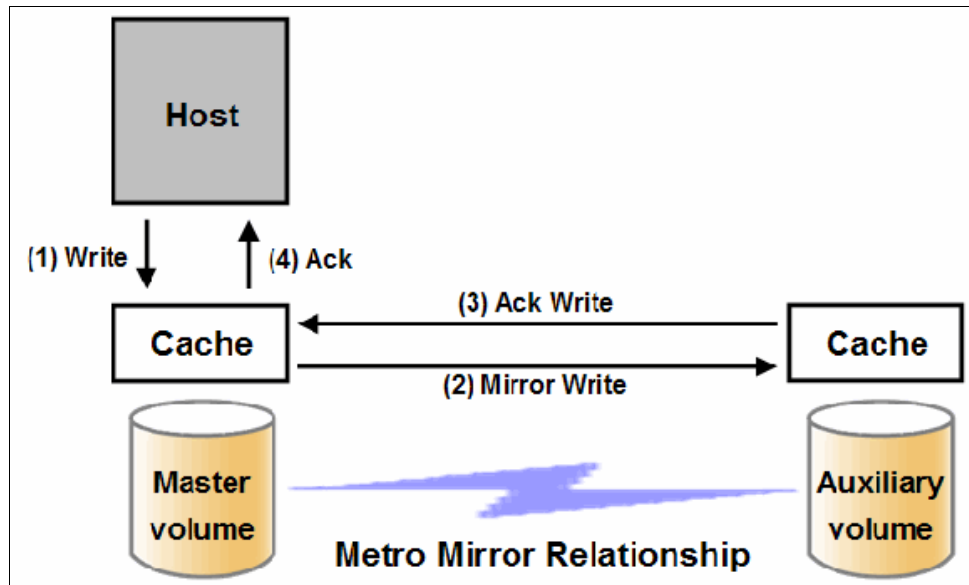


Figure 10-88 Write on volume in Metro Mirror relationship

However, this process also means that the application is exposed to the latency and bandwidth limitations (if any) of the communication link between the master and auxiliary volumes. This process might lead to unacceptable application performance, particularly when placed under peak load. Therefore, the use of traditional FC MM has distance limitations that are based on your performance requirements. IBM Spectrum Virtualize does not support more than 300 km (186.4 miles).

## 10.6.7 Metro Mirror features

The IBM Spectrum Virtualize MM function supports the following features:

- ▶ Synchronous RC of volumes that are dispersed over metropolitan distances.
- ▶ The MM relationships between volume pairs, with each volume in a pair that is managed by an IBM Storwize V7000 system or IBM SAN Volume Controller system (requires V6.3.0 or later).
- ▶ Supports intracluster MM where both volumes belong to the same system (and I/O group).
- ▶ IBM Spectrum Virtualize supports intercluster MM where each volume belongs to a separate system. You can configure a specific system for partnership with another system. All intercluster MM processing occurs between two IBM Spectrum Virtualize systems that are configured in a partnership.
- ▶ Intercluster and intracluster MM can be used concurrently.
- ▶ IBM Spectrum Virtualize does not require that a control network or fabric is installed to manage MM. For intercluster MM, the system maintains a control link between two systems. This control link is used to control the state and coordinate updates at either end. The control link is implemented on top of the same FC fabric connection that the system uses for MM I/O.
- ▶ IBM Spectrum Virtualize implements a configuration model that maintains the MM configuration and state through major events, such as failover, recovery, and resynchronization, to minimize user configuration action through these events.

IBM Spectrum Virtualize supports the resynchronization of changed data so that write failures that occur on the master or auxiliary volumes do not require a complete resynchronization of the relationship.

## 10.6.8 Metro Mirror attributes

The MM function in IBM Spectrum Virtualize features the following attributes:

- ▶ A partnership is created between two IBM Spectrum Virtualize systems that are operating in the replication layer (for intercluster MM).
- ▶ An MM relationship is created between two volumes of the same size.
- ▶ To manage multiple MM relationships as one entity, relationships can be made part of an MM Consistency Group, which ensures data consistency across multiple MM relationships and provides ease of management.
- ▶ When an MM relationship is started and when the background copy completes, the relationship becomes consistent and synchronized.
- ▶ After the relationship is synchronized, the auxiliary volume holds a copy of the production data at the primary, which can be used for DR.
- ▶ The auxiliary volume is in read-only mode when relationship is active.
- ▶ To access the auxiliary volume, the MM relationship must be stopped with the access option enabled before write I/O is allowed to the auxiliary.
- ▶ The remote host server is mapped to the auxiliary volume, and the disk is available for I/O.

## 10.6.9 Practical use of Metro Mirror

The master volume is the production volume, and updates to this copy are mirrored in real time to the auxiliary volume. The contents of the auxiliary volume that existed when the relationship was created are deleted.

**Switching copy direction:** The copy direction for an MM relationship can be switched so that the auxiliary volume becomes the master, and the master volume becomes the auxiliary, which is similar to the FlashCopy restore option. However, although the FlashCopy target volume can operate in read/write mode, the target volume of the started RC is always in read-only mode.

While the MM relationship is active, the auxiliary volume is not accessible for host application write I/O at any time. The IBM Storwize V7000 allows read-only access to the auxiliary volume when it contains a consistent image. IBM Storwize allows boot time operating system discovery to complete without an error, so that any hosts at the secondary site can be ready to start the applications with minimum delay, if required.

For example, many operating systems must read logical block address (LBA) zero to configure a logical unit (LU). Although read access is allowed at the auxiliary in practice, the data on the auxiliary volumes cannot be read by a host because most operating systems write a “dirty bit” to the file system when it is mounted. Because this write operation is not allowed on the auxiliary volume, the volume cannot be mounted.

This access is provided only where consistency can be ensured. However, coherency cannot be maintained between reads that are performed at the auxiliary and later write I/Os that are performed at the master.

To enable access to the auxiliary volume for host operations, you must stop the MM relationship by specifying the **-access** parameter. While access to the auxiliary volume for host operations is enabled, the host must be instructed to mount the volume before the application can be started, or instructed to perform a recovery process.

For example, the MM requirement to enable the auxiliary copy for access differentiates it from third-party mirroring software on the host, which aims to emulate a single, reliable disk regardless of what system is accessing it. MM retains the property that there are two volumes in existence, but it suppresses one volume while the copy is being maintained.

The use of an auxiliary copy demands a conscious policy decision by the administrator that a failover is required, and that the tasks to be performed on the host that is involved in establishing the operation on the auxiliary copy are substantial. The goal is to make this copy rapid (much faster when compared to recovering from a backup copy) but not seamless.

The failover process can be automated through failover management software. The IBM Storwize V7000 provides SNMP traps and programming (or scripting) for the CLI to enable this automation.

## 10.6.10 Global Mirror overview

This section describes the GM copy service, which is an asynchronous RC service. This service provides and maintains a consistent mirrored copy of a source volume to a target volume.

GM function establishes a GM relationship between two volumes of equal size. The volumes in a GM relationship are referred to as the *master* (source) volume and the *auxiliary* (target) volume, which is the same as MM. Consistency groups can be used to maintain data integrity for dependent writes, which is similar to FlashCopy consistency groups.

GM writes data to the auxiliary volume asynchronously, which means that host writes to the master volume provide the host with confirmation that the write is complete before the I/O completes on the auxiliary volume.

GM has the following characteristics:

- ▶ Near-zero RPO
- ▶ Asynchronous
- ▶ Production application performance that is affected by I/O sequencing preparation time

### ***Intracluster Global Mirror***

Although GM is available for intracluster, it has no functional value for production use. Intracluster MM provides the same capability with less processor use. However, leaving this functionality in place simplifies testing and supports client experimentation and testing (for example, to validate server failover on a single test system). As with Intracluster MM, you must consider the increase in the license requirement because source and target exist on the same IBM Spectrum Virtualize system.

### ***Intercluster Global Mirror***

Intercluster GM operations require a pair of IBM Spectrum Virtualize systems that are connected by several intercluster links. The two systems must be defined in a partnership to establish a fully functional GM relationship.

**Limit:** When a local fabric and a remote fabric are connected for GM purposes, the ISL hop count between a local node and a remote node must not exceed seven hops.

## 10.6.11 Asynchronous Remote Copy

GM is an asynchronous RC technique. In asynchronous RC, the write operations are completed on the primary site and the write acknowledgment is sent to the host before it is received at the secondary site. An update of this write operation is sent to the secondary site at a later stage, which provides the capability to perform RC over distances that exceed the limitations of synchronous RC.

The GM function provides the same function as MM RC, but over long-distance links with higher latency without requiring the hosts to wait for the full round-trip delay of the long-distance link.

Figure 10-89 shows that a write operation to the master volume is acknowledged back to the host that is issuing the write before the write operation is mirrored to the cache for the auxiliary volume.

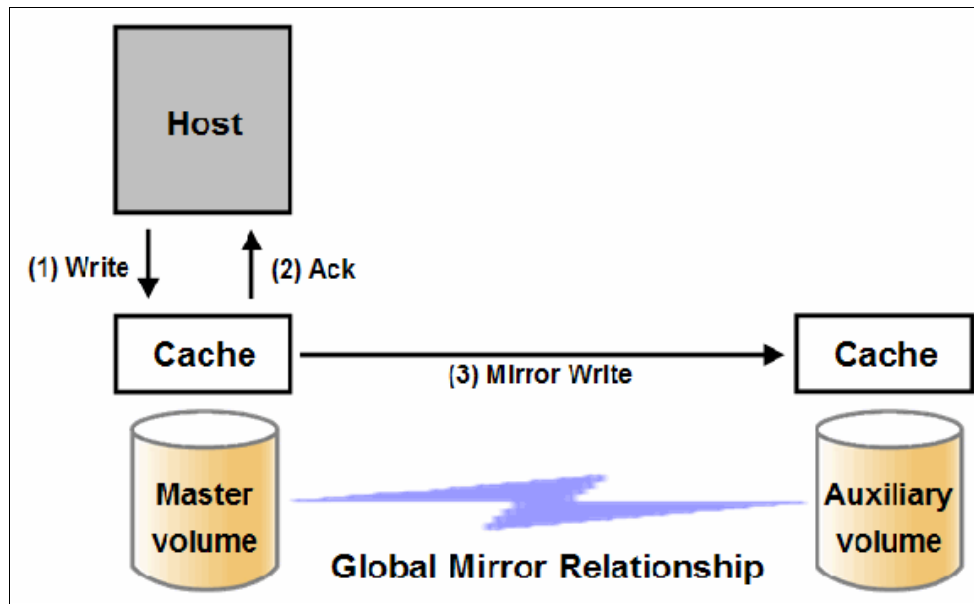


Figure 10-89 Global Mirror write sequence

The GM algorithms maintain a consistent image on the auxiliary. They achieve this consistent image by identifying sets of I/Os that are active concurrently at the master, assigning an order to those sets, and applying those sets of I/Os in the assigned order at the secondary. As a result, GM maintains the features of Write Ordering and Read Stability.

The multiple I/Os within a single set are applied concurrently. The process that marshals the sequential sets of I/Os operates at the secondary system. Therefore, the process is not subject to the latency of the long-distance link. These two elements of the protocol ensure that the throughput of the total system can be grown by increasing system size while maintaining consistency across a growing data set.

GM write I/O from production system to a secondary system requires serialization and sequence-tagging before being sent across the network to a remote site (to maintain a write-order consistent copy of data).

To avoid affecting the production site, IBM Spectrum Virtualize supports more parallelism in processing and managing GM writes on the secondary system by using the following methods:

- ▶ Secondary system nodes store replication writes in new redundant non-volatile cache
- ▶ Cache content details are shared between nodes
- ▶ Cache content details are batched together to make node-to-node latency less of an issue
- ▶ Nodes intelligently apply these batches in parallel as soon as possible
- ▶ Nodes internally manage and optimize GM secondary write I/O processing

In a failover scenario where the secondary site must become the master source of data, specific updates might be missing at the secondary site. Therefore, any applications that use this data must have an external mechanism for recovering the missing updates and reapplying them, such as a transaction log replay.



GM is supported over FC, Fibre Channel over IP (FCIP), Fibre Channel over Ethernet (FCoE), and native IP connections. The maximum distance cannot exceed 80 ms round trip, which is approximately 4000 km (2485.48 miles) between mirrored systems. However, starting with IBM Spectrum Virtualize V7.4, this distance was significantly increased for certain IBM Storwize Gen2 and IBM SAN Volume Controller configurations. Figure 10-90 shows the current supported distances for GM RC.

System Hardware	Partnership		
	FC	1Gbps – IP	10 Gbps – IP
SVC DH8	250ms	80ms	10ms
SVC SV1	250ms	80ms	10ms

Figure 10-90 Supported Global Mirror distances

### 10.6.12 Global Mirror features

IBM Spectrum Virtualize GM supports the following features:

- ▶ Asynchronous RC of volumes that are dispersed over metropolitan-scale distances.
- ▶ IBM Spectrum Virtualize implements the GM relationship between a volume pair, with each volume in the pair being managed by an IBM Spectrum Virtualize system.
- ▶ IBM Spectrum Virtualize supports intracluster GM where both volumes belong to the same system (and I/O group).
- ▶ An IBM Spectrum Virtualize system can be configured for partnership with 1 - 3 other systems. For more information about IP partnership restrictions, see 10.8.2, “IP partnership limitations” on page 615.
- ▶ Intercluster and intracluster GM can be used concurrently, but not for the same volume.
- ▶ IBM Spectrum Virtualize does not require a control network or fabric to be installed to manage GM. For intercluster GM, the system maintains a control link between the two systems. This control link is used to control the state and to coordinate the updates at either end. The control link is implemented on top of the same FC fabric connection that the system uses for GM I/O.
- ▶ IBM Spectrum Virtualize implements a configuration model that maintains the GM configuration and state through major events, such as failover, recovery, and resynchronization, to minimize user configuration action through these events.
- ▶ IBM Spectrum Virtualize implements flexible resynchronization support, enabling it to resynchronize volume pairs that experienced write I/Os to both disks, and to resynchronize only those regions that changed.
- ▶ An optional feature for GM is a delay simulation to be applied on writes that are sent to auxiliary volumes. It is useful in intracluster scenarios for testing purposes.

#### Colliding writes

The GM algorithm requires that only a single write is active on a volume. I/Os that overlap an active I/O are sequential, which is called *colliding writes*. If another write is received from a host while the auxiliary write is still active, the new host write is delayed until the auxiliary write is complete. This rule is needed if a series of writes to the auxiliary must be tried again and is called *reconstruction*. Conceptually, the data for reconstruction comes from the master volume.

If multiple writes are allowed to be applied to the master for a sector, only the most recent write gets the correct data during reconstruction. If reconstruction is interrupted for any reason, the intermediate state of the auxiliary is inconsistent. Applications that deliver such write activity do not achieve the performance that GM is intended to support. A volume statistic is maintained about the frequency of these collisions.

An attempt is made to allow multiple writes to a single location to be outstanding in the GM algorithm. Master writes must still be sequential, and the intermediate states of the master data must be kept in a non-volatile journal while the writes are outstanding to maintain the correct write ordering during reconstruction. Reconstruction must never overwrite data on the auxiliary with an earlier version. The volume statistic that is monitoring colliding writes is now limited to those writes that are not affected by this change.

Figure 10-91 shows a colliding write sequence example.

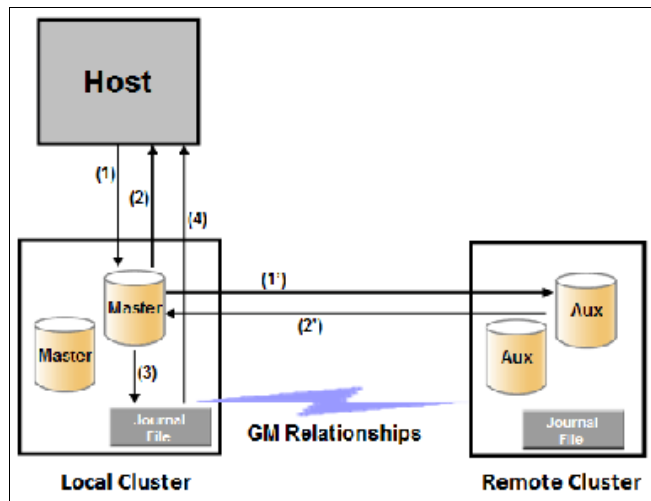


Figure 10-91 Colliding writes example

The following numbers correspond to the numbers that are shown in Figure 10-91:

- ▶ (1) The first write is performed from the host to LBA X.
- ▶ (2) The host is provided acknowledgment that the write completed, even though the mirrored write to the auxiliary volume is not yet complete.
- ▶ (1') and (2') occur asynchronously with the first write.
- ▶ (3) The second write is performed from the host also to LBA X. If this write occurs before (2'), the write is written to the journal file.
- ▶ (4) The host is provided acknowledgment that the second write is complete.

### Delay simulation

GM provides a feature that enables a delay simulation to be applied on writes that are sent to the auxiliary volumes. With this feature, tests can be done to detect colliding writes. It also provides the capability to test an application before the full deployment. The feature can be enabled separately for each of the intracluster or intercluster GMs.

By running the `chsystem` command, the delay setting can be set up and the delay can be checked by running the `lssystem` command. The `gm_intra_cluster_delay_simulation` field expresses the amount of time that intracluster auxiliary I/Os are delayed. The `gm_inter_cluster_delay_simulation` field expresses the amount of time that intercluster auxiliary I/Os are delayed. A value of zero disables the feature.

**Tip:** If you are experiencing repeated problems with the delay on your link, ensure that the delay simulator was correctly disabled.

### 10.6.13 Using Global Mirror with change volumes

GM is designed to achieve an RPO as low as possible so that data is as up-to-date as possible. This design places several strict requirements on your infrastructure. In certain situations with low network link quality, congested hosts, or overloaded hosts, you might be affected by multiple 1920 congestion errors.

Congestion errors occur in the following primary situations:

- ▶ At the source site through the host or network
- ▶ In the network link or network path
- ▶ At the target site through the host or network

GM has functionality that is designed to address the following conditions, which might negatively affect certain GM implementations:

- ▶ The estimation of the bandwidth requirements tends to be complex.
- ▶ Ensuring that the latency and bandwidth requirements can be met is often difficult.
- ▶ Congested hosts on the source or target site can cause disruption.
- ▶ Congested network links can cause disruption with only intermittent peaks.

To address these issues, change volumes were added as an option for GM relationships. Change volumes use the FlashCopy functionality, but they cannot be manipulated as FlashCopy volumes because they are for a special purpose only. Change volumes replicate PiT images on a cycling period. The default is 300 seconds.

Your change rate must include only the condition of the data at the PiT that the image was taken, rather than all the updates during the period. The use of this function can provide significant reductions in replication volume.

GMCV has the following characteristics:

- ▶ Larger RPO
- ▶ PiT copies
- ▶ Asynchronous
- ▶ Possible system performance resource requirements because PiT copies are created locally

Figure 10-92 shows a simple Global Mirror relationship without change volumes.

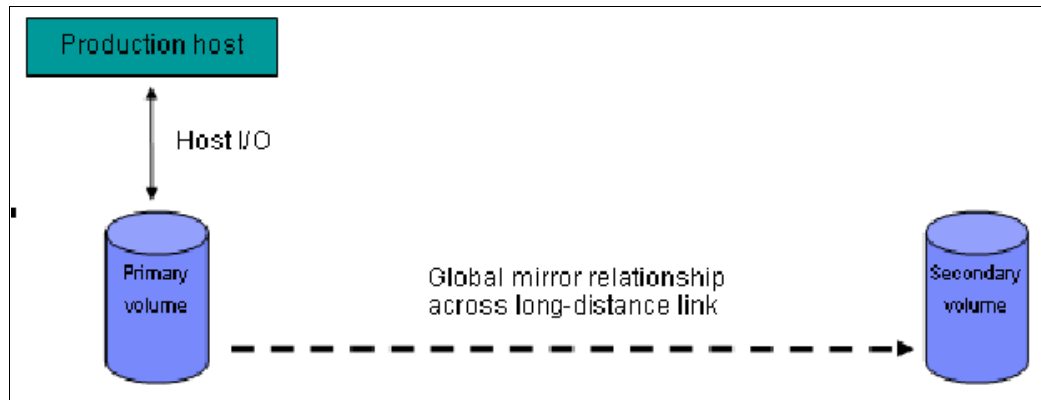


Figure 10-92 Global Mirror without change volumes

With change volumes, this environment looks as it is shown in Figure 10-93.

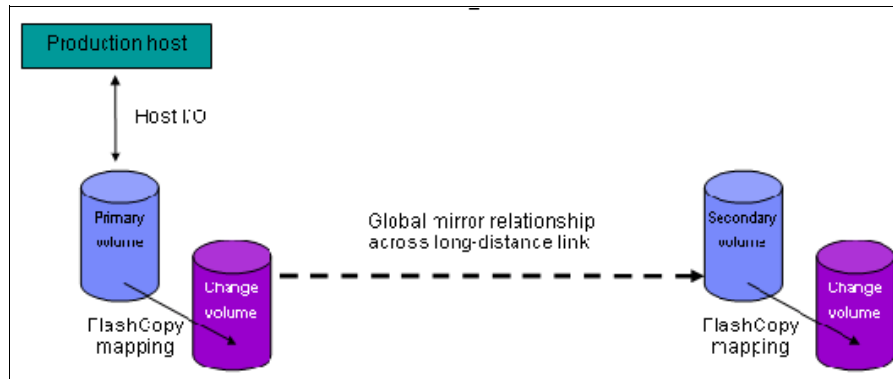


Figure 10-93 Global Mirror with Change Volumes

With change volumes, a FlashCopy mapping exists between the primary volume and the primary change volume. The mapping is updated in the cycling period (60 seconds - 1 day). The primary change volume is then replicated to the secondary GM volume at the target site, which is then captured in another change volume on the target site. This approach provides an always consistent image at the target site and protects your data from being inconsistent during resynchronization.

For more information about IBM FlashCopy, see 10.1, “IBM FlashCopy” on page 492.

You can adjust the cycling period by running the `chrcrelationship -cycleperiodseconds <60 - 86400>` command from the CLI. The default value is 300 seconds. If a copy does not complete in the cycle period, the next cycle does not start until the prior cycle completes. For this reason, the use of change volumes gives you the following possibilities for RPO:

- ▶ If your replication completes in the cycling period, your RPO is twice the cycling period.
- ▶ If your replication does not complete within the cycling period, RPO is twice the completion time. The next cycling period starts immediately after the prior cycling period is finished.

Carefully consider your business requirements versus the performance of GMCV. GMCV increases the intercluster traffic for more frequent cycling periods. Therefore, selecting the shortest cycle periods possible is not always the answer. In most cases, the default must meet requirements and perform well.

**Important:** When you create your Global Mirror volumes with change volumes, ensure that you remember to select the change volume on the auxiliary (target) site. Failure to do so leaves you exposed during a resynchronization operation.

### 10.6.14 Distribution of work among nodes

For the best performance, MM/GM volumes must have their preferred nodes evenly distributed among the nodes of the systems. Each volume within an I/O group has a preferred node property that can be used to balance the I/O load between nodes in that group. MM/GM also uses this property to route I/O between systems.

If this preferred practice is not maintained, such as if source volumes are assigned to only one node in the I/O group, you can change the preferred node for each volume to distribute volumes evenly between the nodes. You can also change the preferred node for volumes that are in an RC relationship without affecting the host I/O to a particular volume.

The RC relationship type does not matter. The RC relationship type can be MM, GM, or GMCV. You can change the preferred node both to the source and target volumes that are participating in the RC relationship.

### 10.6.15 Background copy performance

The background copy performance is subject to sufficient RAID controller bandwidth. Performance is also subject to other potential bottlenecks, such as the intercluster fabric, and possible contention from host I/O for the IBM Spectrum Virtualize system bandwidth resources.

Background copy I/O is scheduled to avoid bursts of activity that might have an adverse effect on system behavior. An entire grain of tracks on one volume is processed at around the same time, but not as a single I/O.

Double buffering is used to try to use sequential performance within a grain. However, the next grain within the volume might not be scheduled for some time. Multiple grains might be copied simultaneously, and might be enough to satisfy the requested rate, unless the available resources cannot sustain the requested rate.

GM paces the rate at which background copy is performed by the appropriate relationships. Background copy occurs on relationships that are in the `InconsistentCopying` state with a status of `Online`.

The quota of background copy (configured on the intercluster link) is divided evenly between all nodes that are performing background copy for one of the eligible relationships. This allocation is made irrespective of the number of disks for which the node is responsible. Each node in turn divides its allocation evenly between the multiple relationships that are performing a background copy.

The default value of the background copy is 25 MBps, per volume.

**Important:** The background copy value is a system-wide parameter that can be changed dynamically, but only on a *per-system* basis and not on a *per-relationship* basis. Therefore, the copy rate of all relationships changes when this value is increased or decreased. In systems with many RC relationships, increasing this value might affect overall system or intercluster link performance. The background copy rate can be changed to 1 - 1000 MBps.

### 10.6.16 Thin-provisioned background copy

MM/GM relationships preserve the space-efficiency of the master. Conceptually, the background copy process detects a deallocated region of the master and sends a special *zero buffer* to the auxiliary.

If the auxiliary volume is thin-provisioned and the region is deallocated, the special buffer prevents a write and therefore, an allocation. If the auxiliary volume is not thin-provisioned or the region in question is an allocated region of a thin-provisioned volume, a buffer of “real” zeros is synthesized on the auxiliary and written as normal.

### 10.6.17 Methods of synchronization

This section describes two methods that can be used to establish a synchronized relationship.

#### Full synchronization after creation

The full synchronization after creation method is the default method. It is the simplest method in that it requires no administrative activity apart from running the necessary commands. However, in certain environments, the available bandwidth can make this method unsuitable.

Run the following command sequence for a single relationship:

- ▶ Run `mkrcrelationship` without specifying the `-sync` option.
- ▶ Run `starttrcrelationship` without specifying the `-clean` option.

#### Synchronized before creation

In this method, the administrator must ensure that the master and auxiliary volumes contain identical data before creating the relationship by using the following technique:

- ▶ Both disks are created with the security delete feature to make all data zero.
- ▶ A complete tape image (or other method of moving data) is copied from one disk to the other disk.

With this technique, do not allow I/O on the master or auxiliary before the relationship is established. Then, the administrator must run the following commands:

- ▶ Run `mkrcrelationship` with the `-sync` flag.
- ▶ Run `starttrcrelationship` without the `-clean` flag.

**Important:** Failure to perform these steps correctly can cause MM/GM to report the relationship as consistent when it is not. This use can cause loss of a data or data integrity exposure for hosts that are accessing data on the auxiliary volume.

## 10.6.18 Practical use of Global Mirror

The practical use of GM is similar to MM, as described in 10.6.9, “Practical use of Metro Mirror” on page 588. The main difference between the two RC modes is that GM and GMCV are mostly used on much larger distances than MM. Weak link quality or insufficient bandwidth between the primary and secondary sites can also be a reason to prefer asynchronous GM over synchronous MM. Otherwise, the use cases for MM/GM are the same.

## 10.6.19 IBM Spectrum Virtualize HyperSwap topology

The IBM HyperSwap topology is based on IBM Spectrum Virtualize RC mechanisms. It is also referred to as an “active-active relationship” in this document.

You can create an HyperSwap topology system configuration where each I/O group in the system is physically on a different site. These configurations can be used to maintain access to data on the system when power failures or site-wide outages occur.

In a HyperSwap configuration, each site is defined as an independent failure domain. If one site experiences a failure, the other site can continue to operate without disruption. You must also configure a third site to host a quorum device or IP quorum application that provides an automatic tie-break in case of a link failure between the two main sites. The main site can be in the same room or across rooms in the data center, buildings on the same campus, or buildings in different cities. Different kinds of sites protect against different types of failures.

For more information about HyperSwap implementation and best practices, see *IBM Storwize V7000, Spectrum Virtualize, HyperSwap, and VMware Implementation*, SG24-8317.

## 10.6.20 Consistency Protection for Global Mirror and Metro Mirror

MM, GM, GMCV, and HyperSwap Copy Services functions create RC or remote replication relationships between volumes or consistency groups. If the secondary volume in a Copy Services relationship becomes unavailable to the primary volume, the system maintains the relationship. However, the data might become out of sync when the secondary volume becomes available.

Since V7.8, it is possible to create a FlashCopy mapping (change volume) for an RC target volume to maintain a consistent image of the secondary volume. The system recognizes it as a *Consistency Protection* and a link failure or an offline secondary volume event is handled differently now.

When Consistency Protection is configured, the relationship between the primary and secondary volumes does not stop if the link goes down or the secondary volume is offline. The relationship does not go in to the consistent stopped status. Instead, the system uses the secondary change volume to automatically copy the previous consistent state of the secondary volume. The relationship automatically moves to the consistent copying status as the system resynchronizes and protects the consistency of the data. The relationship status changes to `consistent_synchronized` when the resynchronization process completes. The relationship automatically resumes replication after the temporary loss of connectivity.

Change volumes that are used for Consistency Protection are not visible and manageable from the GUI because they are used for Consistency Protection internal behavior only.

It is not required to configure a secondary change volume on a MM/GM (without cycling) relationship. However, if the link goes down or the secondary volume is offline, the relationship goes in to the Consistent stopped status. If write operations occur on the primary or secondary volume, the data is no longer synchronized (Out of sync).

Consistency protection must be enabled on all relationships in a consistency group. Every relationship in a consistency group must be configured with a secondary change volume. If a secondary change volume is not configured on one relationship, the entire consistency group stops with a 1720 error if host I/O is processed when the link is down or any secondary volume in the consistency group is offline. All relationships in the consistency group are unable to retain a consistent copy during resynchronization.

The option to add consistency protection is selected by default when MM/GM relationships are created. The option must be cleared to create MM/GM relationships without consistency protection.

### 10.6.21 Valid combinations of FlashCopy, Metro Mirror, and Global Mirror

Table 10-11 lists the combinations of FlashCopy and MM/GM functions that are valid for a single volume.

Table 10-11 Valid combination for a single volume

FlashCopy	MM or GM source	MM or GM target
FlashCopy Source	Supported	Supported
FlashCopy Target	Supported	Not supported

### 10.6.22 Remote Copy configuration limits

Table 10-12 lists the MM/GM configuration limits.

Table 10-12 Metro Mirror configuration limits

Parameter	Value
Number of Metro Mirror or GM Consistency Groups per system	256
Number of Metro Mirror or GM relationships per system	10000
Number of Metro Mirror or GM relationships per Consistency Group	10000
Total Volume size per I/O group	A per I/O group limit of 1024 TB exists on the quantity of master and auxiliary volume address spaces that can participate in Metro Mirror and GM relationships. This maximum configuration uses all 512 MiB of bitmap space for the I/O group and allows 10 MiB of space for all remaining copy services features.



## 10.6.23 Remote Copy states and events

This section describes the various states of a MM/GM relationship and the conditions that cause them to change. In Figure 10-94 shows an overview of the status that can apply to a MM/GM relationship in a connected state.

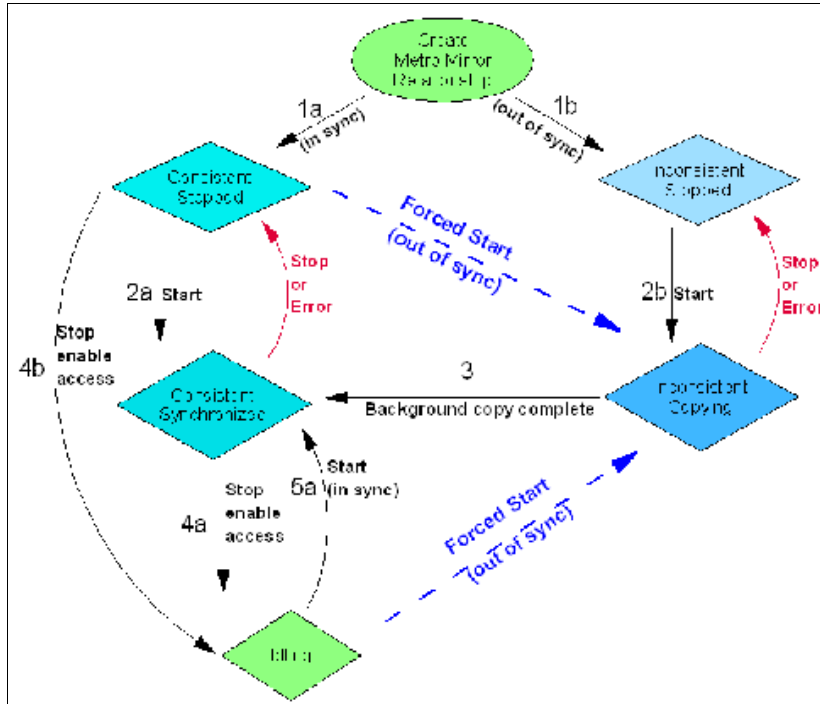


Figure 10-94 Metro Mirror or Global Mirror mapping state diagram

When the MM/GM relationship is created, you can specify whether the auxiliary volume is in sync with the master volume, and the background copy process is then skipped. This capability is useful when MM/GM relationships are established for volumes that were created with the format option.

The following step identifiers are shown in Figure 10-94:

- ▶ Step 1:
  - a. The MM/GM relationship is created with the `-sync` option, and the MM/GM relationship enters the `ConsistentStopped` state.
  - b. The MM/GM relationship is created without specifying that the master and auxiliary volumes are in sync, and the MM/GM relationship enters the `InconsistentStopped` state.
- ▶ Step 2:
  - a. When an MM/GM relationship is started in the `ConsistentStopped` state, the MM/GM relationship enters the `ConsistentSynchronized` state. Therefore, no updates (write I/O) were performed on the master volume while in the `ConsistentStopped` state. Otherwise, the `-force` option must be specified, and the MM/GM relationship then enters the `InconsistentCopying` state while the background copy is started.
  - b. When an MM/GM relationship is started in the `InconsistentStopped` state, the MM/GM relationship enters the `InconsistentCopying` state while the background copy is started.

- ▶ Step 3
 

When the background copy completes, the MM/GM relationship changes from the `InconsistentCopying` state to the `ConsistentSynchronized` state.
- ▶ Step 4:
  - a. When a MM/GM relationship is stopped in the `ConsistentSynchronized` state, the MM/GM relationship enters the `Idling` state when you specify the `-access` option, which enables write I/O on the auxiliary volume.
  - b. When an MM/GM relationship is stopped in the `ConsistentSynchronized` state without an `-access` parameter, the auxiliary volumes remain read-only and the state of the relationship changes to `ConsistentStopped`.
  - c. To enable write I/O on the auxiliary volume, when the MM/GM relationship is in the `ConsistentStopped` state, run the `svctask stopprcrelationship` command, which specifies the `-access` option, and the MM/GM relationship enters the `Idling` state.
- ▶ Step 5:
  - a. When an MM/GM relationship is started from the `Idling` state, you must specify the `-primary` argument to set the copy direction. If no write I/O was performed (to the master or auxiliary volume) while in the `Idling` state, the MM/GM relationship enters the `ConsistentSynchronized` state.
  - b. If write I/O was performed to the master or auxiliary volume, the `-force` option must be specified and the MM/GM relationship then enters the `InconsistentCopying` state while the background copy is started. The background process copies only the data that changed on the primary volume while the relationship was stopped.

## Stop on Error

When a MM/GM relationship is stopped (intentionally, or because of an error), the state changes. For example, the MM/GM relationships in the `ConsistentSynchronized` state enter the `ConsistentStopped` state, and the MM/GM relationships in the `InconsistentCopying` state enter the `InconsistentStopped` state.

If the connection is broken between the two systems that are in a partnership, all (intercluster) MM/GM relationships enter a `Disconnected` state. For more information, see “Connected versus disconnected” on page 600.

**Common states:** Stand-alone relationships and consistency groups share a common configuration and state model. All MM/GM relationships in a consistency group have the same state as the consistency group.

## State overview

The following sections provide an overview of the various MM/GM states.

### ***Connected versus disconnected***

Under certain error scenarios (for example, a power failure at one site that causes one complete system to disappear), communications between two systems in an MM/GM relationship can be lost. Alternatively, the fabric connection between the two systems might fail, which leaves the two systems that are running but cannot communicate with each other.

When the two systems can communicate, the systems and the relationships that spans them are described as *connected*. When they cannot communicate, the systems and the relationships spanning them are described as *disconnected*.

In this state, both systems are left with fragmented relationships and are limited regarding the configuration commands that can be performed. The disconnected relationships are portrayed as having a changed state. The new states describe what is known about the relationship and the configuration commands that are permitted.

When the systems can communicate again, the relationships are reconnected. MM/GM automatically reconciles the two state fragments and considers any configuration or other event that occurred while the relationship was disconnected. As a result, the relationship can return to the state that it was in when it became disconnected, or it can enter a new state.

Relationships that are configured between volumes in the same IBM Storwize V7000 system (intracluster) are never described as being in a disconnected state.

### ***Consistent versus inconsistent***

Relationships that contain volumes that are operating as secondaries can be described as being consistent or inconsistent. Consistency groups that contain relationships can also be described as being consistent or inconsistent. The consistent or inconsistent property describes the relationship of the data on the auxiliary to the data on the master volume. It can be considered a property of the auxiliary volume.

An auxiliary volume is described as *consistent* if it contains data that can be read by a host system from the master if power failed at an imaginary point while I/O was in progress, and power was later restored. This imaginary point is defined as the *recovery point*.

The requirements for consistency are expressed regarding activity at the master up to the recovery point. The auxiliary volume contains the data from all of the writes to the master for which the host received successful completion and that data was not overwritten by a subsequent write (before the recovery point).

Consider writes for which the host did not receive a successful completion (that is, it received bad completion or no completion at all). If the host then performed a read from the master of that data that returned successful completion and no later write was sent (before the recovery point), the auxiliary contains the same data as the data that was returned by the read from the master.

From the point of view of an application, consistency means that an auxiliary volume contains the same data as the master volume at the recovery point (the time at which the imaginary power failure occurred). If an application is designed to cope with an unexpected power failure, this assurance of consistency means that the application can use the auxiliary and begin operation as though it was restarted after the hypothetical power failure. Again, maintaining the application write ordering is the key property of consistency.

For more information about dependent writes, see 10.1.13, “FlashCopy and image mode Volumes” on page 516.

If a relationship (or set of relationships) is inconsistent and an attempt is made to start an application by using the data in the secondaries, the following outcomes are possible:

- ▶ The application might decide that the data is corrupted and crash or exit with an event code.
- ▶ The application might fail to detect that the data is corrupted and return erroneous data.
- ▶ The application might work without a problem.

Because of the risk of data corruption, and in particular undetected data corruption, MM/GM strongly enforces the concept of consistency and prohibits access to inconsistent data.

Consistency as a concept can be applied to a single relationship or a set of relationships in a consistency group. Write ordering is a concept that an application can maintain across several disks that are accessed through multiple systems. Therefore, consistency must operate across all of those disks.

When you are deciding how to use consistency groups, the administrator must consider the scope of an application's data and consider all of the interdependent systems that communicate and exchange information.

If two programs or systems communicate and store details as a result of the information that is exchanged, either of the following actions might occur:

- ▶ All of the data that is accessed by the group of systems must be placed into a single consistency group.
- ▶ The systems must be recovered independently (each within its own consistency group). Then, each system must perform recovery with the other applications to become consistent with them.

### ***Consistent versus synchronized***

A copy that is consistent and up-to-date is described as *synchronized*. In a synchronized relationship, the master and auxiliary volumes differ only in regions where writes are outstanding from the host.

Consistency does not mean that the data is up-to-date. A copy can be consistent and yet contain data that was frozen at a point in the past. Write I/O might continue to a master but not be copied to the auxiliary. This state arises when it becomes impossible to keep data up-to-date and maintain consistency. An example is a loss of communication between systems when you are writing to the auxiliary.

When communication is lost for an extended period and Consistency Protection was not enabled, MM/GM tracks the changes that occurred on the master, but not the order or the details of such changes (write data). When communication is restored, it is impossible to synchronize the auxiliary without sending write data to the auxiliary out of order. Therefore, consistency is lost.

**Note:** MM/GM relationships with Consistency Protection enabled use a PiT copy mechanism (FlashCopy) to keep a consistent copy of the auxiliary. The relationships stay in a consistent state, although not synchronized, even if communication is lost. For more information about Consistency Protection, see 10.6.20, "Consistency Protection for Global Mirror and Metro Mirror" on page 597.

### **Detailed states**

The following sections describe the states that are portrayed to the user for consistency groups or relationships. Also described is the information that is available in each state. The major states are designed to provide guidance about the available configuration commands.

#### ***InconsistentStopped***

*InconsistentStopped* is a connected state. In this state, the master is accessible for read and write I/O, but the auxiliary is not accessible for read or write I/O. A copy process must be started to make the auxiliary consistent. This state is entered when the relationship or consistency group was *InconsistentCopying* and suffered a persistent error or received a **stop** command that caused the copy process to stop.

A **start** command causes the relationship or consistency group to move to the *InconsistentCopying* state. A **stop** command is accepted, but has no effect.

If the relationship or consistency group becomes disconnected, the auxiliary side makes the transition to `InconsistentDisconnected`. The master side changes to `IdlingDisconnected`.

### ***InconsistentCopying***

`InconsistentCopying` is a connected state. In this state, the master is accessible for read and write I/O, but the auxiliary is not accessible for read or write I/O. This state is entered after a **start** command is issued to an `InconsistentStopped` relationship or a consistency group.

It is also entered when a forced start is issued to an `Idling` or `ConsistentStopped` relationship or consistency group. In this state, a background copy process runs that copies data from the master to the auxiliary volume.

In the absence of errors, an `InconsistentCopying` relationship is active, and the copy progress increases until the copy process completes. In certain error situations, the copy progress might freeze or even regress.

A persistent error or **stop** command places the relationship or consistency group into an `InconsistentStopped` state. A **start** command is accepted but has no effect.

If the background copy process completes on a stand-alone relationship or on all relationships for a consistency group, the relationship or consistency group changes to the `ConsistentSynchronized` state.

If the relationship or consistency group becomes disconnected, the auxiliary side changes to `InconsistentDisconnected`. The master side changes to `IdlingDisconnected`.

### ***ConsistentStopped***

`ConsistentStopped` is a connected state. In this state, the auxiliary contains a consistent image, but it might be out-of-date in relation to the master. This state can arise when a relationship was in a `ConsistentSynchronized` state and experienced an error that forces a Consistency Freeze. It can also arise when a relationship is created with a `CreateConsistentFlag` set to `TRUE`.

Normally, write activity that follows an I/O error causes updates to the master, and the auxiliary is no longer synchronized. In this case, consistency must be given up for a period to reestablish synchronization. You must run a **start** command with the **-force** option to acknowledge this condition, and the relationship or consistency group changes to `InconsistentCopying`. Enter this command only after all outstanding events are repaired.

In the unusual case where the master and the auxiliary are still synchronized (perhaps following a user stop, and no further write I/O was received), a **start** command takes the relationship to `ConsistentSynchronized`. No **-force** option is required. Also, in this case, you can run a **switch** command that moves the relationship or consistency group to `ConsistentSynchronized` and reverses the roles of the master and the auxiliary.

If the relationship or consistency group becomes disconnected, the auxiliary changes to `ConsistentDisconnected`. The master changes to `IdlingDisconnected`.

An informational status log is generated whenever a relationship or consistency group enters the `ConsistentStopped` state with a status of `Online`. You can configure this event to generate an SNMP trap that can be used to trigger automation or manual intervention to run a **start** command after a loss of synchronization.

### ***ConsistentSynchronized***

ConsistentSynchronized is a connected state. In this state, the master volume is accessible for read and write I/O, and the auxiliary volume is accessible for read-only I/O. Writes that are sent to the master volume are also sent to the auxiliary volume. Successful completion must be received for both writes, the write must be failed to the host, or a state must change out of the ConsistentSynchronized state before a write is completed to the host.

A **stop** command takes the relationship to the ConsistentStopped state. A **stop** command with the **-access** parameter takes the relationship to the Idling state.

A **switch** command leaves the relationship in the ConsistentSynchronized state, but it reverses the master and auxiliary roles (it switches the direction of replicating data). A **start** command is accepted, but has no effect.

If the relationship or consistency group becomes disconnected, the same changes are made as for ConsistentStopped.

### ***Idling***

Idling is a connected state. Both master and auxiliary volumes operate in the master role. Therefore, both master and auxiliary volumes are accessible for write I/O.

In this state, the relationship or consistency group accepts a **start** command. MM/GM maintains a record of regions on each disk that received write I/O while they were idling. This record is used to determine what areas must be copied following a **start** command.

The **start** command must specify the new copy direction. A **start** command can cause a loss of consistency if either volume in any relationship received write I/O, which is indicated by the Synchronized status. If the **start** command leads to loss of consistency, you must specify the **-force** parameter.

Following a **start** command, the relationship or consistency group changes to ConsistentSynchronized if there is no loss of consistency, or to InconsistentCopying if a loss of consistency occurs.

Also, the relationship or consistency group accepts a **-clean** option on the **start** command while in this state. If the relationship or consistency group becomes disconnected, both sides change their state to IdlingDisconnected.

### ***IdlingDisconnected***

IdlingDisconnected is a disconnected state. The target volumes in this half of the relationship or consistency group are all in the master role and accept read or write I/O.

The priority in this state is to recover the link to restore the relationship or consistency.

No configuration activity is possible (except for deletes or stops) until the relationship becomes connected again. At that point, the relationship changes to a connected state. The exact connected state that is entered depends on the state of the other half of the relationship or consistency group, which depends on the following factors:

- ▶ The state when it became disconnected
- ▶ The write activity since it was disconnected
- ▶ The configuration activity since it was disconnected

If both halves are IdlingDisconnected, the relationship becomes Idling when it is reconnected.

While `IdlingDisconnected`, if a write I/O is received that causes the loss of synchronization (synchronized attribute transitions from `true` to `false`) and the relationship was not already stopped (through a user stop or a persistent error), an event is raised to notify you of the condition. This same event also is raised when this condition occurs for the `ConsistentSynchronized` state.

### ***InconsistentDisconnected***

`InconsistentDisconnected` is a disconnected state. The target volumes in this half of the relationship or consistency group are all in the auxiliary role, and do not accept read *or* write I/O. Except for deletes, no configuration activity is permitted until the relationship becomes connected again.

When the relationship or consistency group becomes connected again, the relationship becomes `InconsistentCopying` automatically unless either of the following conditions are true:

- ▶ The relationship was `InconsistentStopped` when it became disconnected.
- ▶ The user issued a **stop** command while disconnected.

In either case, the relationship or consistency group becomes `InconsistentStopped`.

### ***ConsistentDisconnected***

`ConsistentDisconnected` is a disconnected state. The target volumes in this half of the relationship or consistency group are all in the auxiliary role, and accept read I/O but *not* write I/O.

This state is entered from `ConsistentSynchronized` or `ConsistentStopped` when the auxiliary side of a relationship becomes disconnected.

In this state, the relationship or consistency group displays an attribute of `FreezeTime`, which is the point when consistency was frozen. When it is entered from `ConsistentStopped`, it retains the time that it had in that state. When it is entered from `ConsistentSynchronized`, the `FreezeTime` shows the last time at which the relationship or consistency group was known to be consistent. This time corresponds to the time of the last successful heartbeat to the other system.

A **stop** command with the `-access` flag set to `true` transitions the relationship or consistency group to the `IdlingDisconnected` state. This state allows write I/O to be performed to the auxiliary volume and is used as part of a DR scenario.

When the relationship or consistency group becomes connected again, the relationship or consistency group becomes `ConsistentSynchronized` only if this action does not lead to a loss of consistency. The following conditions must be true:

- ▶ The relationship was `ConsistentSynchronized` when it became disconnected.
- ▶ No writes received successful completion at the master while disconnected.

Otherwise, the relationship becomes `ConsistentStopped`. The `FreezeTime` setting is retained.

### ***Empty***

This state applies only to consistency groups. It is the state of a consistency group that has no relationships and no other state information to show.

It is entered when a consistency group is first created. It is exited when the first relationship is added to the consistency group, at which point the state of the relationship becomes the state of the consistency group.

## 10.7 Remote Copy commands

This section presents commands that must be issued to create and operate RC services.

### 10.7.1 Remote Copy process

The MM/GM process includes the following steps:

1. A system partnership is created between two IBM Spectrum Virtualize systems (for intercluster MM/GM).
2. A MM/GM relationship is created between two volumes of the same size.
3. To manage multiple MM/GM relationships as one entity, the relationships can be made part of a MM/GM consistency group to ensure data consistency across multiple MM/GM relationships, or for ease of management.
4. The MM/GM relationship is started. When the background copy completes, the relationship is consistent and synchronized. When synchronized, the auxiliary volume holds a copy of the production data at the master that can be used for DR.
5. To access the auxiliary volume, the MM/GM relationship must be stopped with the access option enabled before write I/O is submitted to the auxiliary.

Following these steps, the remote host server is mapped to the auxiliary volume and the disk is available for I/O.

The command set for MM/GM contains the following broad groups:

- ▶ Commands to create, delete, and manipulate relationships and consistency group
- ▶ Commands to cause state changes

If a configuration command affects more than one system, MM/GM coordinates configuration activity between the systems. Specific configuration commands can be run only when the systems are connected, and fail with no effect when they are disconnected.

Other configuration commands are permitted, even if the systems are disconnected. The state is reconciled automatically by MM/GM when the systems become connected again.

For any command (with one exception), a single system receives the command from the administrator. This design is significant for defining the context for a CreateRelationship **mkrcrelationship** or CreateConsistencyGroup **mkrcconsistgrp** command. In this case, the system that is receiving the command is called the *local system*.

The exception is a command that sets systems into a MM/GM partnership. The **mkfcpartnership** and **mkippartnership** commands must be issued on both the local and remote systems.

The commands in this section are described as an abstract command set, and are implemented by using one of the following methods:

- ▶ CLI can be used for scripting and automation.
- ▶ GUI can be used for one-off tasks.



## 10.7.2 Listing available system partners

Run the `lspartnershipcandidate` command to list the systems that are available for setting up a two-system partnership. This command is a prerequisite for creating MM/GM relationships.

**Note:** This command is not supported on IP partnerships. Use `mkippartnership` for IP connections.

## 10.7.3 Changing the system parameters

When you want to change system parameters specific to any RC or GM only, use the `chsystem` command. The `chsystem` command features the following parameters for MM/GM:

► **-relationshipbandwidthlimit** *cluster\_relationship\_bandwidth\_limit*

This parameter controls the maximum rate at which any one RC relationship can synchronize. The default value for the relationship bandwidth limit is 25 MBps, but this value can now be specified as 1 - 100,000 MBps.

The partnership overall limit is controlled at a partnership level by the `chpartnership -linkbandwidthhbits` command, and must be set on each involved system.

**Important:** Do not set this value higher than the default without first establishing that the higher bandwidth can be sustained without affecting the host's performance. The limit must never be higher than the maximum that is supported by the infrastructure connecting the remote sites, regardless of the compression rates that you might achieve.

► **-gmlinktolerance** *link\_tolerance*

This parameter specifies the maximum period that the system tolerates delay before stopping GM relationships. Specify values of 60 - 86,400 seconds in increments of 10 seconds. The default value is 300. Do not change this value except under the direction of IBM Support.

► **-gmmaxhostdelay** *max\_host\_delay*

This parameter specifies the maximum time delay, in milliseconds, at which the GM link tolerance timer starts counting down. This threshold value determines the extra effect that GM operations can add to the response times of the GM source volumes. You can use this parameter to increase the threshold from the default value of 5 milliseconds.

► **-maxreplicationdelay** *max\_replication\_delay*

This parameter sets a maximum replication delay in seconds. The value must be a number 0 - 360 (0 being the default value, no delay). This feature sets the maximum number of seconds to be tolerated to complete a single I/O. If I/O cannot complete within the *max\_replication\_delay*, the 1920 event is reported. This setting is system-wide and applies to MM/GM relationships.

Run the `chsystem` command to adjust these values, as shown in the following example:

```
chsystem -gmlinktolerance 300
```

You can view all of these parameter values by running the `lssystem <system_name>` command.

Focus on the **gm1inktolerance** parameter in particular. If poor response extends past the specified tolerance, a 1920 event is logged and one or more GM relationships automatically stop to protect the application hosts at the primary site. During normal operations, application hosts experience a minimal effect from the response times because the GM feature uses asynchronous replication.

However, if GM operations experience degraded response times from the secondary system for an extended period, I/O operations queue at the primary system. This queue results in an extended response time to application hosts. In this situation, the **gm1inktolerance** feature stops GM relationships, and the application host's response time returns to normal.

After a 1920 event occurs, the GM auxiliary volumes are no longer in the `consistent_synchronized` state. Fix the cause of the event and restart your GM relationships. For this reason, ensure that you monitor the system to track when these 1920 events occur.

You can disable the **gm1inktolerance** feature by setting the **gm1inktolerance** value to 0 (zero). However, the **gm1inktolerance** feature cannot protect applications from extended response times if it is disabled. It might be appropriate to disable the **gm1inktolerance** feature under the following circumstances:

- ▶ During SAN maintenance windows in which degraded performance is expected from SAN components, and application hosts can stand extended response times from GM volumes.
- ▶ During periods when application hosts can tolerate extended response times and it is expected that the **gm1inktolerance** feature might stop the GM relationships. For example, if you test by using an I/O generator that is configured to stress the back-end storage, the **gm1inktolerance** feature might detect the high latency and stop the GM relationships.

Disabling the **gm1inktolerance** feature prevents this result at the risk of exposing the test host to extended response times.

A 1920 event indicates that one or more of the SAN components cannot provide the performance that is required by the application hosts. This situation can be temporary (for example, a result of a maintenance activity) or permanent (for example, a result of a hardware failure or an unexpected host I/O workload).

If 1920 events are occurring, you might need to use a performance monitoring and analysis tool, such as the IBM Spectrum Control, to help identify and resolve the problem.

## 10.7.4 System partnership

To create a partnership, run one of the following commands, depending on the connection type:

- ▶ The **mkfcpartnership** command to establish a one-way MM/GM partnership between the local system and a remote system that are linked over an FC (or FCoE) connection.
- ▶ The **mkippartnership** command to establish a one-way MM/GM partnership between the local system and a remote system that are linked over an IP connection.

To establish a fully functional MM/GM partnership, you must run either of these commands (depending on the connection type) on both of the systems that are included of the partnership. This step is a prerequisite for creating MM/GM relationships between volumes on the IBM Spectrum Virtualize systems.

When creating the partnership, you must specify the Link Bandwidth and can specify the Background Copy Rate:

- ▶ The Link Bandwidth, which is expressed in Mbps, is the amount of bandwidth that can be used for the FC or IP connection between the systems within the partnership.
- ▶ The Background Copy Rate is the maximum percentage of the Link Bandwidth that can be used for background copy operations. The default value is 50%.

### **Background copy bandwidth effect on foreground I/O latency**

The combination of the Link Bandwidth value and the Background Copy Rate percentage is referred to as the *Background Copy bandwidth*. It must be at least 8 Mbps. For example, if the Link Bandwidth is set to 10000 and the Background Copy Rate is set to 20, the resulting Background Copy bandwidth that is used for background operations is 200 Mbps.

The background copy bandwidth determines the rate at which the background copy is attempted for MM/GM. The background copy bandwidth can affect foreground I/O latency in one of the following ways:

- ▶ The following results can occur if the background copy bandwidth is set too high compared to the MM/GM intercluster link capacity:
  - The background copy I/Os can back up on the MM/GM intercluster link.
  - There is a delay in the synchronous auxiliary writes of foreground I/Os.
  - The foreground I/O latency increases as perceived by applications.
- ▶ If the background copy bandwidth is set too high for the storage at the primary site, background copy read I/Os overload the primary storage and delay foreground I/Os.
- ▶ If the background copy bandwidth is set too high for the storage at the secondary site, background copy writes at the secondary site overload the auxiliary storage, and again delay the synchronous secondary writes of foreground I/Os.

To set the background copy bandwidth optimally, ensure that you consider all three resources: Primary storage, intercluster link bandwidth, and auxiliary storage. Provision the most restrictive of these three resources between the background copy bandwidth and the peak foreground I/O workload.

Perform this provisioning by calculation or by determining experimentally how much background copy can be allowed before the foreground I/O latency becomes unacceptable. Then, reduce the background copy to accommodate peaks in workload.

### **The `chpartnership` command**

To change the bandwidth that is available for background copy in the system partnership, run the `chpartnership -backgroundcopyrate <percentage_of_link_bandwidth>` command to specify the percentage of whole link capacity to be used by the background copy process.

## **10.7.5 Creating a Metro Mirror/Global Mirror consistency group**

Run the `mkrconsistgrp` command to create an empty MM/GM Consistency Group.

The MM/GM consistency group name must be unique across all consistency groups that are known to the systems owning this consistency group. If the consistency group involves two systems, the systems must be in communication throughout the creation process.

The new consistency group does not contain any relationships and is in the Empty state. You can add MM/GM relationships to the group (upon creation or afterward) by running the `chrelationship` command.

## 10.7.6 Creating a Metro Mirror/Global Mirror relationship

Run the `mkrcrelationship` command to create a MM/GM relationship. This relationship persists until it is deleted.

**Optional parameter:** If you do not use the `-global` optional parameter, an MM relationship is created rather than a GM relationship.

The auxiliary volume must be equal in size to the master volume or the command fails. If both volumes are in the same system, they must be in the same I/O group. The master and auxiliary volume cannot be in a relationship, and they cannot be the target of a FlashCopy mapping. This command returns the new relationship (`relationship_id`) when successful.

When the MM/GM relationship is created, you can add it to a Consistency Group, or it can be a stand-alone MM/GM relationship.

### The `lsrcrelationshipcandidate` command

Run the `lsrcrelationshipcandidate` command to list the volumes that are eligible to form an MM/GM relationship.

When the command is issued, you can specify the master volume name and auxiliary system to list the candidates that comply with the prerequisites to create a MM/GM relationship. If the command is issued with no parameters, all of the volumes that are not disallowed by another configuration state, such as being a FlashCopy target, are listed.

## 10.7.7 Changing Metro Mirror/Global Mirror relationship

Run the `chrcrelationship` command to modify the following properties of an MM/GM relationship:

- ▶ Change the name of an MM/GM relationship.
- ▶ Add a relationship to a group.
- ▶ Remove a relationship from a group by using the `-force` flag.

**Adding an MM/GM relationship:** When an MM/GM relationship is added to a consistency group that is not empty, the relationship must have the same state and copy direction as the group to be added to it.

## 10.7.8 Changing Metro Mirror/Global Mirror consistency group

Run the `chrconsistgrp` command to change the name of an MM/GM consistency group.

## 10.7.9 Starting Metro Mirror/Global Mirror relationship

Run the `startrcrelationship` command to start the copy process of an MM/GM relationship.

When the command is run, you can set the copy direction if it is undefined. Optionally, you can mark the auxiliary volume of the relationship as clean. The command fails if it is used as an attempt to start a relationship that is a part of a consistency group.

You can run this command only to a relationship that is connected. For a relationship that is idling, this command assigns a copy direction (master and auxiliary roles) and begins the copy process. Otherwise, this command restarts a previous copy process that was stopped by a **stop** command or by an I/O error.

If the resumption of the copy process leads to a period when the relationship is inconsistent, you must specify the **-force** parameter when the relationship is restarted. This situation can arise if, for example, the relationship was stopped and then further writes were performed on the original master of the relationship.

The use of the **-force** parameter here is a reminder that the data on the auxiliary becomes inconsistent while resynchronization (background copying) occurs. Therefore, this data is unusable for DR purposes before the background copy completes.

In the `Idling` state, you must specify the master volume to indicate the copy direction. In other connected states, you can provide the **-primary** argument, but it must match the existing setting.

### 10.7.10 Stopping Metro Mirror/Global Mirror relationship

Run the **stopcrelationship** command to stop the copy process for a relationship. You can also use this command to enable write access to a consistent auxiliary volume by specifying the **-access** parameter.

This command applies to a stand-alone relationship. It is rejected if it is addressed to a relationship that is part of a consistency group. You can issue this command to stop a relationship that is copying from master to auxiliary.

If the relationship is in an inconsistent state, any copy operation stops and does not resume until you run a **startcrelationship** command. Write activity is no longer copied from the master to the auxiliary volume. For a relationship in the `ConsistentSynchronized` state, this command causes a Consistency Freeze.

When a relationship is in a consistent state (that is, in the `ConsistentStopped`, `ConsistentSynchronized`, or `ConsistentDisconnected` state), you can use the **-access** parameter with the **stopcrelationship** command to enable write access to the auxiliary volume.

### 10.7.11 Starting Metro Mirror/Global Mirror consistency group

Run the **startcrconsistgrp** command to start an MM/GM consistency group. You can issue this command only to a consistency group that is connected.

For a consistency group that is idling, this command assigns a copy direction (master and auxiliary roles) and begins the copy process. Otherwise, this command restarts a previous copy process that was stopped by a **stop** command or by an I/O error.

### 10.7.12 Stopping Metro Mirror/Global Mirror consistency group

Run the **stopcrconsistgrp** command to stop the copy process for an MM/GM consistency group. You can also use this command to enable write access to the auxiliary volumes in the group if the group is in a consistent state.

If the consistency group is in an inconsistent state, any copy operation stops and does not resume until you run the **startrcconsistgrp** command. Write activity is no longer copied from the master to the auxiliary volumes that belong to the relationships in the group. For a consistency group in the ConsistentSynchronized state, this command causes a Consistency Freeze.

When a consistency group is in a consistent state (for example, in the ConsistentStopped, ConsistentSynchronized, or ConsistentDisconnected state), you can use the **-access** parameter with the **stoprcconsistgrp** command to enable write access to the auxiliary volumes within that group.

### 10.7.13 Deleting Metro Mirror/Global Mirror relationship

Run the **rmmrcrelationship** command to delete the relationship that is specified. Deleting a relationship deletes only the logical relationship between the two volumes. It does not affect the volumes.

If the relationship is disconnected at the time that the command is issued, the relationship is deleted on only the system on which the command is being run. When the systems reconnect, the relationship is automatically deleted on the other system.

Alternatively, if the systems are disconnected and you still want to remove the relationship on both systems, you can run the **rmmrcrelationship** command independently on both of the systems.

A relationship cannot be deleted if it is part of a consistency group. You must first remove the relationship from the consistency group.

If you delete an inconsistent relationship, the auxiliary volume becomes accessible, even though it is still inconsistent. This situation is the one case in which MM/GM does not inhibit access to inconsistent data.

### 10.7.14 Deleting Metro Mirror/Global Mirror consistency group

Run the **rmmrcconsistgrp** command to delete an MM/GM consistency group. This command deletes the specified consistency group.

If the consistency group is disconnected at the time that the command is issued, the consistency group is deleted on only the system on which the command is being run. When the systems reconnect, the consistency group is automatically deleted on the other system.

Alternatively, if the systems are disconnected and you still want to remove the consistency group on both systems, you can run the **rmmrcconsistgrp** command separately on both of the systems.

If the consistency group is not empty, the relationships within it are removed from the consistency group before the group is deleted. These relationships then become stand-alone relationships. The state of these relationships is not changed by the action of removing them from the consistency group.

### 10.7.15 Reversing Metro Mirror/Global Mirror relationship

Run the **switchrcrelationship** command to reverse the roles of the master volume and the auxiliary volume when a stand-alone relationship is in a consistent state. When the command is issued, the wanted master must be specified.

## 10.7.16 Reversing Metro Mirror/Global Mirror consistency group

Run the `switchrconsistgrp` command to reverse the roles of the master volume and the auxiliary volume when a consistency group is in a consistent state. This change is applied to all of the relationships in the consistency group. When the command is issued, the wanted master must be specified.

**Important:** By reversing the roles, your current source volumes become targets, and target volumes become source. Therefore, you lose write access to your current primary volumes.

## 10.8 Native IP replication

IBM Spectrum Virtualize can implement RC services by using FC connections or IP connections. This section describes the IBM Spectrum Virtualize IP replication technology and implementation.

**Demonstration:** The IBM Client Demonstration Center shows how data is replicated by using GMCV (cycling mode set to `multiple`). This configuration perfectly fits the new IP replication functionality because it is well-designed for links with high latency, low bandwidth, or both.

For more information, see this [web page](#) (log in required).

### 10.8.1 Native IP replication technology

Remote Mirroring over IP communication is supported on the IBM SAN Volume Controller and IBM Storwize Family systems by using Ethernet communication links. The IBM Spectrum Virtualize Software IP replication uses innovative Bridgeworks SANSlide technology to optimize network bandwidth and utilization. This function enables the use of a lower-speed and lower-cost networking infrastructure for data replication.

Bridgeworks SANSlide technology, which is integrated into the IBM Spectrum Virtualize Software, uses artificial intelligence (AI) to help optimize network bandwidth use and adapt to changing workload and network conditions.

This technology can improve remote mirroring network bandwidth usage up to three times. Improved bandwidth usage can enable clients to deploy a less costly network infrastructure, or speed up remote replication cycles to enhance DR effectiveness.

With an Ethernet network data flow, the data transfer can slow down over time. This condition occurs because of the latency that is caused by waiting for the acknowledgment of each set of packets that is sent. The next packet set cannot be sent until the previous packet is acknowledged, as shown in Figure 10-95 on page 614.

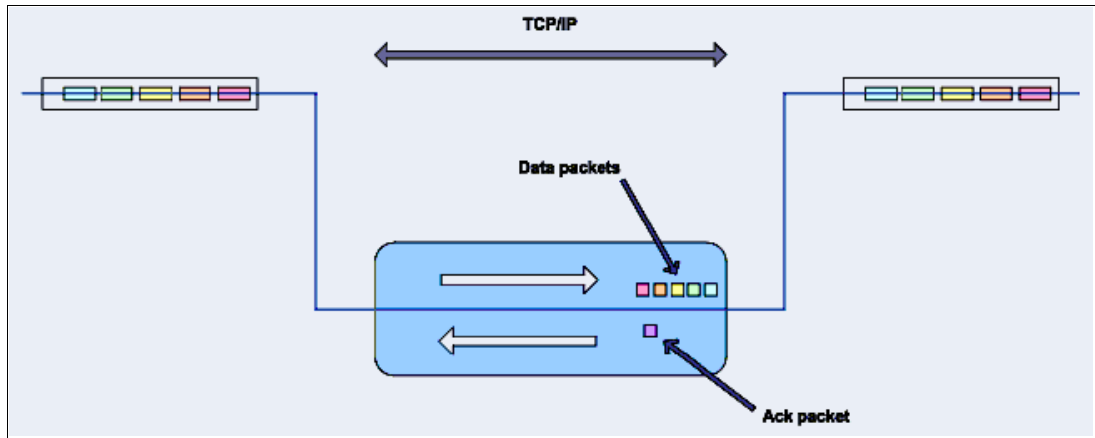


Figure 10-95 Typical Ethernet network data flow

However, by using the embedded IP replication, this behavior can be eliminated with the enhanced parallelism of the data flow by using multiple virtual connections (VC) that share IP links and addresses. The AI engine can dynamically adjust the number of VCs, receive window size, and packet size to maintain optimum performance. While the engine is waiting for one VC's ACK, it sends more packets across other VCs. If packets are lost from any VC, data is automatically retransmitted, as shown in Figure 10-96.

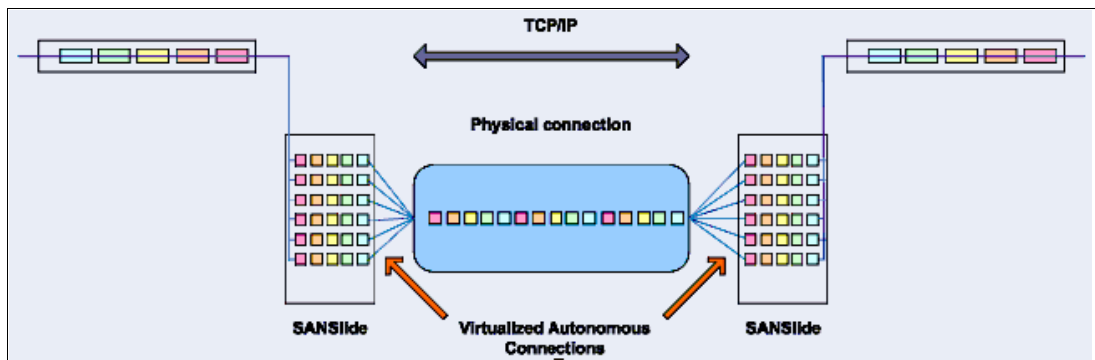


Figure 10-96 Optimized network data flow by using Bridgeworks SANSlide technology

For more information about this technology, see *IBM Storwize V7000 and SANSlide Implementation*, REDP-5023.

With native IP partnership, the following Copy Services features are supported:

- ▶ MM

Referred to as *synchronous replication*, MM provides a consistent copy of a source volume on a target volume. Data is written to the target volume synchronously after it is written to the source volume so that the copy is continuously updated.

- ▶ GM and GMCV

Referred to as *asynchronous replication*, GM provides a consistent copy of a source volume on a target volume. Data is written to the target volume asynchronously so that the copy is continuously updated. However, the copy might not contain the last few updates if a DR operation is performed. An added extension to GM is GMCV. GMCV is the preferred method for use with native IP replication.



**Note:** For IP partnerships, generally use the GMCV method of copying (asynchronous copy of changed grains only). This method can include performance benefits. Also, GM and MM might be more susceptible to the loss of synchronization.

## 10.8.2 IP partnership limitations

The following prerequisites and assumptions must be considered before IP partnership between two IBM Spectrum Virtualize systems can be established:

- ▶ The IBM Spectrum Virtualize systems are successfully installed with V7.2 or later code levels.
- ▶ The systems must have the necessary licenses that enable RC partnerships to be configured between two systems. No separate license is required to enable IP partnership.
- ▶ The storage SANs are configured correctly and the correct infrastructure to support the IBM Spectrum Virtualize systems in RC partnerships over IP links is in place.
- ▶ The two systems must be able to ping each other and perform the discovery.
- ▶ TCP ports 3260 and 3265 are used by systems for IP partnership communications. Therefore, these ports must be open.
- ▶ The maximum number of partnerships between the local and remote systems, including both IP and FC partnerships, is limited to the current maximum that is supported, which is three partnerships (four systems total).
- ▶ Only a single partnership over IP is supported.
- ▶ A system can have simultaneous partnerships over FC and IP, but with separate systems. The FC zones between two systems must be removed before an IP partnership is configured.
- ▶ IP partnerships are supported on both 10 Gbps links and 1 Gbps links. However, the intermix of both on a single link is not supported.
- ▶ The maximum supported round-trip time (RTT) is 80 ms for 1 Gbps links.
- ▶ The maximum supported RTT is 10 ms for 10 Gbps links.
- ▶ The inter-cluster heartbeat traffic uses 1 Mbps per link.
- ▶ Only nodes from two I/O groups can have ports that are configured for an IP partnership.
- ▶ Migrations of RC relationships directly from FC-based partnerships to IP partnerships are not supported.
- ▶ IP partnerships between the two systems can be over IPv4 or IPv6 only, but not both.
- ▶ Virtual local area network (VLAN) tagging of the IP addresses that are configured for RC is supported starting with V7.4.
- ▶ Management IP and internet Small Computer Systems Interface (iSCSI) IP on the same port can be in a different network starting with V7.4.
- ▶ An added layer of security is provided by using Challenge Handshake Authentication Protocol (CHAP) authentication.
- ▶ TCP ports 3260 and 3265 are used for IP partnership communications. Therefore, these ports must be open in firewalls between the systems.
- ▶ Only a single RC data session per physical link can be established. It is intended that only one connection (for sending/receiving RC data) is made for each independent physical link between the systems.

**Note:** A physical link is the physical IP link between the two sites: A (local) and B (remote). Multiple IP addresses on local system A might be connected (by Ethernet switches) to this physical link. Similarly, multiple IP addresses on remote system B might be connected (by Ethernet switches) to the same physical link. At any time, only a single IP address on cluster A can form an RC data session with an IP address on cluster B.

- ▶ The maximum throughput is restricted based on the use of 1 Gbps, 10 Gbps, or 25 Gbps Ethernet ports. It varies based on distance (for example, round-trip latency) and quality of communication link (for example, packet loss):
  - One 1 Gbps port can transfer up to 110 MBps unidirectional, 190 MBps bidirectional
  - Two 1 Gbps ports can transfer up to 220 MBps unidirectional, 325 MBps bidirectional
  - One 10 Gbps port can transfer up to 240 MBps unidirectional, 350 MBps bidirectional
  - Two 10 Gbps port can transfer up to 440 MBps unidirectional, 600 MBps bidirectional

**Note:** IP Replication is supported by 25 Gbps Mellanox and Chelsio adapters, but be aware there is no performance benefit or advantage for IP Replication with these adapters. However, for the purpose of consolidation where these cards are used for other purposes, such as iSCSI Extensions for RDMA (iSER) Host Attach or iSCSI Host Attach/Backend Virtualization, they can be used for IP replication.

The minimum supported link bandwidth is 10 Mbps. However, this requirement scales up with the amount of host I/O that you choose to do. Figure 10-97 shows scaling host I/O.

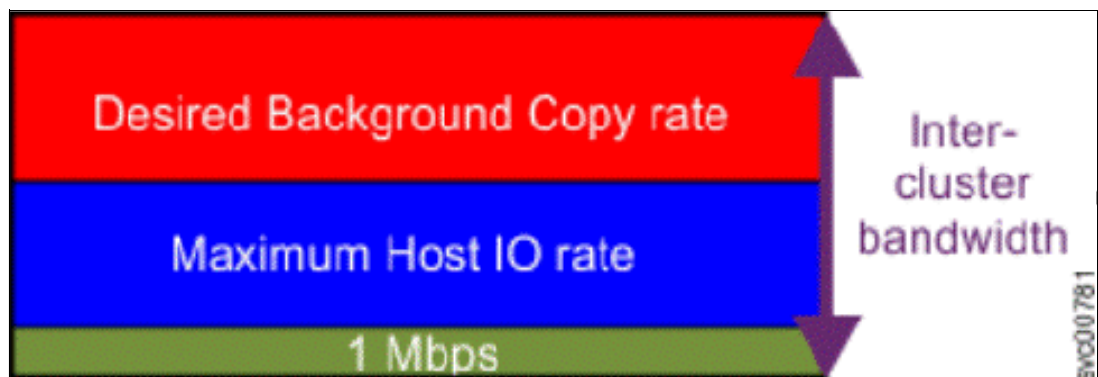


Figure 10-97 Scaling of host I/O

The following equation describes the approximate minimum bandwidth that is required between two systems with < 5 ms RTT and errorless link:

$$\text{Minimum intersite link bandwidth in Mbps} > \text{Required Background Copy in Mbps} + \text{Maximum Host I/O in Mbps} + 1 \text{ Mbps heartbeat traffic}$$

Increasing latency and errors results in a higher requirement for minimum bandwidth.

**Note:** The Bandwidth setting definition when the IP partnerships are created changed in V7.7. Previously, the bandwidth setting defaulted to 50 MiB, and was the maximum transfer rate from the primary site to the secondary site for initial sync/resyncs of volumes.

The Link Bandwidth setting is now configured by using megabits (Mb) not MB. You set the Link Bandwidth setting to a value that the communication link can sustain, or to what is allocated for replication. The Background Copy Rate setting is now a percentage of the Link Bandwidth. The Background Copy Rate setting determines the available bandwidth for the initial sync and resyncs or for GMCV.

### 10.8.3 IP Partnership and data compression

When creating an IP partnership between two systems, you can specify whether you want to use the data compression. When enabled, IP partnership compression compresses the data that is sent from a local system to the remote system and potentially uses less bandwidth than with uncompressed data. It is also used to decompress data that is received by a local system from a remote system.

Data compression is supported for IPv4 or IPv6 partnerships. To enable data compression, both systems in an IP partnership must be running a software level that supports IP partnership compression (V7.7 or later).

To fully enable compression in an IP partnership, each system must enable compression. When compression is enabled on the local system, data sent to the remote system is compressed so it needs to be decompressed on the remote system, and vice versa.

Although IP compression uses the same IBM Real-time Compression (RtC) algorithm as for volumes, a RtC license is not needed on any of the local or remote system.

Replicated volumes by using IP partnership compression can be compressed or uncompressed on the system because no link exists between volumes compression and IP Replication compression. Consider the following points:

- ▶ Read operation decompresses the data
- ▶ Decompressed data is transferred to the RC code
- ▶ Data is compressed before being sent over the IP link
- ▶ Remote system RC code decompresses the received data
- ▶ Write operation on a volume compresses the data

### 10.8.4 VLAN support

Starting with V7.4, VLAN tagging is supported for iSCSI host attachment and IP replication. Hosts and remote-copy operations can connect to the system through Ethernet ports. Each traffic type has different bandwidth requirements, which can interfere with each other if they share IP connections. VLAN tagging creates two separate connections on the same IP network for different types of traffic. The system supports VLAN configuration on both IPv4 and IPv6 connections.

When the VLAN ID is configured for IP addresses that is used for iSCSI host attach or IP replication, the VLAN settings on the Ethernet network and servers must be configured correctly to avoid connectivity issues. After the VLANs are configured, changes to the VLAN settings disrupt iSCSI and IP replication traffic to and from the partnerships.

During the VLAN configuration for each IP address, the VLAN settings for the local and failover ports on two nodes of an I/O group can differ. To avoid any service disruption, switches must be configured so that the failover VLANs are configured on the local switch ports and the failover of IP addresses from a failing node to a surviving node succeeds. If failover VLANs are not configured on the local switch ports, no paths are available to the IBM Spectrum Virtualize system nodes during a node failure and the replication fails.

Consider the following requirements and procedures when implementing VLAN tagging:

- ▶ VLAN tagging is supported for IP partnership traffic between two systems.
- ▶ VLAN provides network traffic separation at the layer 2 level for Ethernet transport.
- ▶ VLAN tagging by default is disabled for any IP address of a node port. You can use the CLI or GUI to optionally set the VLAN ID for port IPs on both systems in the IP partnership.
- ▶ When a VLAN ID is configured for the port IP addresses that are used in RC port groups, appropriate VLAN settings on the Ethernet network must also be configured to prevent connectivity issues.

Setting VLAN tags for a port is disruptive. Therefore, VLAN tagging requires that you stop the partnership first before you configure VLAN tags. Restart the partnership after the configuration is complete.

## 10.8.5 IP partnership and terminology

The IP partnership terminology and abbreviations that are used are listed in Table 10-13.

Table 10-13 Terminology for IP partnership

IP partnership terminology	Description
RC group or RC port group	<p>The following numbers group a set of IP addresses that are connected to the same physical link. Therefore, only IP addresses that are part of the same RC group can form RC connections with the partner system:</p> <ul style="list-style-type: none"> <li>▶ 0: Ports that are not configured for RC</li> <li>▶ 1: Ports that belong to RC port group 1</li> <li>▶ 2: Ports that belong to RC port group 2</li> </ul> <p>Each IP address can be shared for iSCSI host attach and RC functionality. Therefore, appropriate settings must be applied to each IP address.</p>
IP partnership	Two systems that are partnered to perform RC over native IP links.
FC partnership	Two systems that are partnered to perform RC over native FC links.
Failover	Failure of a node within an I/O group causes the volume access to go through the surviving node. The IP addresses fail over to the surviving node in the I/O group. When the configuration node of the system fails, management IPs also fail over to an alternative node.
Failback	When the failed node rejoins the system, all failed over IP addresses are failed back from the surviving node to the rejoined node, and volume access is restored through this node.
linkbandwidthmbits	Aggregate bandwidth of all physical links between two sites in Mbps.

IP partnership terminology	Description
IP partnership or partnership over native IP links	These terms are used to describe the IP partnership feature.
Discovery	<p>Process by which two IBM Spectrum Virtualize systems exchange information about their IP address configuration. For IP-based partnerships, only IP addresses configured for RC are discovered.</p> <p>For example, the first Discovery takes place when the user is running the <code>mkippartnership</code> CLI command. Subsequent Discoveries can take place as a result of user activities (configuration changes) or as a result of hardware failures (for example, node failure, ports failure, and so on).</p>

## 10.8.6 States of IP partnership

The different partnership states in IP partnership are listed in Table 10-14.

Table 10-14 States of IP partnership

State	Systems connected	Support for active RC I/O	Comments
Partially_Configured_Local	No	No	This state indicates that the initial discovery is complete.
Fully_Configured	Yes	Yes	Discovery successfully completed between two systems, and the two systems can establish RC relationships.
Fully_Configured_Stopped	Yes	Yes	The partnership is stopped on the system.
Fully_Configured_Remote_Stopped	Yes	No	The partnership is stopped on the remote system.
Not_Present	Yes	No	The two systems cannot communicate with each other. This state is also seen when data paths between the two systems are not established.
Fully_Configured_Exceeded	Yes	No	There are too many systems in the network, and the partnership from the local system to remote system is disabled.
Fully_Configured_Excluded	No	No	The connection is excluded because of too many problems, or either system cannot support the I/O work load for the MM and GM relationships.

The process to establish two systems in the IP partnerships includes the following steps:

1. The administrator configures the CHAP secret on both the systems. This step is not mandatory, and users can choose to not configure the CHAP secret.
2. The administrator configures the system IP addresses on both local and remote systems so that they can discover each other over the network.
3. If you want to use VLANs, configure your local area network (LAN) switches and Ethernet ports to use VLAN tagging.

4. The administrator configures the systems ports on each node in both of the systems by using the GUI (or the `cfgport ip` CLI command), and completes the following steps:
  - a. Configure the IP addresses for RC data.
  - b. Add the IP addresses in the respective RC port group.
  - c. Define whether the host access on these ports over iSCSI is allowed.
5. The administrator establishes the partnership with the remote system from the local system where the partnership state then changes to `Partially_Configured_Local`.
6. The administrator establishes the partnership from the remote system with the local system. If this process is successful, the partnership state then changes to the `Fully_Configured`, which implies that the partnerships over the IP network were successfully established. The partnership state momentarily remains `Not_Present` before moving to the `Fully_Configured` state.
7. The administrator creates MM, GM, and GMCV relationships.

**Partnership consideration:** When the partnership is created, no master or auxiliary status is defined or implied. The partnership is equal. The concepts of *master or auxiliary* and *primary or secondary* apply to volume relationships only, not to system partnerships.

## 10.8.7 Remote Copy groups

This section describes RC groups (or RC port groups) and different ways to configure the links between the two remote systems. The two IBM Spectrum Virtualize systems can be connected to each other over one link or at most, two links. To address the requirement to enable the systems to know about the physical links between the two sites, the concept of RC port groups was introduced.

RC port group ID is a numerical tag that is associated with an IP port of an IBM Spectrum Virtualize system to indicate to which physical IP link it is connected. Multiple nodes might be connected to the same physical long-distance link, and must therefore share RC port group ID.

In scenarios with two physical links between the local and remote clusters, two RC port group IDs must be used to designate which IP addresses are connected to which physical link. This configuration must be done by the system administrator by using the GUI or running the `cfgport ip` CLI command.

**Remember:** IP ports on both partners must be configured with identical RC port group IDs for the partnership to be established correctly.

The IBM Spectrum Virtualize system IP addresses that are connected to the same physical link are designated with identical RC port groups. The system supports three RC groups: 0, 1, and 2.

The systems' IP addresses are, by default, in RC port group 0. Ports in port group 0 are not considered for creating RC data paths between two systems. For partnerships to be established over IP links directly, IP ports must be configured in RC group 1 if a single inter-site link exists, or in RC groups 1 and 2 if two inter-site links exist.

You can assign one IPv4 address and one IPv6 address to each Ethernet port on the system platforms. Each of these IP addresses can be shared between iSCSI host attach and the IP partnership. The user must configure the required IP address (IPv4 or IPv6) on an Ethernet port with a RC port group.

The administrator might want to use IPv6 addresses for RC operations and use IPv4 addresses on that same port for iSCSI host attach. This configuration also implies that for two systems to establish an IP partnership, both systems must have IPv6 addresses that are configured.

Administrators can choose to dedicate an Ethernet port for IP partnership only. In that case, host access must be specifically disabled for that IP address and any other IP address that is configured on that Ethernet port.

**Note:** To establish an IP partnership, each IBM SAN Volume Controller node must have only a single RC port group that is configured: 1 or 2. The remaining IP addresses must be in RC port group 0.

### 10.8.8 Supported configurations

**Note:** For explanation purposes, this section shows a node with two ports available: 1 and 2. This number generally increments when IBM SAN Volume Controller nodes model DH8 or SV1 are used.

The following supported configurations for IP partnership that were in the first release are described in this section:

- ▶ Two 2-node systems in IP partnership over a single inter-site link, as shown in Figure 10-98 (configuration 1).

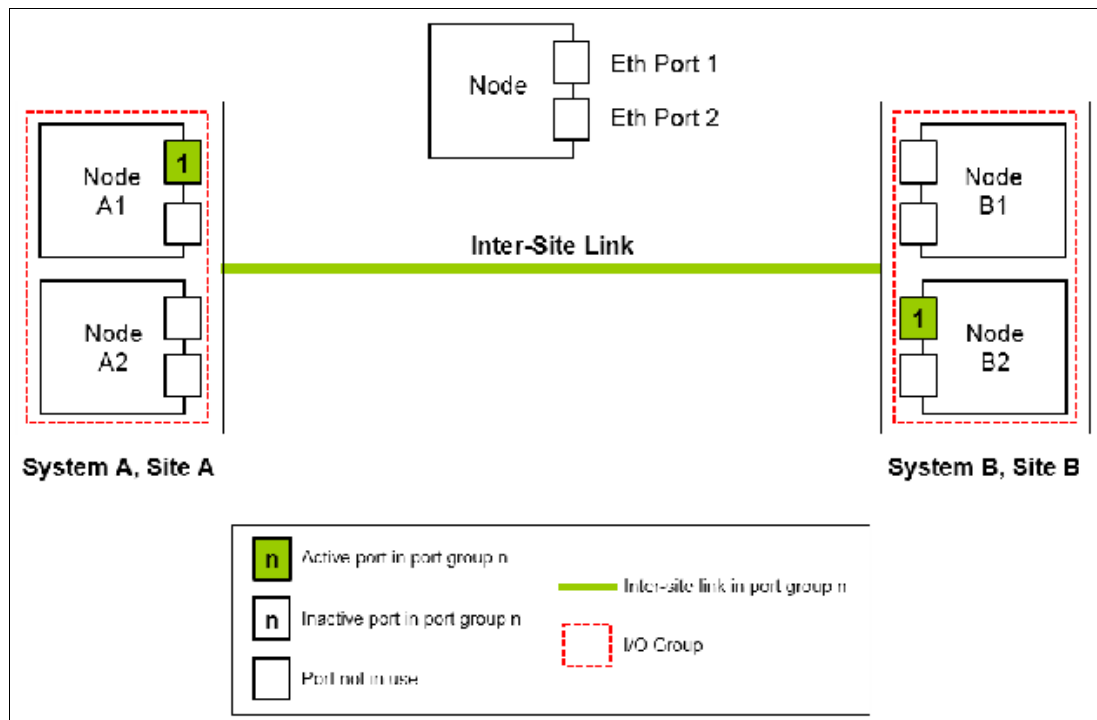


Figure 10-98 Single link with only one Remote Copy port group configured in each system

As shown in Figure 10-98, two systems are available:

- System A
- System B

A single RC port group 1 is created on Node A1 on System A and on Node B2 on System B because only a single inter-site link is used to facilitate the IP partnership traffic. An administrator might choose to configure the RC port group on Node B1 on System B rather than Node B2.

At any time, only the IP addresses that are configured in RC port group 1 on the nodes in System A and System B participate in establishing data paths between the two systems after the IP partnerships are created. In this configuration, no failover ports are configured on the partner node in the same I/O group.

This configuration has the following characteristics:

- Only one node in each system has a RC port group that is configured, and no failover ports are configured.
  - If the Node A1 in System A or the Node B2 in System B encounter a failure, the IP partnership stops and enters the Not\_Present state until the failed nodes recover.
  - After the nodes recover, the IP ports fail back, the IP partnership recovers, and the partnership state goes to the Fully\_Configured state.
  - If the inter-site system link fails, the IP partnerships change to the Not\_Present state.
  - This configuration is not recommended because it is not resilient to node failures.
- Two 2-node systems in IP partnership over a single inter-site link (with failover ports configured), as shown in Figure 10-99 (configuration 2).

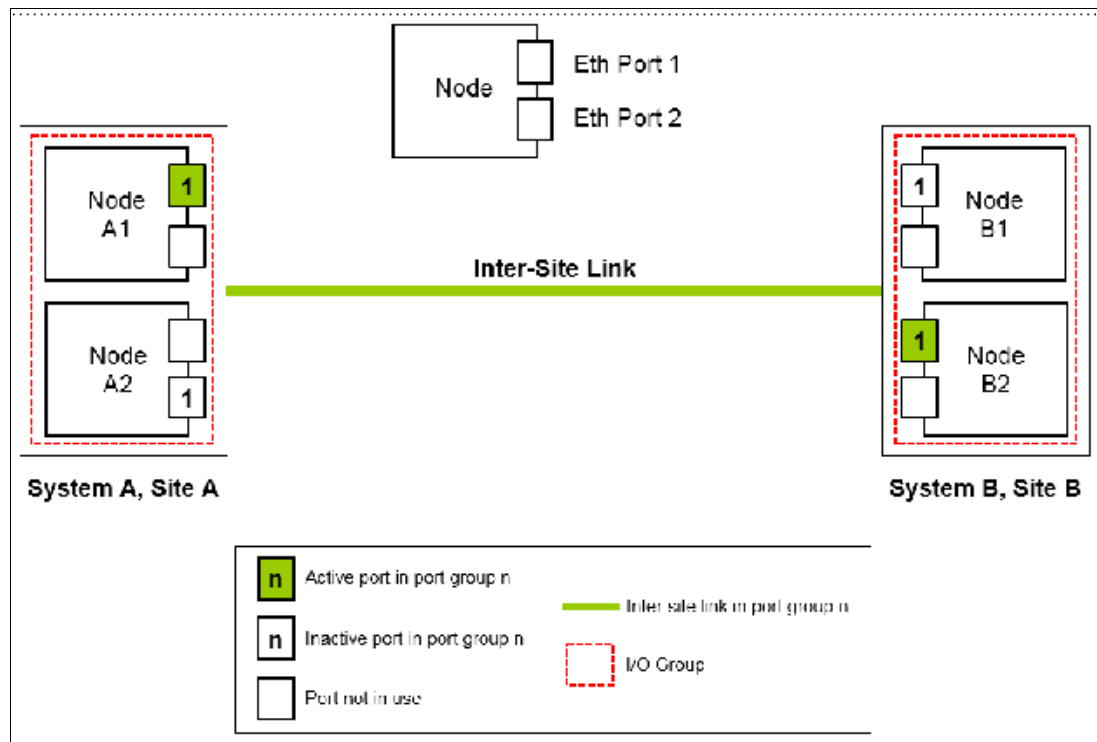


Figure 10-99 One Remote Copy group on each system and nodes with failover ports configured

As shown in Figure 10-99, two systems are available:

- System A
- System B



A single RC port group 1 is configured on two Ethernet ports, one each on Node A1 and Node A2 on System A. Similarly, a single RC port group is configured on two Ethernet ports on Node B1 and Node B2 on System B.

Although two ports on each system are configured for RC port group 1, only one Ethernet port in each system actively participates in the IP partnership process. This selection is determined by a path configuration algorithm that is designed to choose data paths between the two systems to optimize performance.

The other port on the partner node in the I/O group behaves as a standby port that is used if a node fails. If Node A1 fails in System A, IP partnership continues servicing replication I/O from Ethernet Port 2 because a failover port is configured on Node A2 on Ethernet Port 2.

However, it might take some time for discovery and path configuration logic to reestablish paths post failover. This delay can cause partnerships to change to Not\_Present for that time. The details of the particular IP port that is actively participating in IP partnership is provided in the `1sport ip` output (reported as used).

This configuration has the following characteristics:

- Each node in the I/O group has the same RC port group that is configured. However, only one port in that RC port group is active at any time at each system.
  - If the Node A1 in System A or the Node B2 in System B fails in the respective systems, IP partnerships rediscovery is triggered and continues servicing the I/O from the failover port.
  - The discovery mechanism that is triggered because of failover might introduce a delay where the partnerships momentarily change to the Not\_Present state and recover.
- Two 4-node systems in IP partnership over a single inter-site link (with failover ports configured), as shown in Figure 10-100 on page 624 (configuration 3).

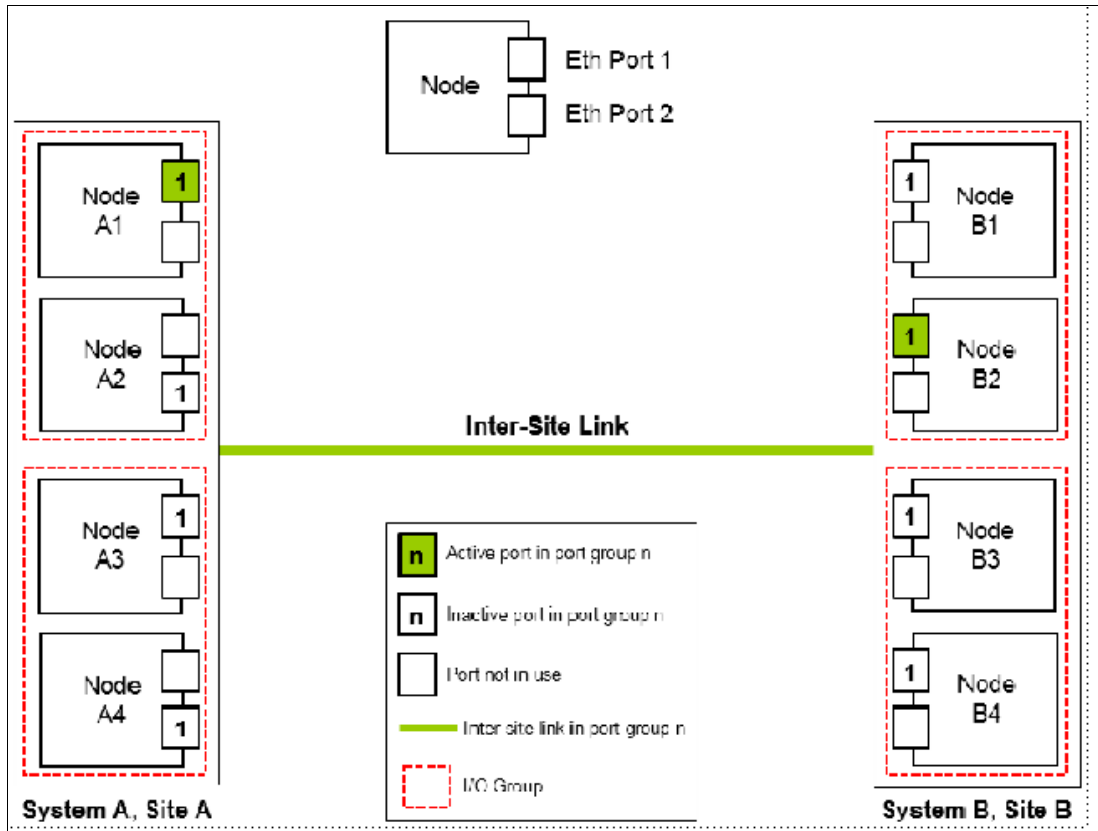


Figure 10-100 Multinode systems single inter-site link with only one RC port group

As shown in Figure 10-100, two 4-node systems are available:

- System A
- System B

A single RC port group 1 is configured on nodes A1, A2, A3, and A4 on System A, Site A; and on nodes B1, B2, B3, and B4 on System B, Site B. Although four ports are configured for RC group 1, only one Ethernet port in each RC port group on each system actively participates in the IP partnership process.

Port selection is determined by a path configuration algorithm. The other ports play the role of standby ports.

If Node A1 fails in System A, the IP partnership selects one of the remaining ports that is configured with RC port group 1 from any of the nodes from either of the two I/O groups in System A. However, it might take some time (generally seconds) for discovery and path configuration logic to reestablish paths post failover. This process can cause partnerships to change to the Not\_Present state.

This result causes RC relationships to stop. The administrator might need to manually verify the issues in the event log and start the relationships or RC consistency groups, if they do not autorecover. The details of the particular IP port actively participating in the IP partnership process is provided in the **1sportip** view (reported as used).

This configuration has the following characteristics:

- Each node has the RC port group that is configured in both I/O groups. However, only one port in that RC port group remains active and participates in IP partnership on each system.

- If the Node A1 in System A or the Node B2 in System B were to encounter some failure in the system, IP partnerships discovery is triggered and it continues servicing the I/O from the failover port.
  - The discovery mechanism that is triggered because of failover might introduce a delay wherein the partnerships momentarily change to the Not\_Present state and then recover.
  - The bandwidth of the single link is used completely.
- Eight-node system in IP partnership with four-node system over single inter-site link, as shown in Figure 10-101 (configuration 4).

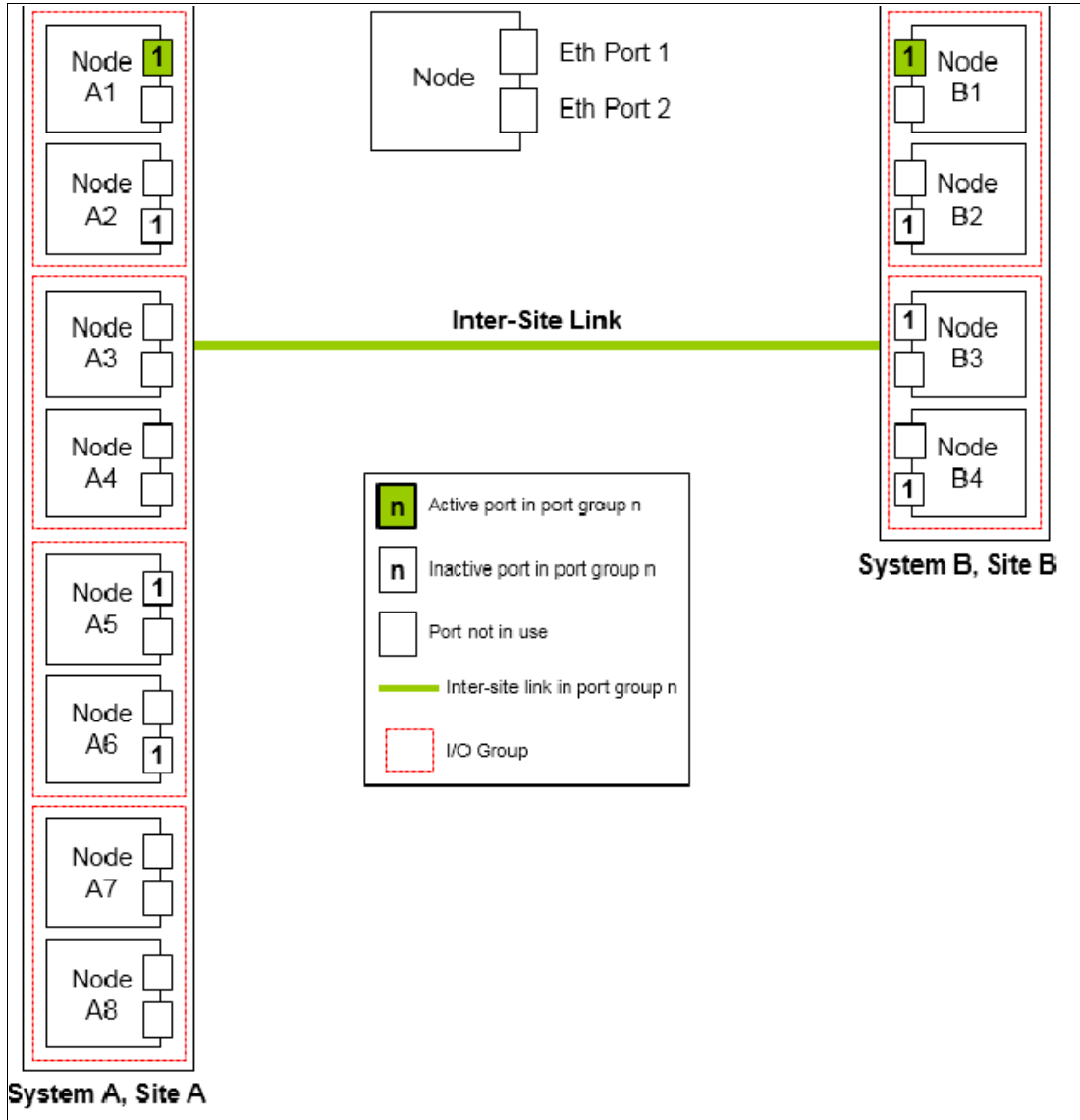


Figure 10-101 Multinode systems single inter-site link with only one Remote Copy port group

As shown in Figure 10-101, an eight-node system (System A in Site A) and a four-node system (System B in Site B) are used. A single RC port group 1 is configured on nodes A1, A2, A5, and A6 on System A at Site A. Similarly, a single RC port group 1 is configured on nodes B1, B2, B3, and B4 on System B.

Although four I/O groups (eight nodes) are in System A, any two I/O groups at maximum are supported to be configured for IP partnerships. If Node A1 fails in System A, IP partnership continues by using one of the ports that is configured in RC port group from any of the nodes from either of the two I/O groups in System A.

However, it might take some time for discovery and path configuration logic to reestablish paths post-failover. This delay might cause partnerships to change to the Not\_Present state.

This process can lead to RC relationships stopping, and the administrator must manually start them if the relationships do not auto-recover. The details of which particular IP port is actively participating in IP partnership process is provided in **lspport ip** output (reported as used).

This configuration features the following characteristics:

- ▶ Each node has the RC port group that is configured in both the I/O groups that are identified for participating in IP Replication. However, only one port in that RC port group remains active on each system and participates in IP Replication.
- ▶ If the Node A1 in System A or the Node B2 in System B fails in the system, the IP partnerships trigger discovery and continue servicing the I/O from the failover ports.
- ▶ The discovery mechanism that is triggered because of failover might introduce a delay wherein the partnerships momentarily change to the Not\_Present state and then recover.
- ▶ The bandwidth of the single link is used completely.
- ▶ Two 2-node systems with two inter-site links, as shown in Figure 10-102 on page 627 (configuration 5).

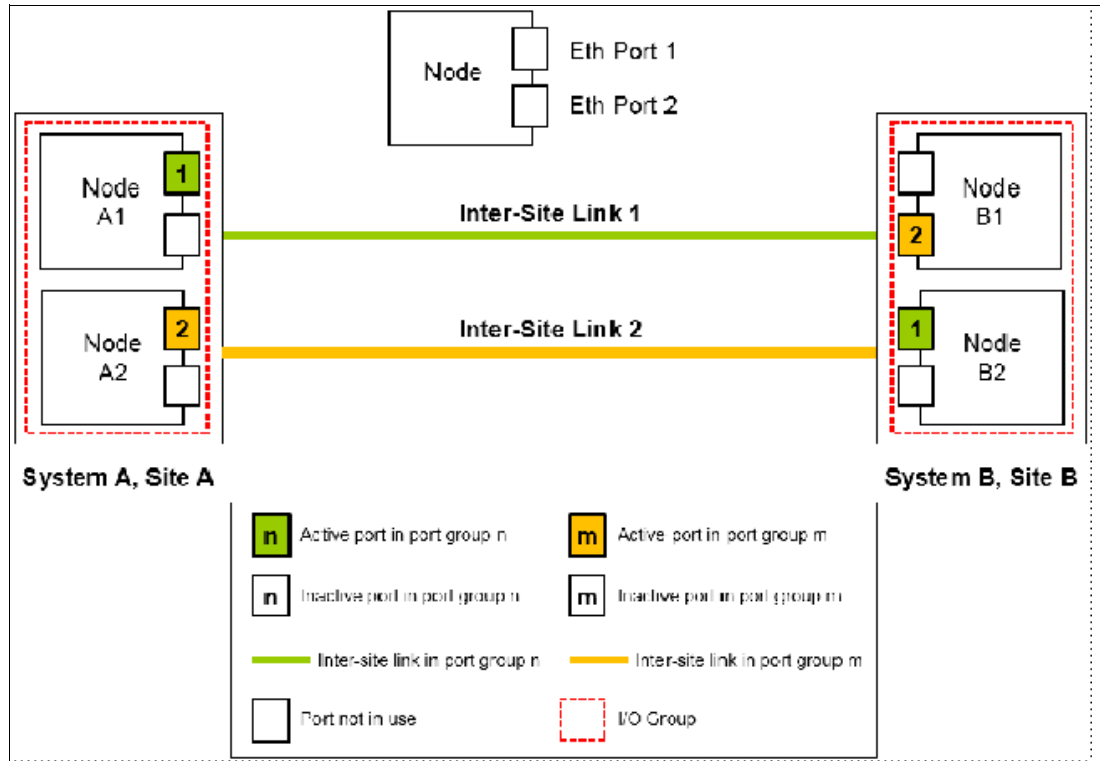


Figure 10-102 Dual links with two Remote Copy groups on each system configured

As shown in Figure 10-102, RC port groups 1 and 2 are configured on the nodes in System A and System B because two inter-site links are available. In this configuration, the failover ports are not configured on partner nodes in the I/O group. Instead, the ports are maintained in different RC port groups on both of the nodes. They remain active and participate in IP partnership by using both of the links.

However, if either of the nodes in the I/O group fail (that is, if Node A1 on System A fails), the IP partnership continues only from the available IP port that is configured in RC port group 2. Therefore, the effective bandwidth of the two links is reduced to 50% because only the bandwidth of a single link is available until the failure is resolved.

This configuration has the following characteristics:

- Two inter-site links and two RC port groups are configured.
- Each node has only one IP port in RC port group 1 or 2.
- Both the IP ports in the two RC port groups participate simultaneously in IP partnerships. Therefore, both of the links are used.
- During node failure or link failure, the IP partnership traffic continues from the other available link and the port group. Therefore, if two links of 10 Mbps each are available and you have 20 Mbps of effective link bandwidth, bandwidth is reduced to 10 Mbps only during a failure.
- After the node failure or link failure is resolved and failback occurs, the entire bandwidth of both of the links is available as before.

- ▶ Two 4-node systems in IP partnership with dual inter-site links, as shown in Figure 10-103 (configuration 6).

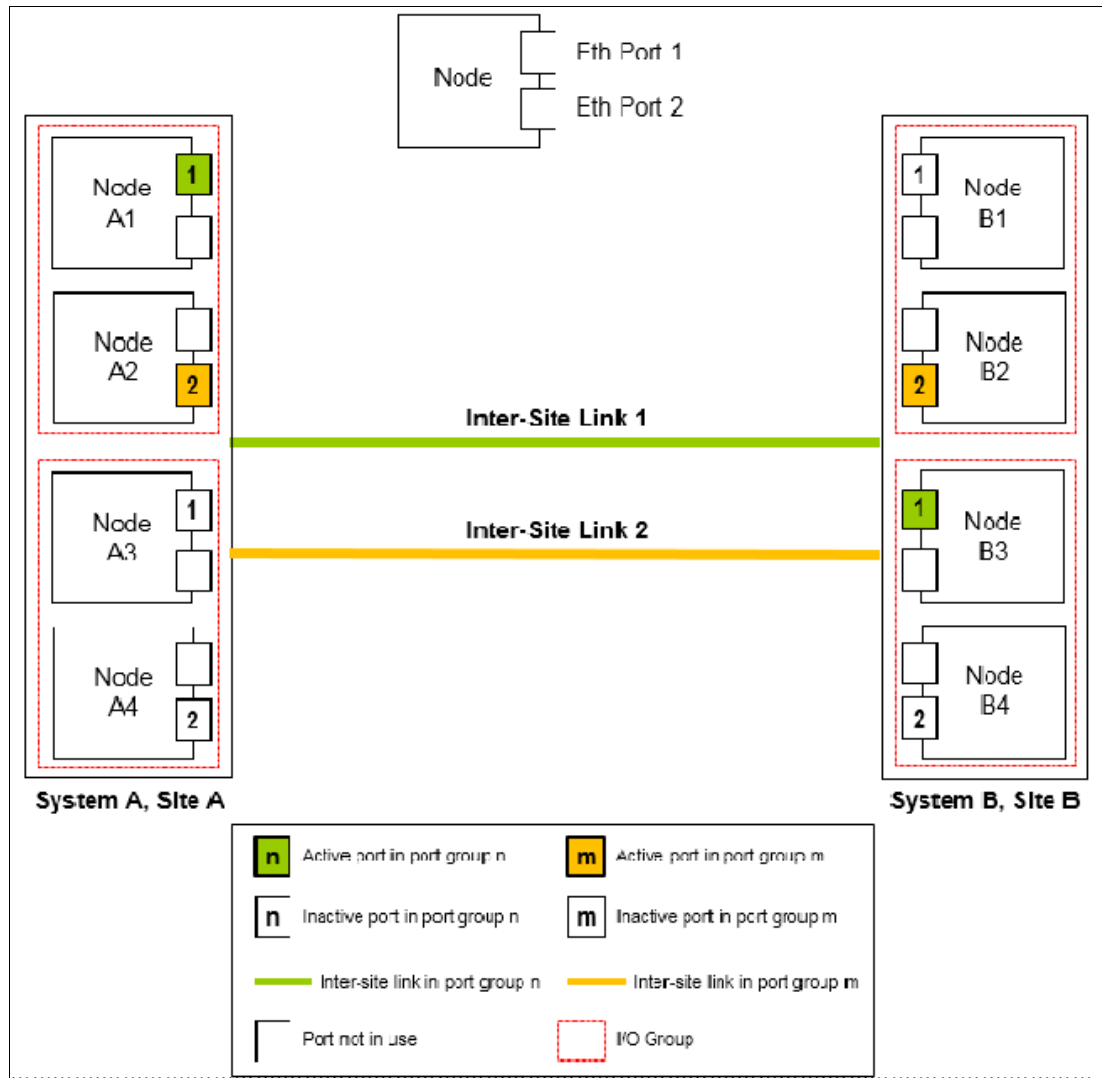


Figure 10-103 Multinode systems with dual inter-site links between the two systems

As shown in Figure 10-103, two 4-node systems are used:

- System A
- System B

This configuration is an extension of Configuration 5 to a multinode multi-I/O group environment. This configuration has two I/O groups, and each node in the I/O group has a single port that is configured in RC port groups 1 or 2.

Although two ports are configured in RC port groups 1 and 2 on each system, only one IP port in each RC port group on each system actively participates in IP partnership. The other ports that are configured in the same RC port group act as standby ports in the event of failure. Which port in a configured RC port group participates in IP partnership at any moment is determined by a path configuration algorithm.

In this configuration, if Node A1 fails in System A, IP partnership traffic continues from Node A2 (that is, RC port group 2) and at the same time the failover also causes discovery in RC port group 1.

Therefore, the IP partnership traffic continues from Node A3 on which RC port group 1 is configured. The details of the particular IP port that is actively participating in IP partnership process is provided in the `1sport ip` output (reported as used).

This configuration has the following characteristics:

- Each node has the RC port group that is configured in the I/O groups 1 or 2. However, only one port per system in both RC port groups remains active and participates in IP partnership.
  - Only a single port per system from each configured RC port group participates simultaneously in IP partnership. Therefore, both of the links are used.
  - During node failure or port failure of a node that is actively participating in IP partnership, IP partnership continues from the alternative port because another port is in the system in the same RC port group but in a different I/O group.
  - The pathing algorithm can start discovery of available ports in the affected RC port group in the second I/O group and pathing is reestablished, which restores the total bandwidth, so both of the links are available to support IP partnership.
- Eight-node system in IP partnership with a four-node system over dual inter-site links, as shown in Figure 10-104 on page 630 (configuration 7).

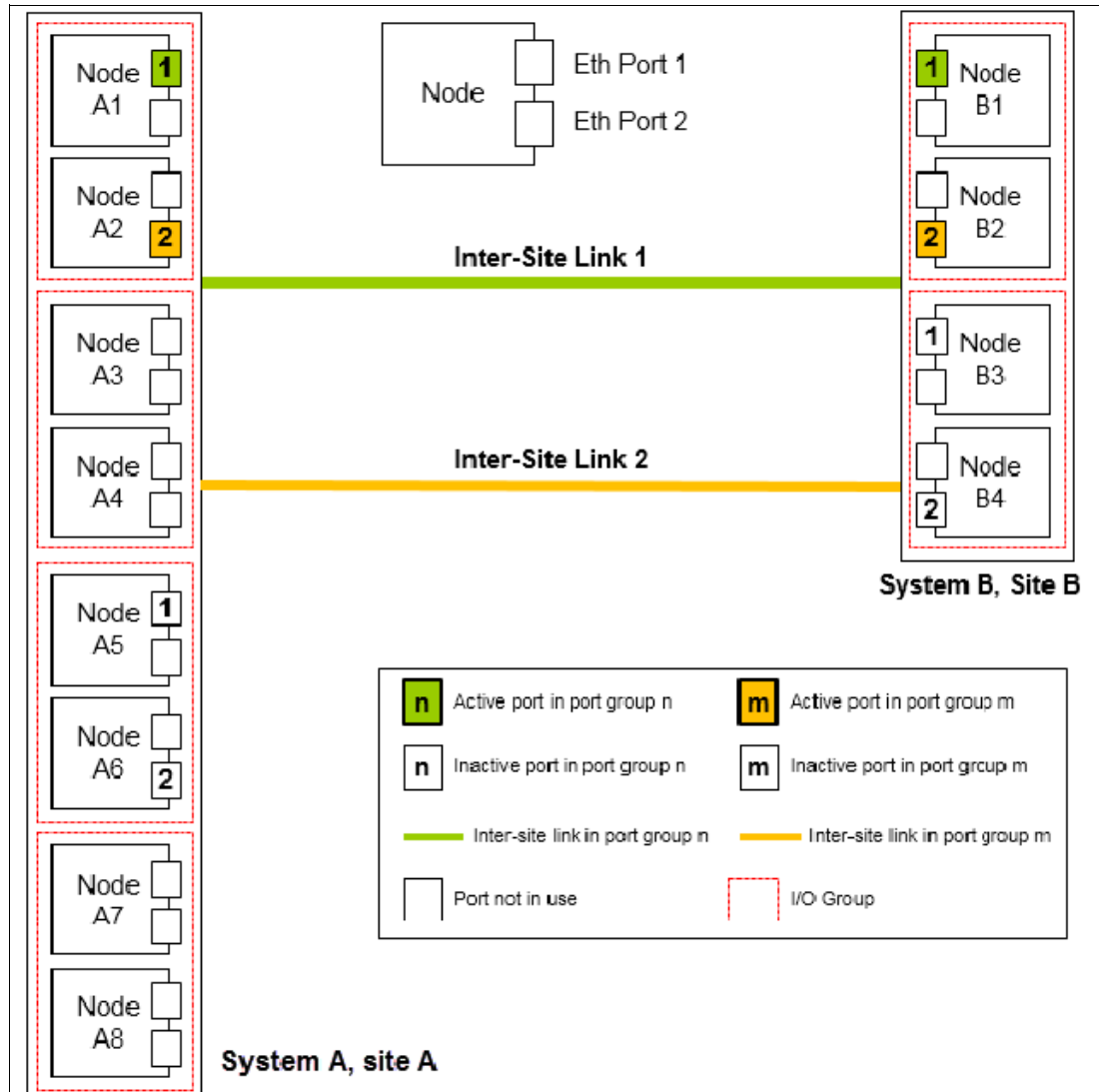


Figure 10-104 Multinode systems (two I/O groups on each system) with dual inter-site links between the two systems

As shown in Figure 10-104, an eight-node System A in Site A and a four-node System B in Site B is used. Because a maximum of two I/O groups in IP partnership is supported in a system, although four I/O groups (eight nodes) exist, nodes from only two I/O groups are configured with RC port groups in System A. The remaining or all of the I/O groups can be configured to be RC partnerships over FC.

In this configuration, two links and two I/O groups are configured with RC port groups 1 and 2, but path selection logic is managed by an internal algorithm. Therefore, this configuration depends on the pathing algorithm to decide which of the nodes actively participates in IP partnership. Even if Node A5 and Node A6 are configured with RC port groups properly, active IP partnership traffic on both of the links might be driven from Node A1 and Node A2 only.

If Node A1 fails in System A, IP partnership traffic continues from Node A2 (that is, RC port group 2). The failover also causes IP partnership traffic to continue from Node A5 on which RC port group 1 is configured. The details of the particular IP port actively participating in IP partnership process is provided in the `1sport ip` output (reported as used).



This configuration has the following characteristics:

- Two I/O groups with nodes in those I/O groups are configured in two RC port groups because two inter-site links are used for participating in IP partnership. However, only one port per system in a particular RC port group remains active and participates in IP partnership.
  - One port per system from each RC port group participates in IP partnership simultaneously. Therefore, both of the links are used.
  - If a node or port on the node that is actively participating in IP partnership fails, the RC data path is established from that port because another port is available on an alternative node in the system with the same RC port group.
  - The path selection algorithm starts discovery of available ports in the affected RC port group in the alternative I/O groups and paths are reestablished, which restores the total bandwidth across both links.
  - The remaining or all of the I/O groups can be in RC partnerships with other systems.
- An example of an *unsupported* configuration for a single inter-site link is shown in Figure 10-105 (configuration 8).

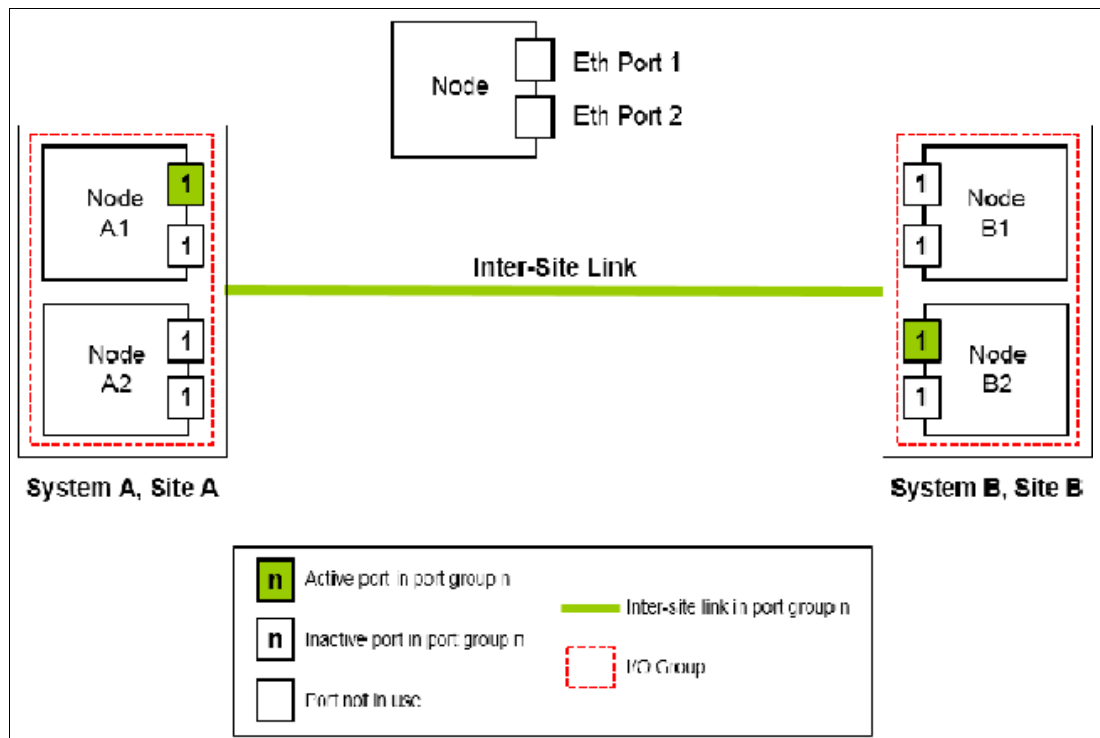


Figure 10-105 Two node systems with single inter-site link and Remote Copy port groups configured

As shown in Figure 10-105, this configuration is similar to Configuration 2, but differs because each node now has the same RC port group that is configured on more than one IP port.

On any node, only one port at any time can participate in IP partnership. Configuring multiple ports in the same RC group on the same node is *not supported*.

- An example of an *unsupported* configuration for a dual inter-site link is shown in Figure 10-106 (configuration 9).

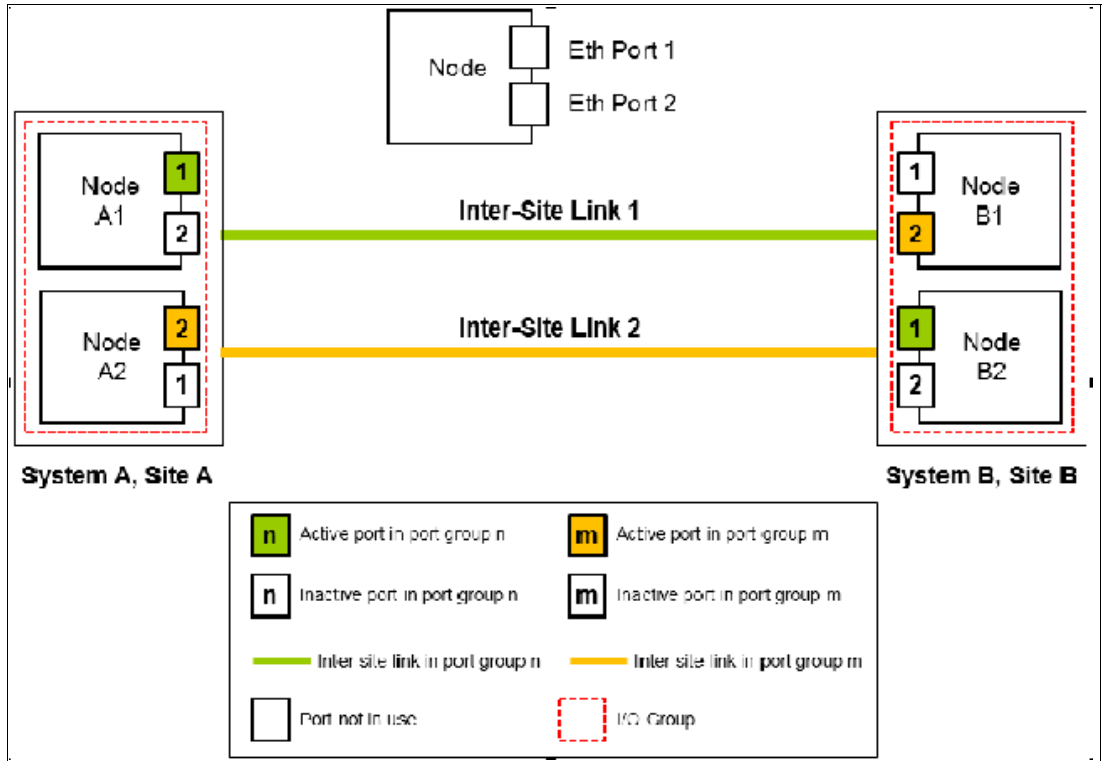


Figure 10-106 Dual Links with two Remote Copy Port Groups with failover Port Groups configured

As shown in Figure 10-106, this configuration is similar to Configuration 5, but differs because each node now also has two ports that are configured with RC port groups. In this configuration, the path selection algorithm can select a path that might cause partnerships to change to the Not\_Present state and then recover, which results in a configuration restriction. The use of this configuration is not recommended until the configuration restriction is lifted in future releases.

- An example deployment for configuration 2 with a dedicated inter-site link is shown in Figure 10-107 (configuration 10).

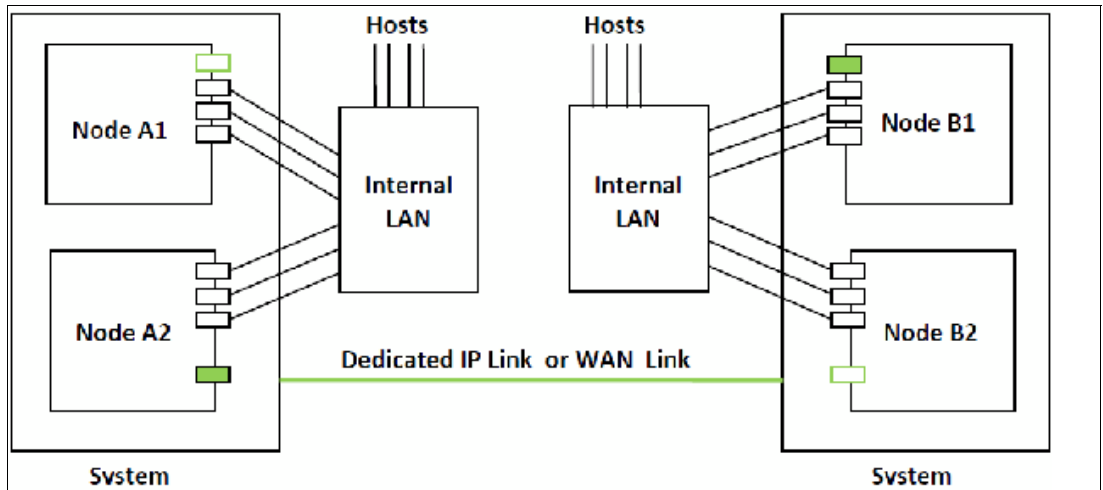


Figure 10-107 Deployment example

In this configuration, one port on each node in System A and System B is configured in RC group 1 to establish IP partnership and support RC relationships. A dedicated inter-site link is used for IP partnership traffic, and iSCSI host attach is disabled on those ports.

The following configuration steps are used:

- a. Configure system IP addresses properly. As such, they can be reached over the inter-site link.
  - b. Qualify if the partnerships must be created over IPv4 or IPv6, and then assign IP addresses and open firewall ports 3260 and 3265.
  - c. Configure IP ports for RC on both the systems by using the following settings:
    - Remote copy group: 1
    - Host: No
    - Assign IP address
  - d. Check that the maximum transmission unit (MTU) levels across the network meet the requirements as set (default MTU is 1500).
  - e. Establish IP partnerships from both of the systems.
  - f. After the partnerships are in the Fully\_Configured state, you can create the RC relationships.
- Figure 10-107 on page 632 is an example deployment for the configuration that is shown in Figure 10-101 on page 625. Ports that are shared with host access are shown in Figure 10-108 (configuration 11).

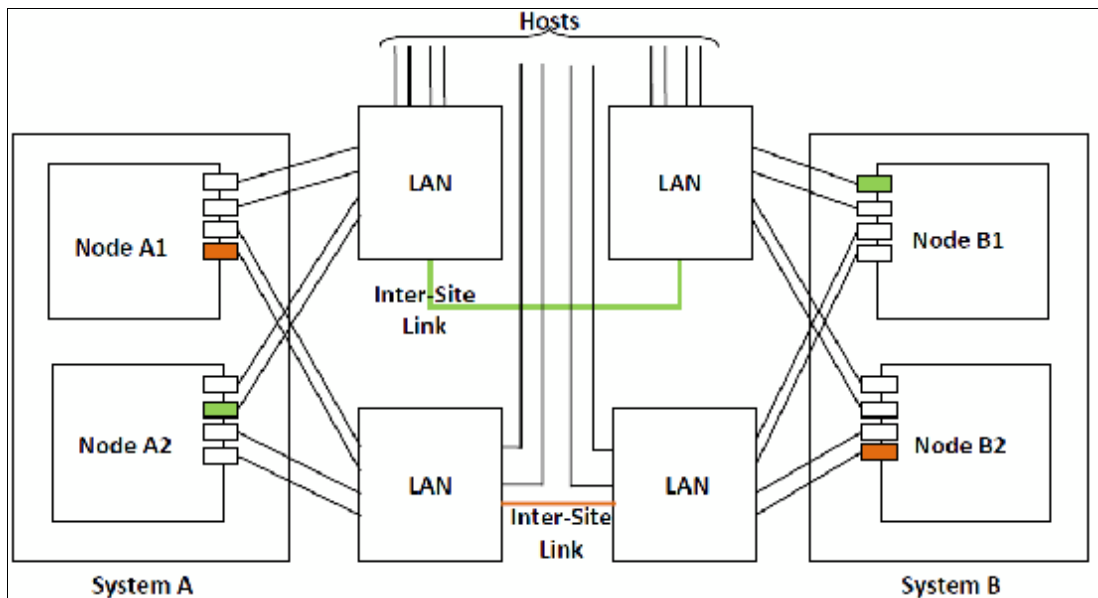


Figure 10-108 Deployment example

In this configuration, IP ports are to be shared by both iSCSI hosts and for IP partnership.

The following configuration steps are used:

- a. Configure System IP addresses properly so that they can be reached over the inter-site link.
- b. Qualify if the partnerships must be created over IPv4 or IPv6, and then assign IP addresses and open firewall ports 3260 and 3265.

- c. Configure IP ports for RC on System A1 by using the following settings:
  - Node 1:
    - Port 1, remote copy port group 1
    - Host: Yes
    - Assign IP address
  - Node 2:
    - Port 4, remote copy port group 2
    - Host: Yes
    - Assign IP address
- d. Configure IP ports for RC on System B1 by using the following settings:
  - Node 1:
    - Port 1, remote copy port group 1
    - Host: Yes
    - Assign IP address
  - Node 2:
    - Port 4, remote copy port group 2
    - Host: Yes
    - Assign IP address
- e. Check the MTU levels across the network as set (default MTU is 1500 on SAN Volume Controller and Storwize V7000).
- f. Establish IP partnerships from both systems.
- g. After the partnerships are in the Fully\_Configured state, you can create the RC relationships.

## 10.9 Managing Remote Copy by using the GUI

It is often easier to control MM/GM with the GUI if you have few mappings. When many mappings are used, run your commands by using the CLI. This section describes the tasks that you can perform at a RC level.

**Note:** The **Copy Services** → **Consistency Groups** menu relates to FlashCopy consistency groups only, not RC groups.

The following panels are used to visualize and manage your remote copies:

► Remote Copy panel

To open the Remote Copy panel, click **Copy Services** → **Remote Copy** in the main menu, as shown in Figure 10-109.

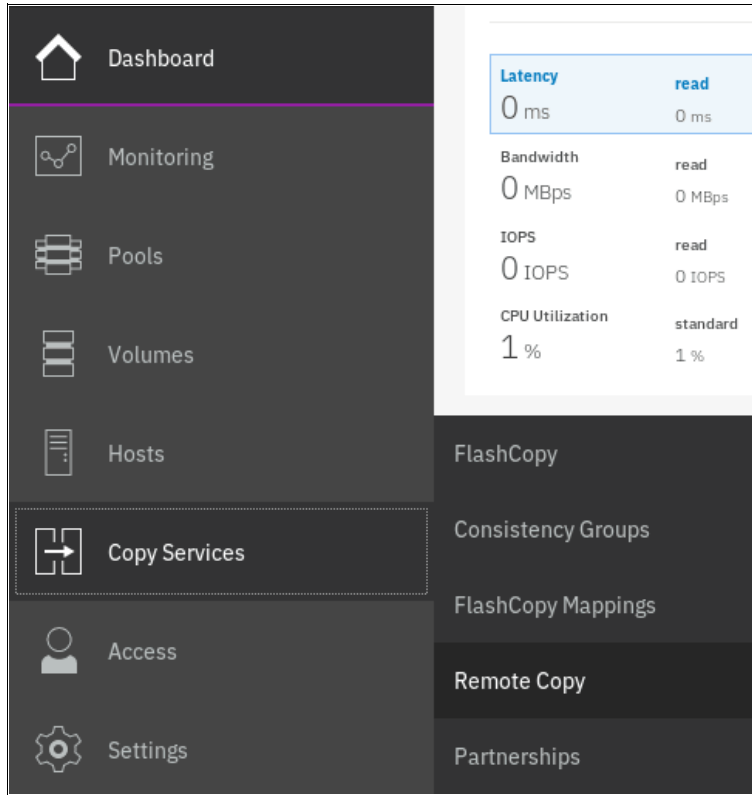


Figure 10-109 Remote Copy menu

The Remote Copy panel is displayed, as shown in Figure 10-110.

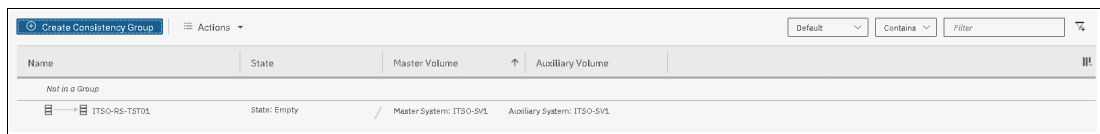


Figure 10-110 Remote Copy panel

► Partnerships panel

To open the Partnership panel, click **Copy Services** → **Partnership** in the main menu, as shown in Figure 10-111.

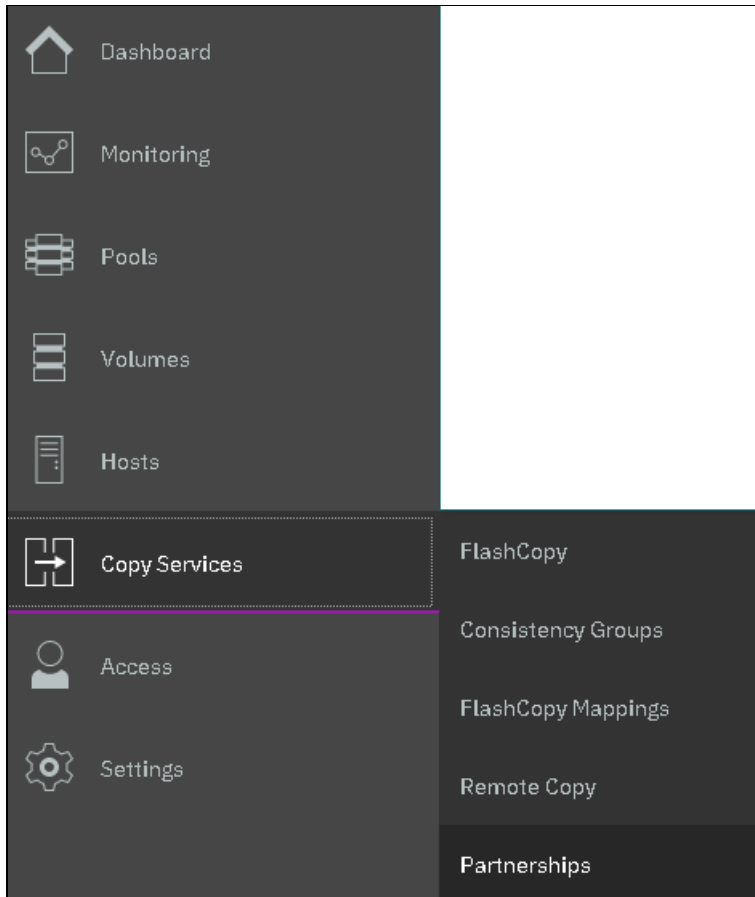


Figure 10-111 Partnership menu

The Partnership panel is displayed, as shown in Figure 10-112.

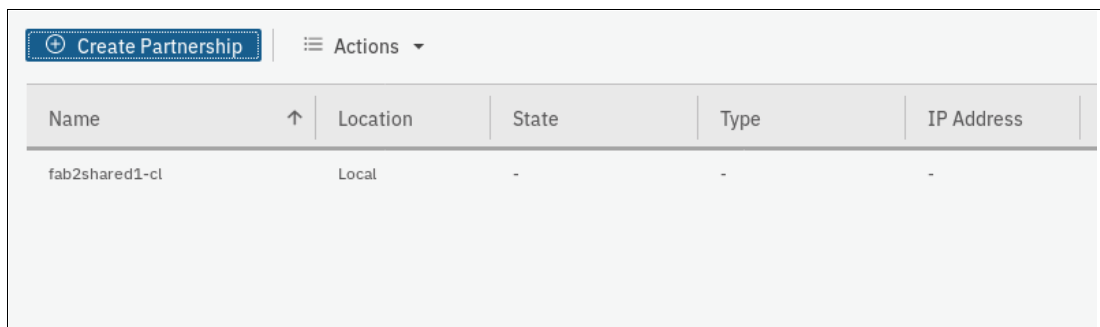


Figure 10-112 Partnership panel

## 10.9.1 Creating Fibre Channel partnership

**Intra-cluster MM:** If you are creating an intra-cluster MM, do not perform this next step to create the MM partnership. Instead, see 10.9.2, “Creating Remote Copy relationships” on page 638.

To create an FC partnership between IBM Spectrum Virtualize systems by using the GUI, open the Partnerships panel and click **Create Partnership** to create a partnership.

In the Create Partnership window, enter the following information, as shown in Figure 10-113:

1. Select the partnership type (**Fibre Channel** or **IP**). If you choose IP partnership, you must provide the IP address of the partner system and the partner system’s CHAP key.
2. If your partnership is based on Fibre Channel Protocol (FCP), select an available partner system from the menu. If no candidate is available, the This system does not have any candidates error message is displayed.
3. Enter a link bandwidth in Mbps that is used by the background copy process between the systems in the partnership.
4. Enter the background copy rate.
5. Click **OK** to confirm the partnership relationship.

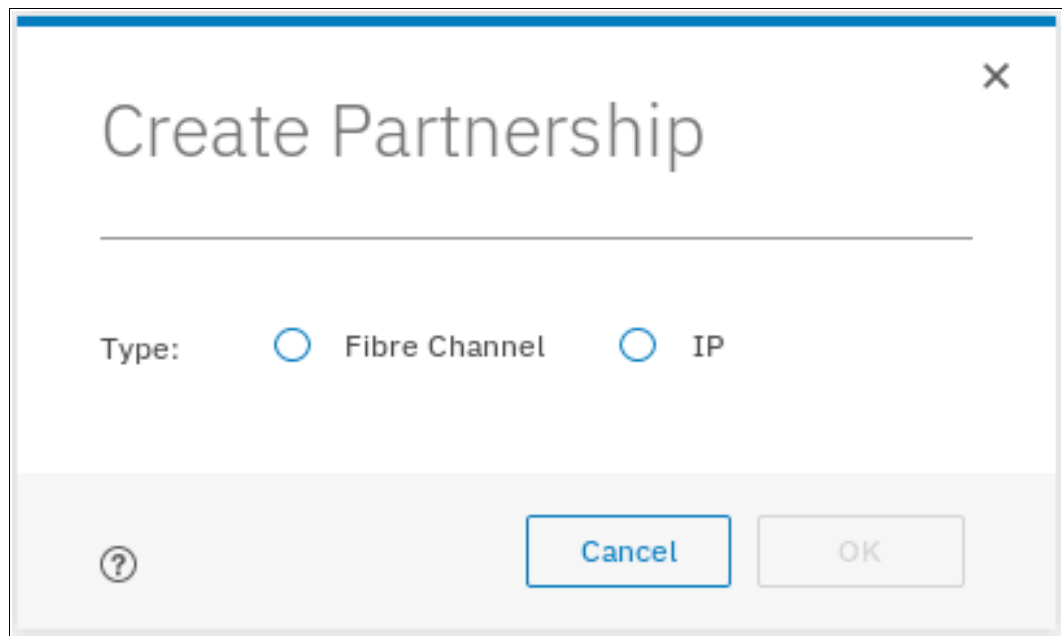


Figure 10-113 Creating a Partnership details

To fully configure the partnership between both systems, perform the same steps on the other system in the partnership. If not configured on the partner system, the partnership is displayed as **Partially Configured**.

When both sides of the system partnership are defined, the partnership is displayed as **Fully Configured** as shown in Figure 10-114 on page 638.

+ Create Partnership		☰ Actions ▾			
Name	↑	Location	State	Type	IP Address
fab2shared1-cl		Local	-	-	-
fab2shared1a-cl		Remote	✓ Fully Configured	Fibre Channel	-

Figure 10-114 Fully configured FC partnership

## 10.9.2 Creating Remote Copy relationships

This section shows how to create RC relationships for volumes with their respective remote targets. Before creating a relationship between a volume on the local system and a volume on a remote system, both volumes must exist and have the same virtual size.

To create a RC relationship, complete the following steps:

1. Open the **Copy Services** → **Remote Copy** panel.
2. Right-click the Consistency Group for which you want to create the relationship and select **Create Relationship**, as shown in Figure 10-115. If you want to create a stand-alone relationship (not in a consistency group), right-click the **Not in a Group** group.

+ Create Consistency Group		☰ Actions ▾	
Name		State	
<i>Not in a Group</i>			
ITSO-RCCG-01		State: Empty	/
		<div style="border: 1px solid #ccc; padding: 5px; margin-bottom: 5px;">Create Relationship</div> <div style="border: 1px solid #ccc; padding: 5px; margin-bottom: 5px;">Rename</div> <div style="border: 1px solid #ccc; padding: 5px; margin-bottom: 5px;">Start</div> <div style="border: 1px solid #ccc; padding: 5px; margin-bottom: 5px;">Stop</div> <div style="border: 1px solid #ccc; padding: 5px; margin-bottom: 5px;">Switch</div> <div style="border: 1px solid #ccc; padding: 5px; margin-bottom: 5px;">Edit Consistency Group</div> <div style="border: 1px solid #ccc; padding: 5px;">Delete</div>	
<i>Selected 1 Remote-copy consistency group</i>			

Figure 10-115 Creating Remote Copy relationships



3. In the Create Relationship window, select one of the following types of relationships that you want to create, as shown in Figure 10-116:

- MM
- GM (with or without Consistency Protection)
- GMCV

Click **Next**.

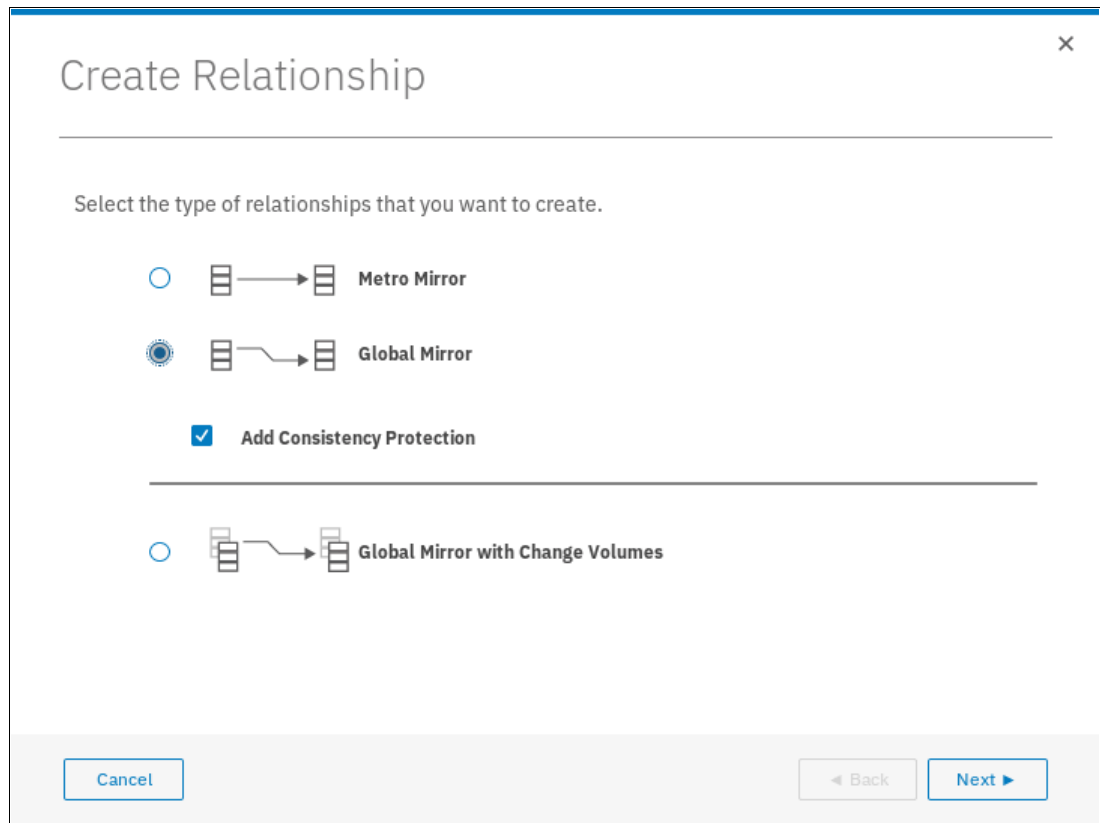


Figure 10-116 Creating a Remote Copy relationship

4. In the next window, select the master and auxiliary volumes, as shown in Figure 10-117. Click **ADD**.

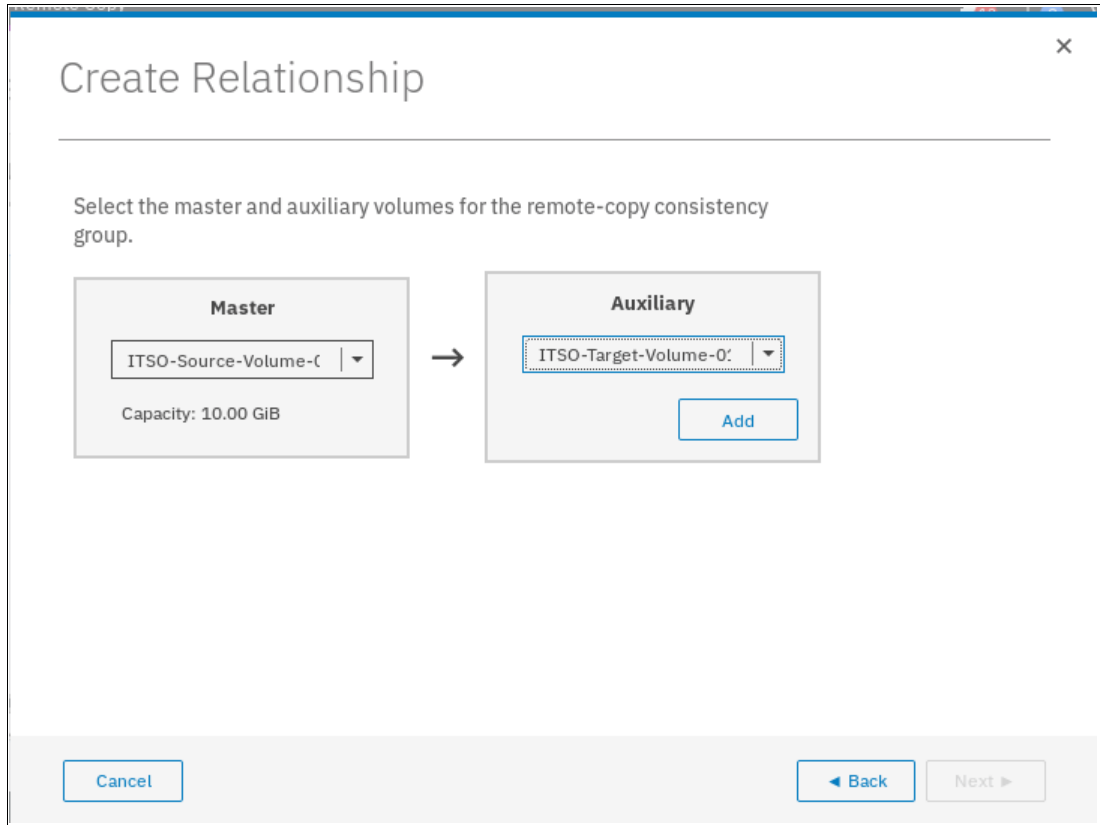


Figure 10-117 Selecting the master and auxiliary volumes

**Important:** The master and auxiliary volumes must be of equal size. Therefore, only the targets with the suitable size are shown in the list for a specific source volume.

5. In the next window, you can add change volumes if needed, as shown in Figure 10-118. Click **Finish**.

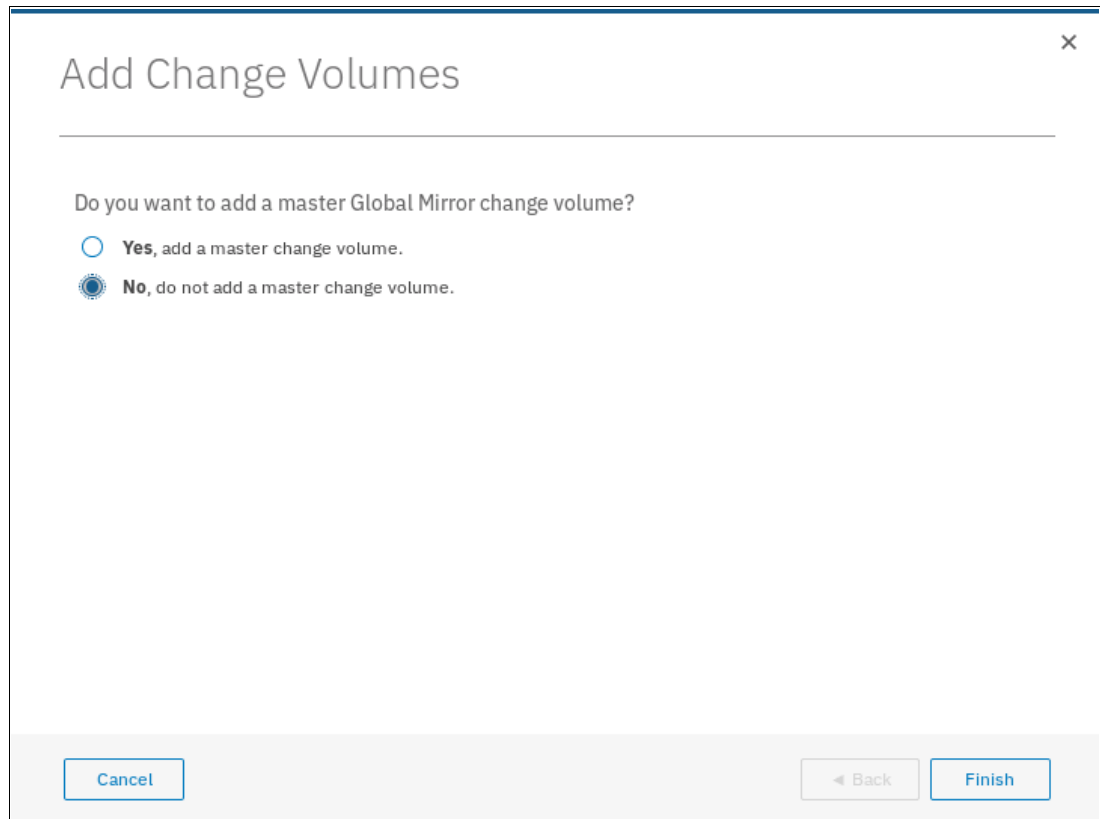


Figure 10-118 Add Change Volumes panel

6. In the next window, check the relationship that was created. If you want to add relationships at the same consistency group, add the new relationships, as shown in Figure 10-119. Then, click **Next**.

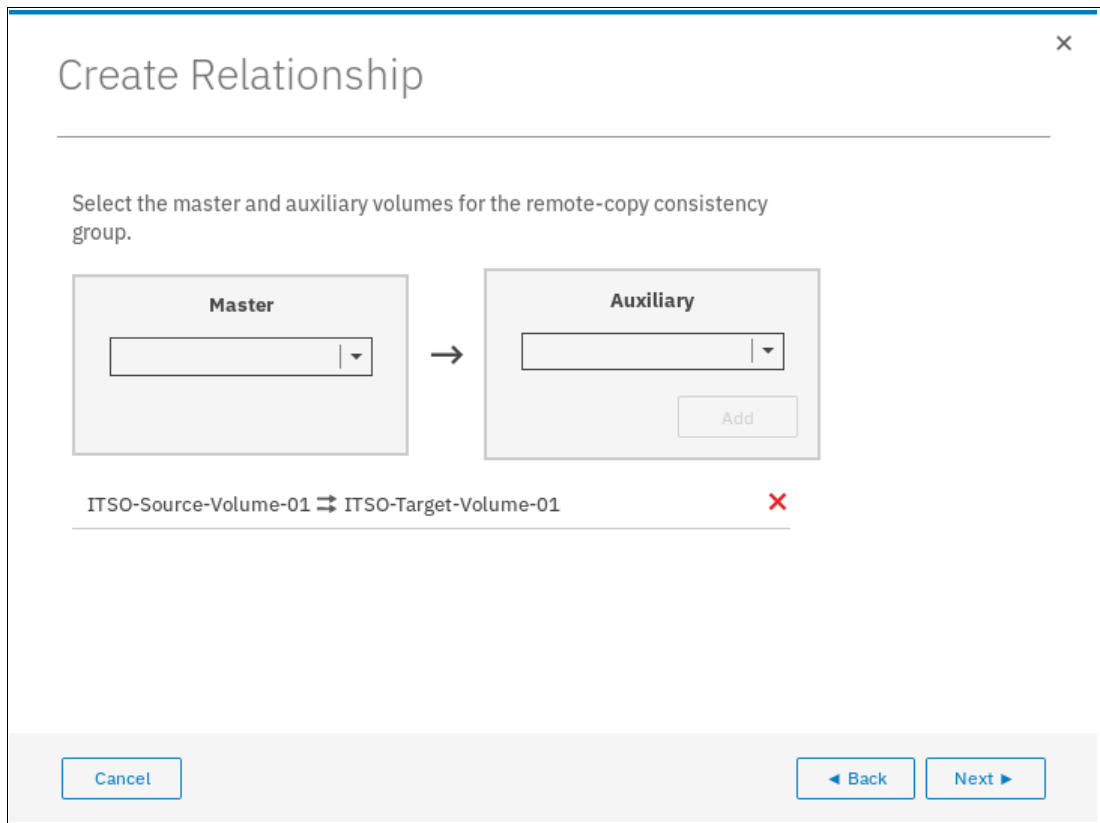


Figure 10-119 Checking and adding the relationship

7. In the next window, select whether the volumes are synchronized so that the relationship is created, as shown in Figure 10-120. Click **Finish**.

Create Relationship

Are the volumes in this relationship already synchronized?

Yes, the volumes are already synchronized.

No, the volumes are not synchronized.

Cancel Back Finish

Figure 10-120 Selecting if volumes are synchronized

**Note:** If the volumes are not synchronized, the initial copy copies the entire source volume to the remote target volume. If you suspect volumes are different or if you have a doubt, synchronize them to ensure consistency on both sides of the relationship.

### 10.9.3 Creating a consistency group

To create a consistency group, complete the following steps:

1. Open the **Copy Services** → **Remote Copy** panel and click **Create Consistency Group**, as shown in Figure 10-121.

Create Consistency Group Actions

Name	State
Not in a Group	
> ITSO-RCCG-01	State: Inconsistent Stopped

Figure 10-121 Creating a Remote Copy consistency group

2. Enter a name for the consistency group and click **Next**, as shown in Figure 10-122.

Create Consistency Group

Enter a name for the consistency group:

Consistency Group Name:

Cancel      < Back      Next >

Figure 10-122 Entering a name for the new consistency group

3. In the next window, select the location of the auxiliary volumes in the consistency group, as shown in Figure 10-123, and click **Next**:
- **On this system**, which means that the volumes are local.
  - **On another system**, which means that you select the remote system from the menu.

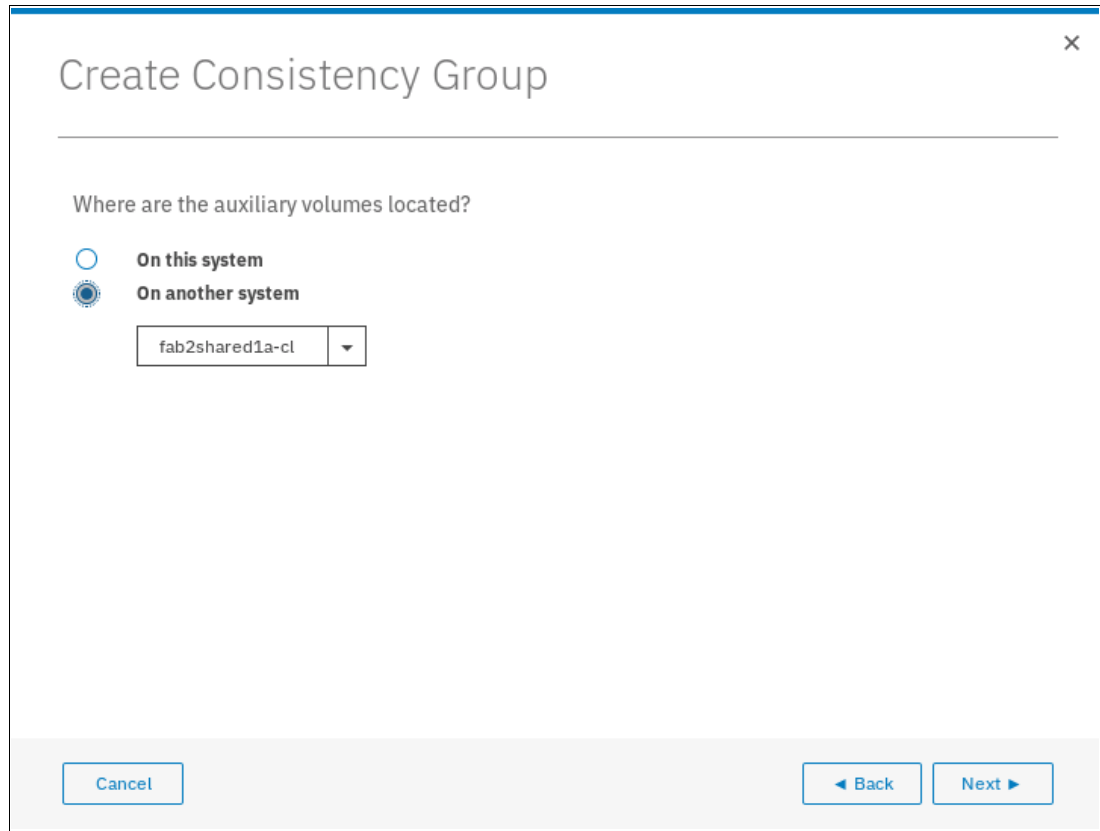


Figure 10-123 Selecting the system with which to create the consistency group

4. Select whether you want to add relationships to this group, as shown in Figure 10-124. The following options are available:
  - If you select **No**, click **Finish** to create an empty consistency group that can be used later.
  - If you select **Yes**, click **Next** to continue the wizard and continue with the next steps.

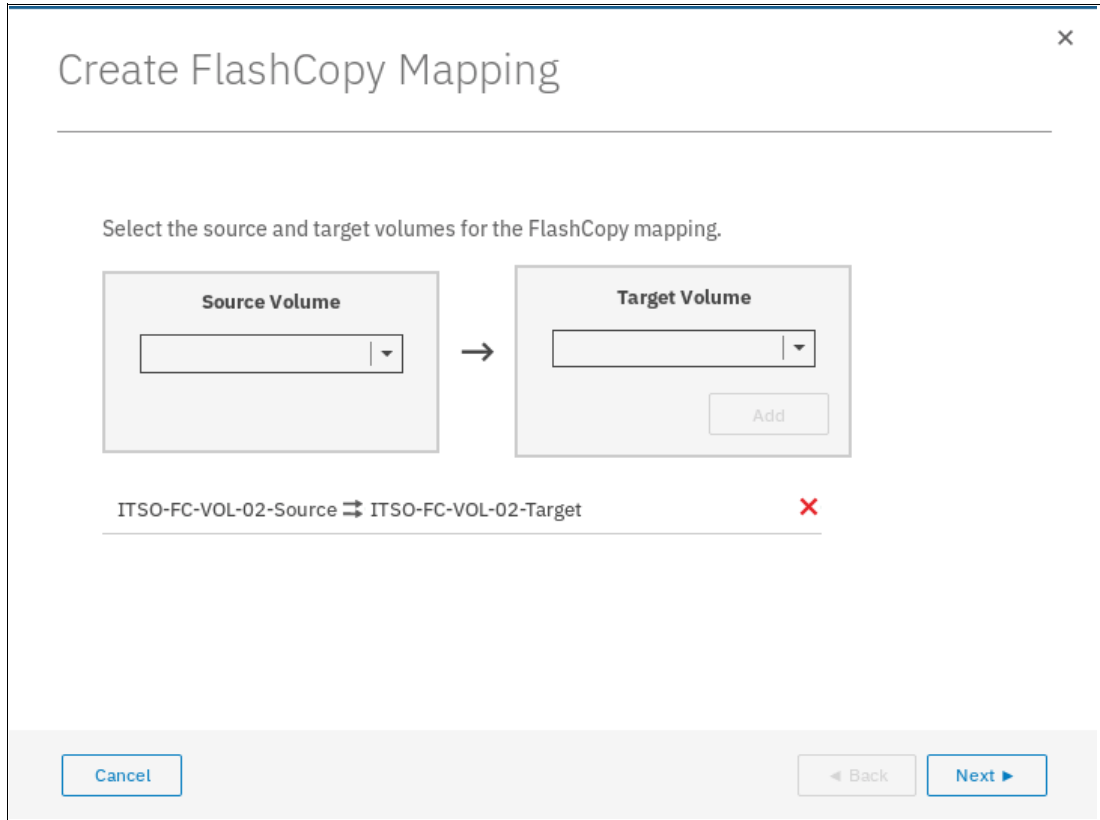


Figure 10-124 Selecting whether relationships should be added to the new consistency group



Select one of the following types of relationships that you want to create or add, as shown in Figure 10-125, and click **Next**:

- MM
- GM (with or without Consistency Protection)
- Global Mirror with Change Volumes

Create Relationship

Select the type of relationships that you want to create.

Metro Mirror

Global Mirror

Add Consistency Protection

Global Mirror with Change Volumes

Cancel      < Back      Next >

Figure 10-125 Selecting the type of remote copy relationships to create or add

5. As shown in Figure 10-126, you can optionally select existing relationships to the group. Click **Next**.

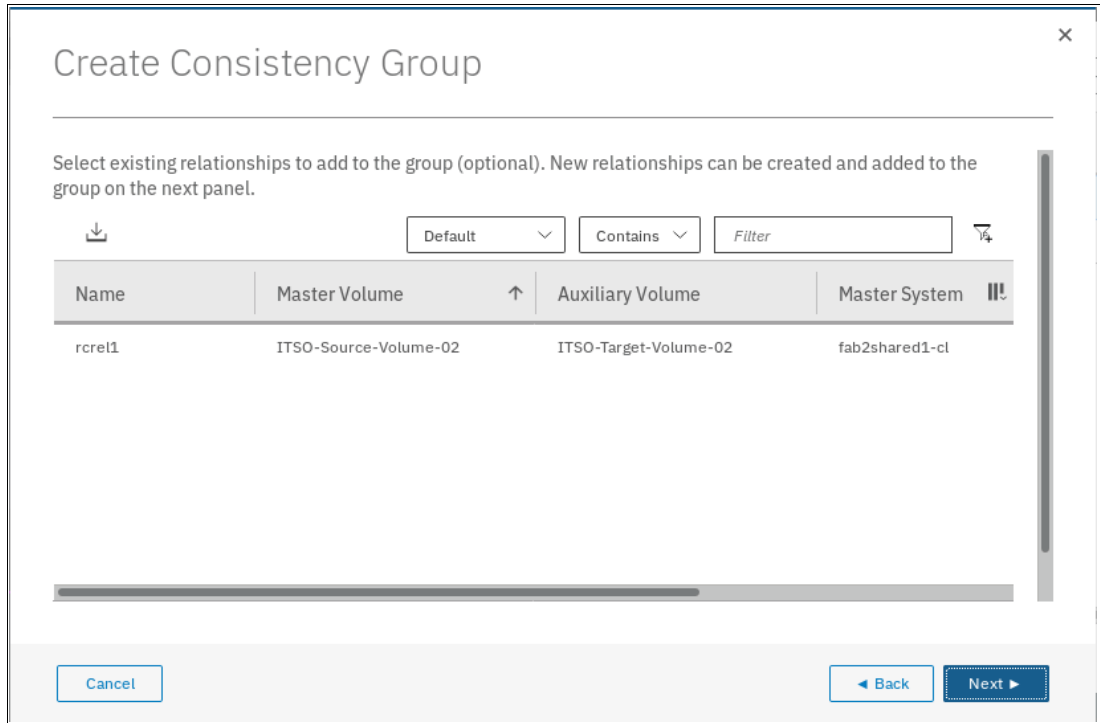


Figure 10-126 Adding Remote Copy relationships to the new consistency group

**Note:** Only relationships of the type that were selected are listed.

6. In the next window, you can create relationships between master volumes and auxiliary volumes to be added to the consistency group that is being created, as shown in Figure 10-127. Click **Add** when both volumes are selected. You can add multiple relationships in this step by repeating the selection.

When all the relationships you need are created, click **Next**.

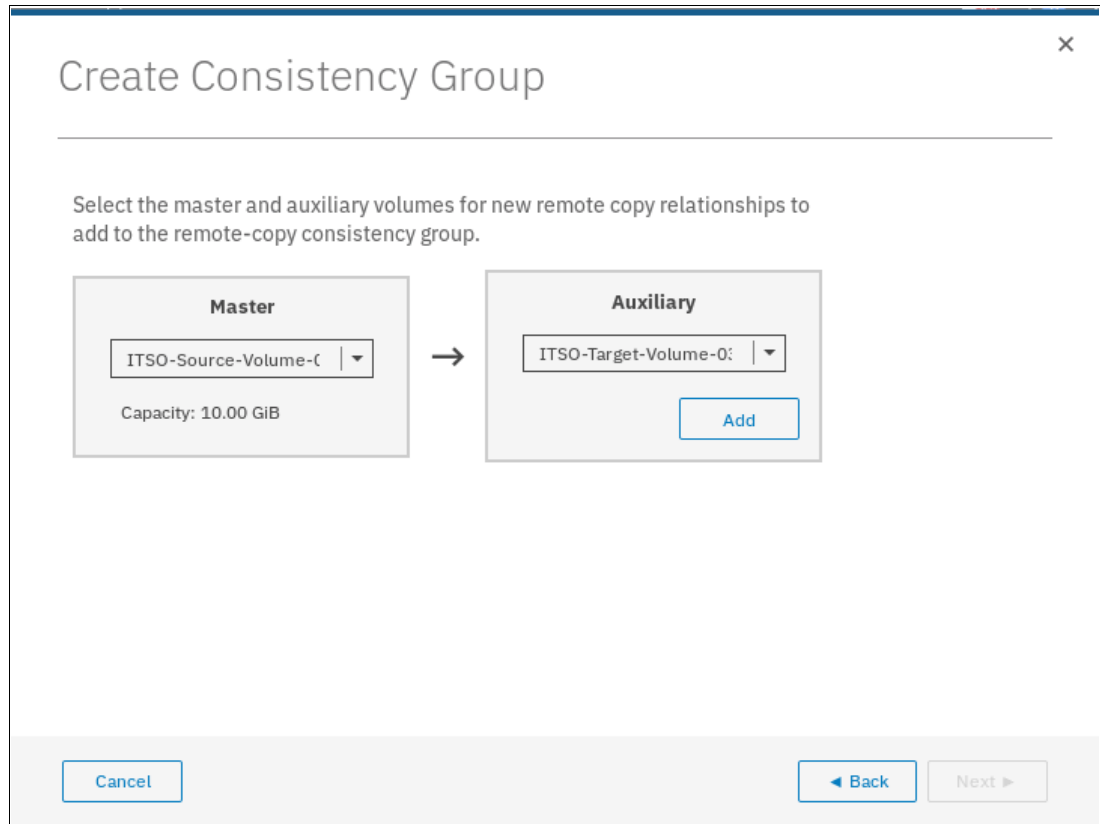
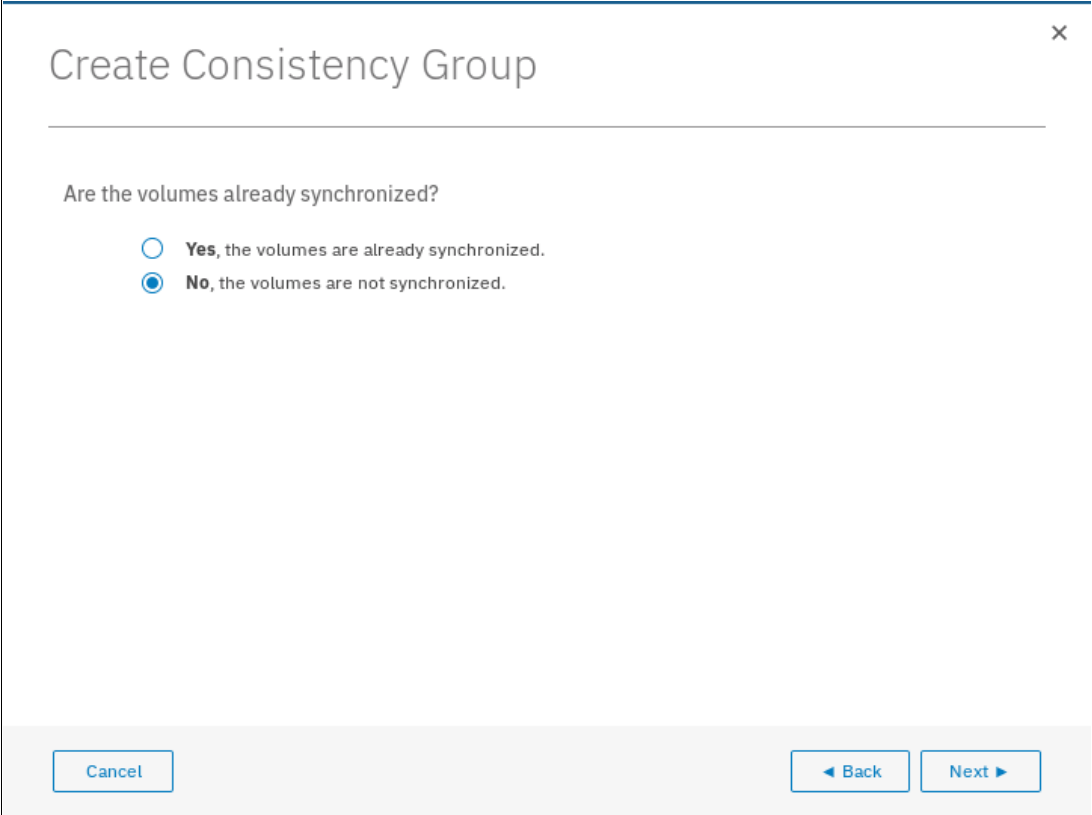


Figure 10-127 Creating new relationships for the new consistency group

**Important:** The master and auxiliary volumes must be of equal size. Therefore, only the targets with the suitable size are shown in the list for a specific source volume.

7. Specify whether the volumes in the consistency group are synchronized, as shown in Figure 10-128. Click **Next**.



Create Consistency Group ×

---

Are the volumes already synchronized?

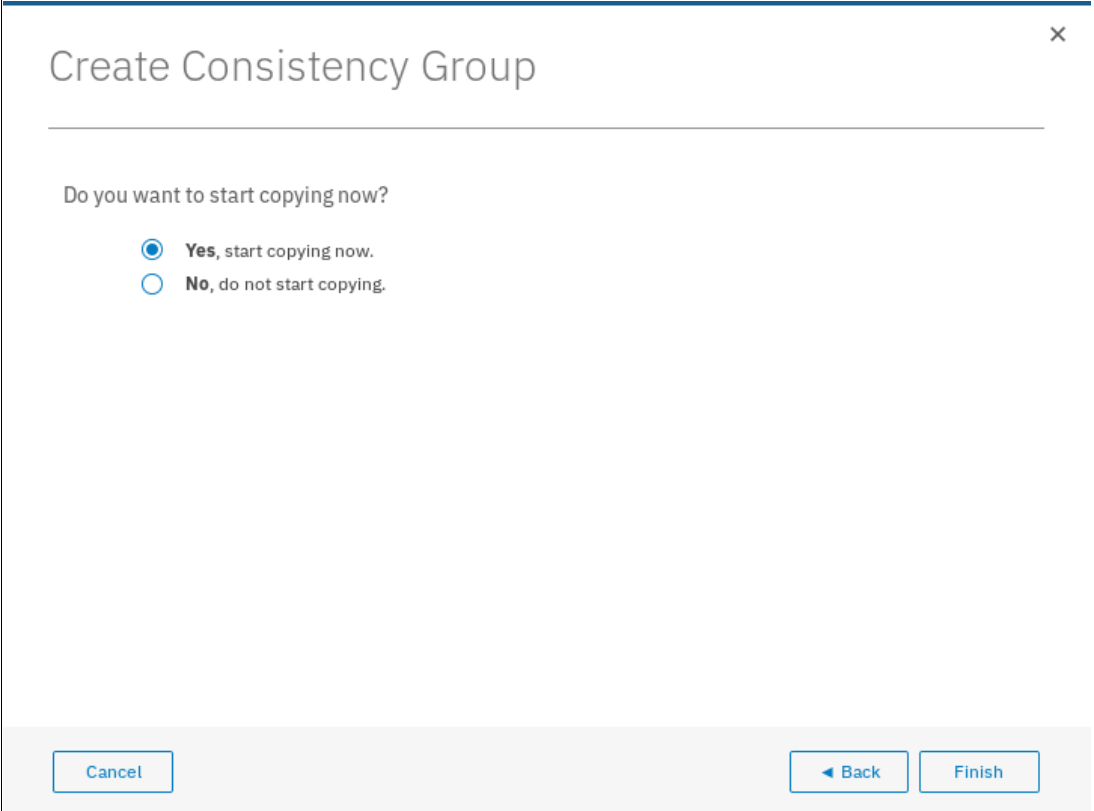
Yes, the volumes are already synchronized.

No, the volumes are not synchronized.

Figure 10-128 Selecting if volumes in the new consistency group are already synchronized or not

**Note:** If the volumes are not synchronized, the initial copy copies the entire source volume to the remote target volume. If you suspect volumes are different or if you have a doubt, synchronize them to ensure consistency on both sides of the relationship.

8. In the last window, select whether you want to start the copy of the consistency group, as shown in Figure 10-129. Click **Finish**.



Create Consistency Group ×

---

Do you want to start copying now?

**Yes**, start copying now.

**No**, do not start copying.

Figure 10-129 Selecting whether copy should start

## 10.9.4 Renaming Remote Copy relationships

To rename one or multiple RC relationships, complete the following steps:

1. Open the **Copy Services** → **Remote Copy** panel.
2. Right-click the relationships to be renamed and select **Rename**, as shown in Figure 10-130.

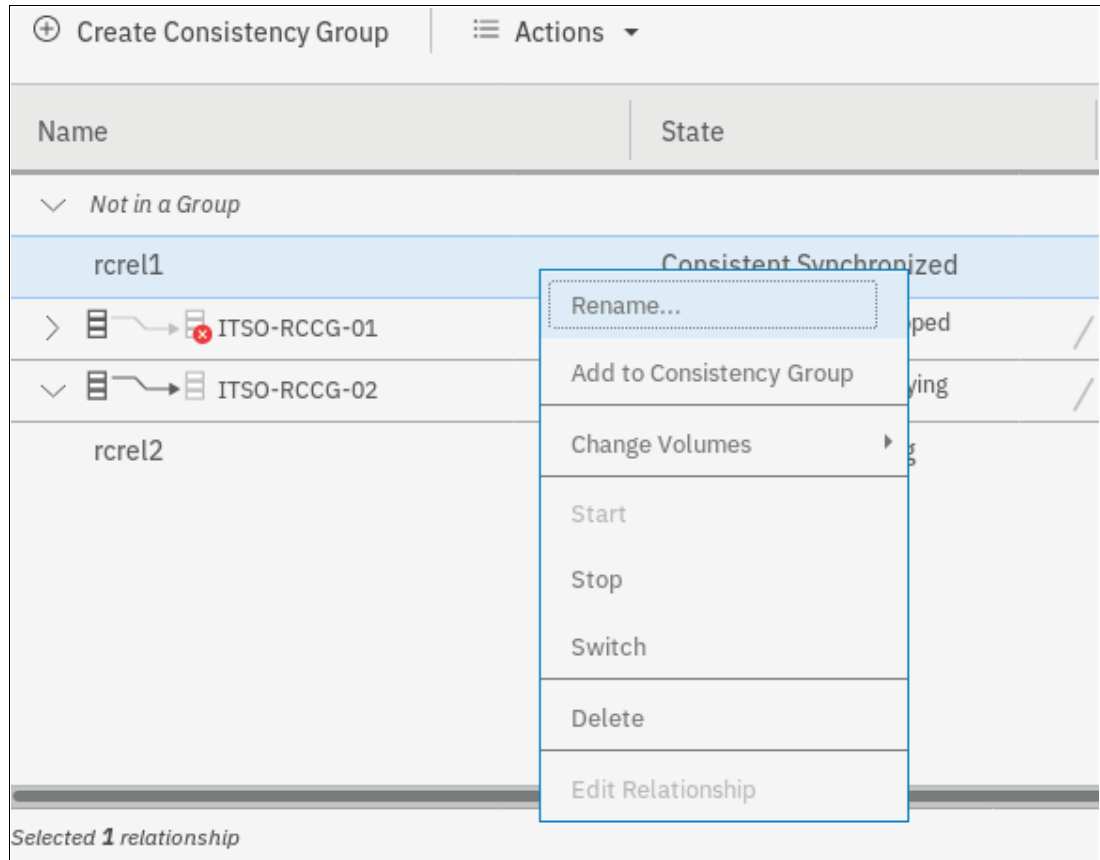


Figure 10-130 Renaming Remote Copy relationships

3. Enter the new name that you want to assign to the relationships and click **Rename**, as shown in Figure 10-131.

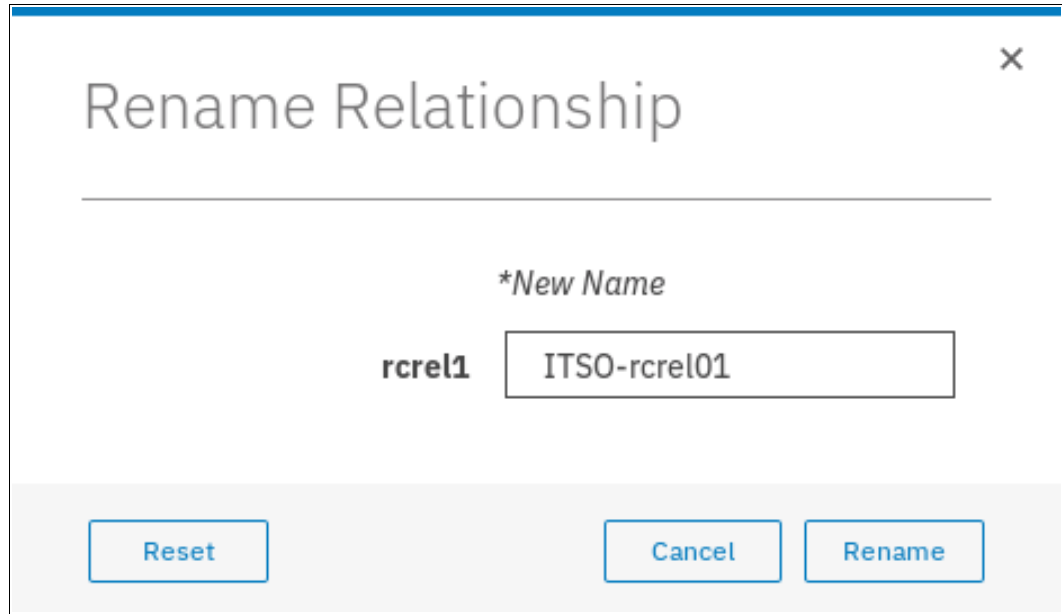


Figure 10-131 Renaming Remote Copy relationships

**RC relationship name:** You can use the letters A - Z and a - z, the numbers 0 - 9, and the underscore ( \_ ) character. The RC name can be 1 - 15 characters. Blanks cannot be used.

### 10.9.5 Renaming a Remote Copy consistency group

To rename a RC consistency group, complete the following steps:

1. Open the **Copy Services** → **Remote Copy** panel.
2. Right-click the consistency group to be renamed and select **Rename**, as shown in Figure 10-132 on page 654.

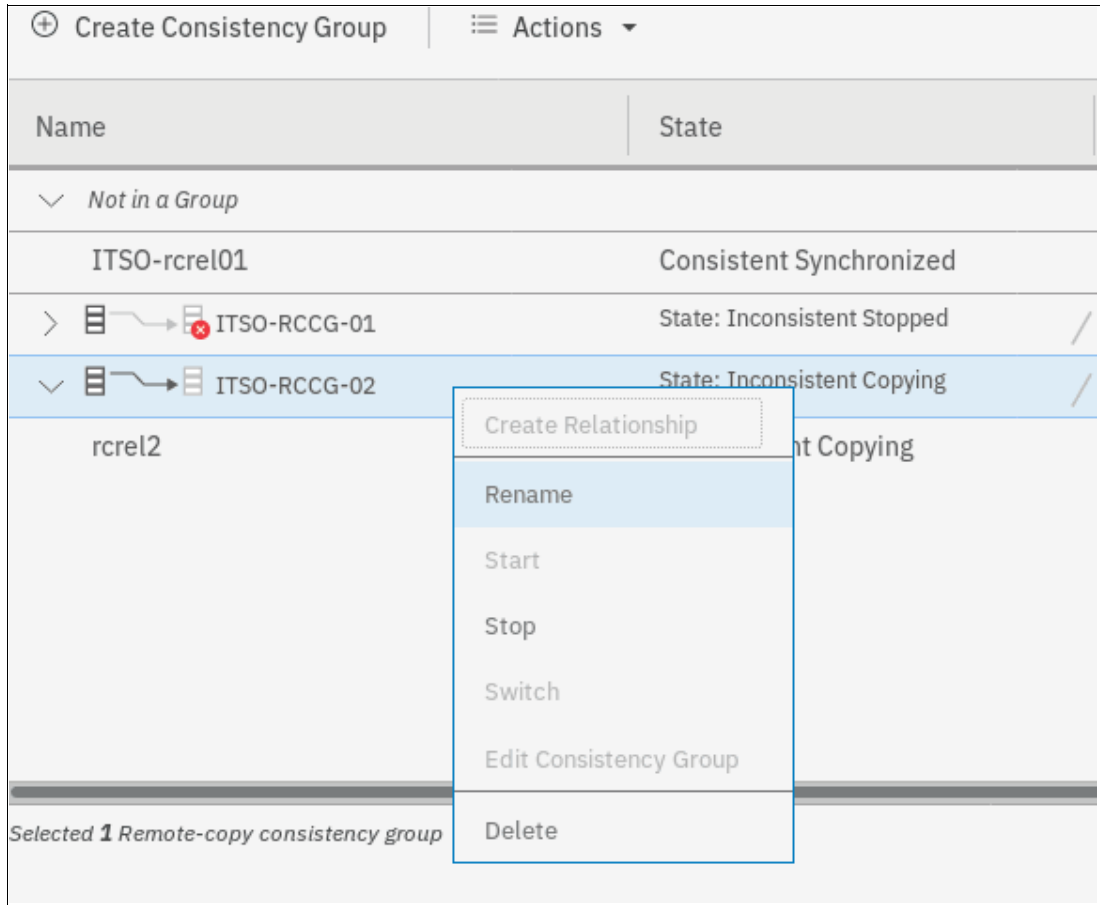


Figure 10-132 Renaming a Remote Copy consistency group

3. Enter the new name that you want to assign to the consistency group and click **Rename**, as shown in Figure 10-133.

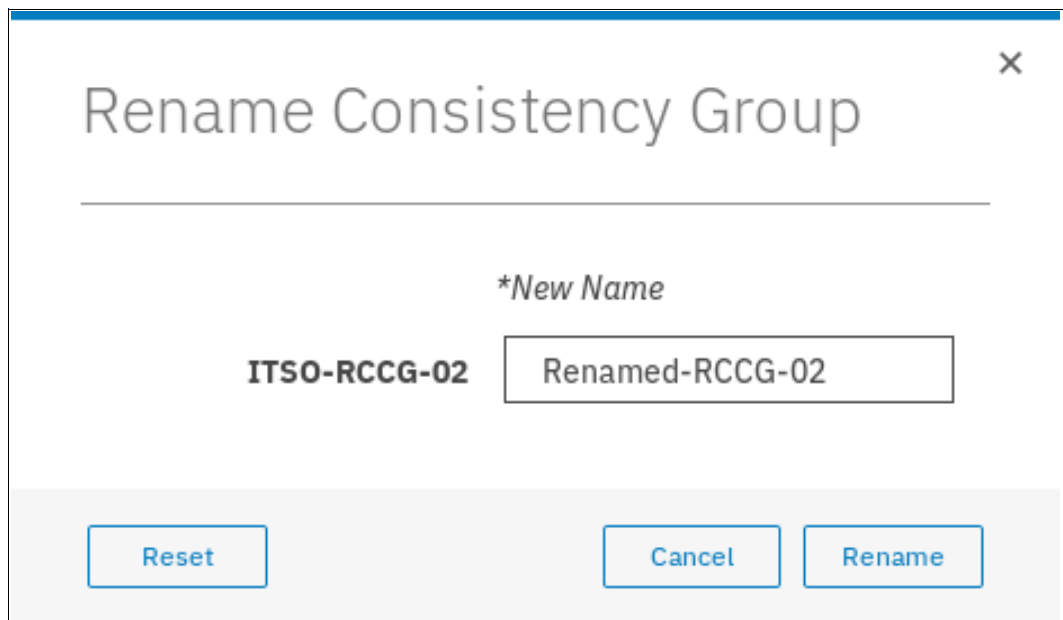


Figure 10-133 Entering new name for consistency group



**RC consistency group name:** You can use the letters A - Z and a - z, the numbers 0 - 9, and the underscore ( \_ ) character. The RC name can be 1 - 15 characters. Blanks cannot be used.

## 10.9.6 Moving stand-alone Remote Copy relationships to consistency group

To add one or multiple stand-alone relationships to a RC consistency group, complete the following steps:

1. Open the **Copy Services** → **Remote Copy** panel.
2. Right-click the relationships to be moved and select **Add to Consistency Group**, as shown in Figure 10-134.

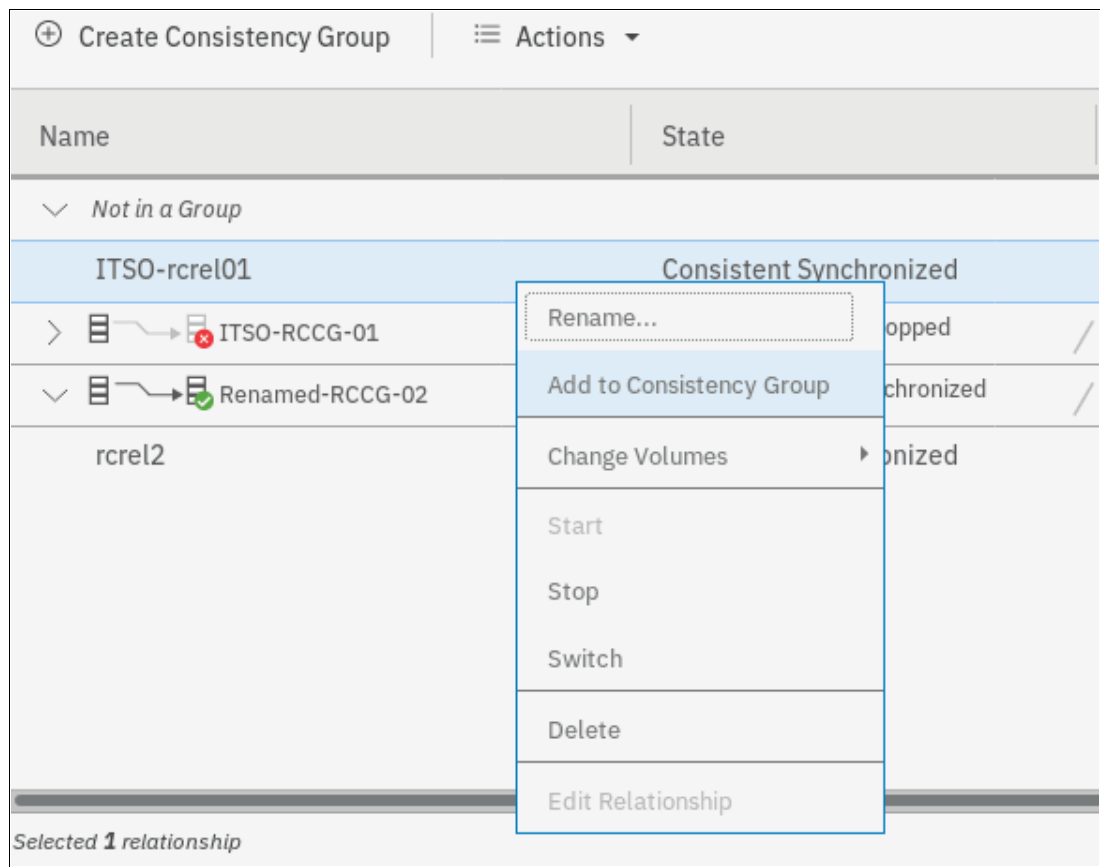


Figure 10-134 Adding relationships to a consistency group

3. Select the consistency group for this RC relationship by using the menu, as shown in Figure 10-135. Click **Add to Consistency Group** to confirm your changes.

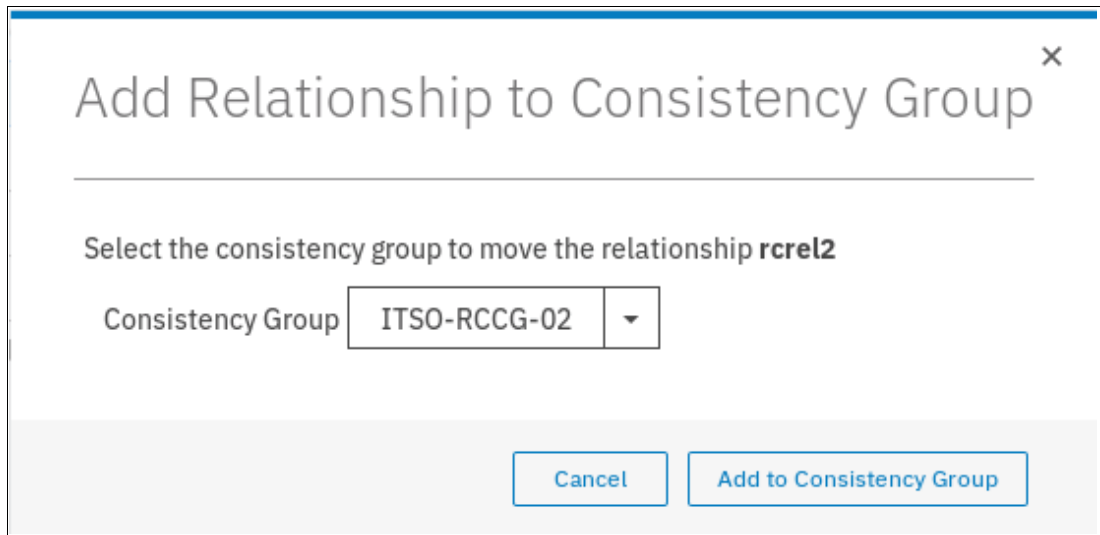


Figure 10-135 Selecting the consistency group to add the relationships to

## 10.9.7 Removing Remote Copy relationships from consistency group

To remove one or multiple relationships from a RC consistency group, complete the following steps:

1. Open the **Copy Services** → **Remote Copy** panel.
2. Right-click the relationships to be removed and select **Remove from Consistency Group**, as shown in Figure 10-136.

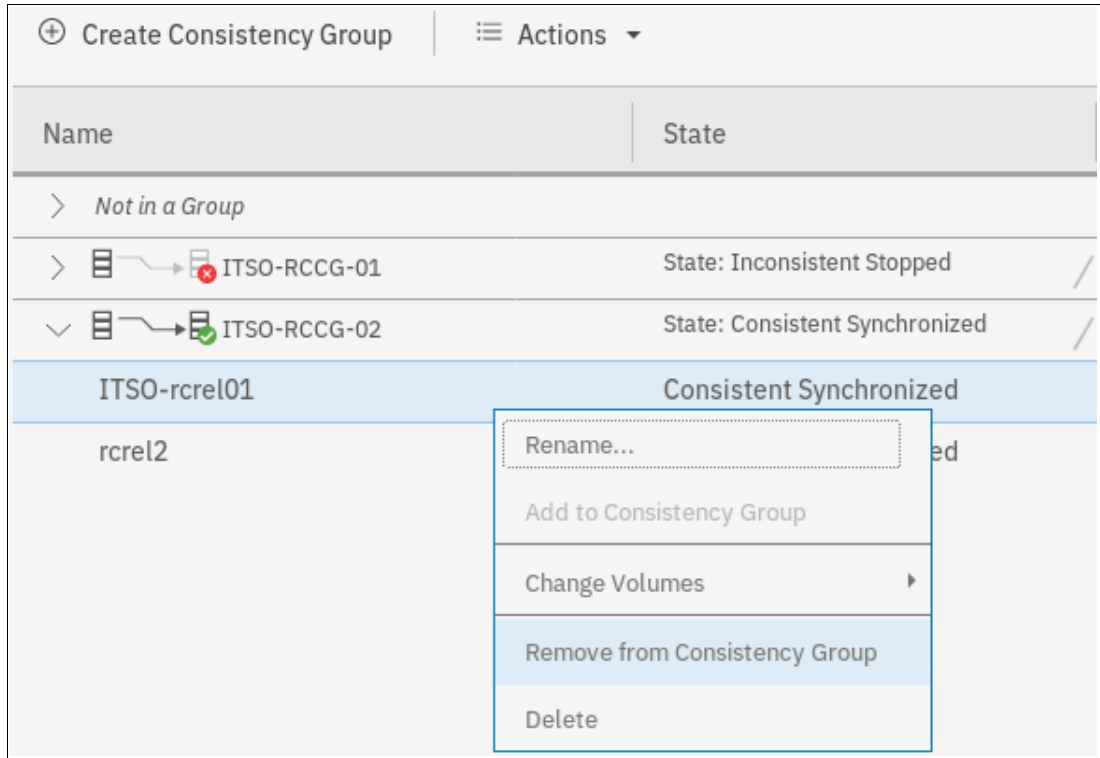


Figure 10-136 Removing relationships from a consistency group

3. Confirm your selection and click **Remove**, as shown in Figure 10-137.

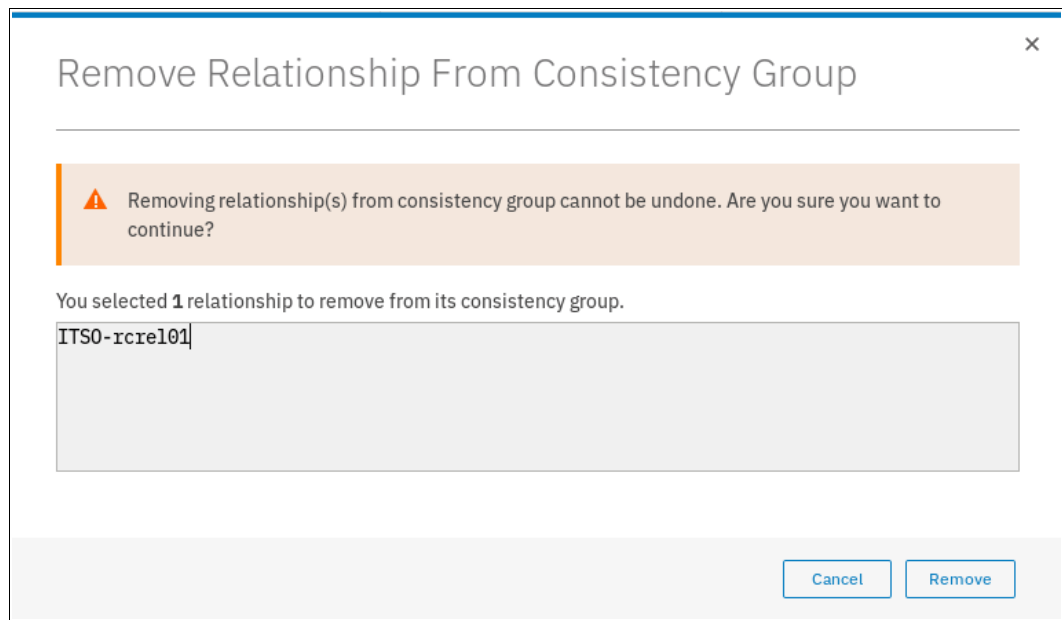


Figure 10-137 Confirm the removal of relationships from a consistency group

## 10.9.8 Starting Remote Copy relationships

When a RC relationship is created, the RC process can be started. Only relationships that are not members of a consistency group, or the only relationship in a consistency group, can be started. In any other case, consider starting the consistency group instead.

To start one or multiple relationships, complete the following steps:

1. Open the **Copy Services** → **Remote Copy** panel.
2. Right-click the relationships to be started and select **Start**, as shown in Figure 10-138.

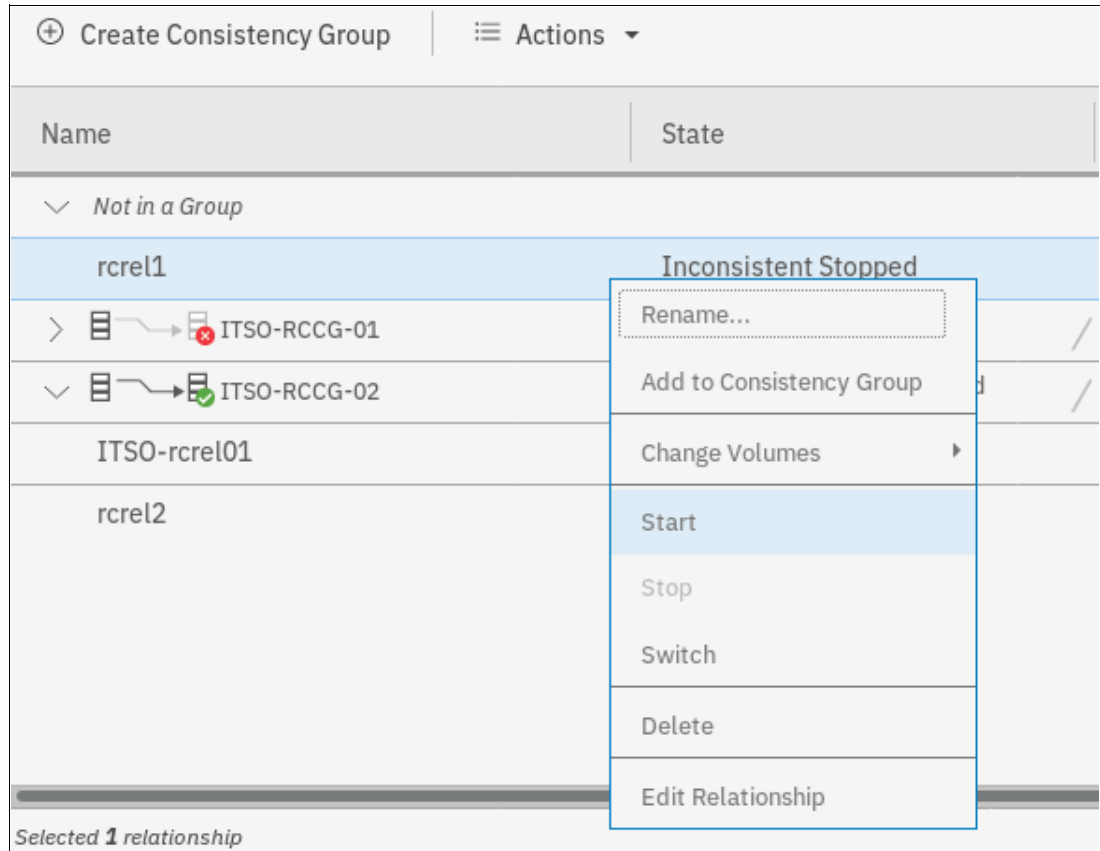


Figure 10-138 Starting Remote Copy relationships

## 10.9.9 Starting a Remote Copy consistency group

When a RC consistency group is created, the RC process can be started for all the relationships that are part of the consistency groups.

To start a consistency group, open the **Copy Services** → **Remote Copy** panel, right-click the consistency group to be started, and select **Start**, as shown in Figure 10-139.

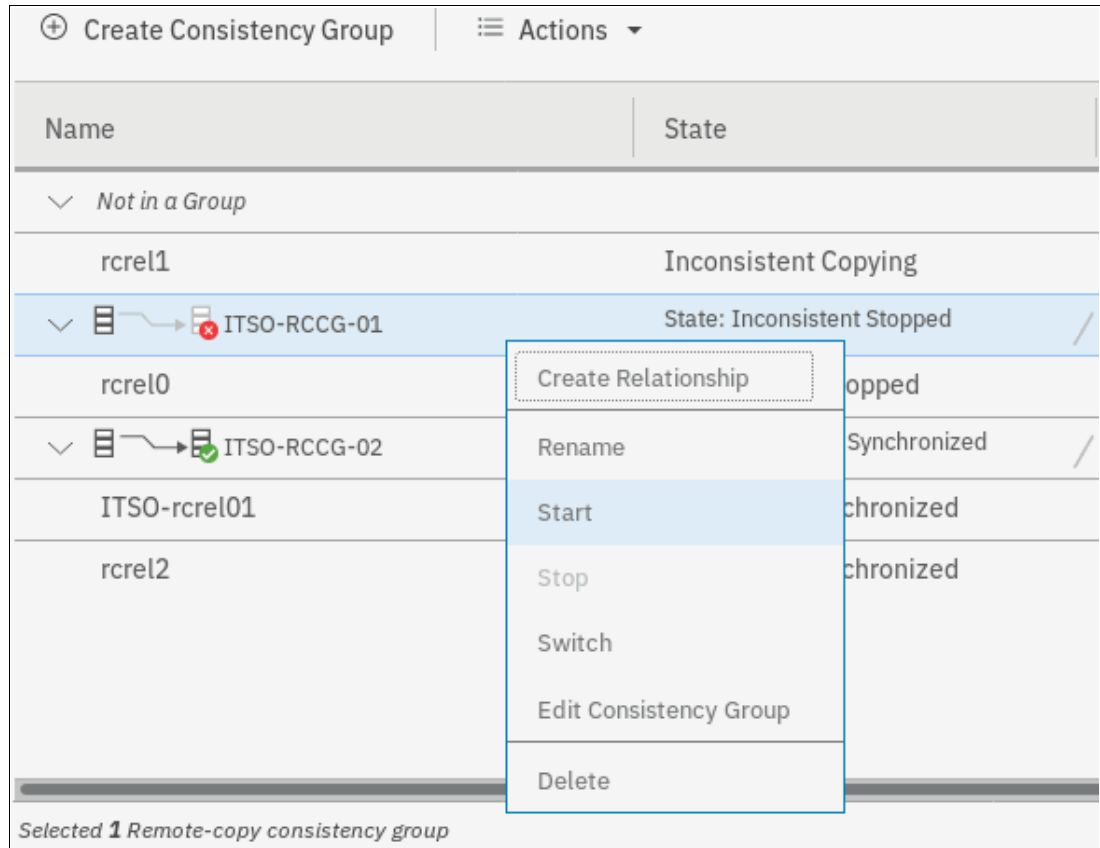


Figure 10-139 Starting a Remote Copy consistency group

## 10.9.10 Switching a relationship copy direction

When a RC relationship is in the Consistent synchronized state, the copy direction for the relationship can be changed. Only relationships that are not a member of a consistency group, or the only relationship in a consistency group, can be switched. In any other case, consider switching the consistency group instead.

**Important:** When the copy direction is switched, it is crucial that no outstanding I/O exists to the volume that changes from primary to secondary because all of the I/O is inhibited to that volume when it becomes the secondary. Therefore, careful planning is required before you switch the copy direction for a relationship.

To switch the direction of a RC relationship, complete the following steps:

1. Open the **Copy Services** → **Remote Copy** panel.
2. Right-click the relationship to be switched and select **Switch**, as shown in Figure 10-140.

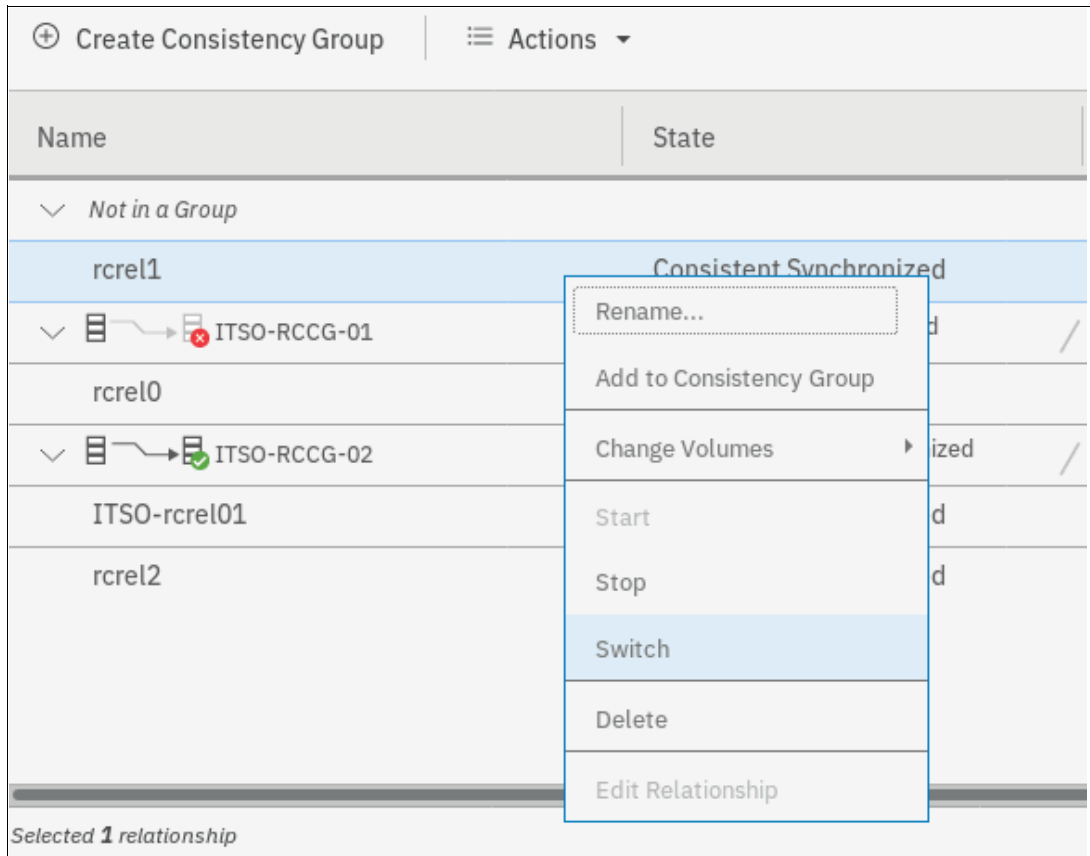


Figure 10-140 Switching Remote Copy relationship direction

3. Because the master-auxiliary relationship direction is reversed, write access is disabled on the new auxiliary volume (former master volume), whereas it is enabled on the new master volume (former auxiliary volume). A warning message is displayed, as shown in Figure 10-141. Click **Yes**.

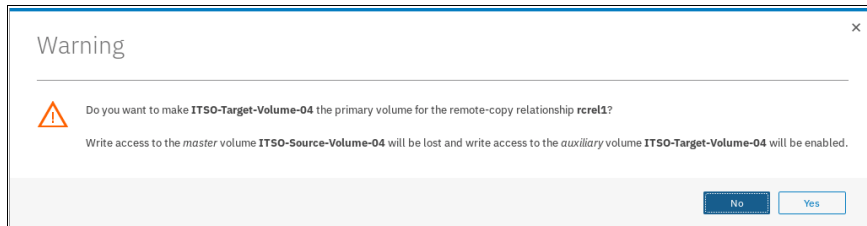


Figure 10-141 Switching master-auxiliary direction of relationships changes the write access

When a RC relationship is switched, an icon is displayed in the Remote Copy panel list, as shown in Figure 10-142.


+ Create Consistency Group		☰ Actions ▾
Name	State	
▾ Not in a Group		
rcrel1	Consistent Synchronized 	

Figure 10-142 Switched Remote Copy Relationship

### 10.9.11 Switching a consistency group direction

When a RC consistency group is in the consistent synchronized state, the copy direction for the consistency group can be changed.

**Important:** When the copy direction is switched, it is crucial that no outstanding I/O exists to the volume that changes from primary to secondary because all of the I/O is inhibited to that volume when it becomes the secondary. Therefore, careful planning is required before you switch the copy direction for a relationship.

To switch the direction of a RC consistency group, complete the following steps:

1. Open the **Copy Services** → **Remote Copy** panel.



- Right-click the consistency group to be switched and select **Switch**, as shown in Figure 10-143.

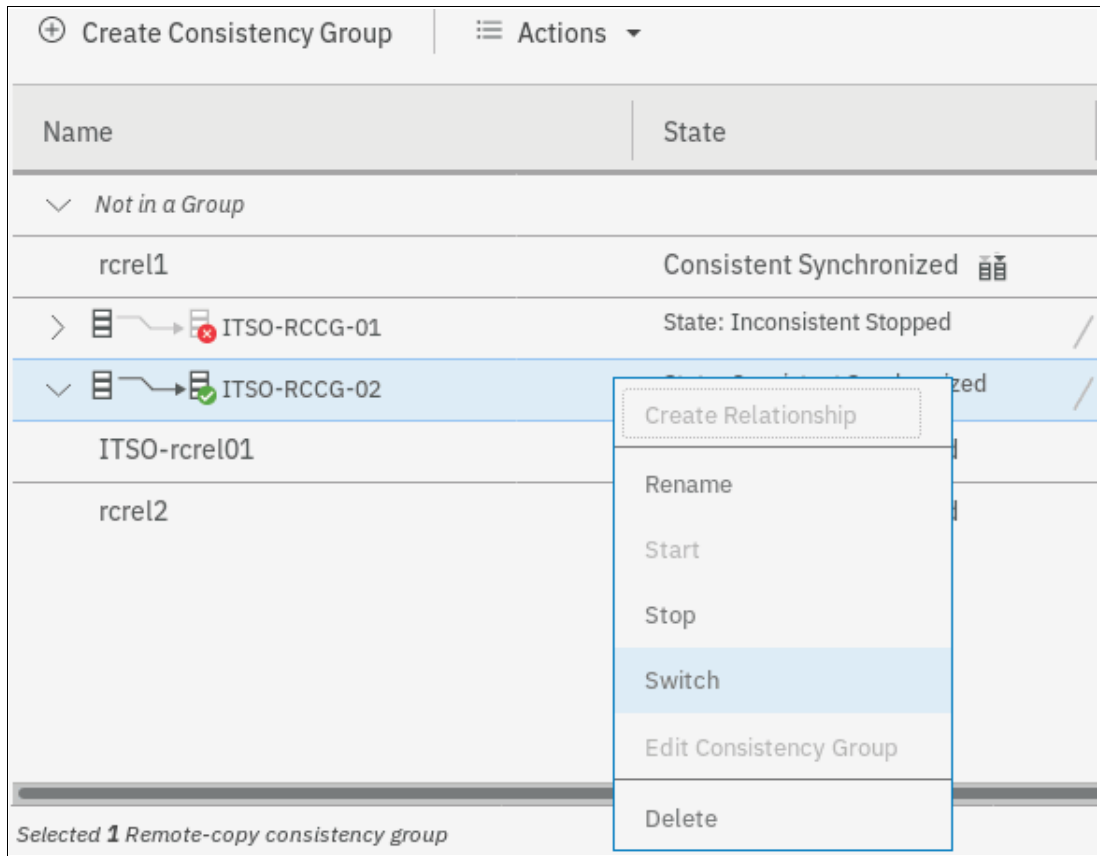


Figure 10-143 Switching a Consistency Group direction

- Because the master-auxiliary relationship direction is reversed, write access is disabled on the new auxiliary volume (former master volume), while it is enabled on the new master volume (former auxiliary volume). A warning message is displayed, as shown in Figure 10-144. Click **Yes**.

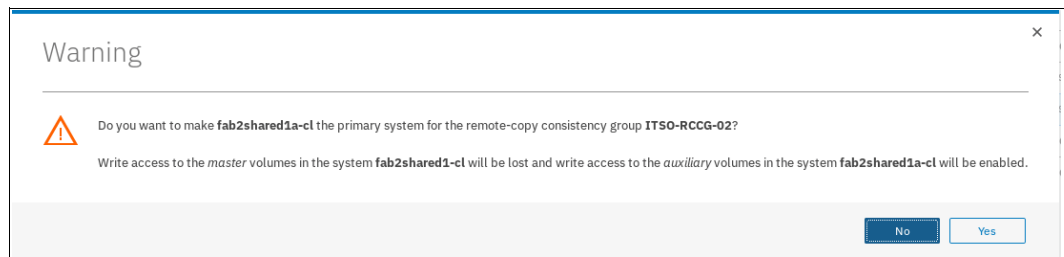


Figure 10-144 Switching direction of Consistency Groups changes the write access

## 10.9.12 Stopping Remote Copy relationships

When a RC relationship is created and started, the RC process can be stopped. Only relationships that are not members of a consistency group, or the only relationship in a consistency group, can be stopped. In any other case, consider stopping the consistency group instead.

To stop one or multiple relationships, complete the following steps:

1. Open the **Copy Services** → **Remote Copy** panel.
2. Right-click the relationships to be stopped and select **Stop**, as shown in Figure 10-145.

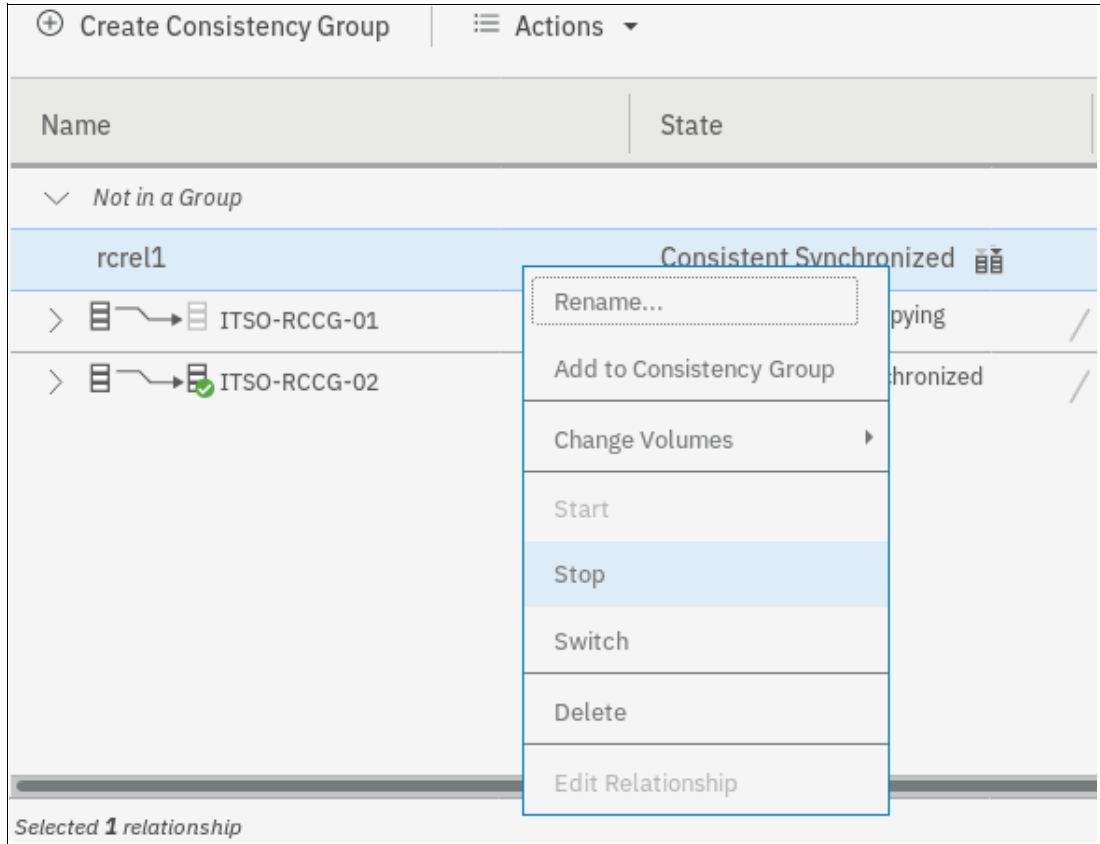


Figure 10-145 Stopping a Remote Copy relationship

3. When a RC relationship is stopped, access to the auxiliary volume can be changed so it can be read and written by a host. A confirmation message is displayed, as shown in Figure 10-146.

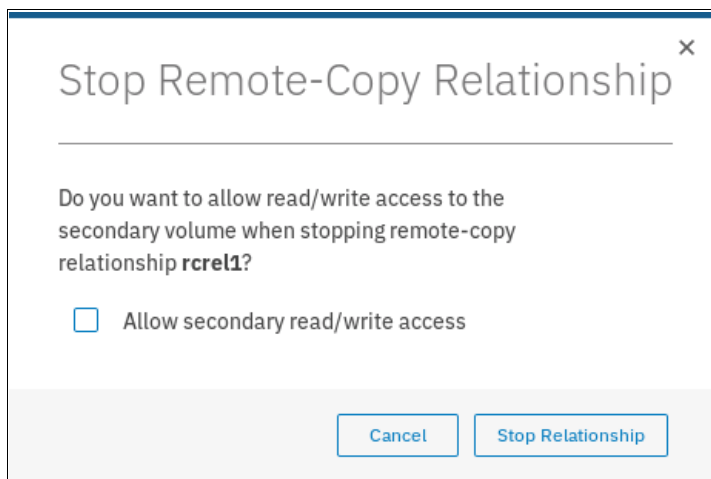


Figure 10-146 Grant access in read and write to the auxiliary volume

### 10.9.13 Stopping a consistency group

When a RC consistency group is created and started, the RC process can be stopped.

To stop a consistency group, complete the following steps:

1. Open the **Copy Services** → **Remote Copy** panel.
2. Right-click the consistency group to be stopped and select **Stop**, as shown in Figure 10-147.

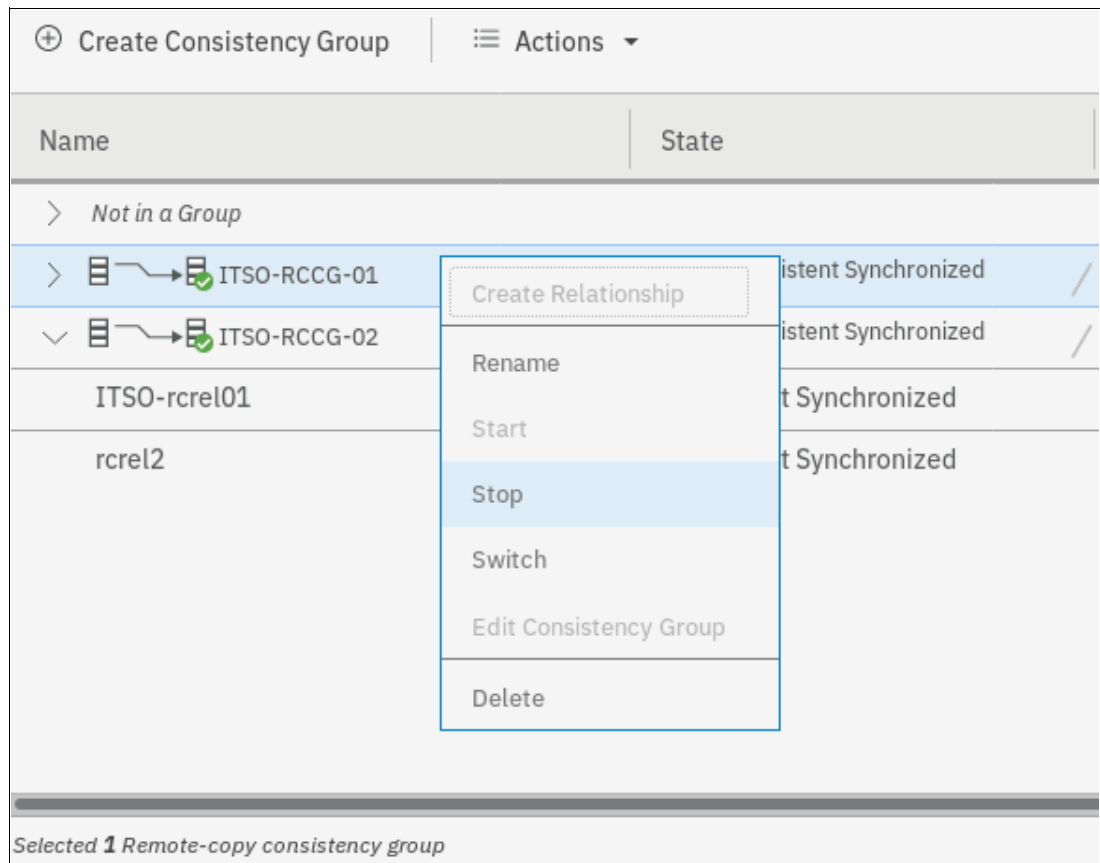


Figure 10-147 Stopping a Consistency Group

3. When a RC consistency group is stopped, access to the auxiliary volumes can be changed so it can be read and written by a host. A confirmation message is displayed, as shown in Figure 10-148 on page 666.

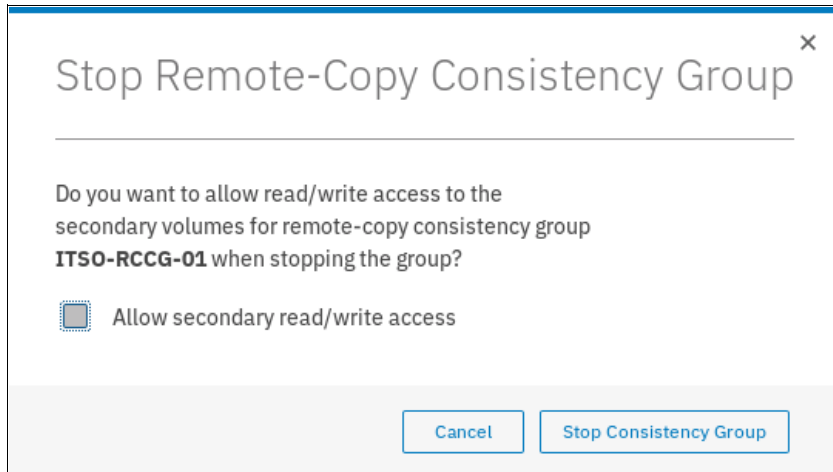


Figure 10-148 Grant access in read and write to the auxiliary volumes

### 10.9.14 Deleting Remote Copy relationships

To delete RC relationships, complete the following steps:

1. Open the **Copy Services** → **Remote Copy** panel.
2. Right-click the relationships that you want to delete and select **Delete**, as shown in Figure 10-149.

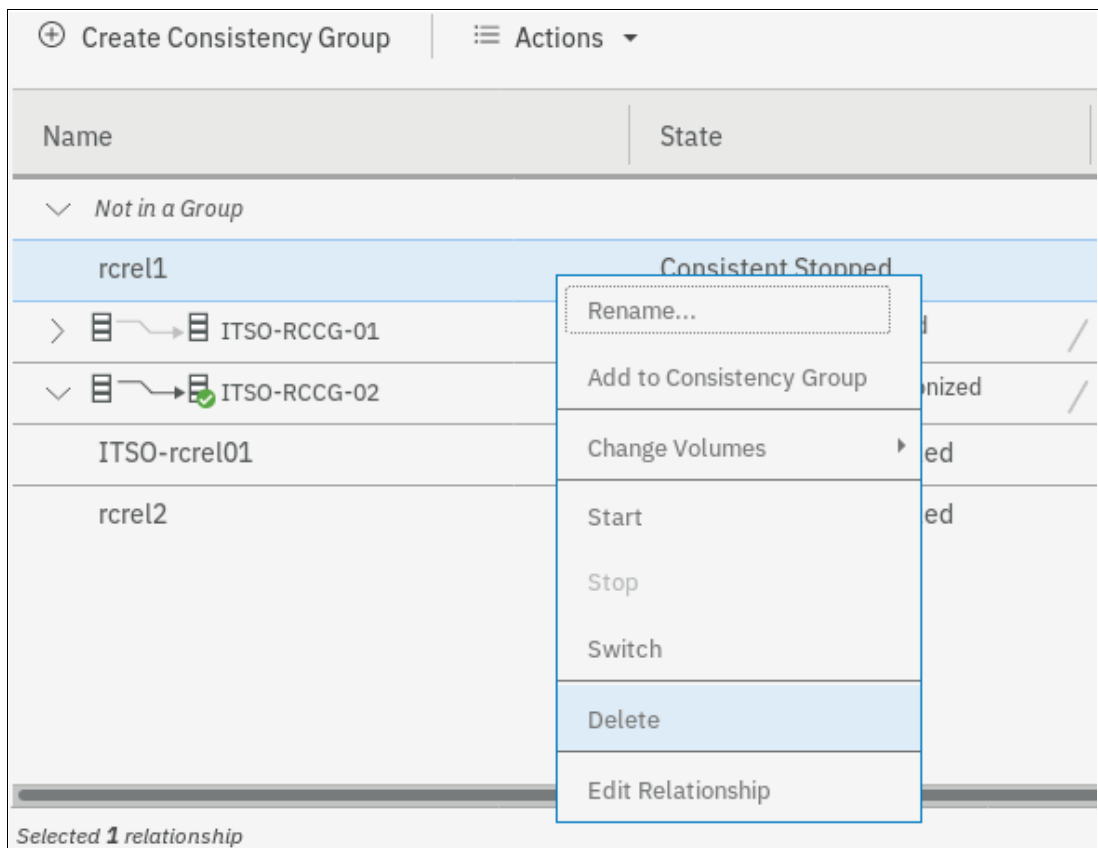


Figure 10-149 Deleting Remote Copy Relationships

3. A confirmation message is displayed that requests that the user enter the number of relationships to be deleted, as shown in Figure 10-150.

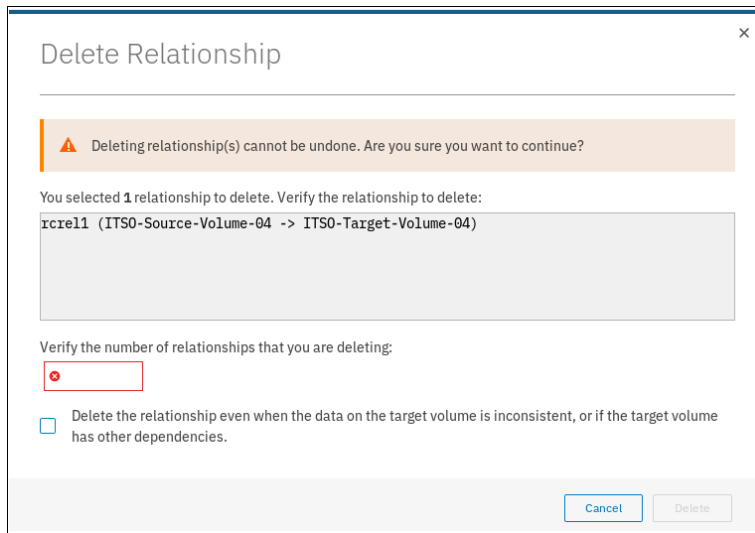


Figure 10-150 Confirmation of relationships deletion

### 10.9.15 Deleting a consistency group

To delete a RC consistency group, complete these steps:

1. Open the **Copy Services** → **Remote Copy** panel.
2. Right-click the consistency group that you want to delete and select **Delete**, as shown in Figure 10-151 on page 668.

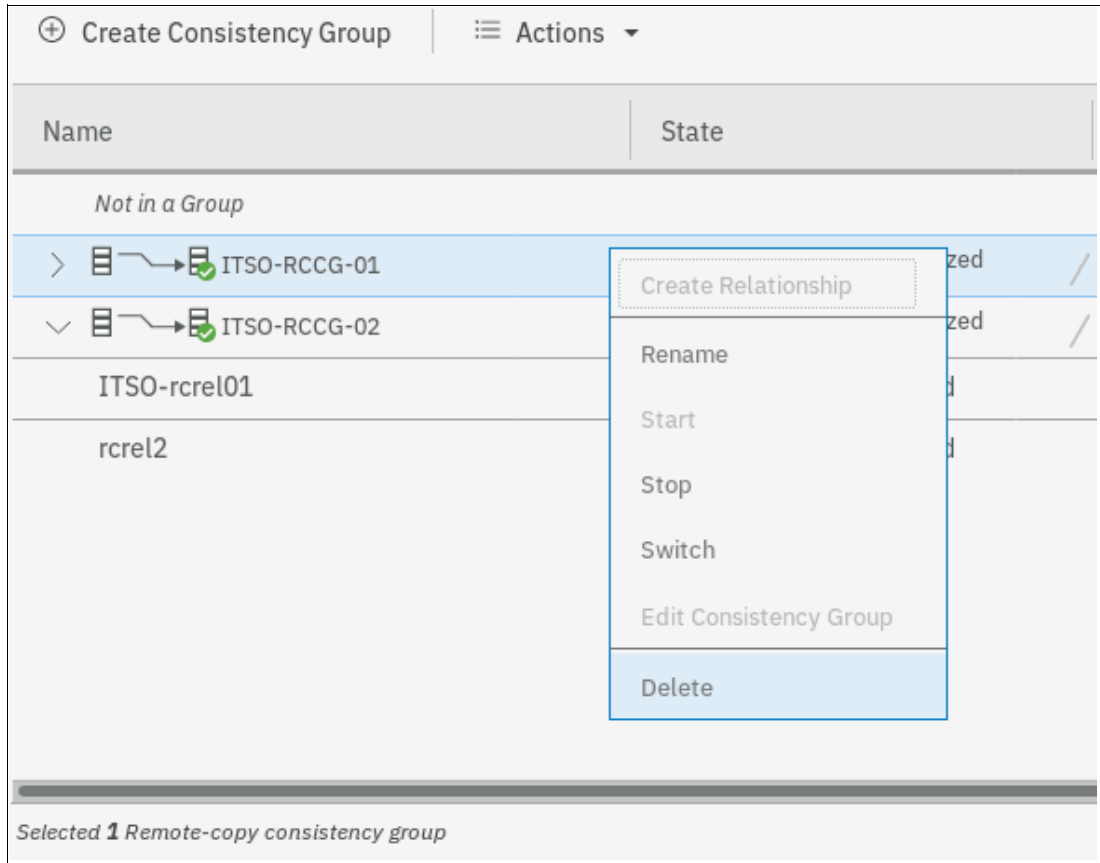


Figure 10-151 Deleting a Consistency Group

3. A confirmation message is displayed, as shown in Figure 10-152. Click **Yes**.

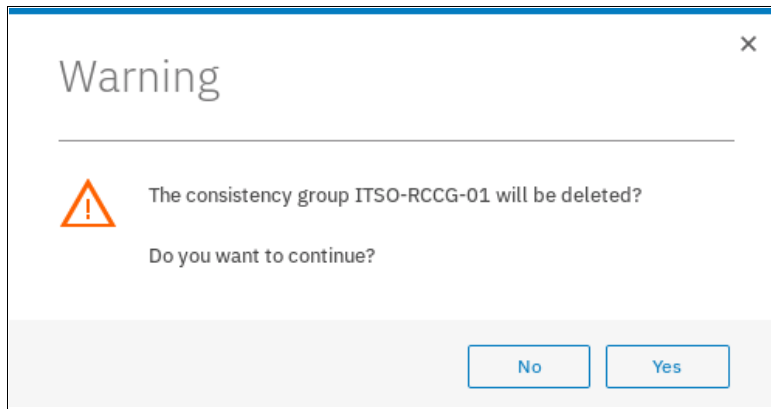


Figure 10-152 Confirmation of a Consistency Group deletion

**Important:** Deleting a consistency group does *not* delete its RC mappings.

## 10.10 Remote Copy memory allocation

Copy Services features require that small amounts of volume cache be converted from cache memory into bitmap memory to allow the functions to operate at an I/O group level. If you do not have enough bitmap space allocated when you try to use one of the functions, the configuration cannot be completed.

The total memory that can be dedicated to these functions is not defined by the physical memory in the system. The memory is constrained by the software functions that use the memory.

For every RC relationship that is created on an IBM Spectrum Virtualize system, a bitmap table is created to track the copied grains. By default, the system allocates 20 MiB of memory for a minimum of 2 TiB of remote copied source volume capacity. Every 1 MiB of memory provides the following volume capacity for the specified I/O group: for 256 KiB grains size, 2 TiB of total MM, GM, or active-active volume capacity.

Review Table 10-15 to calculate the memory requirements and confirm that your system can accommodate the total installation size.

Table 10-15 Memory allocation for FlashCopy services

Minimum allocated bitmap space	Default allocated bitmap space	Maximum allocated bitmap space	Minimum functionality when using the default values <sup>1</sup>
0	20 MiB	512 MiB	40 TiB of remote mirroring volume capacity
<sup>1</sup> RC includes MM, GM, and active-active relationships.			

When you configure change volumes for use with GM, two internal FlashCopy mappings are created for each change volume.

Two bitmaps exist for MM, GM, and HyperSwap active-active relationships. For MM/GM relationships, one is used for the master clustered system and one is used for the auxiliary system because the direction of the relationship can be reversed. For active-active relationships, which are configured automatically when HyperSwap volumes are created, one bitmap is used for the volume copy on each site because the direction of these relationships can be reversed.

MM/GM relationships do not automatically increase the available bitmap space. You might need to run the `chiogrp` command to manually increase the space in one or both of the master and auxiliary systems.

You can modify the resource allocation for each I/O group of an IBM SAN Volume Controller system by selecting **Settings** → **System** and clicking the **Resources** menu, as shown in Figure 10-153.

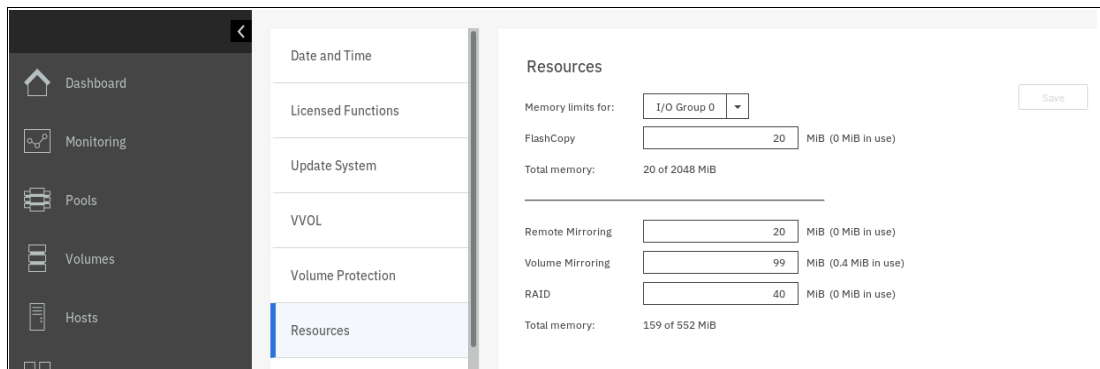


Figure 10-153 Modifying resources allocation

## 10.11 Troubleshooting Remote Copy

RC (MM and GM) features two primary error codes that are displayed:

- ▶ A 1920 error can be considered as a voluntary stop of a relationship by the system when it evaluates the replication causes errors on the hosts. A 1920 is a congestion error. This error means that the source, the link between the source and target, or the target cannot keep up with the requested copy rate. The system then triggers a 1920 error to prevent replication from having undesired effects on hosts.
- ▶ A 1720 error is a heartbeat or system partnership communication error. This error often is more serious because failing communication between your system partners involves extended diagnostic time.

### 10.11.1 1920 error

A 1920 error is deliberately generated by the system and is considered as a control mechanism. It occurs after 985003 (“Unable to find path to disk in the remote cluster (system) within the time-out period”) or 985004 (“Maximum replication delay has been exceeded”) events.

It can have several triggers, including the following probable causes:

- ▶ Primary system or SAN fabric problem (10%)
- ▶ Primary system or SAN fabric configuration (10%)
- ▶ Secondary system or SAN fabric problem (15%)
- ▶ Secondary system or SAN fabric configuration (25%)
- ▶ Intercluster link problem (15%)
- ▶ Intercluster link configuration (25%)

In practice, the most often overlooked cause is latency. GM has a RTT tolerance limit of 80 or 250 milliseconds, depending on the firmware version and the hardware model. A message that is sent from the source IBM Spectrum Virtualize system to the target system and the accompanying acknowledgment must have a total time of 80 or 250 milliseconds round trip. That is, it must have up to 40 or 125 milliseconds latency each way.



The primary component of your RTT is the physical distance between sites. For every 1000 kilometers (621.4 miles), you observe a 5-millisecond delay each way. This delay does not include the time that is added by equipment in the path. Every device adds a varying amount of time, depending on the device, but a good rule is 25 microseconds for pure hardware devices.

For software-based functions (such as compression that is implemented in applications), the added delay tends to be much higher (usually in the millisecond plus range.) The following is an example of a physical delay.

Company A has a production site that is 1900 kilometers (1180.6 miles) away from its recovery site. The network service provider uses a total of five devices to connect the two sites. In addition to those devices, Company A uses a SAN FC router at each site to provide FCIP to encapsulate the FC traffic between sites.

Now, there are seven devices and 1900 kilometers (1180.6 miles) of distance delay. All the devices are adding 200 microseconds of delay each way. The distance adds 9.5 milliseconds each way, for a total of 19 milliseconds. Combined with the device latency, the delay is 19.4 milliseconds of physical latency minimum, which is under the 80-millisecond limit of GM until you realize that this number is the best case number.

The link quality and bandwidth play a large role. Your network provider likely ensures a latency maximum on your network link. Therefore, be sure to stay as far beneath the GM RTT limit as possible. You can easily double or triple the expected physical latency with a lower quality or lower bandwidth network link. Then, you are within the range of exceeding the limit if high I/O occurs that exceeds the bandwidth capacity.

When you get a 1920 event, always check the latency first. The FCIP routing layer can introduce latency if it is not properly configured. If your network provider reports a much lower latency, you might have a problem at your FCIP routing layer. Most FCIP routing devices have built-in tools to enable you to check the RTT. When you are checking latency, remember that TCP/IP routing devices (including FCIP routers) report RTT by using standard 64-byte ping packets.

Effective transit time must be measured only by using packets that are large enough to hold an FC frame, or 2148 bytes (2112 bytes of payload and 36 bytes of header). Allow estimated resource requirements to be a safe amount because various switch vendors have optional features that might increase this size. After you verify your latency by using the proper packet size, proceed with normal hardware troubleshooting.

Before proceeding, look at the second largest component of your RTT, which is *serialization delay*. Serialization delay is the amount of time that is required to move a packet of data of a specific size across a network link of a certain bandwidth. The required time to move a specific amount of data decreases as the data transmission rate increases.

The amount of time in microseconds that is required to transmit a packet across network links of varying bandwidth capacity is compared. The following packet sizes are used:

- ▶ 64 bytes: The size of the common ping packet
- ▶ 1500 bytes: The size of the standard TCP/IP packet
- ▶ 2148 bytes: The size of an FC frame

Finally, your path (MTU) affects the delay that is incurred to get a packet from one location to another location. An MTU might cause fragmentation or be too large and cause too many retransmits when a packet is lost.

For more information, see [IBM Knowledge Center](#).

**Note:** Unlike 1720 errors, 1920 errors are deliberately generated by the system because it evaluated that a relationship can affect the host's response time. The system has no indication about if or when the relationship can be restarted. Therefore, the relationship cannot be restarted automatically and it must be done manually.

## 10.11.2 1720 error

The 1720 error (event ID 050020) is the other problem RC might encounter. The amount of bandwidth that is needed for system-to-system communications varies based on the number of nodes. It is important that it is not zero. When a partner on either side stops communication, a 1720 is displayed in your error log. According to the product documentation, there are no likely field-replaceable unit (FRU) breakages or other causes.

The source of this error is most often a fabric problem or a problem in the network path between your partners. When you receive this error, check your fabric configuration for zoning of more than one host bus adapter (HBA) port for each node per I/O group if your fabric has more than 64 HBA ports zoned. The suggested zoning configuration for fabrics is one port for each node per I/O group per fabric that is associated with the host.

For those fabrics with 64 or more host ports, this suggestion becomes a rule. Therefore, you see four paths to each volume discovered on the host because each host needs to have at least two FC ports from separate HBA cards, each in a separate fabric. On each fabric, each host FC port is zoned to two IBM SAN Volume Controller node ports where each node port comes from a different IBM SAN Volume Controller node. This configuration provides four paths per volume. More than four paths per volume are supported but not recommended.

Improper zoning can lead to SAN congestion, which can inhibit remote link communication intermittently. Checking the zero buffer credit timer with IBM Spectrum Control and comparing against your sample interval reveals potential SAN congestion. If a zero buffer credit timer is more than 2% of the total time of the sample interval, it might cause problems.

Next, always ask your network provider to check the status of the link. If the link is acceptable, watch for repeats of this error. It is possible in a normal and functional network setup to have occasional 1720 errors, but multiple occurrences could indicate a larger problem.

If you receive multiple 1720 errors, recheck your network connection and then check the system partnership information to verify its status and settings. Then, perform diagnostics for every piece of equipment in the path between your two IBM SAN Volume Controller systems. It often helps to have a diagram that shows the path of your replication from both logical and physical configuration viewpoints.

**Note:** With Consistency Protection enabled on the MM/GM relationships, the system tries to resume the replication when possible. Therefore, it is not necessary to manually restart the failed relationship after a 1720 error is triggered.

For more information, see [IBM Knowledge Center](#).

If your investigations fail to resolve your RC problems, contact your IBM Support representative for a more complete analysis.



## Ownership groups

The ownership groups feature, or object-based access control (OBAC), provides a method of implementing a multi-tenant solution on the IBM SAN Volume Controller systems. Its principles of operations and implementation steps are provided in this chapter.

This chapter includes the following topics:

- ▶ 11.1, “Ownership groups principles of operations” on page 674
- ▶ 11.2, “Implementing ownership groups on a new system” on page 675
- ▶ 11.3, “Migrating objects to ownership groups” on page 680

## 11.1 Ownership groups principles of operations

Ownership groups allow the allocation of storage resources to several independent tenants with the assurance that one tenant cannot access resources that are associated with another tenant.

Ownership groups restrict access for users in the ownership group to only those objects that are defined within that ownership group. An owned object can belong to one ownership group. Users in an ownership group are restricted to viewing and managing objects within their ownership group. Users that are not in an ownership group can continue to view or manage all the objects on the system based on their defined user role, including objects within ownership groups.

Only users with Security Administrator roles (for example, superuser) can configure and manage ownership groups.

The system supports several resources that you assign to ownership groups:

- ▶ Child pools
- ▶ Volumes
- ▶ Volume groups
- ▶ Hosts
- ▶ Host clusters
- ▶ Host mappings
- ▶ FlashCopy mappings
- ▶ FlashCopy consistency groups

An owned object can belong to only one ownership group. An owner is a user with an ownership group that can view and manipulate objects within that group.

Before you create ownership groups and assign resources and users, it important to understand the following guidelines:

- ▶ Users can be in only one ownership group at a time (that applies to both local and remotely authenticated users).
- ▶ Objects can be within at most one ownership group.
- ▶ Global resources, such as drives, enclosures, arrays, cannot be assigned to ownership groups.
- ▶ Global users that do not belong to any ownership group can view and manage (depending on their user role) all resources on the system, including those that belong to an ownership group, and users within an ownership group.
- ▶ Users within an ownership group cannot have the Security Administrator role. All Security Administrator role users are global users.
- ▶ Users within an ownership group can view or change resources within the ownership group in which they belong.
- ▶ Users within an ownership group cannot change any objects outside of their ownership group. This restriction includes global resources that are related to resources within the ownership group. For example, a user can change a volume in the ownership group, but not the drive that provides the storage for that volume.
- ▶ Users within an ownership group cannot view or change resources if those resources are assigned to another ownership group or are not assigned to any ownership group. However, users within ownership groups can view and display global resources.

For example, users can display information about drives on the system because drives are a global resource that cannot be assigned to any ownership group.

When a user group is assigned to an ownership group, the users in that user group retain their role but are restricted to only those resources that belong to the same ownership group. The role that is associated with a user group can define the permitted operations on the system and the ownership group can further limit access to individual resources. For example, you can configure a user group with the Copy Operator role, which limits user access to FlashCopy operations. Access to individual resources, such as a specific FlashCopy consistency group, can be further restricted by assigning it to an ownership group.

Child pools is a key requirement for the ownership groups feature. By defining a child pool and assigning it to an ownership group, the system administrator provides capacity for volumes that ownership group users can create or manage. Child pools are supported only with standard pools. You cannot create a child pool for a Data Reduction Pool (DRP).

Depending on the type of resource, the owning group for it can be defined explicitly or inherited from these explicitly defined objects. For example, a child pool needs an ownership group parameter to be set by a system administrator. But volumes that are created in that child pool automatically inherit the ownership group from a child pool. For more information about ownership inheritance, see [IBM Knowledge Center](#).

When the user logs on to the management GUI or command-line interface (CLI), only resources that they can through the ownership group are available. Also, only events and commands that are related to the ownership group in which a user belongs are viewable by those users.

## 11.2 Implementing ownership groups on a new system

This section describes ownership group implementation process for a new system, which has no volumes and users that must be migrated to ownership groups.

### 11.2.1 Creating an ownership group

Complete the following steps to create an ownership group:

1. To create the first ownership group, select **Access** → **Ownership Groups**, as shown in Figure 11-1. Enter the name for it, and click **Create Ownership Group**.

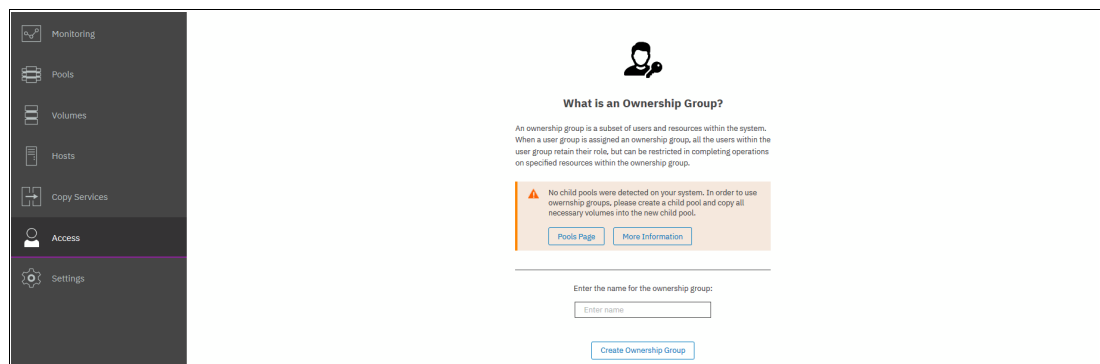


Figure 11-1 Creating a first ownership group

After the first group is created, the pane changes to ownership group mode, as shown in Figure 11-2. The new ownership group has no user groups and no resources assigned to it.

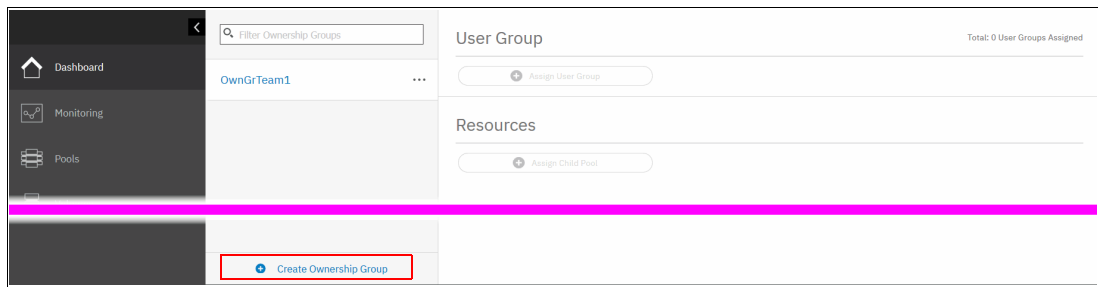


Figure 11-2 Ownership groups management pane

2. To create more ownership groups, click **Create Ownership Group**.

## 11.2.2 Assigning users to an ownership group

To create accounts for users that use ownership group resources, the user group must be created and assigned to the ownership group.

To create a user group, select **Access** → **Users by Group** and click **Create User Group**. The Create User Group window opens, as shown in Figure 11-3. Within it, specify the User Group name and select an ownership group to tie this user group to and select a role for the users in this group.

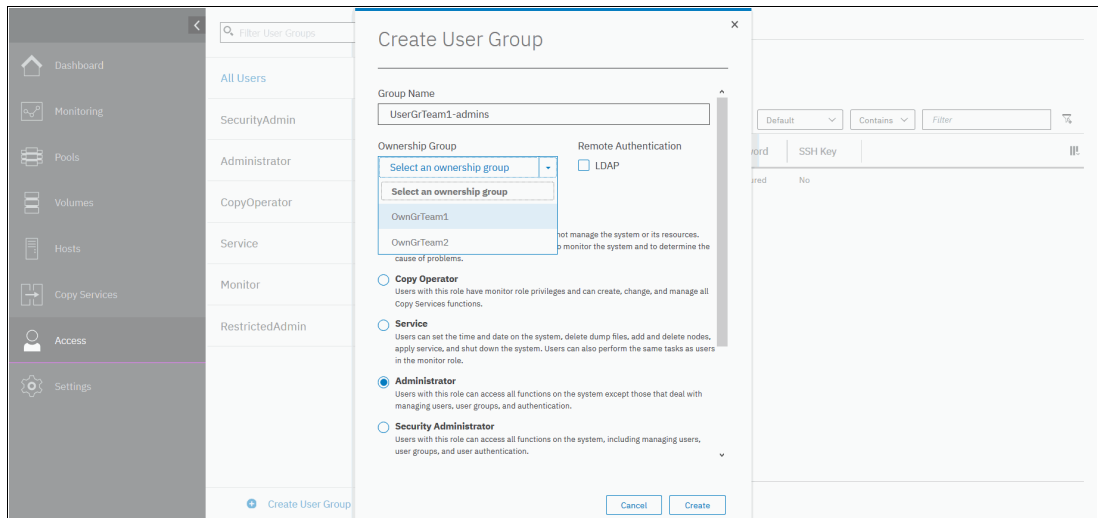


Figure 11-3 Creating and assigning a user group

For more information about user roles, see [IBM Knowledge Center](#).

To create a volume, host, and other objects in an ownership group, users must have an Administrator or Restricted Administrator role. Users with Security Administrator role cannot be assigned to an ownership group.

You can also set up a user group to use remote authentication; that is, Lightweight Directory Access Protocol (LDAP), if it is enabled. To set up remote authentication, select the **LDAP** check box.

**Note:** Users that use LDAP can belong to multiple user groups, but belong to only one ownership group that is associated with one of the user groups.

If remote authentication is not configured, you also must create a user (or users) and assign it to a created user group, as shown in Figure 11-4.

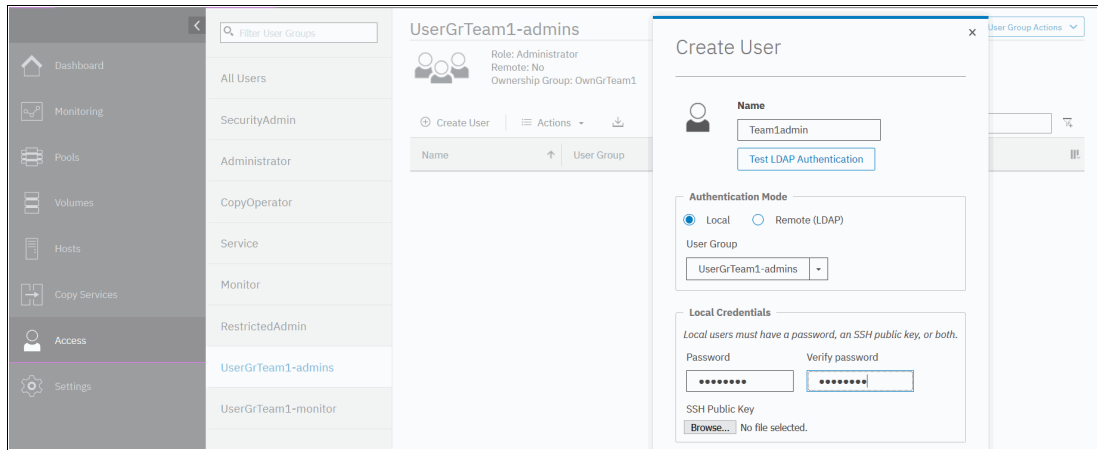


Figure 11-4 Creating a user

You can manage user groups that are assigned to an ownership group by selecting **Access** → **Ownership Groups**, as shown in Figure 11-5. To assign a user group that exists but is not yet assigned to any ownership group, click **Assign User Group**. To unassign, click the “...” icon that is next to the assigned user group name.

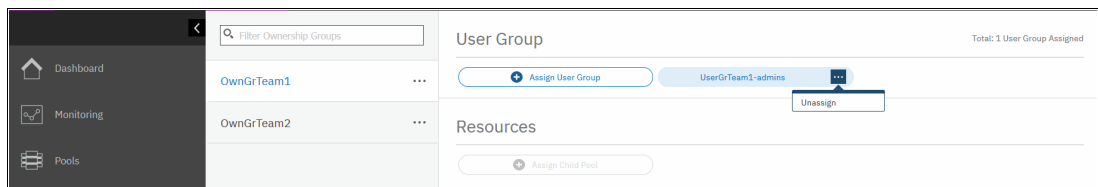


Figure 11-5 Unassigning user groups

Multiple user groups with different user roles can be assigned to one ownership group. For example, you can create and assign a user group with the Monitor role in addition to a group with the Administrator role to have two sets of users with different privilege levels accessing an ownership groups resources.

### 11.2.3 Creating ownership group resources

To create ownership group volumes and other resources, a child pool must be created and assigned to the ownership group.

Select **Pools** → **Pools**, right-click a parent pool that is designated to be a container for child pools, and click **Create Child Pool**, as shown in Figure 11-6 on page 678.

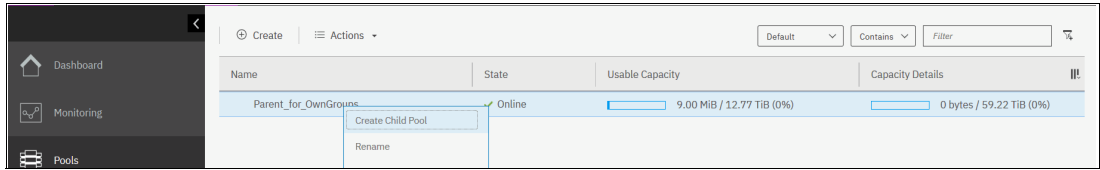


Figure 11-6 Creating a child pool

When creating a child pool, specify an ownership group for it and assign a part of the parent's pool capacity, as shown in Figure 11-7. Ownership group objects can use only capacity that is provisioned for them with the child pool.

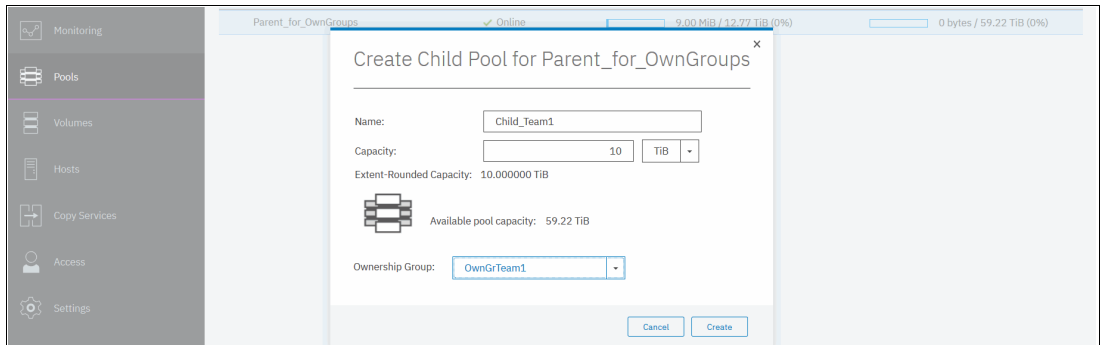


Figure 11-7 Creating a child pool and assigning it to an ownership group

Multiple child pools that are created from the same or different parent pools can be assigned to a single ownership group.

After a child pool is created and assigned, the ownership group management pane at **Access** → **Ownership Groups** changes to show the assigned and available resources, as shown in Figure 11-8.

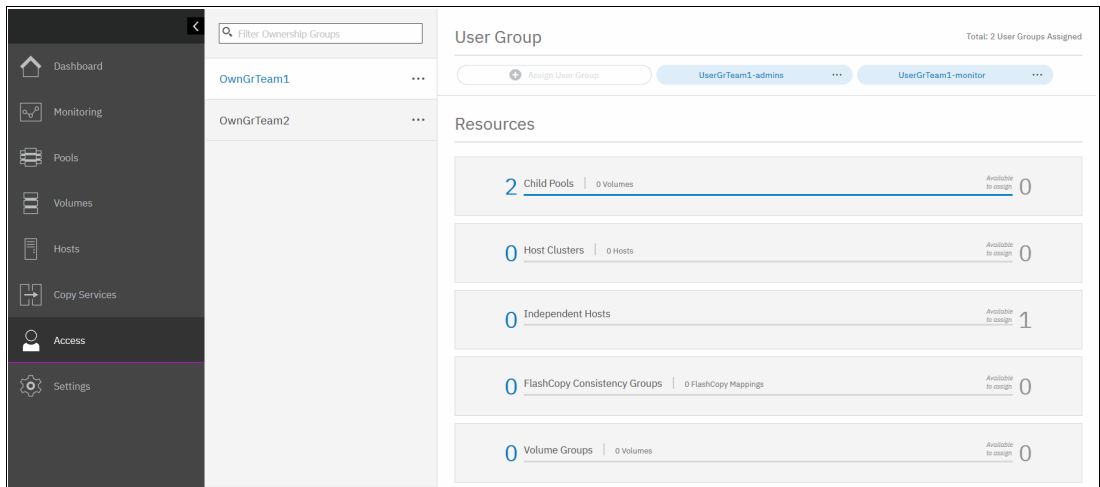


Figure 11-8 Ownership group management pane

Any volumes that are created on a child pool that are assigned to an ownership group inherit ownership from the child pool.



After a child pool and user group are assigned to an ownership group, ownership group administrators can log in with their credentials and start creating volumes, host and host clusters, or FlashCopy mappings.

For more information about creating these objects, see Chapter 6, “Volumes” on page 255, Chapter 7, “Hosts” on page 351, and Chapter 10, “Advanced Copy Services” on page 491.

Although an ownership group administrator can create objects only within the resources that are assigned to them, the system administrator can create, monitor, and assign objects for any ownership group.

## 11.2.4 Listing ownership group resources

By default, the Ownership Group attribute is not enabled in the GUI panes that list volumes and other objects that can be owned. For convenience, the system administrator can enable this attribute, as shown in Figure 11-9.

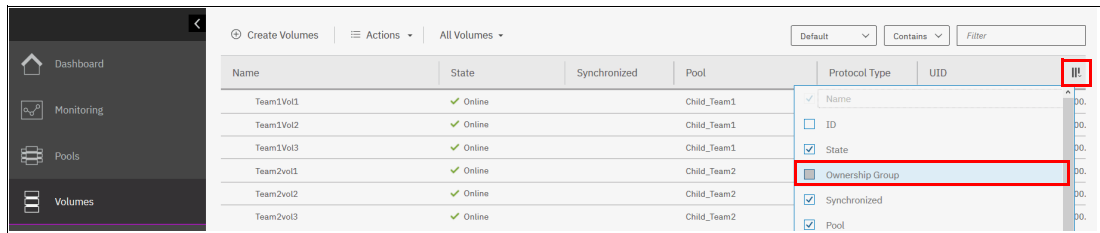


Figure 11-9 Enabling ownership group attribute display

As an example, the volume listing for a system administrator looks as shown in Figure 11-10.

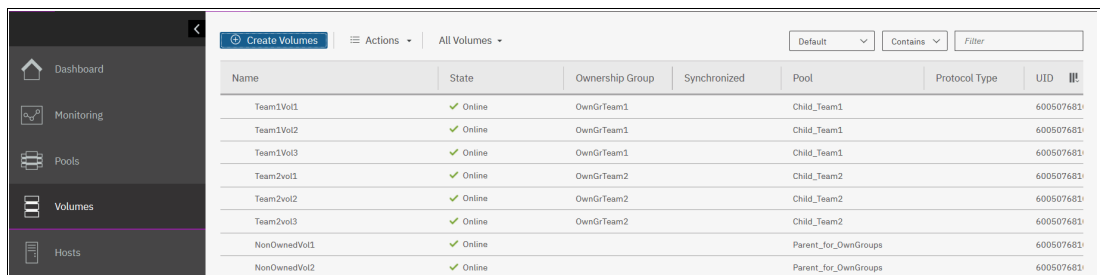


Figure 11-10 Listing volumes for all ownership groups

The system administrator can see and manage resources of all ownership groups, and resources that are not assigned to any of the groups.

When the ownership group user logs in, they can see and manage only resources that are assigned to their group. The example in Figure 11-11 on page 680 shows the dashboard pane for an ownership group user with the Administrator role.

It is different from a dashboard visible to system administrator. Global system performance and capacity parameters are not shown. Only tiles for resources that are assigned to this ownership group are displayed. Out of eight volumes are configured on a system and shown in Figure 11-10; they can see and manage only three volumes that belong to the group.

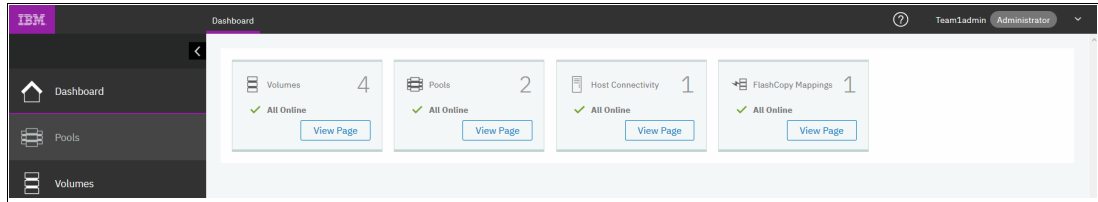


Figure 11-11 Ownership group administrator view

The ownership group user can use the GUI to browse, create, and delete (depending on their user role) resources that are assigned to their group. To see information about the global resources (for example, list managed disks [MDisks] or arrays on the pool), they must use the CLI. Ownership group users cannot manage global resources; they can only view them.

### 11.2.5 Actions on ownership groups

The global system administrator can rename or remove an ownership group. Select **Access** → **Ownership Groups**, click the “...” icon that is next to the group name, and select the required task, as shown in Figure 11-12.

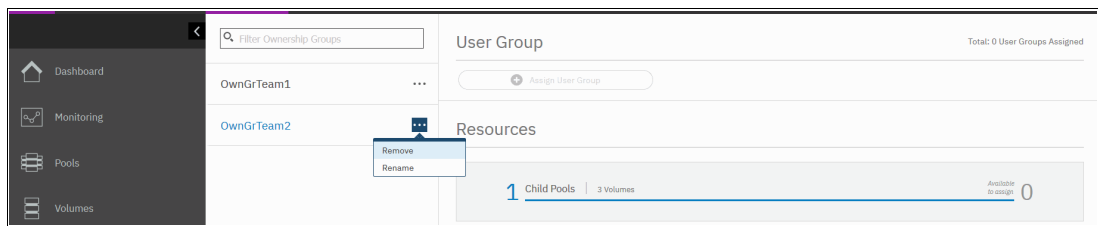


Figure 11-12 Renaming or removing an ownership group

When an ownership group is removed with the GUI, all ownership assignment information for all the objects of the ownership group is removed, but the objects remain configured. Only the system administrator can manage those resources later.

## 11.3 Migrating objects to ownership groups

If you want to use ownership groups for objects on your system, you must reconfigure specific resources if you want to configure ownership groups.

If child pools are on the system, you can define an ownership group to the child pool or child pools. Before you define an ownership group to child pools, determine other related objects that you want to migrate. Any volumes that are in the child pool inherit the ownership group that is defined for the child pool.

If no child pool is on the system, you must create child pools and move any volumes to those child pools before you can assign them to ownership groups. If volumes are in a parent pool, volume mirroring can be used to create copies of the volume within the child pool. Alternatively, volume migration can be used to relocate a volume from a parent pool to a child pool within that parent pool without requiring copying.

Consider the following example scenario in which resources are non-disruptively migrated to an ownership group:

1. Create a child pool. Do not assign it to an ownership group yet.

Migrate volumes that must be assigned to an ownership group to that child pool by using the volume migration or volume mirroring function. The example that is shown in Figure 11-13 shows volume NonOwnedVol1, which is in a parent pool and does not belong to any ownership group. This volume is mapped to a Small Computer System Interface (SCSI) host.

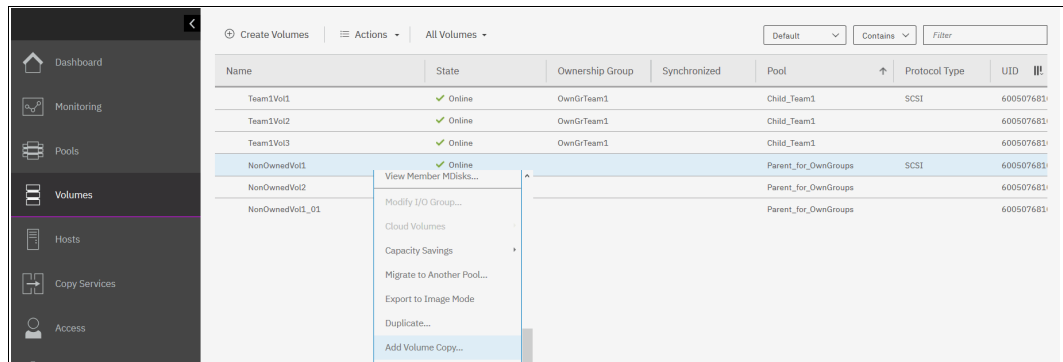


Figure 11-13 Adding volume copy for migration

2. To start migration, right-click the volume and select **Add Volume Copy**. You can also select **Migrate to Another Pool**; however, this method is suitable only if you are migrating from a pool with the same extent size (for example, from a parent pool to a child pool of the same pool), and provides less flexibility.

In the **Volume Copy** window, select the child pool that is to be assigned later to an ownership group, as shown in Figure 11-14.

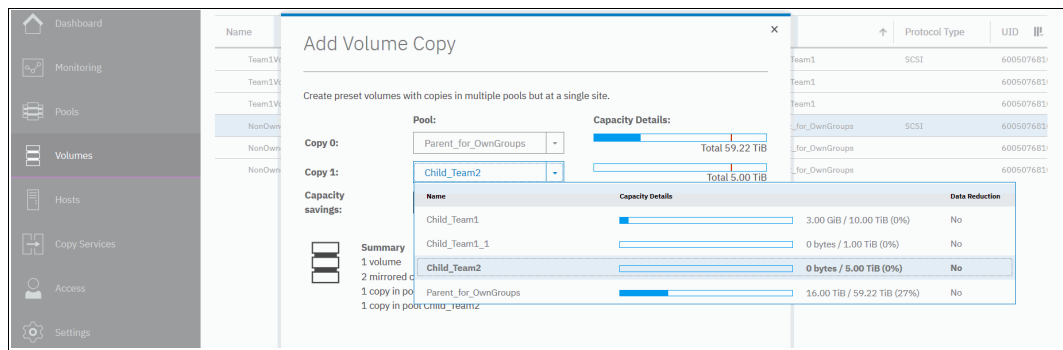


Figure 11-14 Migrating to a child pool

3. Repeat the previous step for all volumes that must belong to an ownership group. Then, remove the source copies.
4. Create an ownership group as described in 11.2.1, “Creating an ownership group” on page 675. Assign a user group to it, as described in 11.2.2, “Assigning users to an ownership group” on page 676.
5. As shown in Figure 11-15 on page 682, select **Access** → **Ownership Groups**. Then, select the wanted ownership group and click **Assign Child Pool**.

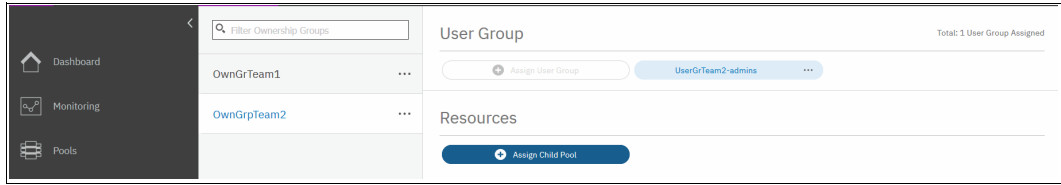


Figure 11-15 Assign child pool to an ownership group

6. Select a child pool to assign, as shown in Figure 11-16.

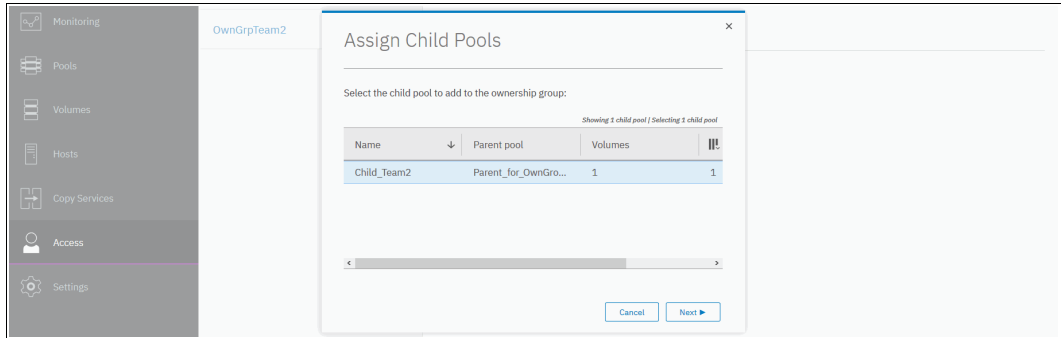


Figure 11-16 Selecting child pool to assign

After you click **Next**, the system notifies you that more resources inherit ownership from a volume and as the volume is mapped to a host, the host becomes an ownership group object, as shown in Figure 11-17.

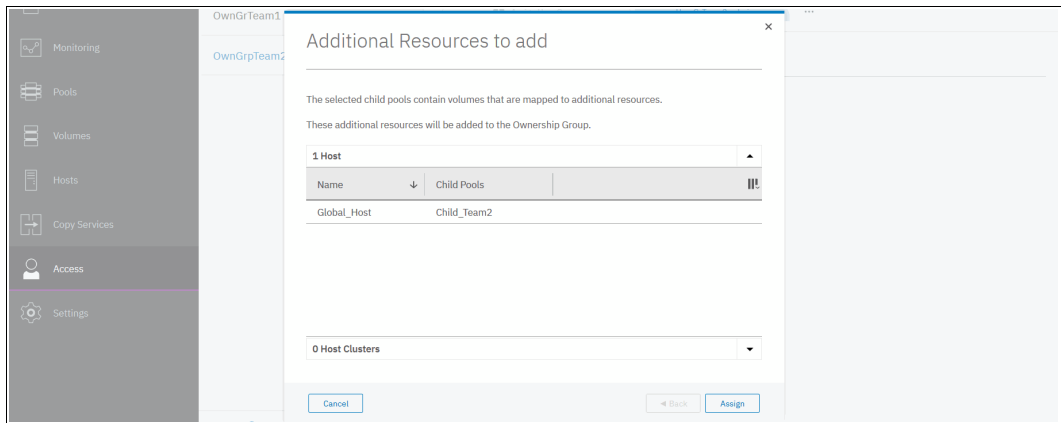


Figure 11-17 More resources to add

As shown in Figure 11-18, a volume and a host both belong to an ownership group. As a host and a volume are in a group, host mapping inherits ownership and also becomes a part of an ownership group.

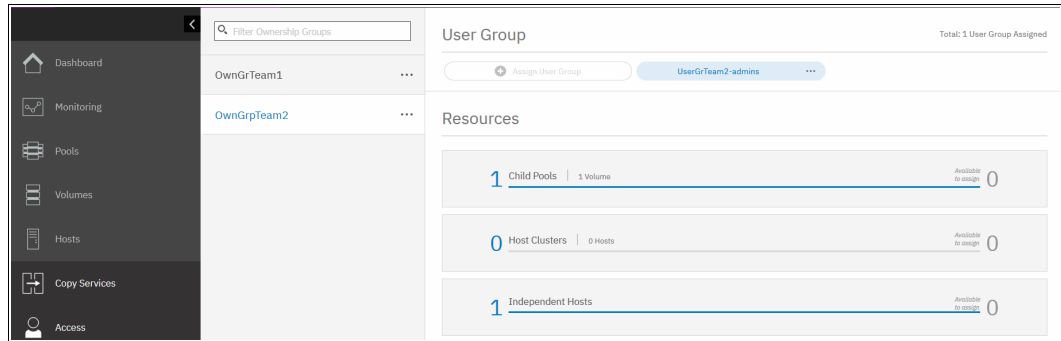


Figure 11-18 Resources of an ownership group

After the steps are completed, a child pool is assigned to an ownership group. If later you must migrate volumes to it, the same process can be used. However, during migration one volume copy is in an owned child pool, and the original copy remains in an unowned parent pool. Such a condition causes Inconsistent ownership, as shown in Figure 11-19.

The screenshot shows the 'All Volumes' table with the following data:

Name	State	Ownership Group	Synchronized	Pool	Protocol Type	UID
Team1Vol1	Online	OwnGrTeam1		Child_Team1	SCSI	600507681
Team1Vol2	Online	OwnGrTeam1		Child_Team1		600507681
Team1Vol3	Online	OwnGrTeam1		Child_Team1		600507681
NonOwnedVol1	Online	OwnGrTeam2		Child_Team2	SCSI	600507681
NonOwnedVol2	Online	Inconsistent		Parent_for_OwnGroups		600507681
Copy 0*	Online		Yes	Parent_for_OwnGroups		600507681
Copy 1	Online		No	Child_Team2		600507681

Figure 11-19 Example of inconsistent volume ownership

Until the inconsistent volume ownership is resolved, the volume does not belong to an ownership group and cannot be seen or managed by an ownership group administrator. To resolve it, delete one of the copies after both are synchronized.





# Encryption

Encryption protects against the potential exposure of sensitive user data that is stored on discarded, lost, or stolen storage devices. IBM SAN Volume Controller and other storage devices driven by IBM Spectrum Virtualize supports optional encryption of data at-rest.

This chapter includes the following topics:

- ▶ 12.2, “Planning for encryption” on page 687
- ▶ 12.3, “Defining encryption of data at-rest” on page 687
- ▶ 12.4, “Activating encryption” on page 692
- ▶ 12.5, “Enabling encryption” on page 703
- ▶ 12.6, “Configuring more providers” on page 729
- ▶ 12.7, “Migrating between providers” on page 733
- ▶ 12.8, “Recovering from a provider loss” on page 737
- ▶ 12.9, “Using encryption” on page 737
- ▶ 12.10, “Rekeying an encryption-enabled system” on page 746
- ▶ 12.11, “Disabling encryption” on page 752

## 12.1 General types of encryption across IBM Spectrum Virtualize

Within IBM Spectrum Virtualize, IBM SAN Volume Controller, and IBM FlashSystems, three different types of encryption are available.

### 12.1.1 Externally virtualized storage

Data is decrypted and encrypted as we issue read/write I/Os to the external storage. You can have an encryption key per storage pool or per child storage pool. Migrating a volume between pools (by using volume mirroring) can be used as a technique for encrypting and decrypting the data.

The key per pool (and in particular allowing different keys for child pools) supports some part of the multi-tenant use case (if you delete a pool then you delete the key and cryptoerase the data). However, all the keys are wrapped and protected by a single master key that is obtained from a USB stick or an external key server.

As a special case, it is possible to turn off encryption for individual managed disks (MDisks) within the storage pool. If an external storage controller supports encryption, you can choose to allow it to encrypt the data instead.

### 12.1.2 Serial-attached SCSI internal storage

Data is decrypted and encrypted by the serial-attached Small Computer System Interface (SCSI) (SAS) controller. An encryption key is available per Redundant Array of Independent Disks (RAID). Normally, all arrays in a storage pool are encrypted to form an encrypted storage pool. Although you can create child storage pools, only one key is available per RAID array. Multi-tenancy is possible only if you have more than one array and storage pool, which usually is not practical.

You can migrate volumes from a non-encrypted storage pool to an encrypted storage pool, or you can add an encrypted array to a storage pool and then delete the unencrypted array (which migrates all of the data automatically) as a way of encrypting data.

### 12.1.3 Non-Volatile Memory Express internal storage

Data is decrypted and encrypted by the Non-Volatile Memory Express (NVMe) drives. Each drive has a media encryption key, but this key is wrapped and protected by an encryption key per RAID array. Therefore, it has the same properties as SAS internal storage.

A storage pool can include a mixture of two or all three types of storage. In this case, the SAS and NVMe internal storage uses a key per RAID array for encryption and the externally virtualized storage uses the pool level key.

Because it is almost impossible to control exactly what storage is used for each volume, from a security viewpoint this is effectively a single key for the whole pool and a cryptographic erase is only possible by deleting the entire storage pool and arrays.



## 12.2 Planning for encryption

Data-at-rest encryption is a powerful tool that can help organizations protect the confidentiality of sensitive information. However, like any other tool, encryption must be used correctly to fulfill its purpose.

Multiple drivers exist for an organization to implement data-at-rest encryption. These drivers can be internal, such as protection of confidential company data, and ease of storage sanitization, or external, such as compliance with legal requirements or contractual obligations.

Therefore, before configuring encryption on the storage, the organization defines its needs and, if it is decided that data-at-rest encryption is required, includes it in the security policy. Without defining the purpose of the particular implementation of data-at-rest encryption, it is difficult or impossible to choose the best approach to implement encryption and verify whether the implementation meets the set of goals.

The following factors are worth considering during the design of a solution that includes data-at-rest encryption:

- ▶ Legal requirements
- ▶ Contractual obligations
- ▶ Organization's security policy
- ▶ Attack vectors
- ▶ Expected resources of an attacker
- ▶ Encryption key management
- ▶ Physical security

Multiple regulations mandate data-at-rest encryption, from processing of Sensitive Personal Information to the guidelines of the Payment Card Industry. If any regulatory or contractual obligations govern the data that is held on the storage system, they often provide a wide and detailed range of requirements and characteristics that must be realized by that system. Apart from mandating data-at-rest encryption, these documents might contain requirements concerning encryption key management.

Another document that should be consulted when planning data-at-rest encryption is the organization's security policy.

The outcome of a data-at-rest encryption planning session answers the following questions:

- ▶ What are the goals that the organization wants to realize by using data-at-rest encryption?
- ▶ How will data-at-rest encryption be implemented?
- ▶ How can it be demonstrated that the proposed solution realizes the set of goals?

## 12.3 Defining encryption of data at-rest

*Encryption* is the process of encoding data so that only authorized parties can read it. Secret keys are used to encode the data according to well-known algorithms.

Encryption of data-at-rest as implemented in IBM Spectrum Virtualize is defined by the following characteristics:

- ▶ *Data-at-rest* means that the data is encrypted on the end device (drives).
- ▶ The algorithm that is used is the Advanced Encryption Standard (AES) US government standard from 2001.

- ▶ Encryption of data at-rest complies with the Federal Information Processing Standard 140-2 (FIPS-140-2) standard.
- ▶ AES 256 is used for master access keys.
- ▶ XTS-AES 256 encryption, is a FIPS 140-2 compliant algorithm.
- ▶ XTS-AES-256 is used for data encryption.
- ▶ The algorithm is public; the only secrets are the keys.
- ▶ A symmetric key algorithm is used. The same key is used to encrypt and decrypt data.

The encryption of system data and metadata is not required; therefore, they are not encrypted.

### 12.3.1 Encryption methods

There are two types of encryption on devices running IBM Spectrum Virtualize: hardware encryption and software encryption. Both methods of encryption protect against the potential exposure of sensitive user data that is stored on discarded, lost, or stolen media. Both can also facilitate the warranty return or disposal of hardware.

Which method is used for encryption is chosen automatically by the system based on the placement of the data:

- ▶ Hardware encryption: Data is encrypted by using SAS hardware. It is used only for internal storage (drives).
- ▶ Software encryption: Data is encrypted by using nodes' CPU (encryption code uses AES-NI CPU instruction set). It is used only for external storage.

**Note:** Software encryption is available in IBM Spectrum Virtualize code V7.6 and later.

Both methods of encryption use the same encryption algorithm, key management infrastructure, and license.

**Note:** The design for encryption is based on the concept that a system is encrypted or not encrypted. Encryption implementation is intended to encourage solutions that contain only encrypted volumes or only unencrypted volumes. For example, after encryption is enabled on the system, all new objects (for example, pools) are created as encrypted by default.

### 12.3.2 Encrypted data

IBM Spectrum Virtualize performs data-at-rest encryption, which is the process of encrypting data that is stored on the end devices, such as physical drives.

Data is encrypted or decrypted when it is written to or read from internal drives (hardware encryption) or external storage systems (software encryption).

Therefore, data is encrypted when transferred across the storage area network (SAN) only between IBM Spectrum Virtualize systems and external storage. Data in transit is *not* encrypted when transferred on SAN interfaces under the following circumstances:

- ▶ Server-to-storage data transfer
- ▶ Remote Copy (RC); for example, Global Mirror (GM) or Metro Mirror (MM)
- ▶ Intracluster communication

**Note:** Only data-at-rest is encrypted. Host-to-storage communication and data that is sent over links that are used for Remote Mirroring are not encrypted.

Figure 12-1 shows an encryption example. Encrypted disks and encrypted data paths are marked in blue. Unencrypted disks and data paths are marked in red. The server sends unencrypted data to an IBM SAN Volume Controller 2145-DH8 system, which stores hardware-encrypted data on internal disks. The data is mirrored to a remote Storwize V7000 Gen1 system by using RC. The data that is flowing through the RC link is not encrypted. Because the Storwize V7000 Gen1 (2076-324) cannot perform any encryption activities, data on the Storwize V7000 Gen1 is not encrypted.

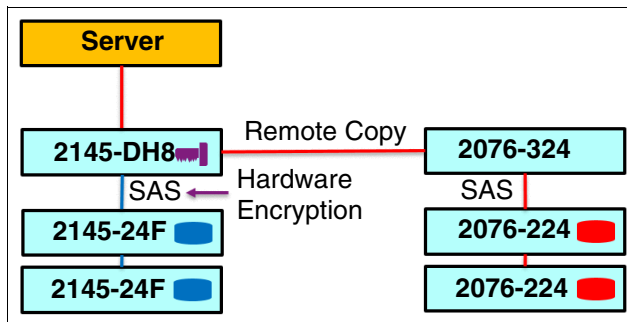


Figure 12-1 Encryption on single site

To enable encryption of both data copies, the Storwize V7000 Gen1 must be replaced by an encryption capable (with optional encryption enabled), IBM Spectrum Virtualize system as shown in Figure 12-2. After such replacement, both copies of data are encrypted, but the RC communication between both sites remains unencrypted.

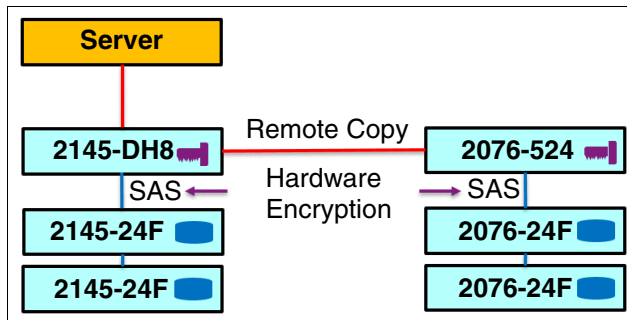


Figure 12-2 Encryption on both sites

Figure 12-3 shows an example configuration that uses software and hardware encryption. Software encryption is used to encrypt an external virtualized storage system (2076-324 in Figure 12-3). Hardware encryption is used for internal, SAS-attached disk drives.

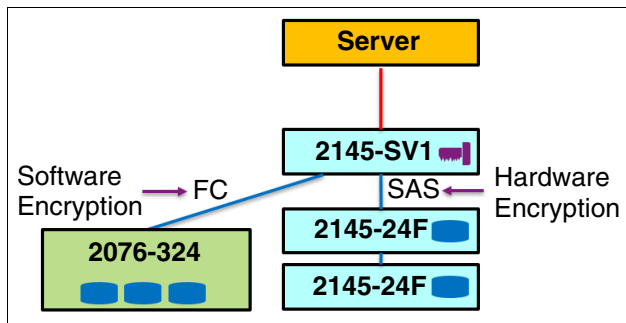


Figure 12-3 Example of software encryption and hardware encryption

Placement of hardware encryption and software encryption in the IBM Spectrum Virtualize code stack is shown in Figure 12-4. Because compression is performed before encryption, it is possible to realize the benefits of compression for the encrypted data.

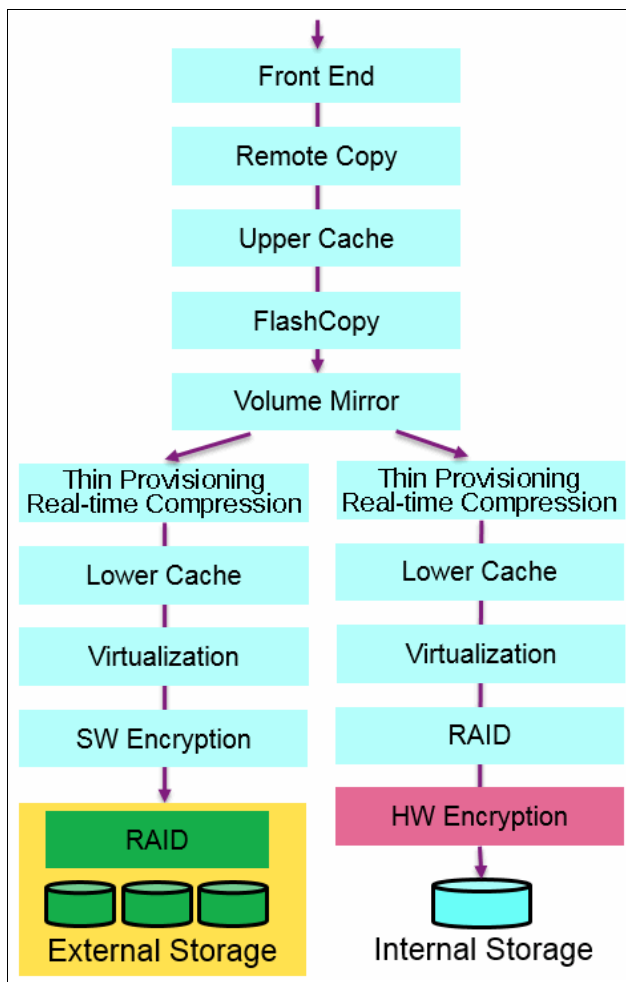


Figure 12-4 Encryption placement in the IBM Spectrum Virtualize software stack

Each volume copy can use different encryption methods (hardware and software). It also can have volume copies with different encryption status (encrypted versus unencrypted). The encryption method depends only on the pool that is used for the specific copy. You can migrate data between different encryption methods by using volume migration or volume mirroring.

### 12.3.3 Encryption keys

Hardware and software encryption use the same encryption key infrastructure. The only difference is the object that is encrypted by using the keys. The following objects can be encrypted:

- ▶ Pools (software encryption)
- ▶ Child pools (software encryption)
- ▶ Arrays (hardware encryption)

Consider the following points regarding encryption keys:

- ▶ Keys are unique for each object, and they are created when the object is created.
- ▶ Two types of keys are defined in the system:
  - Master access key:
    - The master access key is created when encryption is enabled.
    - The master access key can be stored on USB flash drives, key servers, or both. One master access key is created for each enabled encryption key provider.
    - It can be copied or backed up as necessary.
    - It is *not* permanently stored anywhere in the system.
    - It is required at boot time to unlock access to encrypted data.
  - Data encryption keys (one for each encrypted object):
    - Data encryption keys are used to encrypt data. When an encrypted object (such as an array, pool, or child pool) is created, a new data encryption key is generated for this object.
    - MDisk that are not self-encrypting are automatically encrypted by using the data encryption key of the pool or child pool to which they belong.
    - MDisk that are self-encrypting are not reencrypted by using the data encryption key of the pool or child pool they belong to by default. You can override this default by manually configuring the MDisk as not self-encrypting.
    - Data encryption keys are stored in secure memory.
    - During cluster internal communication, data encryption keys are encrypted with the master access key.
    - Data encryption keys cannot be viewed or changed.
    - When an encrypted object is deleted, its data encryption key is discarded (*secure erase*).

**Important:** Consider the following points:

- ▶ If all master access key copies are lost and the system must cold restart, all encrypted data is gone. No method exists, even for IBM, to decrypt the data without the keys. If encryption is enabled and the system cannot access the master access key, all SAS hardware is offline, including unencrypted arrays.
- ▶ A self-encrypting MDisk is an MDisk from an encrypted volume in an external storage system.

### 12.3.4 Encryption licenses

Encryption is a licensed feature that uses key-based licensing. A license must be present for each node in the system before you can enable encryption.

If you add a node to a system that in which encryption is enabled, the node must also be licensed.

No trial licenses for encryption exist on the basis that when the trial runs out, the access to the data is lost. Therefore, you must purchase an encryption license before you activate encryption. Licenses are generated by IBM Data Storage Feature Activation (DSFA) based on the serial number (S/N) and the machine type and model (MTM) of the node.

You can activate an encryption license during the initial system setup (on the Encryption window of the initial setup wizard) or later on, in the running environment.

Contact your IBM marketing representative or IBM Business Partner to purchase an encryption license.

## 12.4 Activating encryption

Encryption is enabled at a system level and all of the following prerequisites must be met to use encryption:

- ▶ You must purchase an encryption license before you activate the function.  
If you did not purchase a license, contact an IBM marketing representative or IBM Business Partner to purchase an encryption license.
- ▶ At least three USB flash drives are required if you plan to *not* use a key management server. They are available as a feature code from IBM.
- ▶ You must activate the license that you purchased.
- ▶ Encryption must be enabled.

Activation of the license can be performed in one of two ways:

- ▶ **Automatic activation:** Used when you have the authorization code and the workstation that is being used to activate the license has access to external network. In this case, you must enter only the authorization code. The license key is automatically obtained from the internet and activated in the IBM Spectrum Virtualize system.
- ▶ **Manual activation:** If you cannot activate the license automatically because any of the requirements are not met, you can follow the instructions that are provided in the GUI to obtain the license key from the web and activate in the IBM Spectrum Virtualize system.

Both methods are available during the initial system setup and when the system is in use.

### 12.4.1 Obtaining an encryption license

You must purchase an encryption license before you activate encryption. If you did not purchase a license, contact an IBM marketing representative or IBM Business Partner to purchase an encryption license.

When you purchase a license, you receive a function authorization document with an authorization code printed on it. This code allows you to proceed by using the automatic activation process.

If the automatic activation process fails or if you prefer the use of the manual activation process, see [this web page](#) to retrieve your license keys.

Ensure that you have the following information:

- ▶ Machine type (MT)
- ▶ Serial number (S/N)
- ▶ Machine signature
- ▶ Authorization code

For more information about how to retrieve the machine signature of a node, see 12.4.5, “Manual license activation” on page 701.

### 12.4.2 Starting the activation process during initial system setup

One of the steps in the initial setup enables encryption license activation. The system asks, “Was the encryption feature purchased for this system?”

To activate encryption at this stage, complete the following steps:

1. Select **Yes**, as shown in Figure 12-5 on page 694.

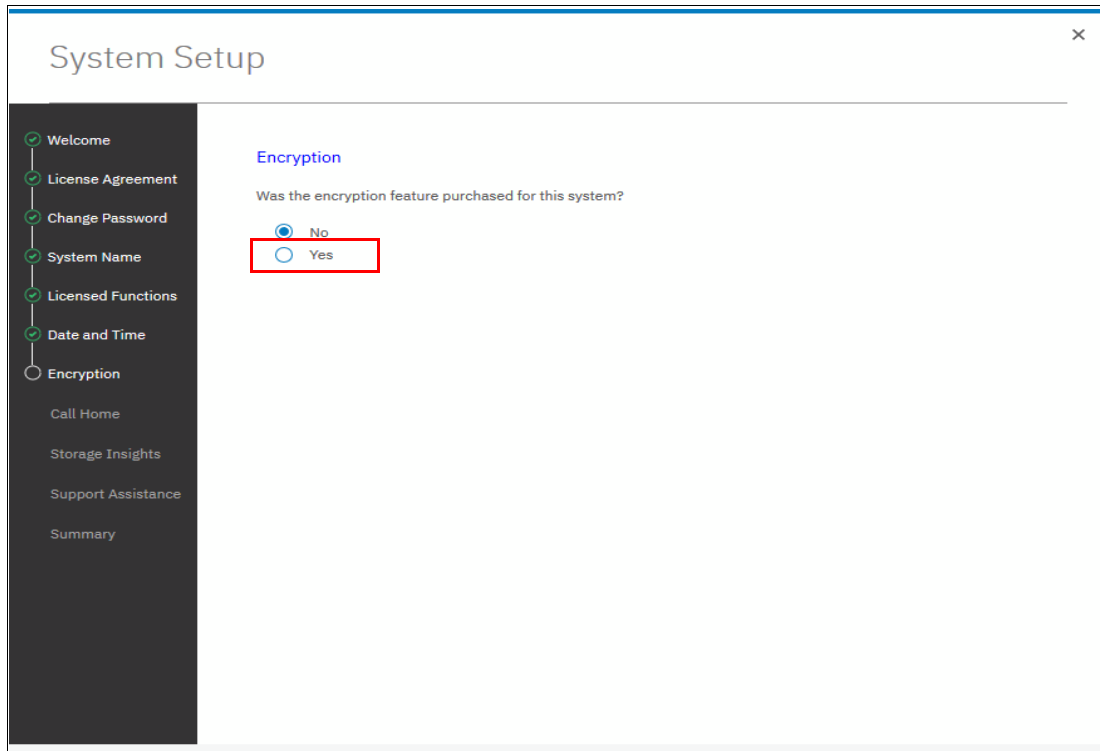


Figure 12-5 Encryption activation during initial system setup

The Encryption window displays information about your storage system, as shown in Figure 12-6.

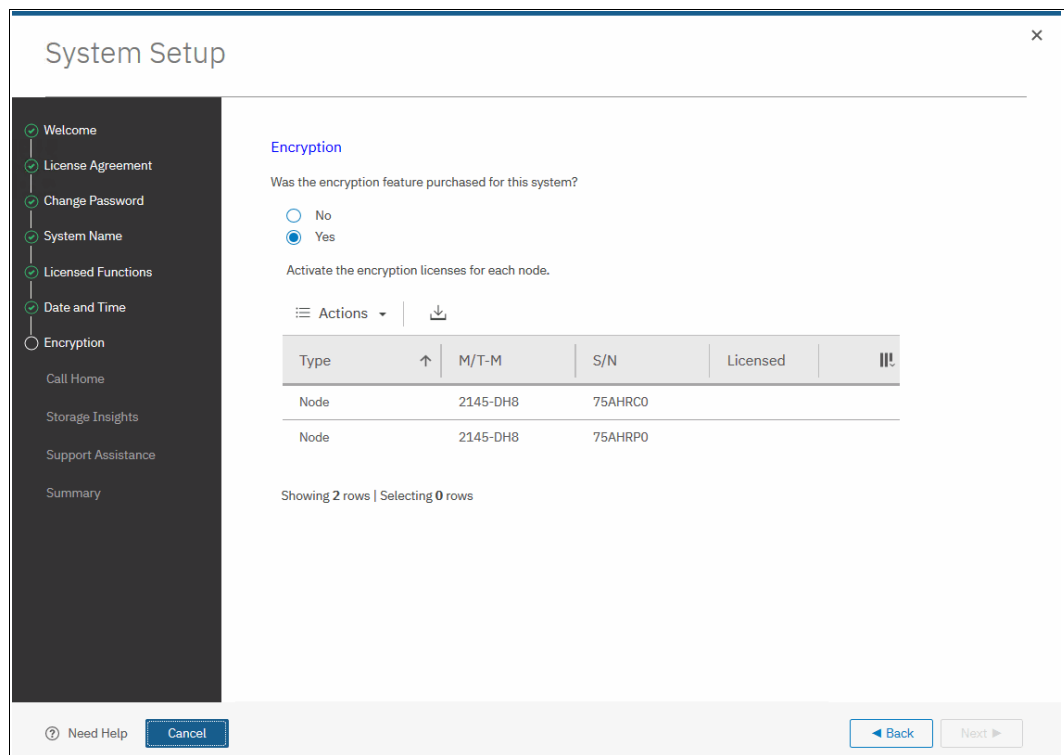


Figure 12-6 Information storage system during initial system setup



2. Right-click the node to open a menu with two license activation options (**Activate License Automatically** and **Activate License Manually**), as shown in Figure 12-7. Use either option to activate encryption. For more information about how to complete the automatic activation process, see 12.4.4, “Activate the license automatically” on page 697. For more information about how to complete a manual activation process, see 12.4.5, “Manual license activation” on page 701.

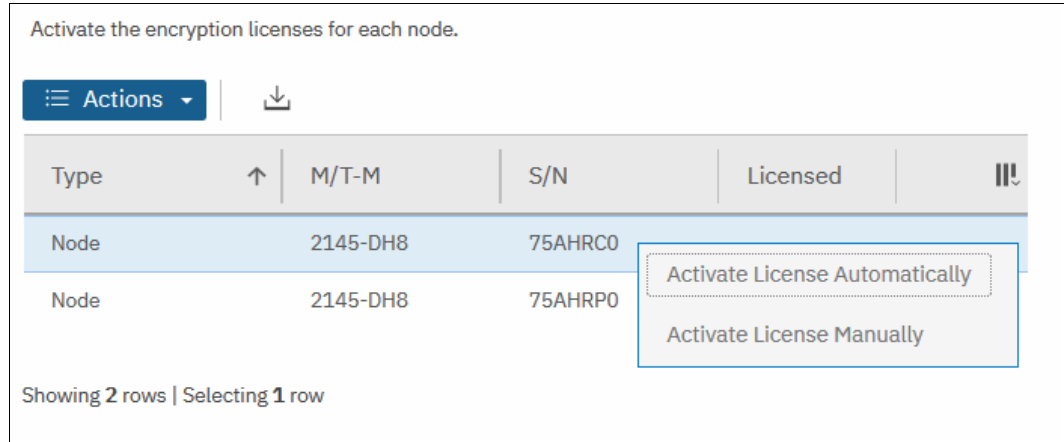


Figure 12-7 Selecting license activation method

3. After either activation process is complete, you can see a green check mark in the column labeled Licensed next to a node for which the license was enabled. You can proceed with the initial system setup by clicking **Next** (see Figure 12-8 on page 696).

**Note:** Every enclosure needs an active encryption license before you can enable encryption on the system. Attempting to add a non-licensed enclosure to an encryption-enabled system fails.

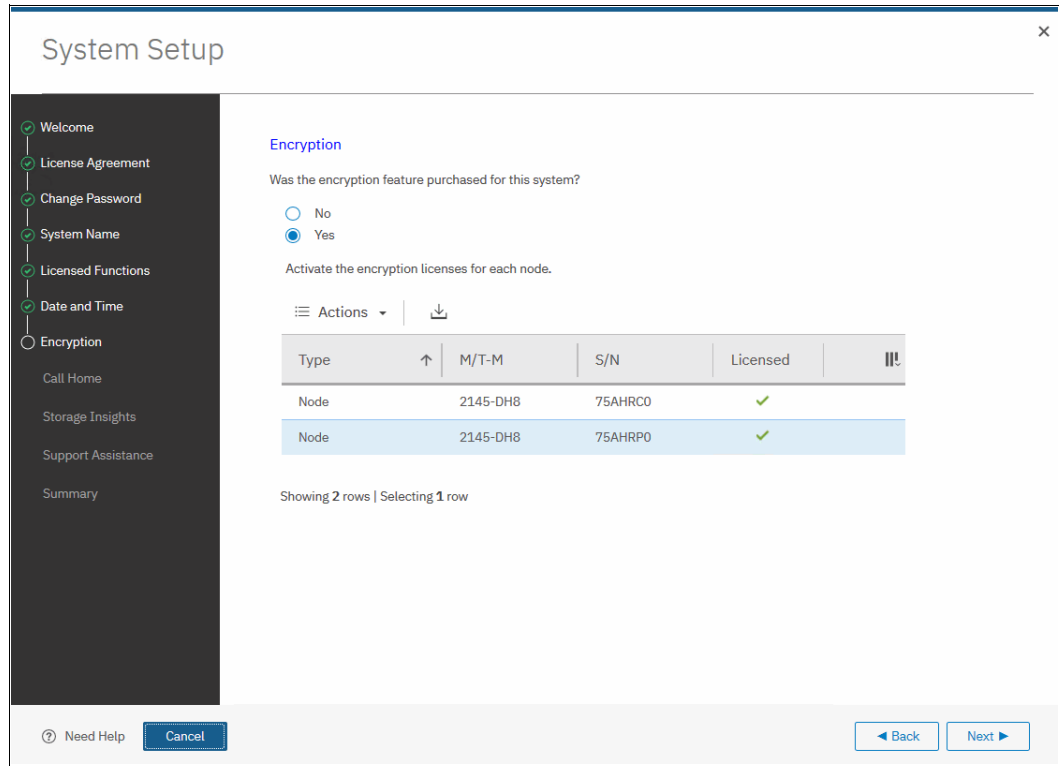


Figure 12-8 Successful encryption license activation during initial system setup

### 12.4.3 Starting the activation process on a running system

To activate encryption on a running system, complete the following steps:

1. Click **Settings** → **System** → **Licensed Functions**.
2. Click **Encryption Licenses**, as shown in Figure 12-9.

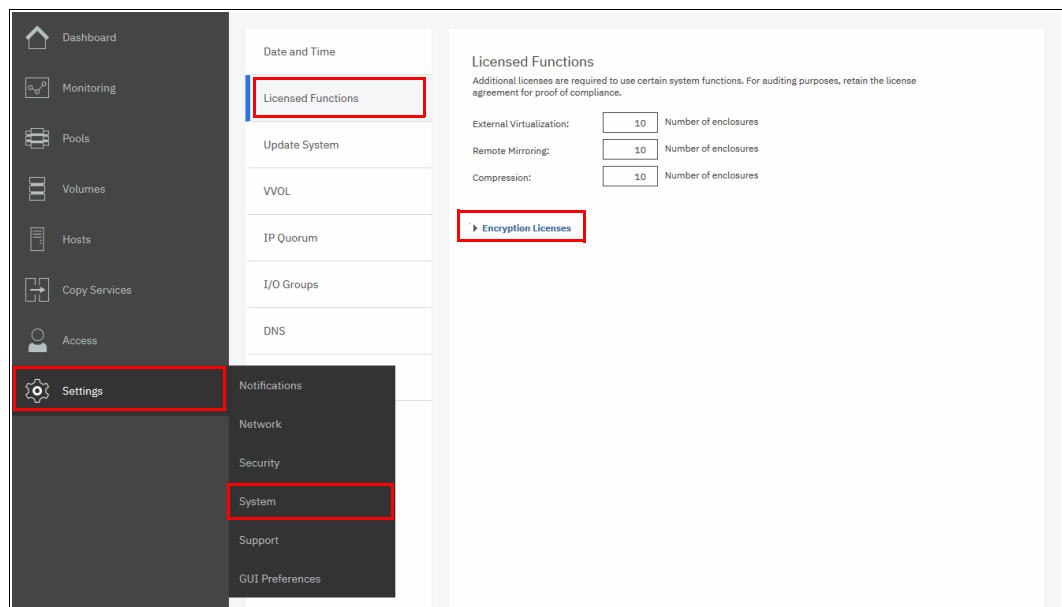


Figure 12-9 Expanding Encryption Licenses section on the Licensed Functions window

- The Encryption Licenses window displays information about your nodes. Right-click the enclosure on which you want to install an encryption license. This action opens a menu with two license activation options (**Activate License Automatically** and **Activate License Manually**), as shown in Figure 12-10. Use either option to activate encryption. For more information about how to complete an automatic activation process, see 12.4.4, “Activate the license automatically” on page 697. For more information about how to complete a manual activation process, see 12.4.5, “Manual license activation” on page 701.

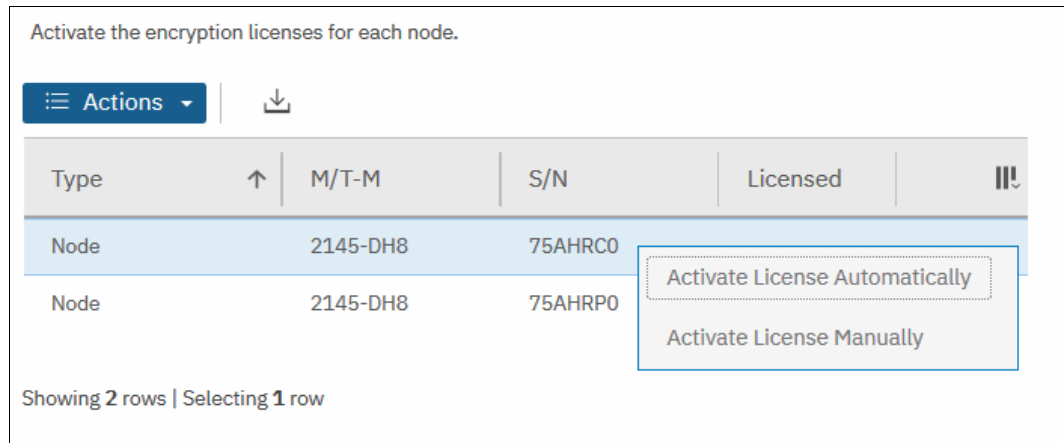


Figure 12-10 Select the node on which you want to enable the encryption

After either activation process is complete, you can see a green check mark in the column labeled Licensed for the node, as shown in Figure 12-11.

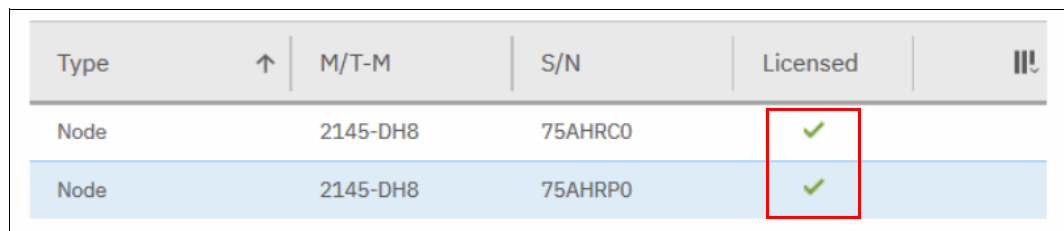


Figure 12-11 Successful encryption license activation on a running system

### 12.4.4 Activate the license automatically

The automatic license activation is the faster method to activate the encryption license for IBM Spectrum Virtualize. You need the authorization code and the workstation that is used to access the GUI that can access the external network.

**Note:** The PC that was used to connect to the GUI and activate the license must connect to the internet.

To activate the encryption license for a node automatically, complete the following steps:

- Select **Activate License Automatically** to open the Activate License Automatically window, as shown in Figure 12-12 on page 698.

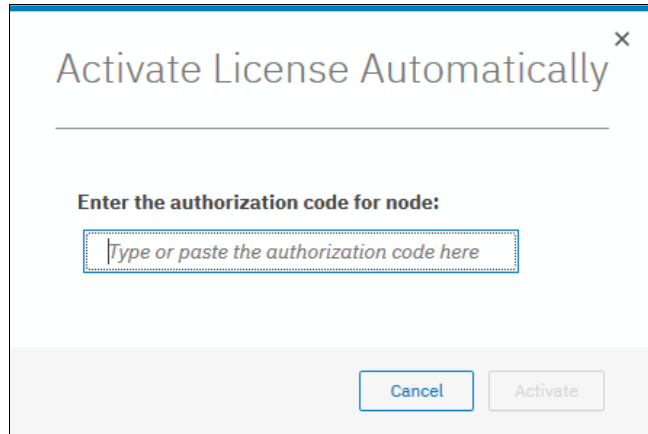


Figure 12-12 Encryption license Activate License Automatically window

2. Enter the authorization code that is specific to the node that you selected, as shown in Figure 12-13. You can now click **Activate**.

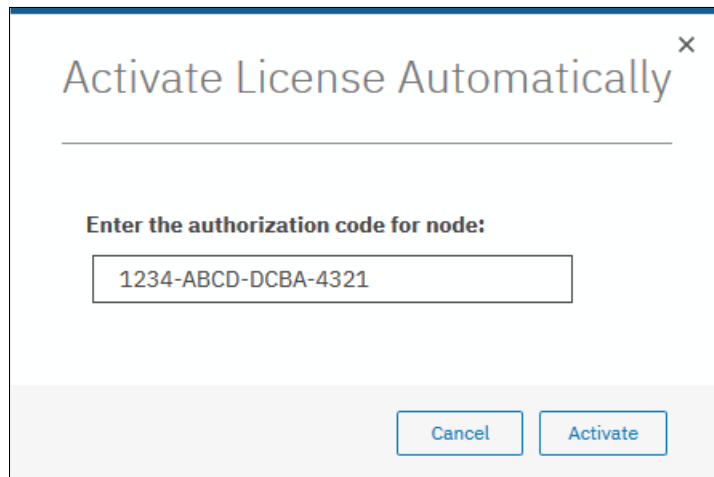


Figure 12-13 Entering an authorization code

The system connects to IBM to verify the authorization code and retrieve the license key. Figure 12-14 on page 699 shows a window that is displayed during this connection. If everything works correctly, the procedure takes less than a minute.

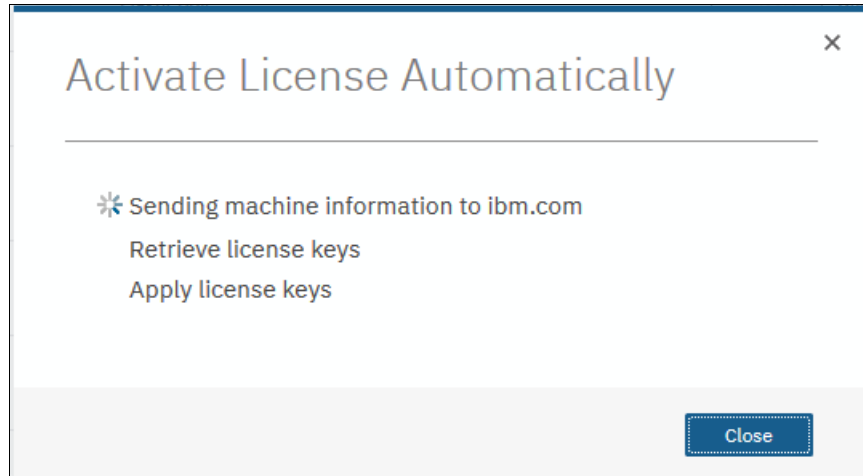


Figure 12-14 Activating encryption

After the license key is retrieved, it is automatically applied, as shown in Figure 12-15.

Type	↑	M/T-M	S/N	Licensed	!!!
Node		2145-DH8	75AHRC0	✓	
Node		2145-DH8	75AHRP0	✓	

Figure 12-15 Successful encryption license activation

### Problems with automatic license activation

If connections problems occur with the automatic license activation procedure, the system times out after 3 minutes with an error.

Check whether the PC that is used to connect to the IBM SAN Volume Controller GUI and activate the license can access the internet. If you cannot complete the automatic activation procedure, use the manual activation procedure that is described in 12.4.5, “Manual license activation” on page 701.

Although authorization codes and encryption license keys use the same format (four groups of four hexadecimal digits), you can only use each of them in the appropriate activation process. If you use a license key when the system expects an authorization code, the system displays an error message, as shown in Figure 12-16.

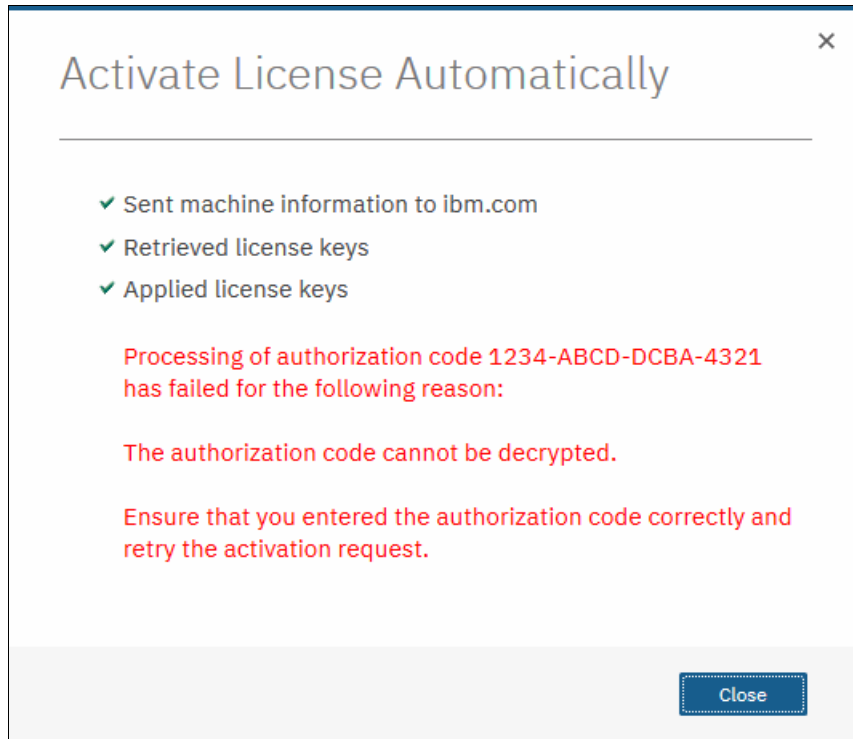


Figure 12-16 Authorization code failure

## 12.4.5 Manual license activation

To manually activate the encryption license for a node, complete the following steps:

1. Select **Activate License Manually** to open the Manual Activation window, as shown in Figure 12-17.

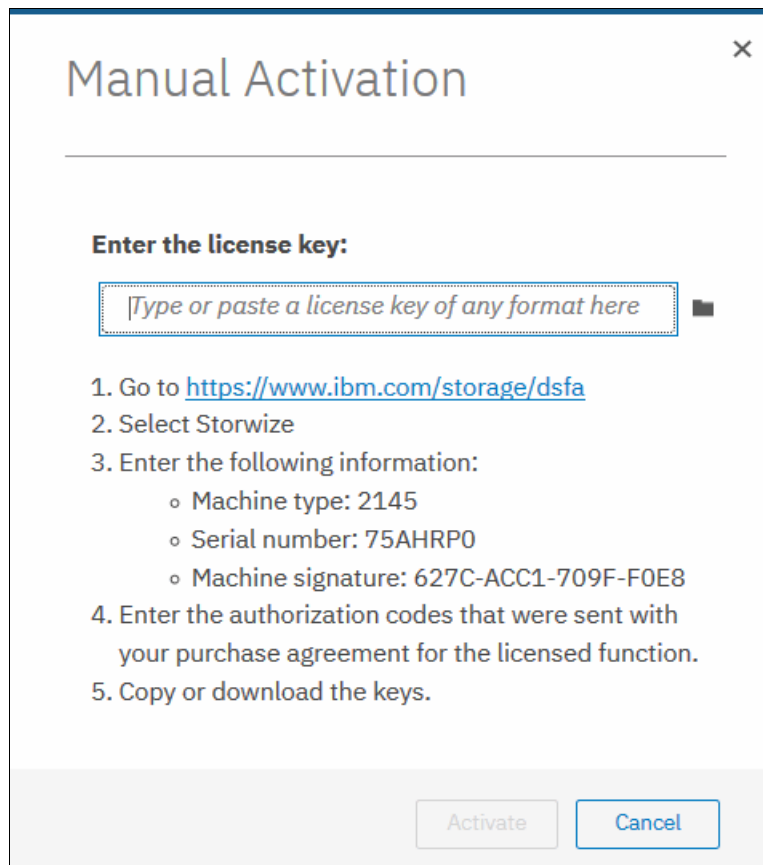


Figure 12-17 Manual encryption license activation window

2. If you have not done so already, obtain the encryption license for the node. The information that is required to obtain the encryption license is displayed in the Manual Activation window. Use this data to follow the instructions in 12.4.1, “Obtaining an encryption license” on page 693.

- You can enter the license key by typing it, pasting it, or clicking the folder icon and uploading the license key file to the storage system that was downloaded from DSFA. In Figure 12-18, the sample key is entered. Click **Activate**.

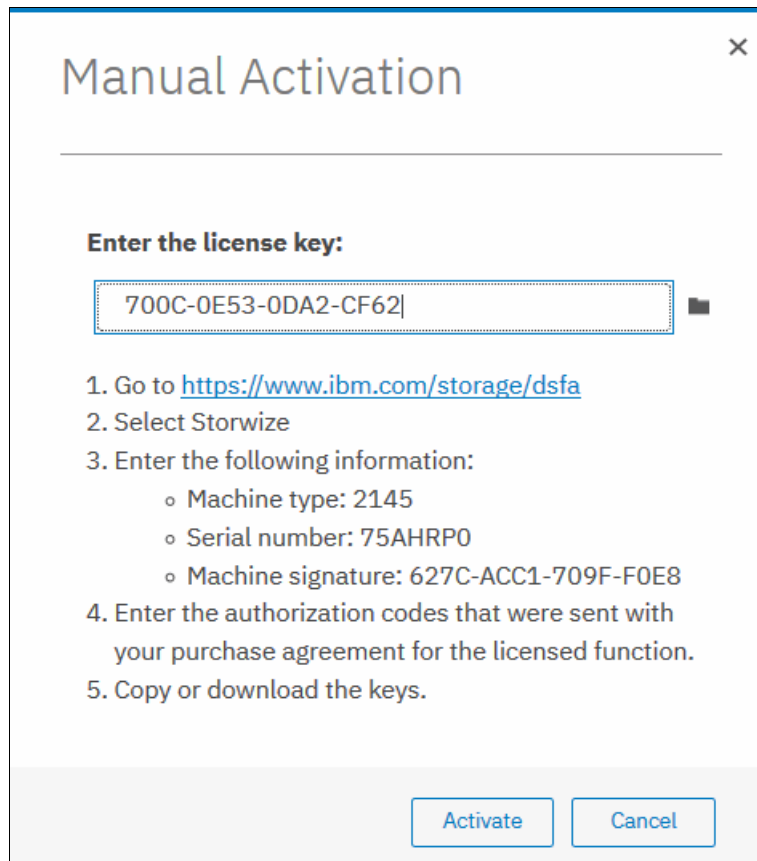


Figure 12-18 Entering an encryption license key

After the task completes successfully, the GUI shows that encryption is licensed for the specified node, as shown in Figure 12-19.

Type	↑	M/T-M	S/N	Licensed	⋮
Node		2145-DH8	75AHRP0	✓	
Node		2145-DH8	75AHRP0	✓	

Figure 12-19 Successful encryption license activation



## Problems with manual license activation

Although authorization codes and encryption license keys use the same format (four groups of four hexadecimal digits), you can only use each of them in the suitable activation process. If you use an authorization code when the system expects a license key, the system displays an error message, as shown in Figure 12-20.

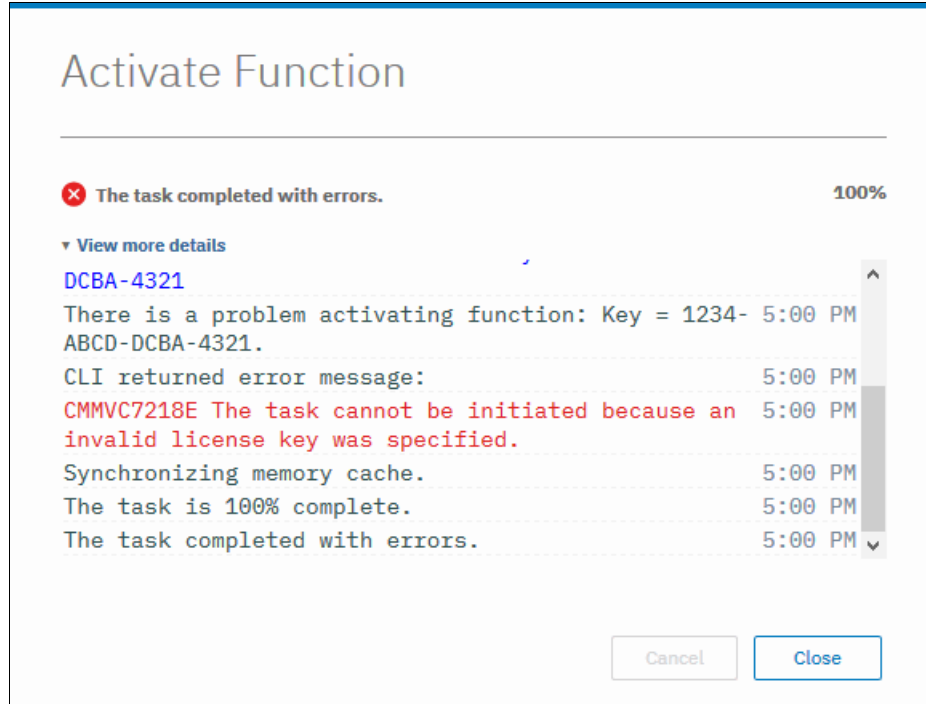


Figure 12-20 License key failure

## 12.5 Enabling encryption

This section describes the process to create and store system master access key copies, also referred to as *encryption keys*. These keys can be stored on any or both of two key providers: USB flash drives or a key server.

Two types of key servers are supported by IBM Spectrum Virtualize:

- ▶ IBM Security Key Lifecycle Manager (SKLM), which was introduced in IBM Spectrum Virtualize V7.8.
- ▶ Gemalto SafeNet KeySecure, introduced in IBM Spectrum Virtualize V8.2.

For more information about supported key servers, see [this IBM Support web page](#).

IBM Spectrum Virtualize code V8.1 introduced the ability to define up to four encryption key servers, which is a preferred configuration because it increases key provider availability. In this version, support for simultaneous use of both USB flash drives and key server was added.

Organizations that use encryption key management servers might consider parallel use of USB flash drives as a backup solution. During normal operation, such drives can be disconnected and stored in a secure location. However, during a catastrophic loss of encryption servers, the USB drives can still be used to unlock the encrypted storage.

The key server and USB flash drive characteristics that are described next might help you to choose the type of encryption key provider that you want to use.

Key servers can have the following characteristics:

- ▶ Physical access to the system is not required to perform a rekey operation.
- ▶ Support for businesses that have security requirements that preclude use of USB ports.
- ▶ Possibility to use hardware security modules (HSMs) for encryption key generation.
- ▶ Ability to replicate keys between servers and perform automatic backups.
- ▶ Implementations follow an open standard (Key Management Interoperability Protocol [KMIP]) that aids in interoperability.
- ▶ Ability to audit operations related to key management.
- ▶ Ability to separately manage encryption keys and physical access to storage systems.

USB flash drives have the following characteristics:

- ▶ Physical access to the system might be required to process a rekey operation.
- ▶ No moving parts with almost no read or write operations to the USB flash drive.
- ▶ Inexpensive to maintain and use.
- ▶ Convenient and easy to have multiple identical USB flash drives available as backups.

**Important:** Maintaining confidentiality of the encrypted data hinges on security of the encryption keys. Pay special attention to ensure secure creation, management, and storage of the encryption keys.

## 12.5.1 Starting the Enable Encryption wizard

After the license activation step is successfully completed on IBM SAN Volume Controller nodes, you can now enable encryption. You can enable encryption after completion of the initial system setup by using the GUI or CLI.

The GUI can be used in two ways to start the Enable Encryption wizard. The first method is by clicking **Run Task** that is next to Enable Encryption on the Suggested Tasks window, as shown in Figure 12-21.

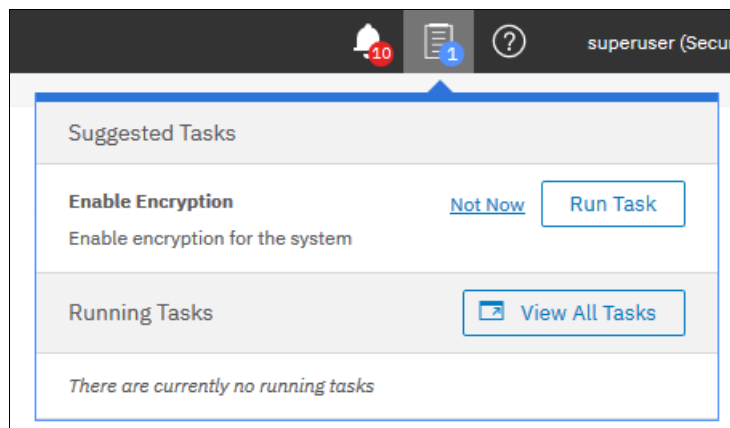


Figure 12-21 Enable Encryption from the Suggested Tasks window

You can also click **Settings** → **Security** → **Encryption** and then, click **Enable Encryption**, as shown in Figure 12-22.

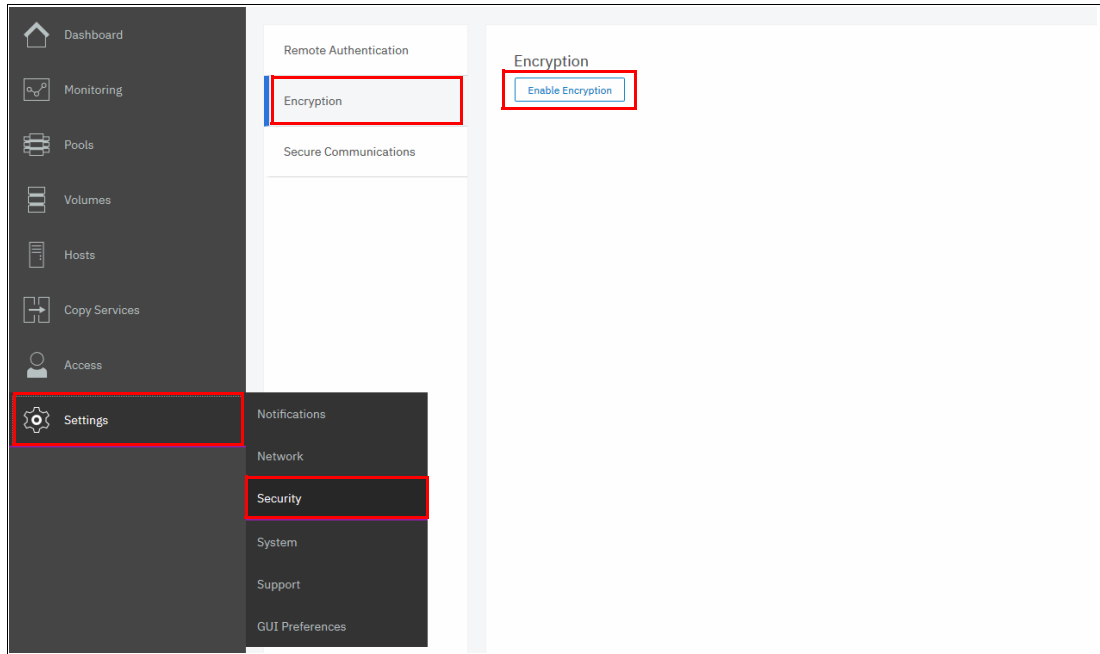


Figure 12-22 Enable Encryption from the Security pane

The Enable Encryption wizard starts by prompting you to select the encryption key provider to use for storing the encryption keys, as shown in Figure 12-23. You can enable either or both providers.

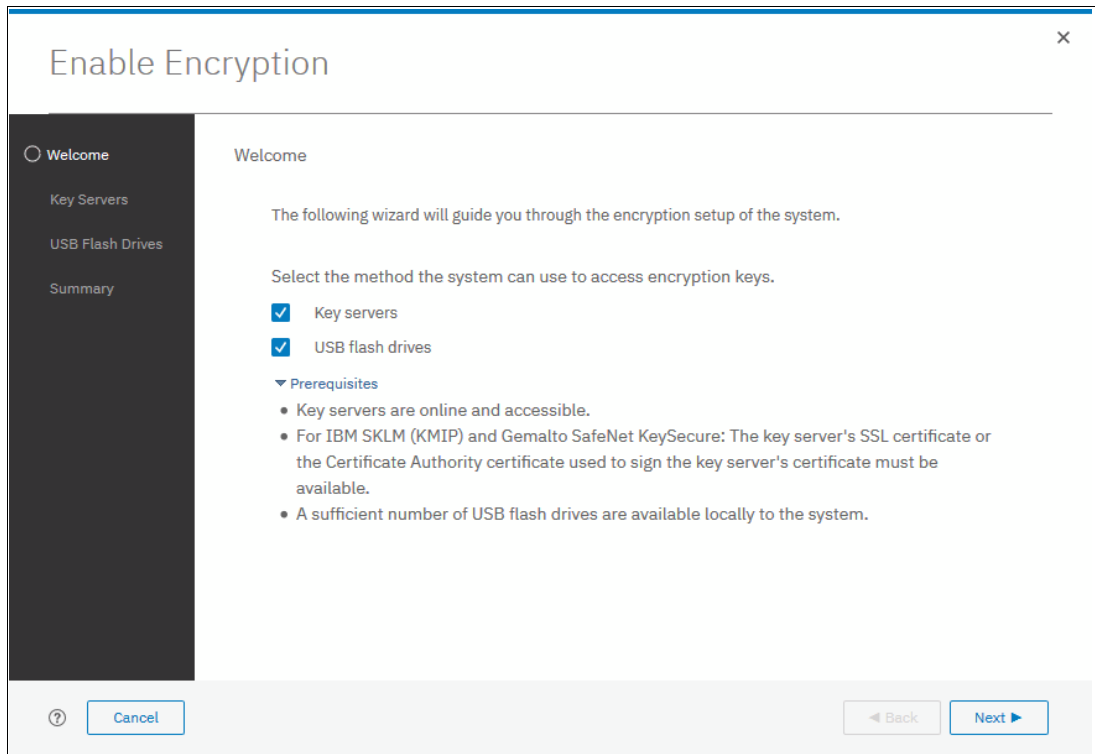


Figure 12-23 Enable Encryption wizard Welcome window

The next section presents a scenario in which both encryption key providers are enabled at the same time. For more information about how to enable encryption by using only USB flash drives, see 12.5.2, “Enabling encryption by using USB flash drives” on page 706.

For more information about how to enable encryption by using key servers as the sole encryption key provider, see 12.5.3, “Enabling encryption by using key servers” on page 711.

## 12.5.2 Enabling encryption by using USB flash drives

**Note:** The system needs at least three USB flash drives to be present before you can enable encryption by using this encryption key provider. IBM USB flash drives are preferred and can be obtained from IBM with the feature name Encryption USB Flash Drives (Four Pack). Other flash drives might also work. You can use any USB ports in any node of the cluster.

Using USB flash drives as the encryption key provider requires a minimum of three USB flash drives to store the generated encryption keys. Because the system attempts to write the encryption keys to any USB key that is inserted into a node port, it is critical to maintain physical security of the system during this procedure.

While the system enables encryption, you are prompted to insert USB flash drives into the system. The system generates and copies the encryption keys to all available USB flash drives.

Ensure that each copy of the encryption key is valid before you write any user data to the system. The system validates any key material on a USB flash drive when it is inserted into the canister. If the key material is invalid, the system logs an error.

If the USB flash drive is unusable or fails, the system does not display it as output. Figure 12-26 on page 709 shows an example where the system detected and validated three USB flash drives.

If your system is in a secure location with controlled access, one USB flash drive for each canister can remain inserted in the system. If a risk of unauthorized access exists, all USB flash drives with the master access keys must be removed from the system and stored in a secure place.

Securely store all copies of the encryption key. For example, any USB flash drives that are holding an encryption key copy that are not left plugged into the system can be locked in a safe. Similar precautions must be taken to protect any other copies of the encryption key that are stored on other media.

**Notes:** Generally, create at least one extra copy on another USB flash drive for storage in a secure location. You can also copy the encryption key from the USB drive and store the data on other media, which can provide extra resilience and mitigate risk that the USB drives used to store the encryption key come from a faulty batch.

Every encryption key copy must be stored securely to maintain confidentiality of the encrypted data.

A minimum of one USB flash drive with the correct master access key is required to unlock access to encrypted data after a system restart such as a system-wide restart or power loss. No USB flash drive is required during a warm restart, such as a node that is exiting service mode or a single node restart. The data center power-on procedure must ensure that USB flash drives containing encryption keys are plugged into the storage system before it is powered on.

During power-on, insert USB flash drives into the USB ports on two supported canisters to safeguard against failure of a node, node's USB port, or USB flash drive during the power-on procedure.

To enable encryption by using USB flash drives as the only encryption key provider, complete the following steps:

1. In the Enable Encryption wizard Welcome tab, select **USB flash drives** and click **Next**, as shown in Figure 12-24.

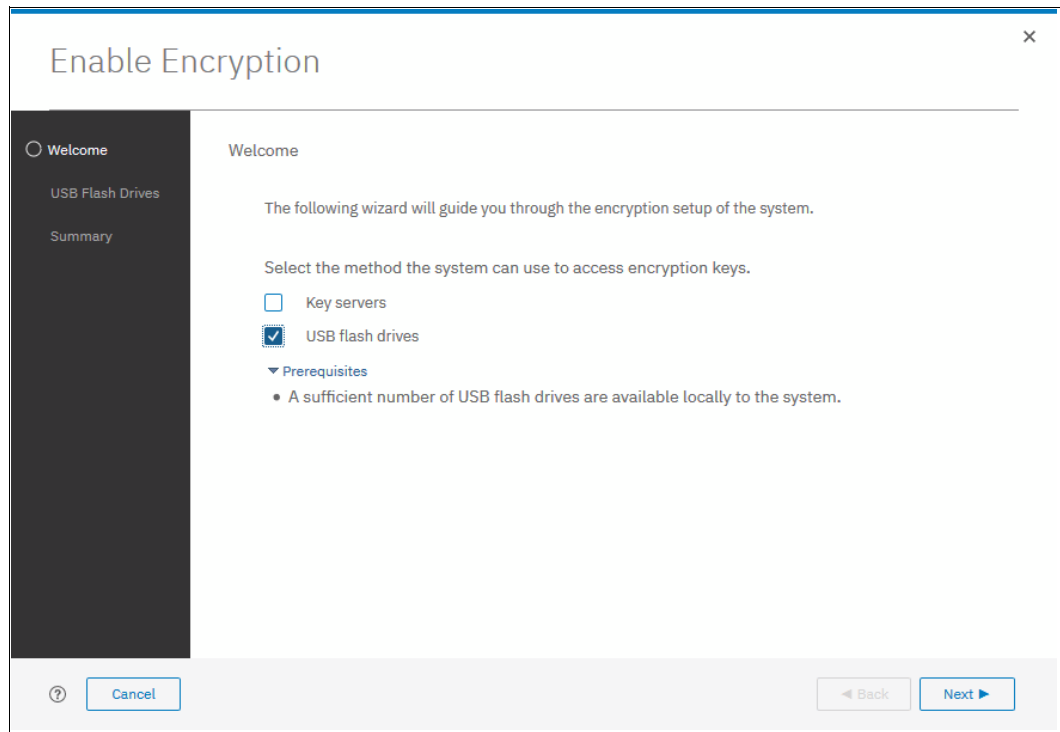


Figure 12-24 Selecting USB flash drives in the Enable Encryption wizard

2. If fewer than three USB flash drives are inserted into the system, you are prompted to insert more drives, as shown in Figure 12-25. The system reports how many more drives must be inserted.

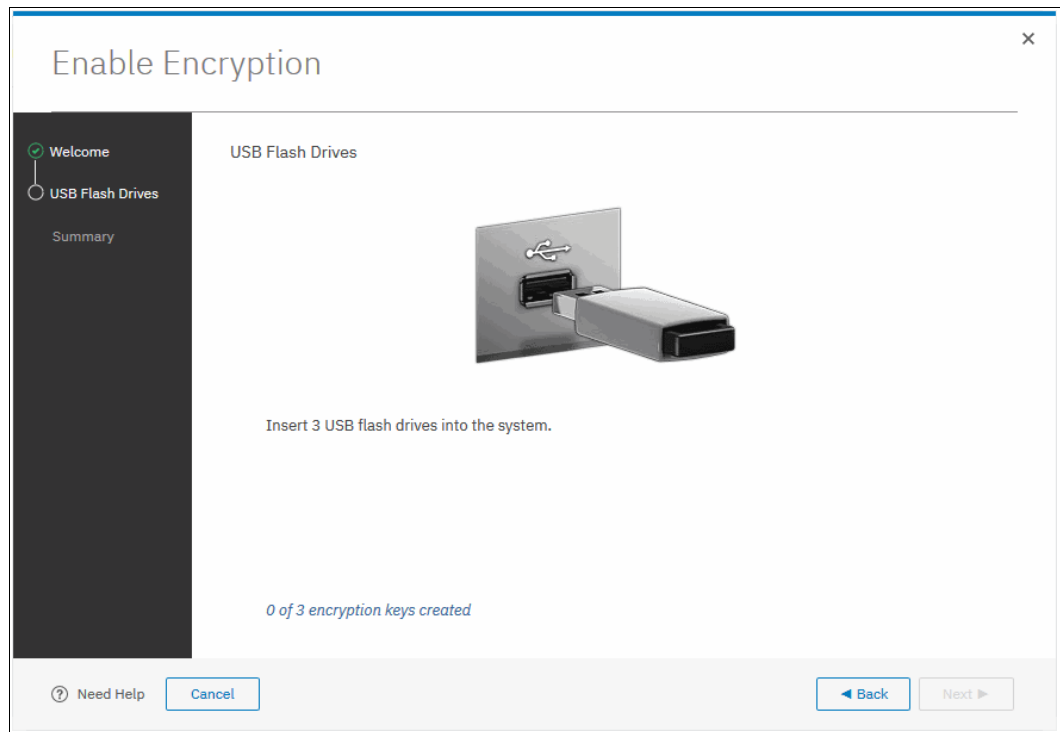


Figure 12-25 Waiting for USB flash drives to be inserted

**Note:** The **Next** option remains disabled until at least three USB flash drives are detected.

3. Insert the USB flash drives into the USB ports as requested.

After the minimum required number of drives is detected, the encryption keys are automatically copied on the USB flash drives, as shown in Figure 12-26.

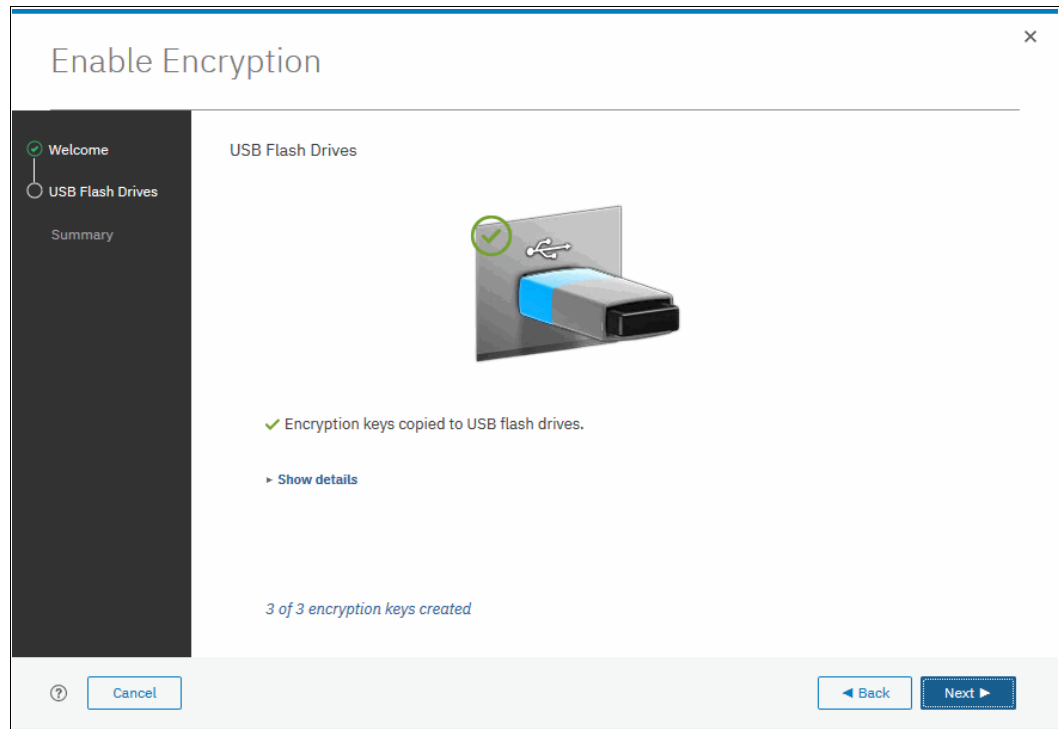


Figure 12-26 Writing the master access key to USB flash drives

You can keep adding USB flash drives or replacing the drives that are plugged in to create new copies. When done, click **Next**.

- The number of keys that were created is shown in the Summary tab, as shown in Figure 12-27. Click **Finish** to finalize the encryption enablement.

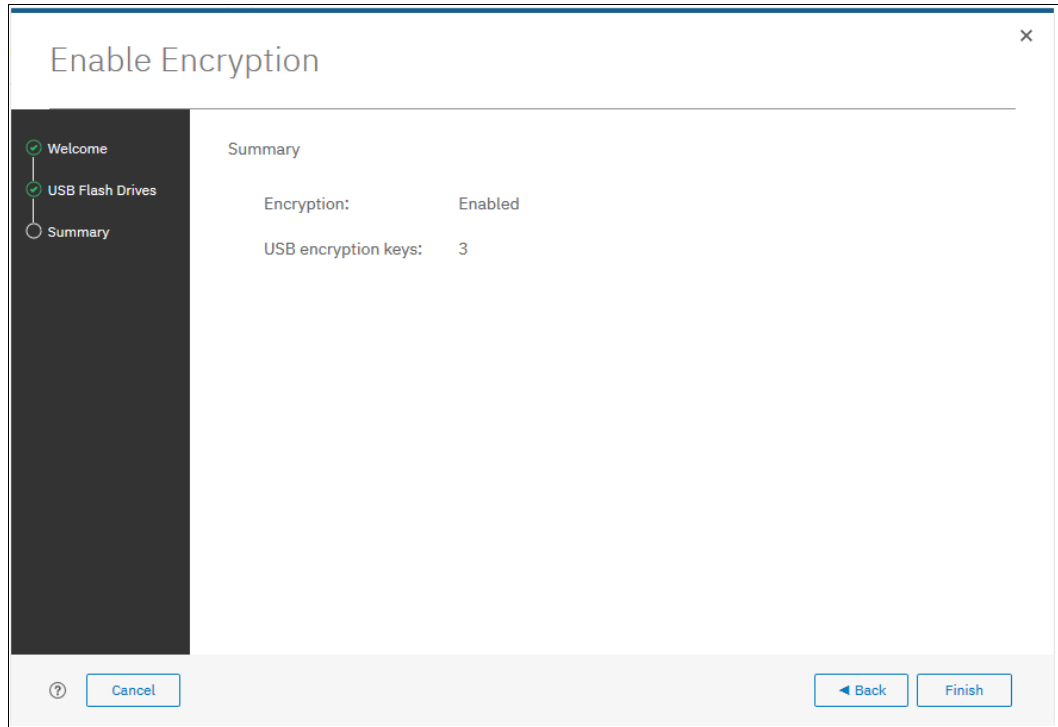


Figure 12-27 Commit the encryption enablement

You receive a message confirming that the encryption is now enabled on the system, as shown in Figure 12-28.



Figure 12-28 Encryption enabled message that uses USB flash drives



5. You can confirm that encryption is enabled and verify which key providers are in use by selecting **Settings** → **Security** → **Encryption**, as shown in Figure 12-29.

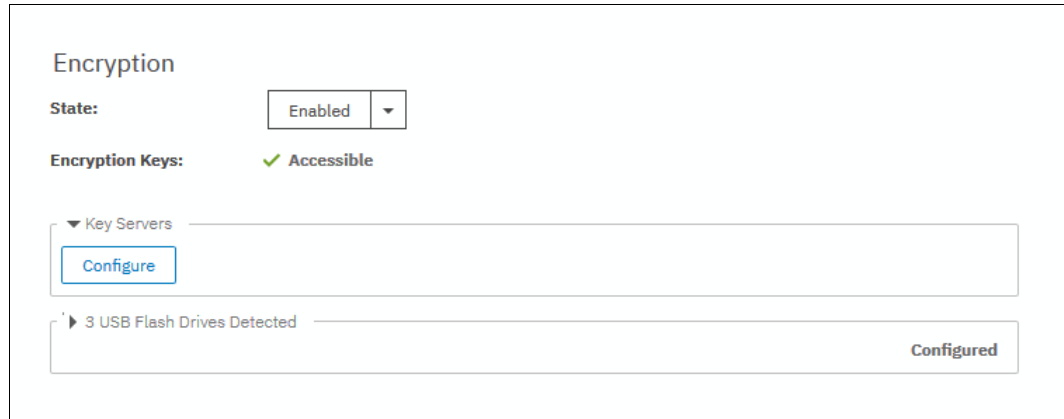


Figure 12-29 Encryption view showing using USB flash drives as the enabled provider

### 12.5.3 Enabling encryption by using key servers

A key server is a centralized system that receives and then distributes encryption keys to its clients, including IBM Spectrum Virtualize systems.

IBM Spectrum Virtualize supports use of the following key servers as encryption key providers:

- ▶ SKLM
- ▶ Gemalto SafeNet KeySecure

**Note:** Support for SKLM was introduced in IBM Spectrum Virtualize V7.8. Support for Gemalto SafeNet KeySecure was introduced in IBM Spectrum Virtualize V8.2.1.

SKLM and SafeNet KeySecure support KMIP, which is a standard for management of cryptographic keys.

**Note:** Make sure that the key management server function is fully independent from encrypted storage that has encryption managed by this key server environment. Failure to observe this requirement might create an encryption deadlock. An encryption deadlock is a situation in which none of key servers in the environment can become operational because some critical part of the data in each server is stored on a storage system that depends on one of the key servers to unlock access to the data.

IBM Spectrum Virtualize code V8.1 and later supports up to four key server objects that are defined in parallel. But, only one key server type (SKLM or KeySecure) can be enabled at one time.

Another characteristic when working with key servers is that it is not possible to migrate from one key server type directly to another. If you want to migrate from one type to another, you first must migrate from your current key server to USB encryption, and then migrate from USB to the other type of key server.

## Enabling encryption by using SKLM

Before you create a key server object in the storage system, the key server must be configured. Ensure that you complete the following tasks on the SKLM server before you enable encryption on the storage system:

- ▶ Configure the SKLM server to use Transport Layer Security version 1.2. The default setting is TLSv1, but IBM Spectrum Virtualize supports only version 1.2. Therefore, set the value of security protocols to SSL\_TLSv2 (which is a set of protocols that includes TLSv1.2) in the SKLM server configuration properties.
- ▶ Ensure that the database service is started automatically on start.
- ▶ Ensure that at least one Secure Sockets Layer (SSL) certificate is available for browser access.
- ▶ Create an IBM Spectrum\_VIRT device group for IBM Spectrum Virtualize systems.

For more information about completing these tasks, see the [SKLM documentation](#) at IBM Knowledge Center.

Access to the key server that is storing the correct master access key is required to enable access to encrypted data in the system after a system restart. System restart can be a system-wide restart or power loss. Access to the key server is not required during a warm restart, such as a node that is exiting service mode or a single node restart.

The data center power-on procedure must ensure key server availability before the storage system that is using encryption is started. If a system with encrypted data is restarted and cannot access the encryption keys, the encrypted storage pools are offline until the encryption keys are detected.

To enable encryption by using an SKLM key server, complete the following steps:

1. Ensure that service IPs are configured on all your nodes.
2. In the Enable Encryption wizard Welcome tab, select **Key servers** and click **Next**, as shown in Figure 12-30.

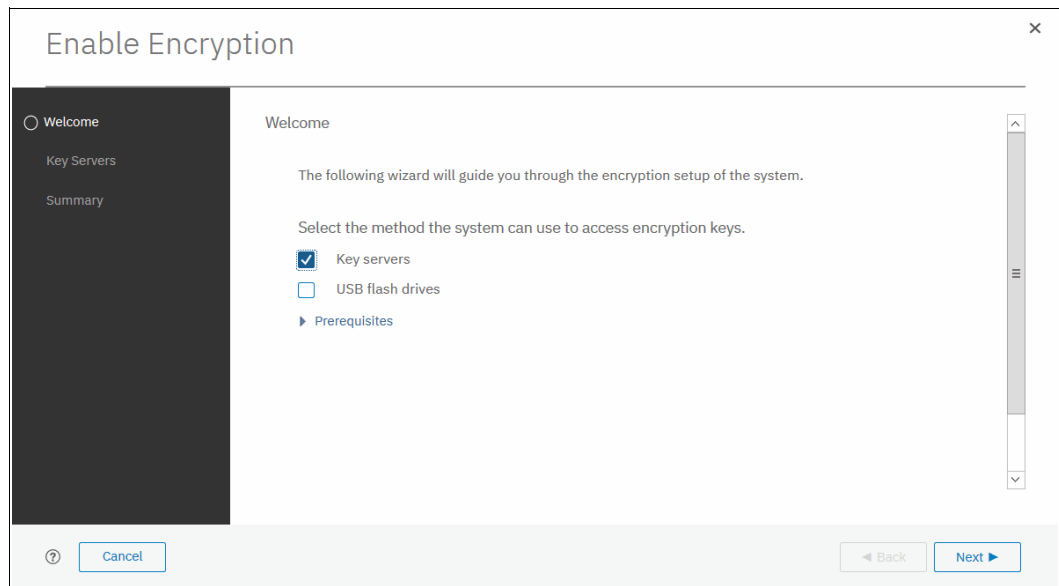


Figure 12-30 Selecting Key server as the only provider in the Enable Encryption wizard

3. Select **IBM SKLM (with KMIP)** as the key server type, as shown in Figure 12-31.

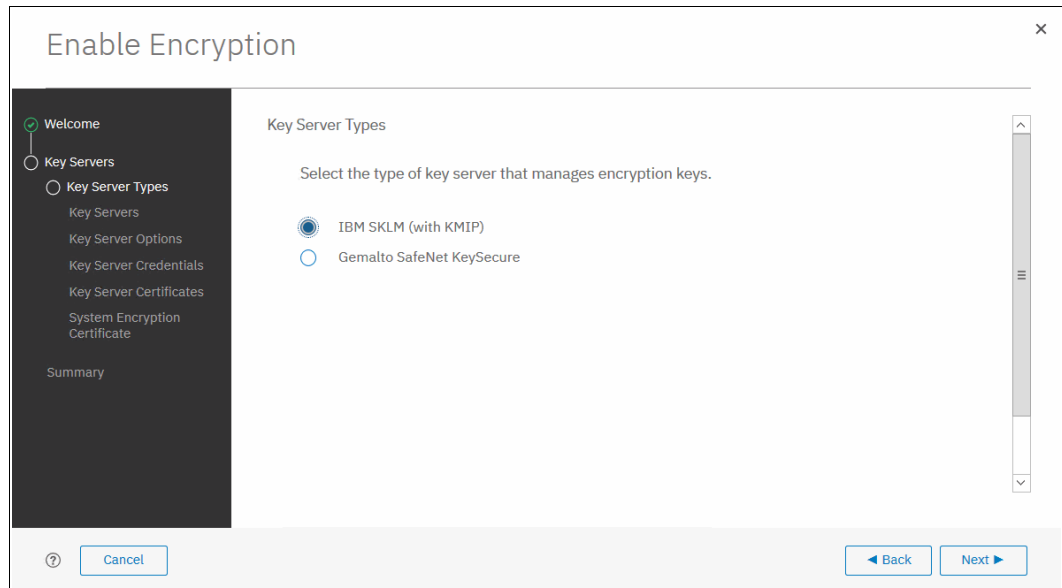


Figure 12-31 Selecting SKLM as key server type

4. The wizard moves to the Key Servers tab, as shown in Figure 12-32 on page 714. Enter the name and IP address of the key servers. The first key server that is specified must be the primary SKLM key server.

**Note:** The supported versions of SKLM (up to V3.0, which was the latest code version available at the time of this writing) differentiate between the primary and secondary key server role. The Primary SKLM server as defined on the Key Servers window of the Enable Encryption wizard must be the server defined as the primary by SKLM administrators.

The key server name serves only as a label. Only the provided IP address is used to contact the server. If the key server's TCP port number differs from the default value for the KMIP protocol (that is, 5696), enter the port number. An example of a complete primary SKLM configuration is shown in Figure 12-32.

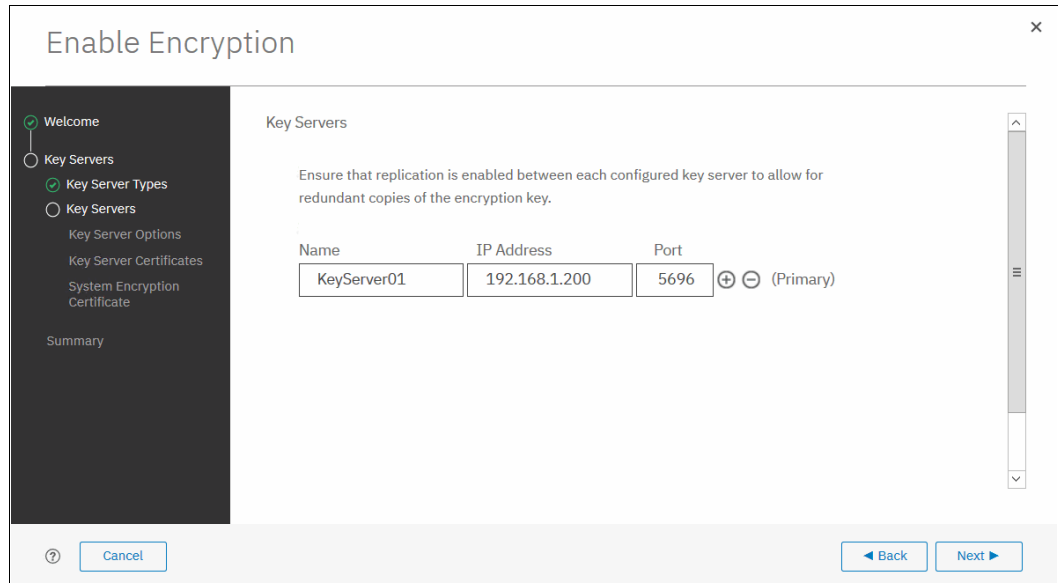


Figure 12-32 Configuration of the primary SKLM server

- If you want to add secondary SKLM servers, click the “+” symbol and enter the data for secondary SKLM servers, as shown on Figure 12-33. You can define up to four SKLM servers. Click **Next** when you are done.

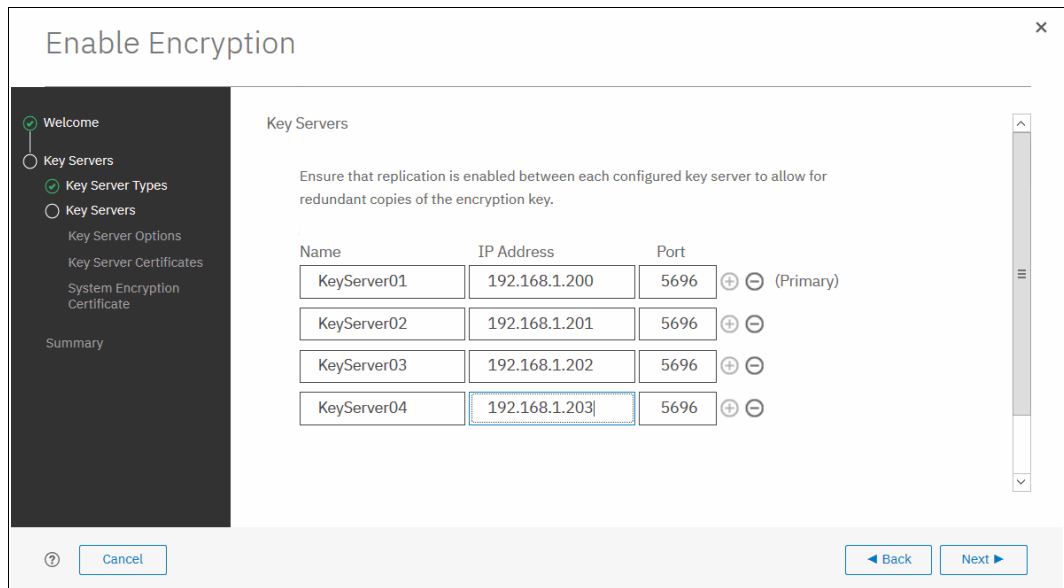


Figure 12-33 Configuring multiple SKLM servers

- The next window in the wizard is a reminder that IBM Spectrum\_VIRT device group that is dedicated for IBM Spectrum Virtualize systems must exist on the SKLM key servers. Make sure that this device group exists and click **Next** to continue, as shown in Figure 12-34.

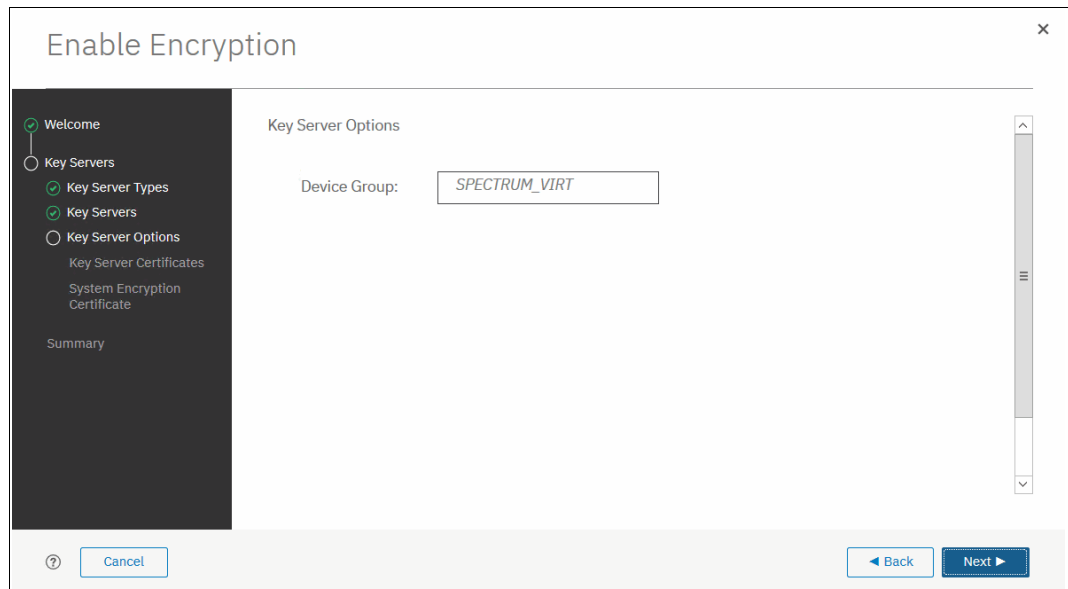


Figure 12-34 Checking key server device group

- Enable secure communication between the IBM Spectrum Virtualize system and the SKLM key servers by uploading the key server certificate from a trusted third-party certificate authority (CA) or by using a self-signed certificate. The self-signed certificate can be obtained from each of key servers directly.

After uploading any of the certificates in the window that is shown in Figure 12-35, click **Next**.

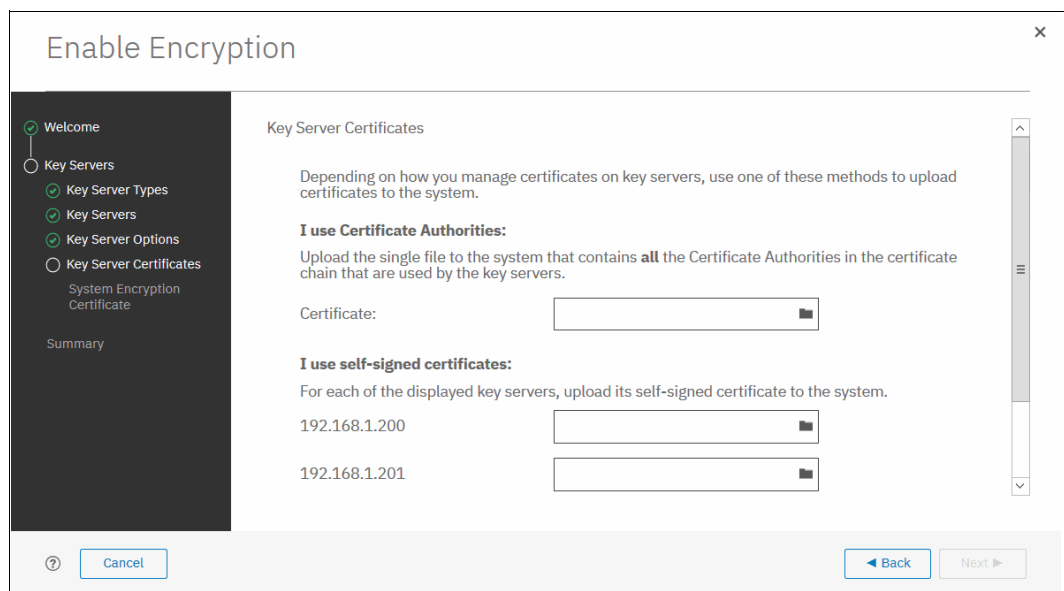


Figure 12-35 Uploading key servers or certificate authority SSL certificate

- Configure the SKLM key server to trust the public key certificate of the IBM Spectrum Virtualize system. You can download the IBM Spectrum Virtualize system public SSL certificate by clicking **Export Public Key**, as shown in Figure 12-36. Install this certificate in the SKLM key server in the IBM Spectrum\_VIRT device group.

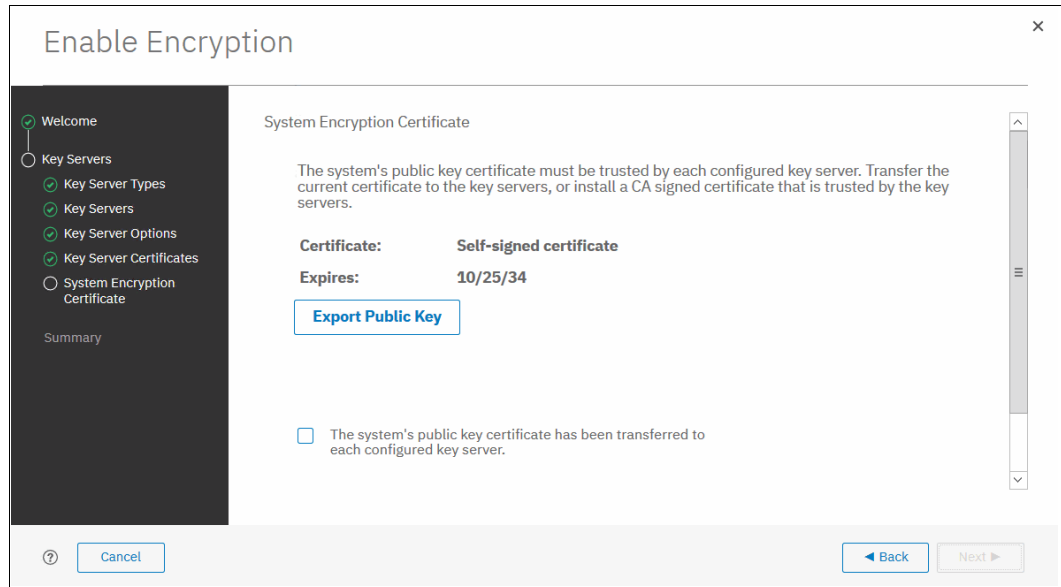


Figure 12-36 Downloading the IBM Spectrum Virtualize SSL certificate

- When the IBM Spectrum Virtualize system public key certificate is installed on the SKLM key servers, acknowledge this installation by selecting the box that is indicated in Figure 12-37 and click **Next**.

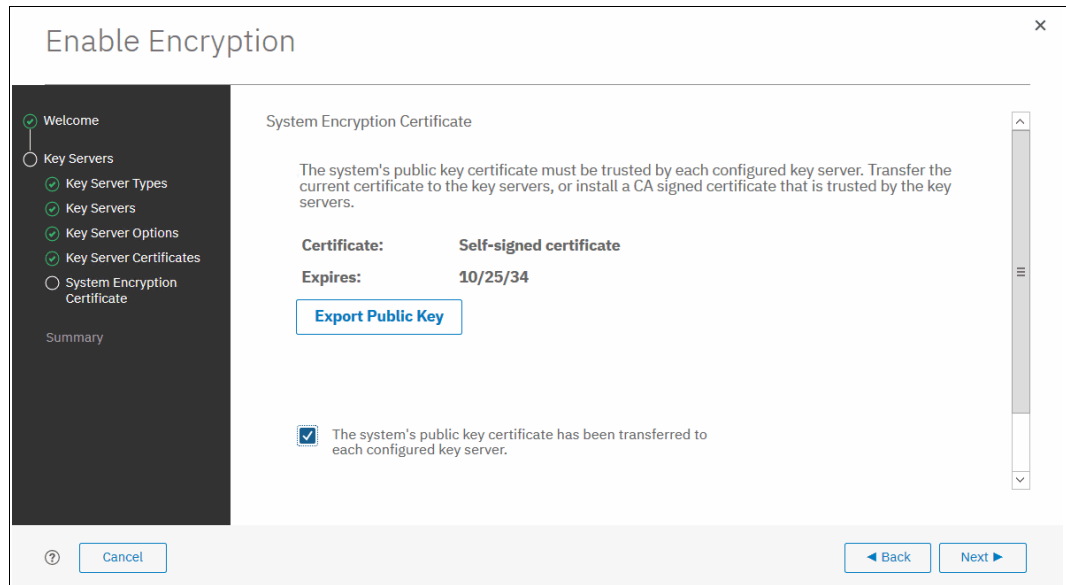


Figure 12-37 Acknowledge IBM Spectrum Virtualize public key certificate transfer

10. The key server configuration is shown in the Summary tab, as shown in Figure 12-38. Click **Finish** to create the key server object and finalize the encryption enablement.

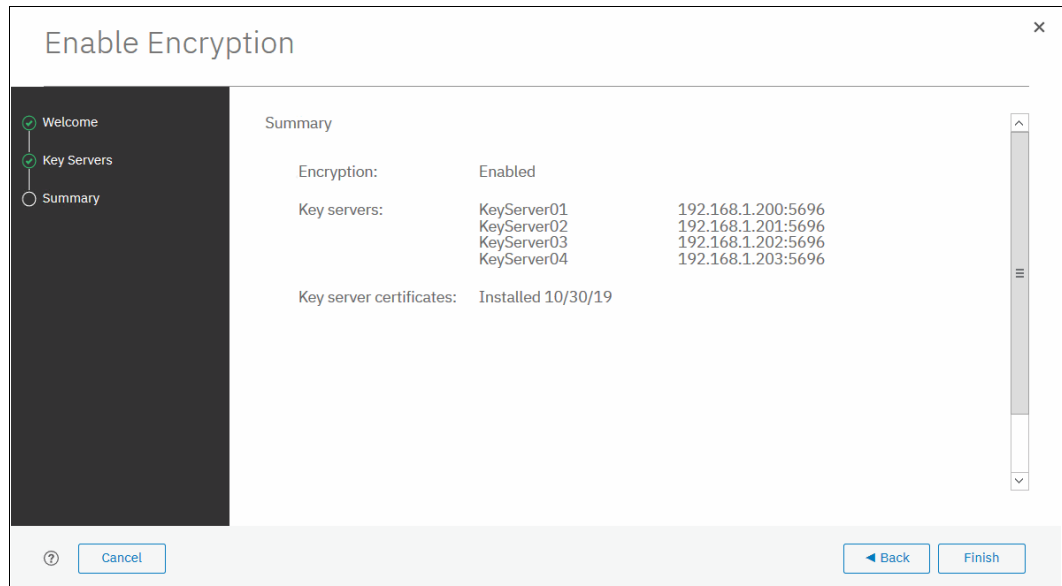


Figure 12-38 Finish enabling encryption using SKLM key servers

11. If no errors occur while the key server object is created, you receive a message that confirms that the encryption is now enabled on the system, as shown in Figure 12-39. Click **Close**.

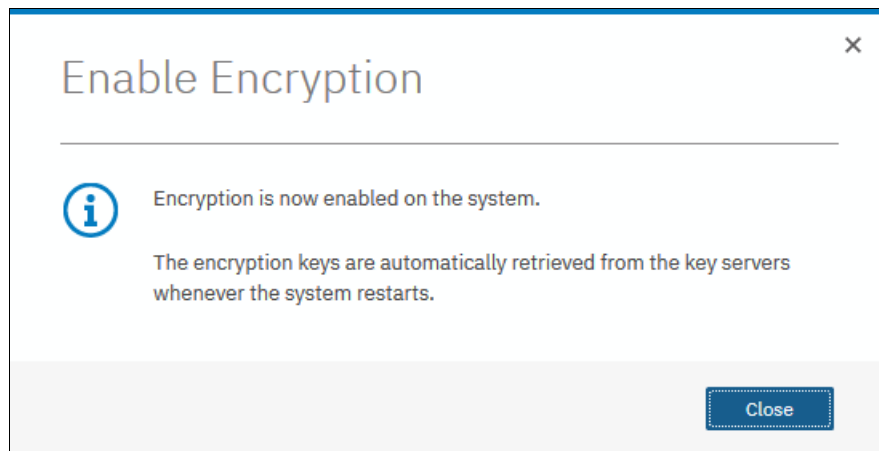


Figure 12-39 Encryption enabled message using an SKLM key server

12. Confirm that encryption is enabled in **Settings** → **Security** → **Encryption**, as shown in Figure 12-40. The Online state indicates which SKLM servers are detected as available by the system.

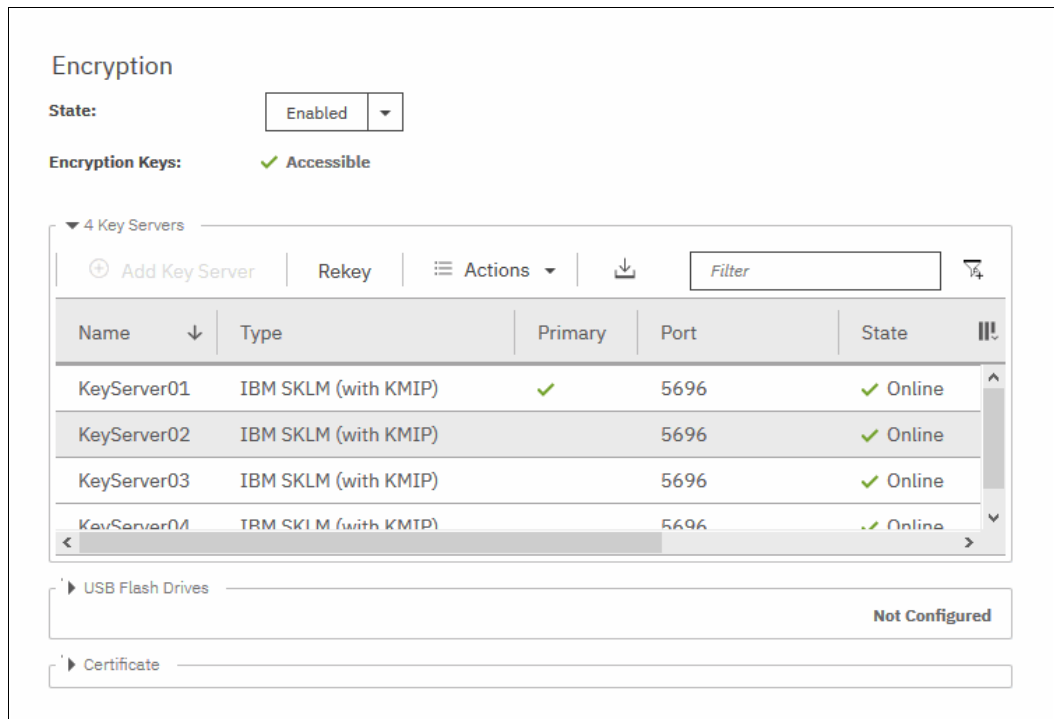


Figure 12-40 Encryption enabled with only SKLM servers as encryption key providers

## Enabling encryption by using SafeNet KeySecure

IBM Spectrum Virtualize V8.2.1 introduced support for Gemalto SafeNet KeySecure, which is a third-party key management server. It can be used as an alternative to SKLM.

IBM Spectrum Virtualize supports Gemalto SafeNet KeySecure version 8.3.0 and later, and by using only KMIP protocol. It is possible to configure up to four SafeNet KeySecure servers in IBM Spectrum Virtualize for redundancy, and they can coexist with USB flash drive encryption.

It is not possible to have both SafeNet KeySecure and SKLM key servers configured at the same time in IBM Spectrum Virtualize. It is also not possible to migrate directly from one type of key server to another (from SKLM to SafeNet KeySecure or vice versa). If you want to migrate from one type to another, first migrate to USB flash drives encryption, and then, migrate to the other type of key servers.

KeySecure uses an active-active clustered model. All changes to one key server are instantly propagated to all other servers in the cluster.

Although KeySecure uses KMIP protocol as IBM SKLM does, an option is available to configure the user name and password for IBM Spectrum Virtualize and KeySecure server authentication, which is not possible when the configuration is performed with SKLM.

The certificate for client authentication in SafeNet KeySecure can be self-signed or signed by a CA.



To enable encryption in IBM Spectrum Virtualize by using a Gemalto SafeNet KeySecure key server, complete the following steps:

1. Ensure that the service IPs are configured on all of your nodes.
2. In the Enable Encryption wizard Welcome tab, select **Key servers** and click **Next**, as shown in Figure 12-41.

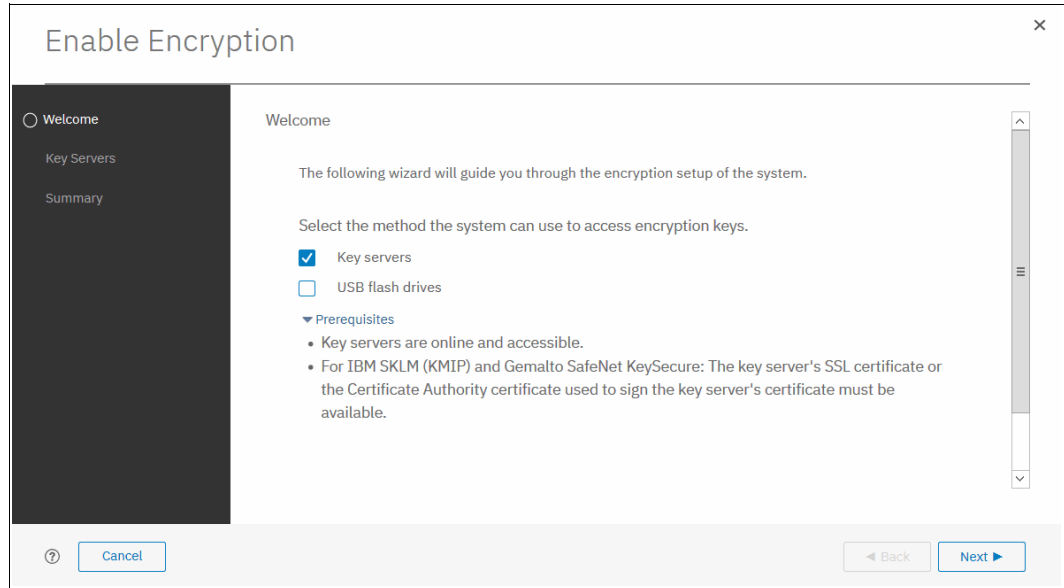


Figure 12-41 Selecting Key servers as the only provider in the Enable Encryption wizard

3. In the next window, you can choose between IBM SKLM or Gemalto SefeNet KeySecure server types, as shown in Figure 12-42. Select **Gemalto SefeNet KeySecure** and click **Next**.

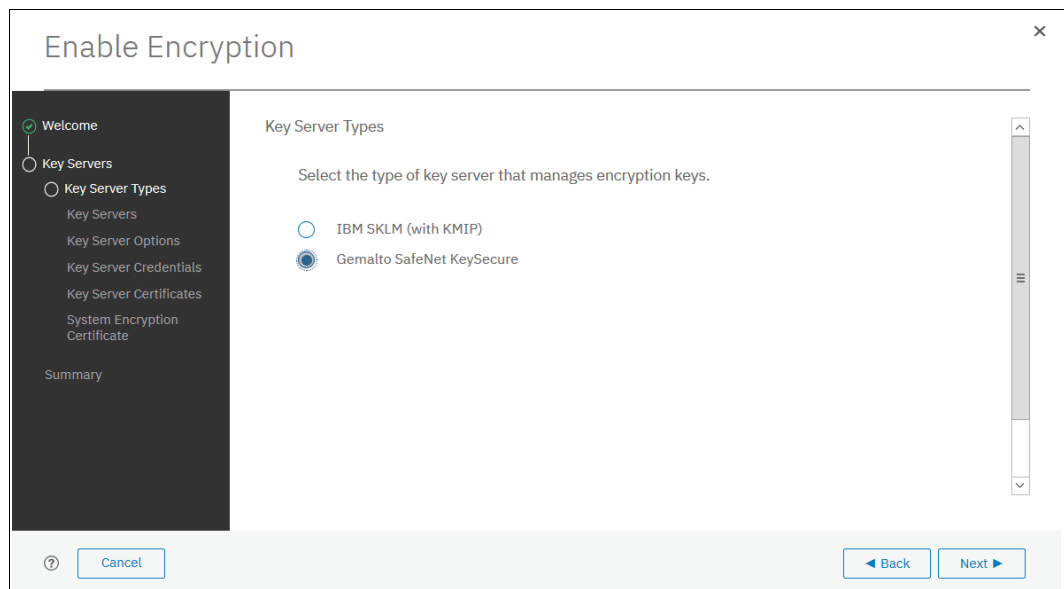


Figure 12-42 Selecting Gemalto SafeNet KeySecure as key server type

4. Add up to four SafeNet KeySecure servers in the next wizard window, as shown in Figure 12-43. For each key server, enter the name, IP address, and TCP port for KMIP protocol (default value is 5696). The server name is only a label, so it does not need to be the real host name of the server.

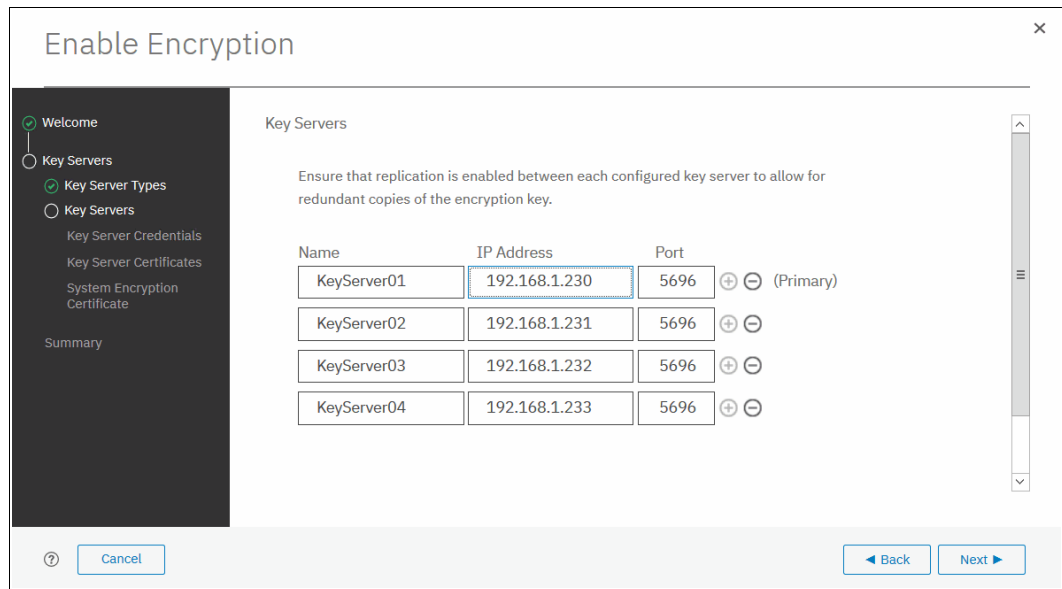


Figure 12-43 Configuring multiple SafeNet KeySecure servers

Although Gemalto SafeNet KeySecure uses an active-active clustered model, IBM Spectrum Virtualize asks for a primary key server. The primary key server represents only the KeySecure server that is used for key create and rekey operations. Therefore, any of the clustered key servers can be selected as the primary.

Selecting a primary key server is beneficial for load balancing. Any four key servers can be used to retrieve the master key.

- The next window in the wizard prompts for key servers credentials (user name and password), as shown in Figure 12-44. This setting is optional because it depends on how SafeNet KeySecure servers are configured.

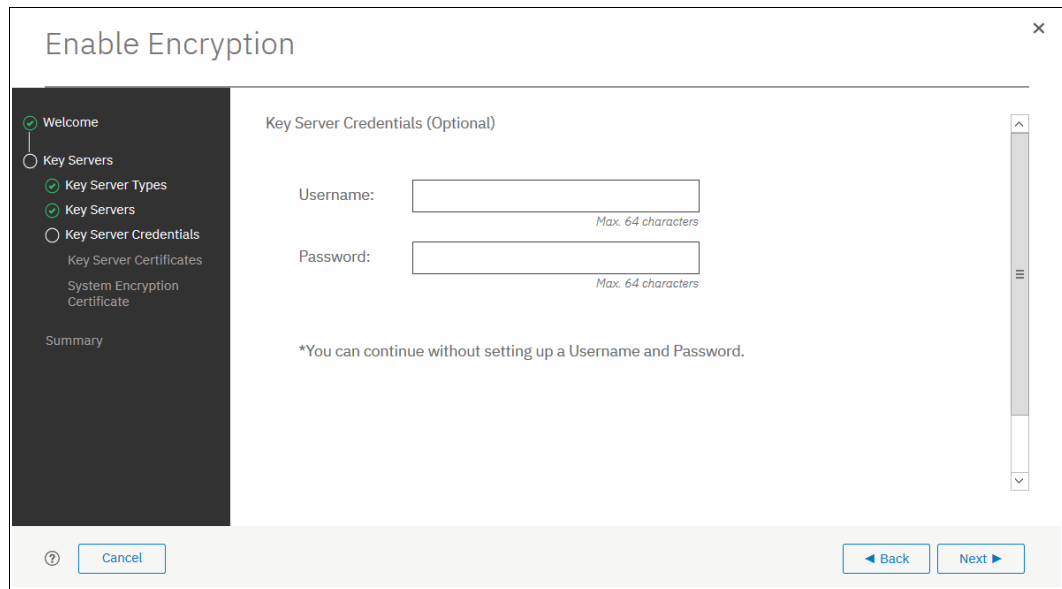


Figure 12-44 Key server credentials input (optional)

- Enable secure communication between the IBM Spectrum Virtualize system and the SafeNet KeySecure key servers by uploading the key server certificate from a trusted third-party CA or by using a self-signed certificate. The self-signed certificate can be obtained from each of key servers directly. After uploading any of the certificates in the window that is shown in Figure 12-45, click **Next**.

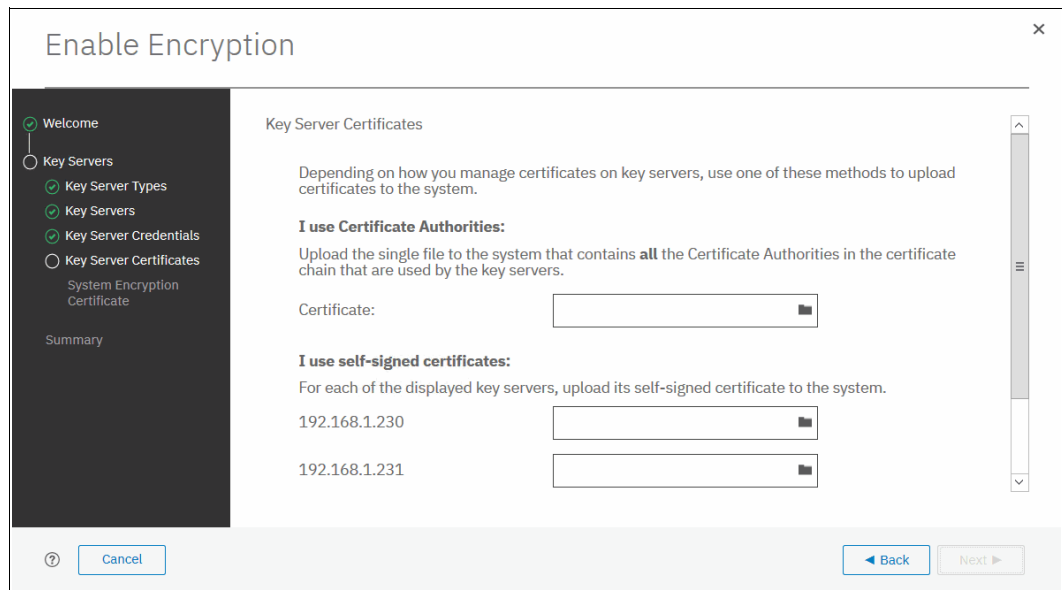


Figure 12-45 Uploading SafeNet KeySecure key servers certificate

- Configure the SafeNet KeySecure key servers to trust the public key certificate of the IBM Spectrum Virtualize system. You can download the IBM Spectrum Virtualize system public SSL certificate by clicking **Export Public Key**, as shown in Figure 12-45 on page 721. After adding the public key certificate to the key servers, select the check box and click **Next**.

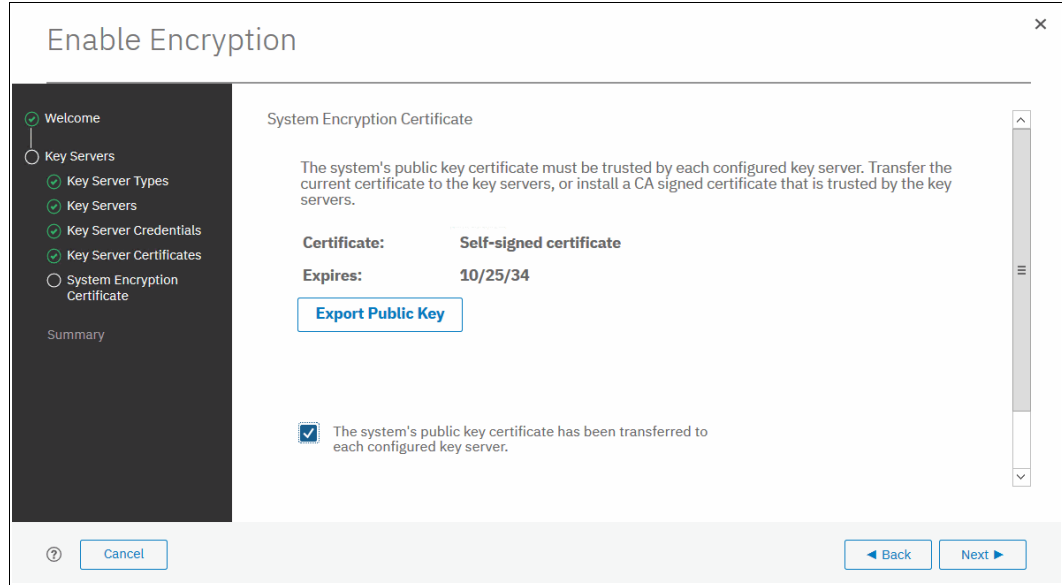


Figure 12-46 Downloading the IBM Spectrum Virtualize SSL certificate

- The key server configuration is shown in the Summary tab, as shown in Figure 12-47. Click **Finish** to create the key server object and finalize the encryption enablement.

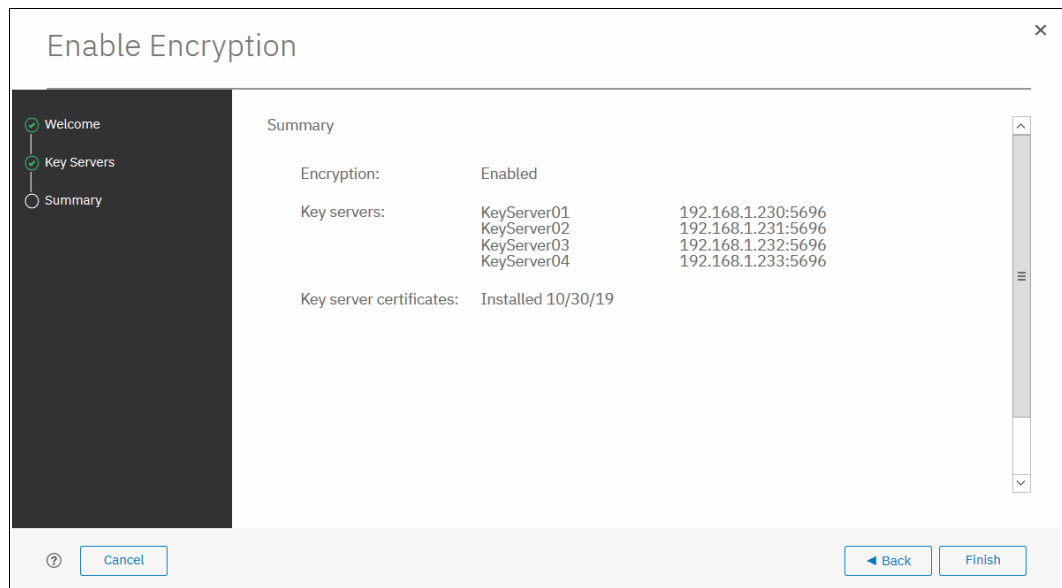


Figure 12-47 Finish enabling encryption using SafeNet KeySecure key servers

- If no errors occurred while creating the key server object, you receive a message that confirms that the encryption is now enabled on the system, as shown in Figure 12-48. Click **Close**.

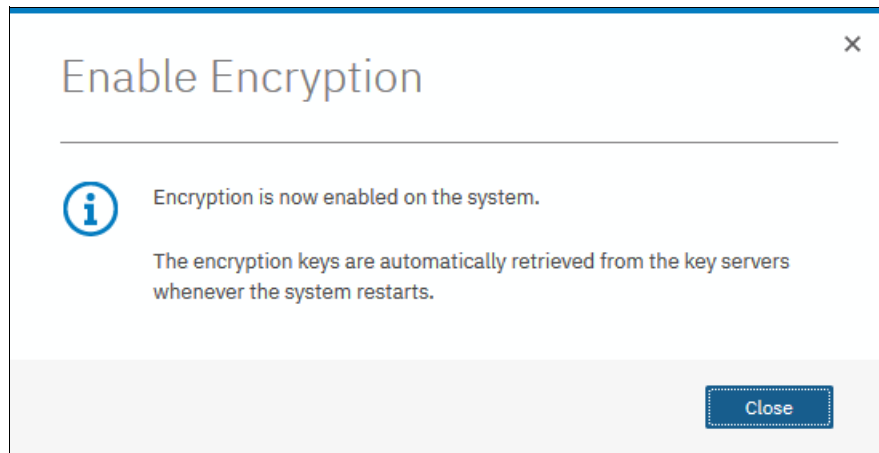


Figure 12-48 Encryption enabled using SafeNet KeySecure key servers

- Confirm that encryption is enabled in **Settings** → **Security** → **Encryption**, as shown in Figure 12-49. Check whether the four servers are shown as online, which indicate that all four SafeNet KeySecure servers are detected as available by the system.

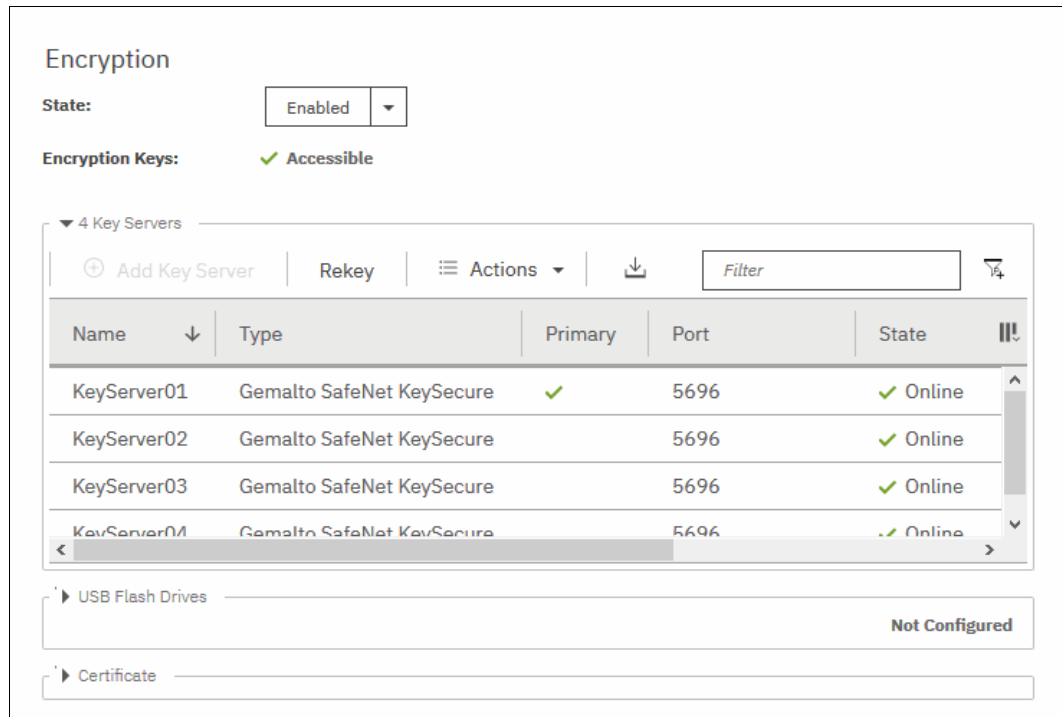


Figure 12-49 Encryption enabled with 4 SafeNet KeySecure key servers

## 12.5.4 Enabling encryption by using both providers

IBM Spectrum Virtualize allows parallel use of both USB flash drive and one type of key server (SKLM or SafeNet KeySecure) as encryption key providers. It is possible to configure both providers in a single run of encryption enable wizard. To perform this configuration process, the system must meet requirements of both key server (SKLM or SafeNet KeySecure) and USB flash drive encryption key providers.

**Note:** Make sure that the key management server functionality is fully independent from an encrypted storage that has encryption managed by this key server environment. Failure to observe this requirement might create an encryption deadlock. An encryption deadlock is a situation in which none of key servers in the environment can become operational because some critical part of the data in each server is stored on an encrypted storage system that depends on one of the key servers to unlock access to the data.

IBM Spectrum Virtualize code V8.1 and later supports up to four key server objects that are defined in parallel.

Before you start to enable encryption by using both USB flash drives and a key servers, confirm the requirements that are described in section 12.5.2, “Enabling encryption by using USB flash drives” on page 706, and 12.5.3, “Enabling encryption by using key servers” on page 711.

To enable encryption by using a key server and USB flash drive, complete the following steps:

1. Ensure that you have service IPs configured on all your nodes.
2. In the Enable Encryption wizard Welcome tab, select **Key servers** and **USB flash drives** and then, click **Next**, as shown in Figure 12-50.

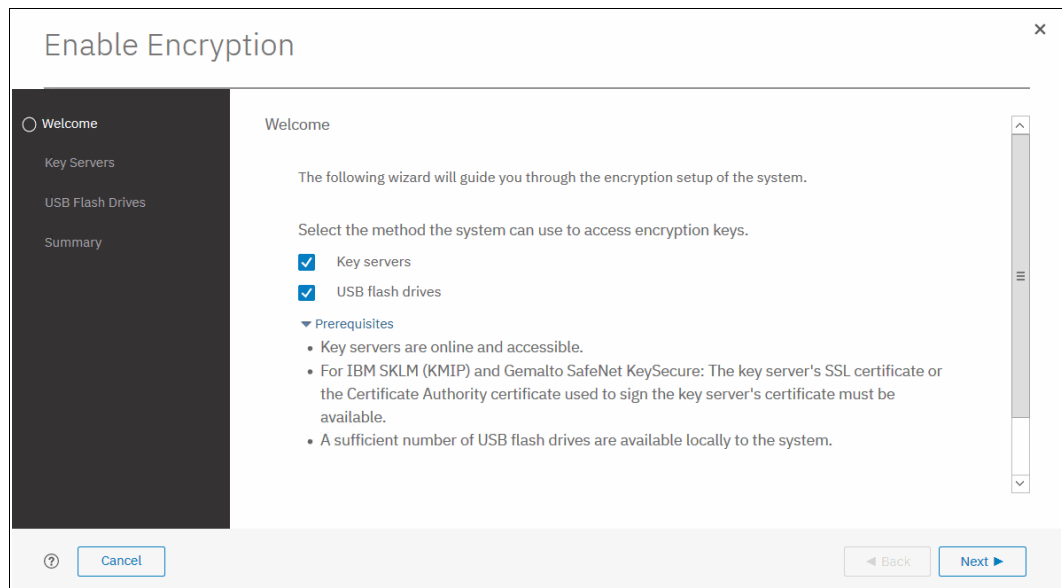


Figure 12-50 Selecting Key servers and USB flash drives in the Enable Encryption wizard

- The wizard moves to the Key Server Types window, as shown in Figure 12-51. Then, select the key server type that manages the encryption keys.

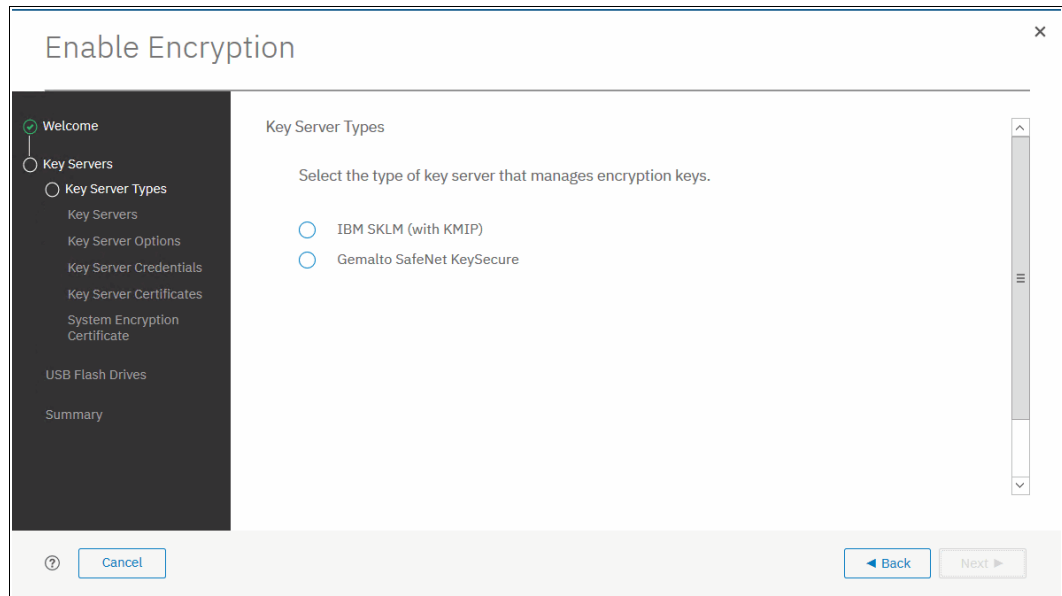


Figure 12-51 Selecting the key server type

The next windows that are displayed are the same as described in 12.5.3, “Enabling encryption by using key servers” on page 711, depending on the type of key server selected.

When the key servers details are entered, the USB flash drive encryption configuration is displayed. In this step, master encryption key copies are stored in the USB flash drives. If fewer than three drives are detected, the system requests plugging in more USB flash drives, as shown on Figure 12-52. You cannot proceed until the required minimum number of USB flash drives is detected by the system.

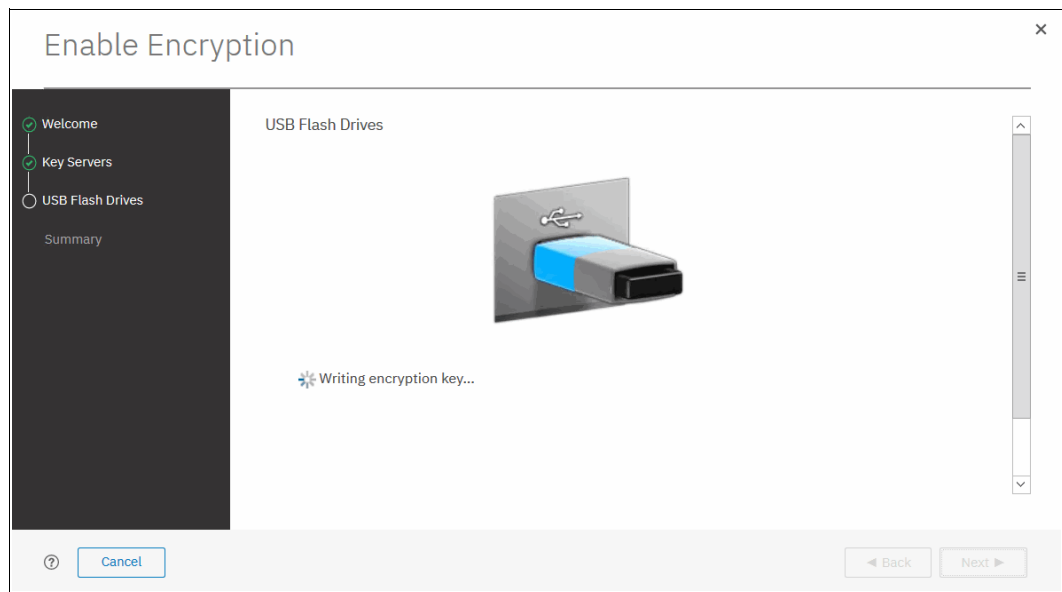


Figure 12-52 Prompt to insert USB flash drives

After at least three USB flash drives are detected, the system writes master access key to each of the drives. Notice that the system attempts to write the encryption key to any flash drive it detects. Therefore, it is crucial to maintain physical security of the system during this procedure. After the keys are successfully copied to at least three USB flash drives, the system displays a window, as shown in Figure 12-53.

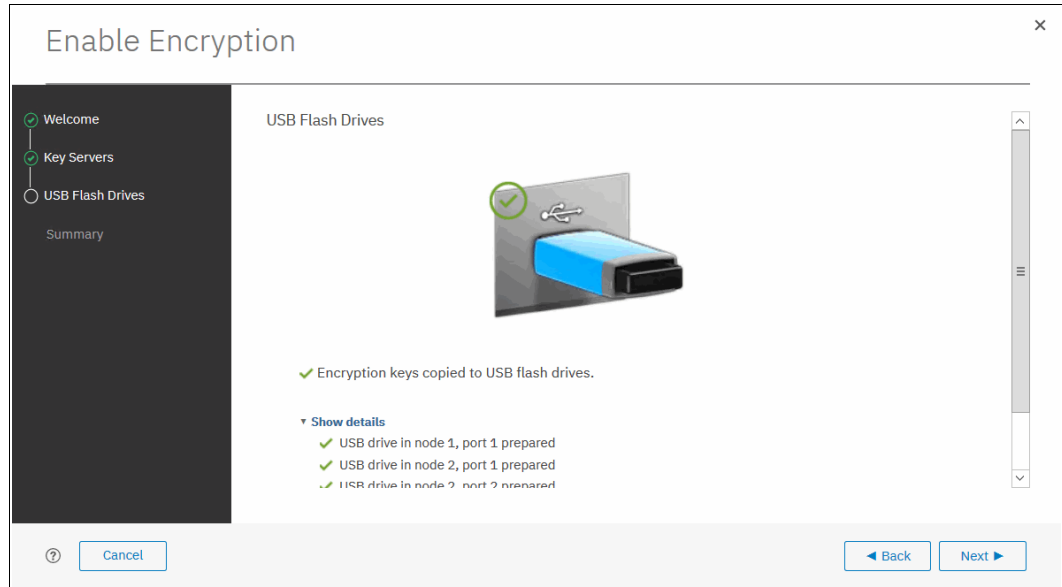


Figure 12-53 Master Access Key successfully copied to USB flash drives

4. After copying encryption keys to USB flash drives, the next window is shown with the summary of the configuration that is implemented on the system (see Figure 12-54). Click **Finish** to create the key server object and finalize the encryption enablement.

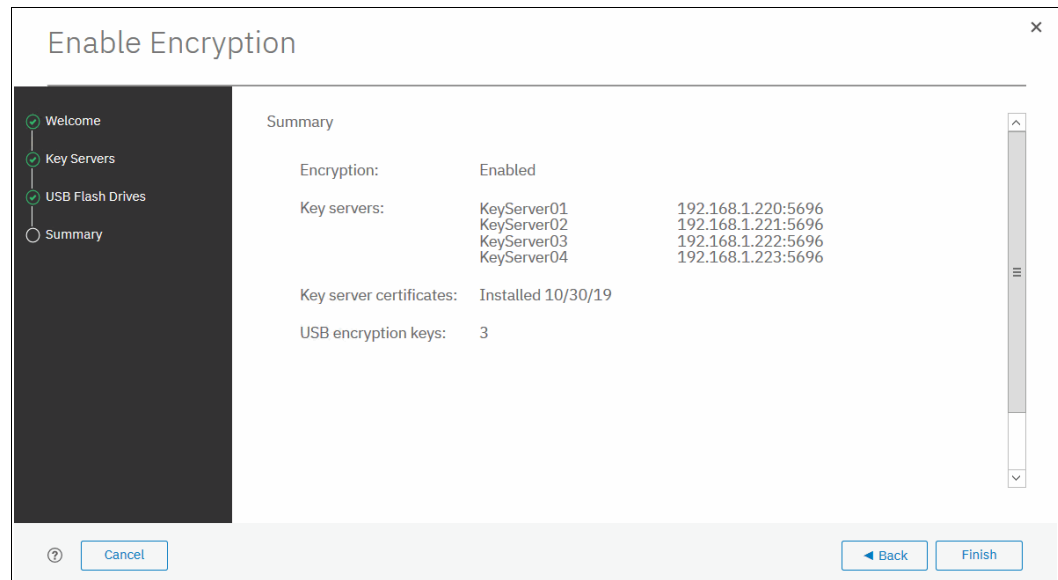
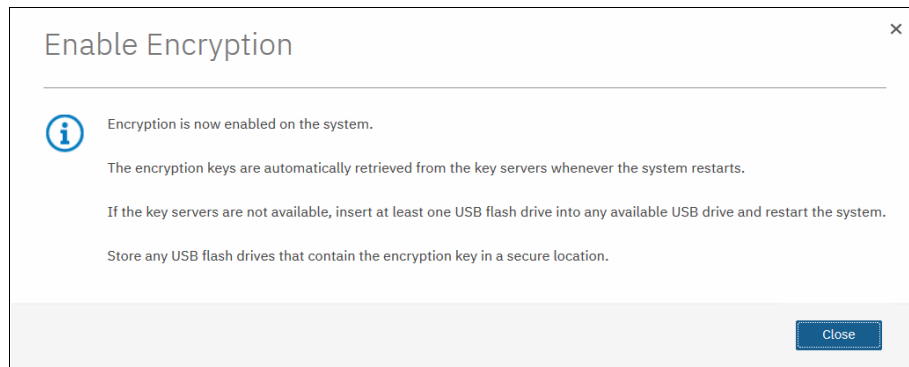


Figure 12-54 Encryption configuration summary in two providers scenario



5. If no errors occur while creating the key server object, the system displays a window that confirms that the encryption is now enabled on the system and that both encryption key providers are enabled (see Figure 12-55).



*Figure 12-55 Encryption enabled message using both encryption key providers*

- You can confirm that encryption is enabled and verify which key providers are in use by selecting **Settings** → **Security** → **Encryption**, as shown in Figure 12-56. Notice the **Online** state of the key servers and **Validated** state of the USB ports where USB flash drives are inserted to ensure that they are correctly configured.

**Encryption**

**State:** Enabled ▼

**Encryption Keys:** ✓ Accessible

▼ 4 Key Servers

+ Add Key Server | Rekey | ☰ Actions ▼ | ↓  ⌵

Name	Type	Primary	Port	State
KeyServer01	Gemalto SafeNet KeySecure	✓	5696	✓ Online
KeyServer02	Gemalto SafeNet KeySecure		5696	✓ Online
KeyServer03	Gemalto SafeNet KeySecure		5696	✓ Online
KeyServer04	Gemalto SafeNet KeySecure		5696	✓ Online

▼ 3 USB Flash Drives Detected

☰ Actions ▼ | Rekey | ↓  ⌵

ID	USB Port	State
0	1	✓ Validated
9	2	✓ Validated

Showing 3 ports | Selecting 0 ports

▶ Certificate

Figure 12-56 Encryption enabled with both USB flash drives and key servers

## 12.6 Configuring more providers

After the system is configured with a single encryption key provider, a second provider can be added.

**Note:** If you set up encryption of your storage system when it was running IBM Spectrum Virtualize code version earlier than V7.8.0, you must rekey the master encryption key before you can enable second encryption provider when you upgrade to code version V8.1 or later.

### 12.6.1 Adding key servers as a second provider

If the storage system is configured with the USB flash drive provider, it is possible to configure SKLM or SafeNet KeySecure servers as a second provider.

To enable key servers as a second provider, complete the following steps:

1. Select **Settings** → **Security** → **Encryption**. Expand the Key Servers section and click **Configure**, as shown in Figure 12-57. To enable key server as a second provider, the system must detect at least one USB flash drive with a current copy of the master access key.

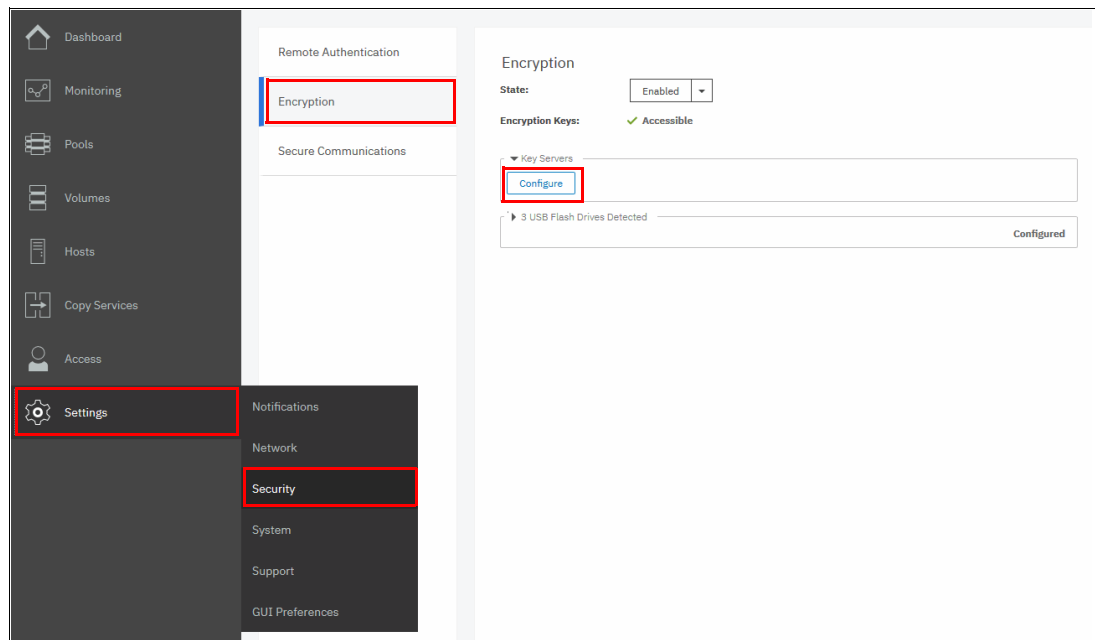


Figure 12-57 Enable key servers as a second provider

2. Complete the steps that are required to configure the key server provider, as described in 12.5.3, “Enabling encryption by using key servers” on page 711. The difference in the process that is described in that section is that the wizard gives you an option to disable USB flash drive encryption, which aims to migrate from the USB flash drive to key server provider.

Select **No** to enable both encryption key providers, as shown in Figure 12-58.

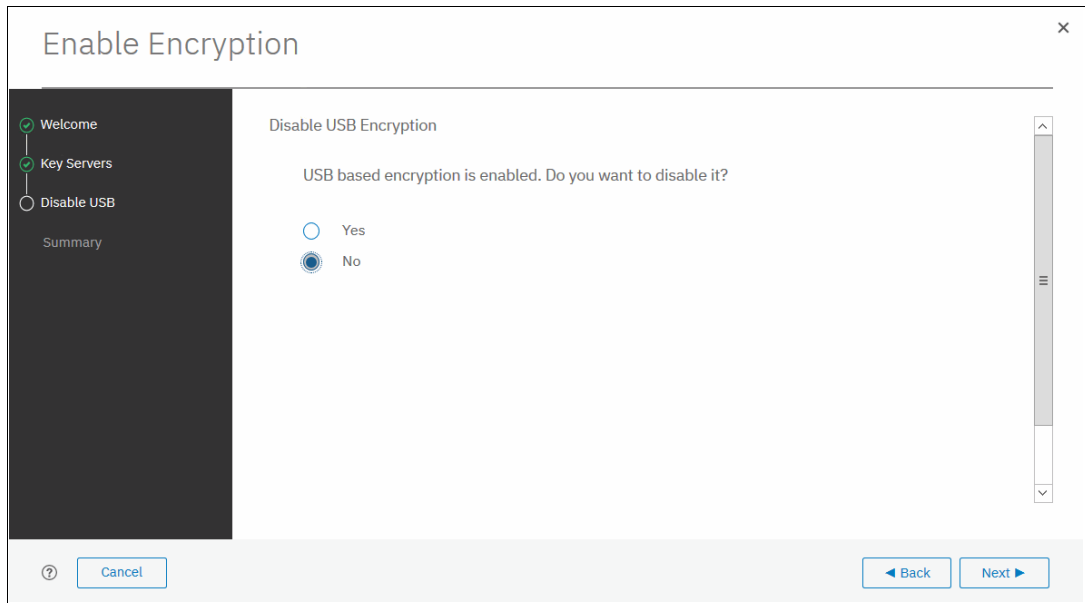


Figure 12-58 Do not disable USB flash drive encryption key provider

This choice is confirmed on the summary window before the configuration is committed, as shown in Figure 12-59.

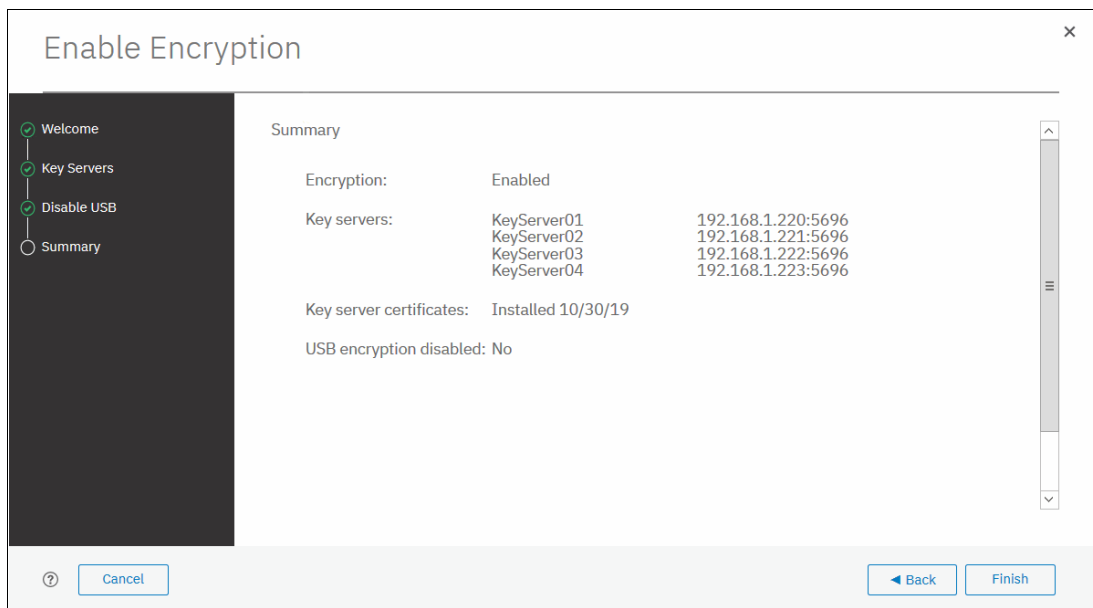


Figure 12-59 Configuration summary before committing

- After you click **Finish**, the system configures keys servers as a second encryption key provider. Successful completion of the task is confirmed by a message, as shown in Figure 12-60. Click **Close**.

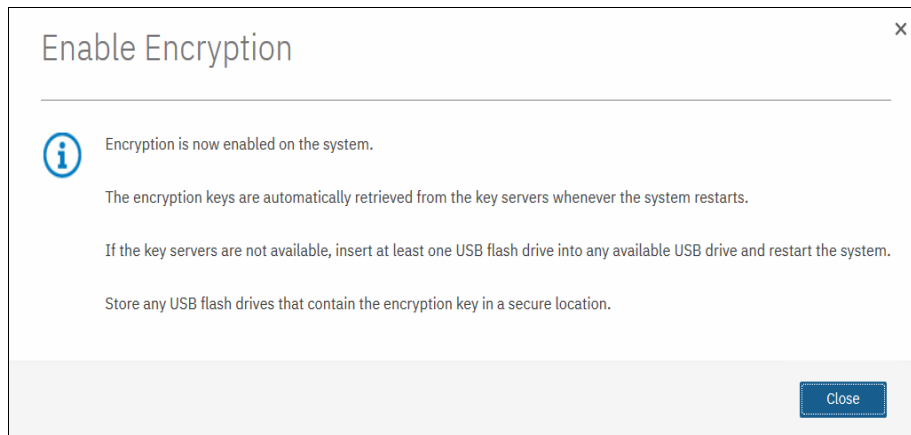


Figure 12-60 Confirmation of successful configuration of two encryption key providers

- You can confirm that encryption is enabled and verify which key providers are in use by selecting **Settings** → **Security** → **Encryption**, as shown in Figure 12-61. Notice the **Online** state of the key servers and the **Validated** state of the USB ports where USB flash drives are inserted to make sure that they are correctly configured.

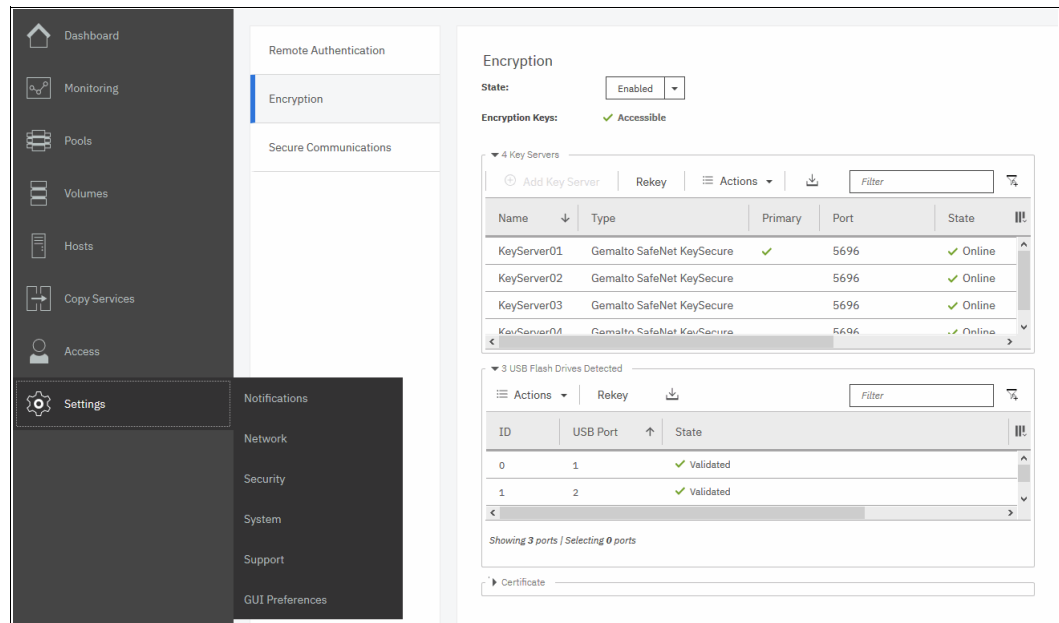


Figure 12-61 Encryption enabled with two key providers available

## 12.6.2 Adding USB flash drives as a second provider

If the storage system is configured with an SKLM or SafeNet KeySecure encryption key provider, it is possible to configure USB flash drives as a second provider. To enable USB flash drives as a second provider, complete the following steps:

1. Select **Settings** → **Security** → **Encryption**. Expand the USB Flash Drives section and click **Configure**, as shown in Figure 12-62. To enable USB flash drives as a second provider, the system must access key servers with the current master access key.

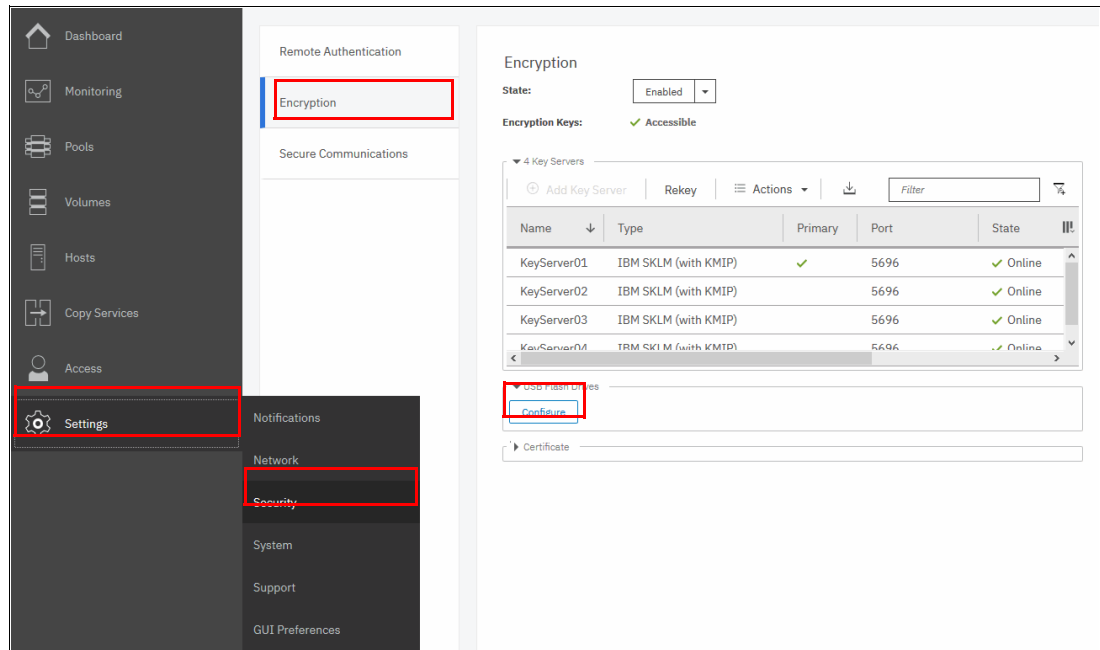


Figure 12-62 Enable USB flash drives as a second encryption key provider

2. After you click **Configure**, you are presented with a wizard that is similar to the one that is described in 12.5.2, “Enabling encryption by using USB flash drives” on page 706. You cannot disable key server providers during this process.

After successful completion of the process, you are presented with a message confirming that both encryption key providers are enabled, as shown in Figure 12-63.

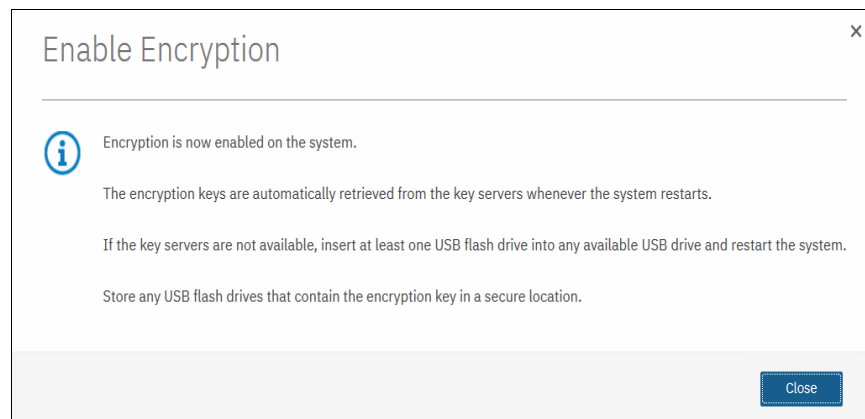


Figure 12-63 Confirmation of successful configuration of two encryption key providers

- You can confirm that encryption is enabled and verify which key providers are in use by selecting **Settings** → **Security** → **Encryption**, as shown in Figure 12-64. Notice the **Online** state of the key servers and the **Validated** state of the USB ports where USB flash drives are inserted to make sure that they are correctly configured.

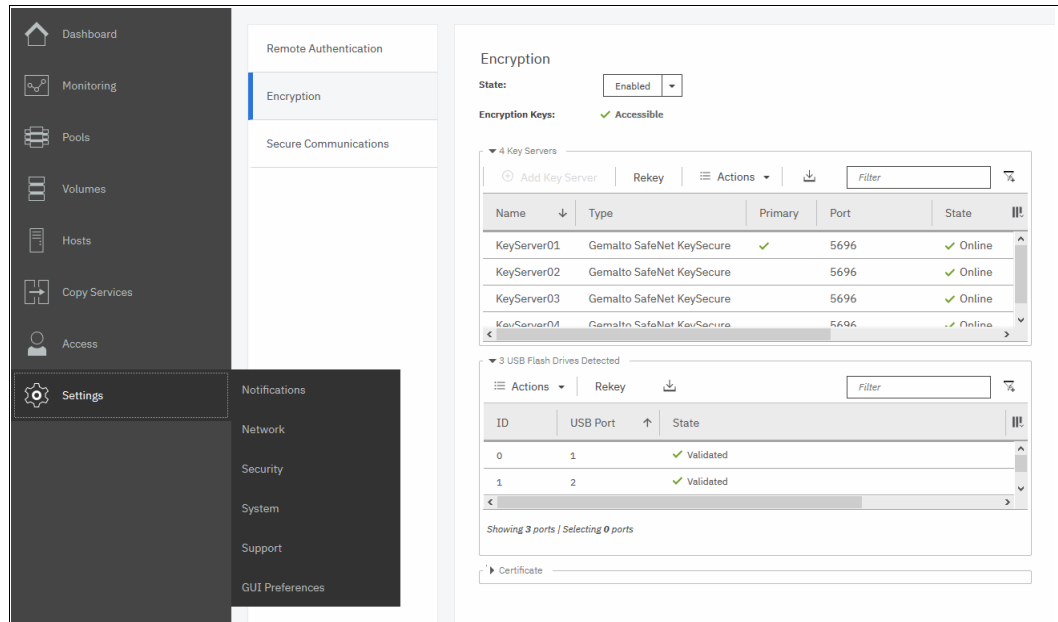


Figure 12-64 Encryption enabled with two key providers available

## 12.7 Migrating between providers

IBM Spectrum Virtualize V8.1 introduced support for simultaneous use of both USB flash drives and a key server as encryption key providers. The system also allows migration from configuration by using only USB flash drive provider to key servers provider, and vice versa.

If you want to migrate from one key server type to another (for example, migrating from SKLM to SafeNet KeySecure or vice versa), direct migration is not possible. In this case, it is required first to migrate from the current key server type to a USB flash drive, and then migrate to the other type of key server.

## 12.7.1 Changing from USB flash drive provider to encryption key server

The system is designed to facilitate changing from USB flash drives encryption key provider to encryption key server provider. If you follow the steps that are described in 12.6.1, “Adding key servers as a second provider” on page 729, but when completing step 2 on page 729, select **Yes** instead of **No** (see Figure 12-65). This action de-activates the USB flash drives provider, and the procedure completes with only key servers configured as key provider.

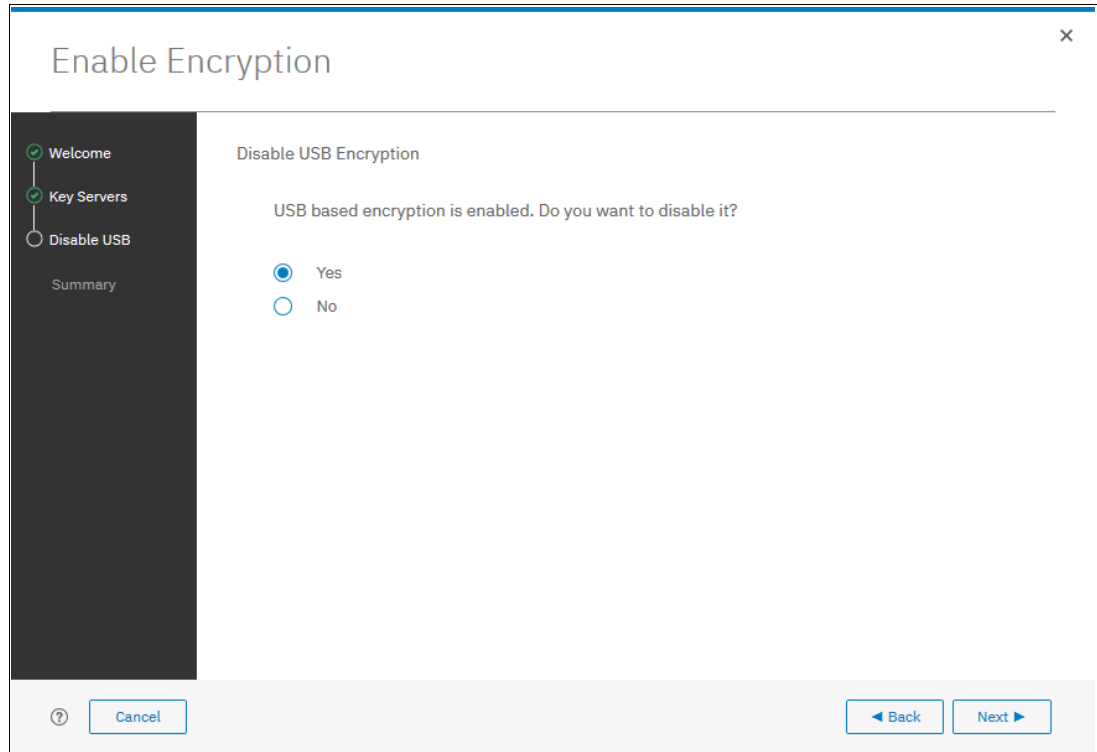


Figure 12-65 Disable USB flash drive provider while changing to SKLM provider

## 12.7.2 Changing from encryption key server to USB flash drive provider

Change in the other direction (that is, from the use of encryption key servers provider to USB flash drives provider), is not possible by using only the GUI.

To change the direction, add USB flash drives as a second provider by completing the steps described in 12.6.2, “Adding USB flash drives as a second provider” on page 732.

Then, run the following command in the CLI:

```
chencryption -usb validate
```

To make sure that USB drives contain the correct master access key, disable the encryption key server provider by running the following command:

```
chencryption -keyserver disable
```

This command disables the encryption key server provider, which effectively migrates your system from encryption key server to USB flash drive provider.



### 12.7.3 Migrating between different key server types

The migration between different key server types cannot be performed directly from one type of key server to another. USB flash drives encryption must be used to facilitate this process.

If you want to migrate from one type of key server to another, you first must migrate from your current key servers to USB encryption, and then, migrate from USB to the other type of key servers.

The procedure to migrate from one key server type to another is shown here. In this example, we migrate an IBM Spectrum Virtualize system that is configured with IBM SKLM key servers (as shown in Figure 12-66) to SafeNet KeySecure servers.

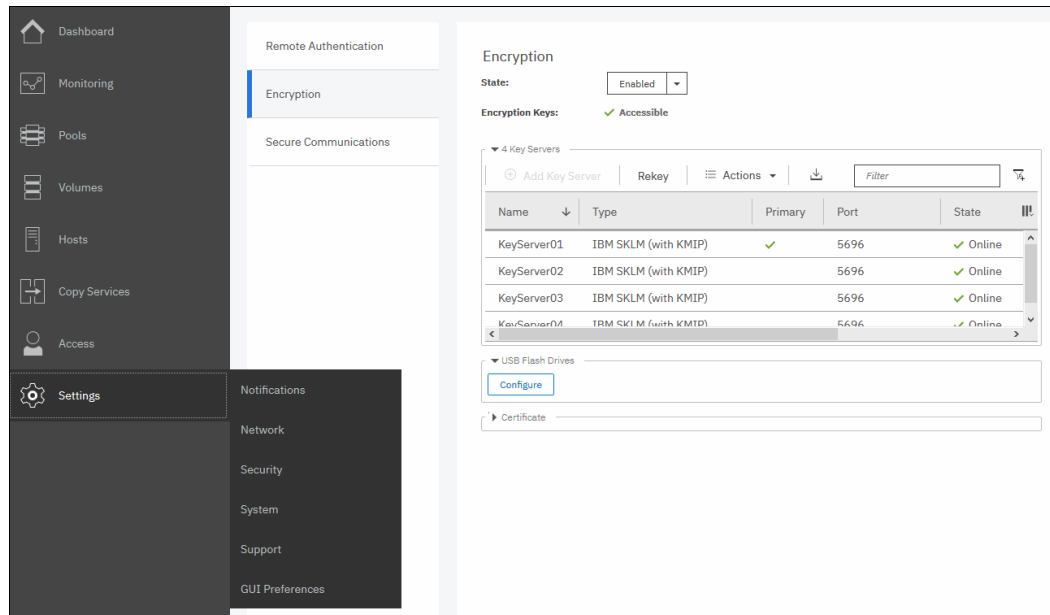


Figure 12-66 IBM Spectrum Virtualize encryption configured with IBM SKLM servers

Complete the following steps to migrate to Gemalto SafeNet KeySecure:

1. Migrate from key server encryption to USB flash drives encryption, as described in 12.7.2, “Changing from encryption key server to USB flash drive provider” on page 734. After this step, only USB flash drives encryption are configured, as shown in Figure 12-67 on page 736.

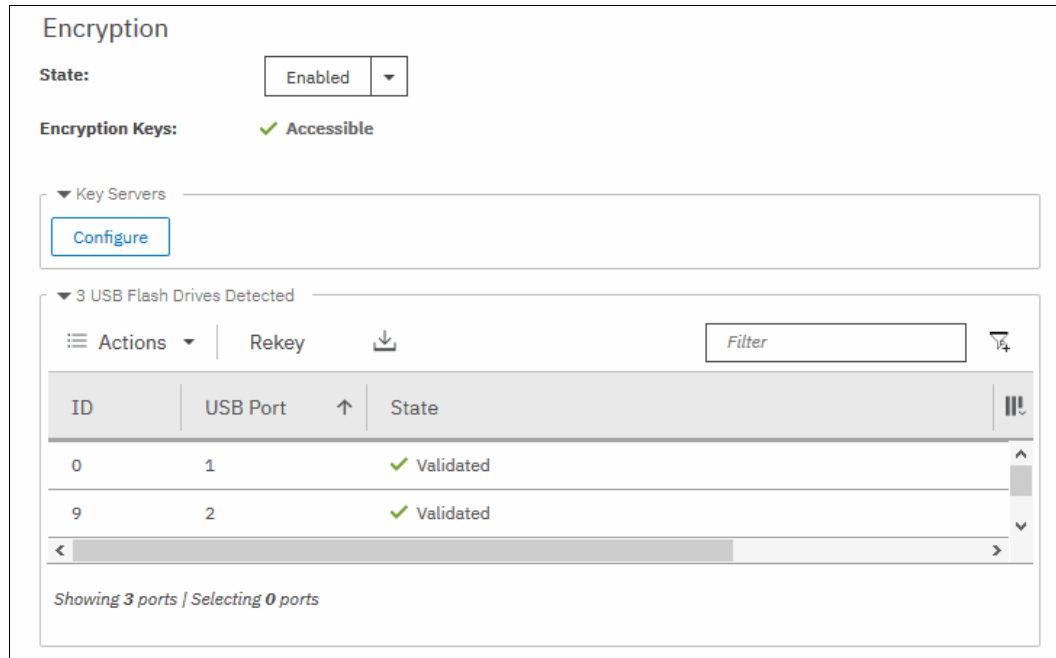


Figure 12-67 IBM SAN Volume Controller encryption configured with USB Flash Drives

2. Migrate from USB flash drives encryption to the other key server type encryption (in this example, Gemalto SafeNet KeySecure), following the steps that are described in 12.7.1, “Changing from USB flash drive provider to encryption key server” on page 734. After completing this step, the other key server type are configured as encryption provider in IBM Spectrum Virtualize, as shown in Figure 12-68.

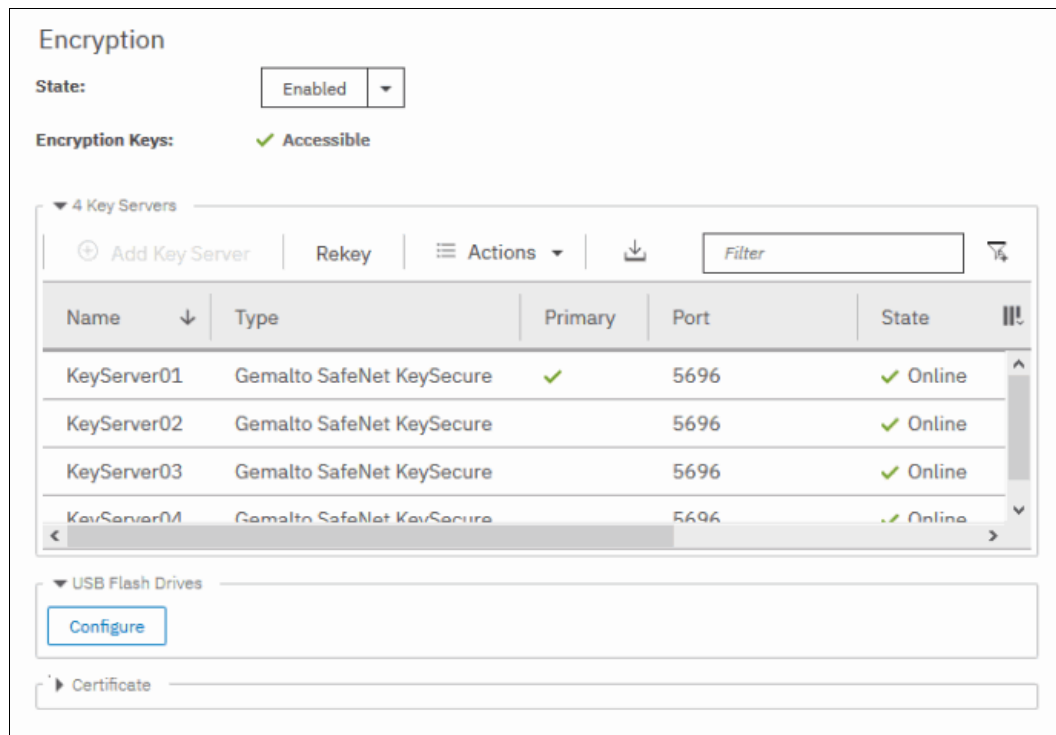


Figure 12-68 IBM SAN Volume Controller encryption configured with SafeNet KeySecure

## 12.8 Recovering from a provider loss

If both encryption key providers are enabled, and you lose one of them (by losing all copies of the encryption key kept on the USB flash drives or by losing all SKLM servers), you can recover from this situation by disabling the provider to which you lost the access. To disable the unavailable provider, you must have access to a valid master access key on the remaining provider.

If you lost access to the encryption key server provider, run the following command:

```
chencryption -keyserver disable
```

If you lost access to the USB flash drives provider, run the following command:

```
chencryption -usb disable
```

If you want to restore the configuration with both encryption key providers, follow the instructions that are described in 12.6, “Configuring more providers” on page 729.

**Note:** If you lose access to all encryption key providers that are defined in the system, no method is available to recover access to the data protected by the master access key.

## 12.9 Using encryption

The design for encryption is based on the concept that a system is fully encrypted or not encrypted. Encryption implementation is intended to encourage solutions that contain only encrypted volumes or only unencrypted volumes. For example, after encryption is enabled on the system, all new objects (for example, pools) are by default created as encrypted.

Some unsupported configurations are actively policed in code. For example, no support exists for creating unencrypted child pools from encrypted parent pools. However, the following exceptions exist:

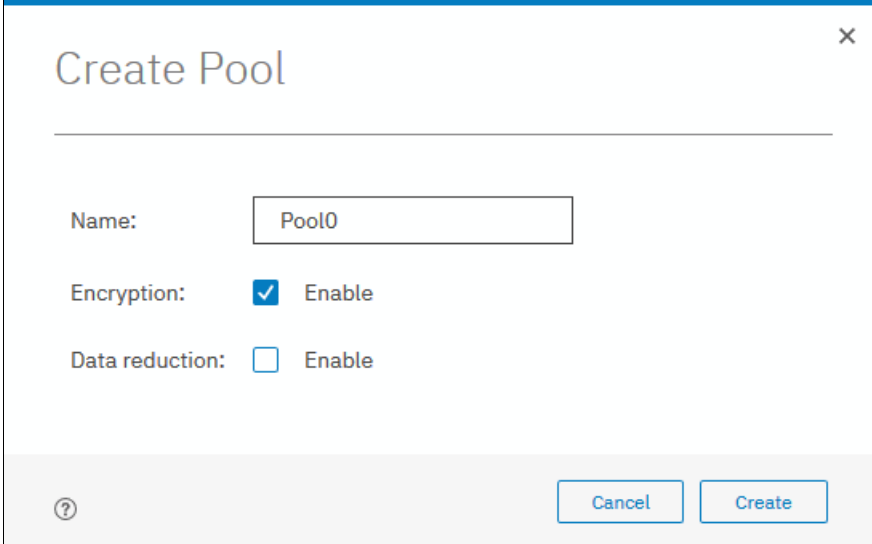
- ▶ During the migration of volumes from unencrypted to encrypted volumes, a system might report encrypted and unencrypted volumes.
- ▶ It is possible to create unencrypted arrays from CLI by manually overriding the default encryption setting.

**Notes:** Encryption support for distributed Redundant Array of Independent Disks (DRAID) is available in IBM Spectrum Virtualize code V7.7 and later.

You must decide whether to encrypt or not encrypt an object when it is created. You cannot change this setting later. To change the encryption state of stored data, you must migrate from an encrypted object (for example, pool) to an unencrypted one, or vice versa. Volume migration is the only way to encrypt any volumes that were created before enabling encryption on the system.

## 12.9.1 Encrypted pools

For more information about how to open the Create Pool window, see Chapter 5, “Storage pools” on page 199. After encryption is enabled, any new pool is created by default as encrypted, as shown in Figure 12-69.



The screenshot shows a 'Create Pool' dialog box. The 'Name' field is filled with 'Pool0'. The 'Encryption' checkbox is checked, and the 'Data reduction' checkbox is unchecked. The 'Create' button is highlighted in blue.

Figure 12-69 Create Pool window basic

You can click **Create** to create an encrypted pool. All storage that is added to this pool is encrypted.

You can customize the Pools view in the management GUI to show pool encryption status. Select **Pools** → **Pools**, and then, select **Actions** → **Customize Columns** → **Encryption**, as shown in Figure 12-70.

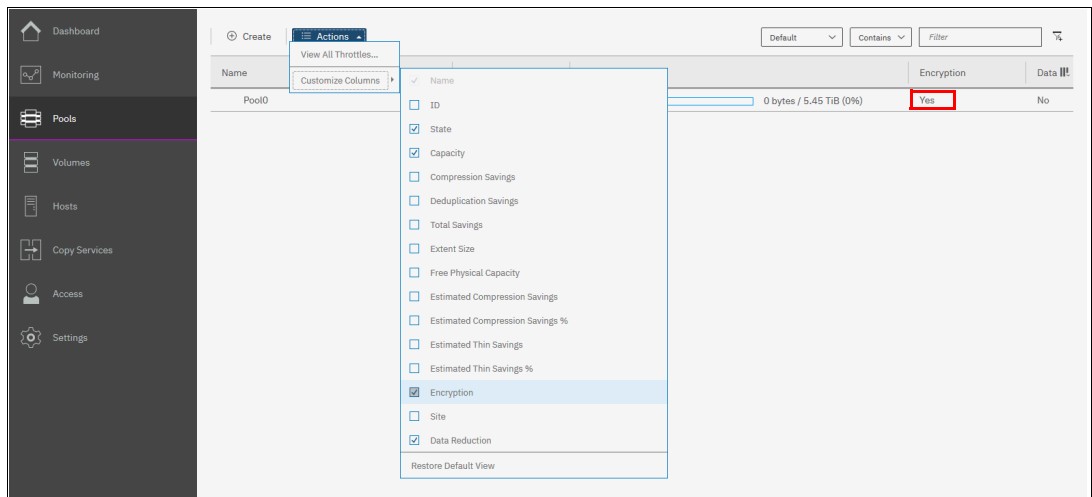


Figure 12-70 Pool encryption state

If you create an unencrypted pool but you add only encrypted arrays or self-encrypting MDisks to the pool, the pool is reported as encrypted because all extents in the pool are encrypted. The pool reverts to the unencrypted state if you add an unencrypted array or MDisk.

More information about how to add encrypted storage to encrypted pools is described next. You can mix and match storage encryption types in a pool. Figure 12-71 shows an example of an encrypted pool that contains storage by using different encryption methods.

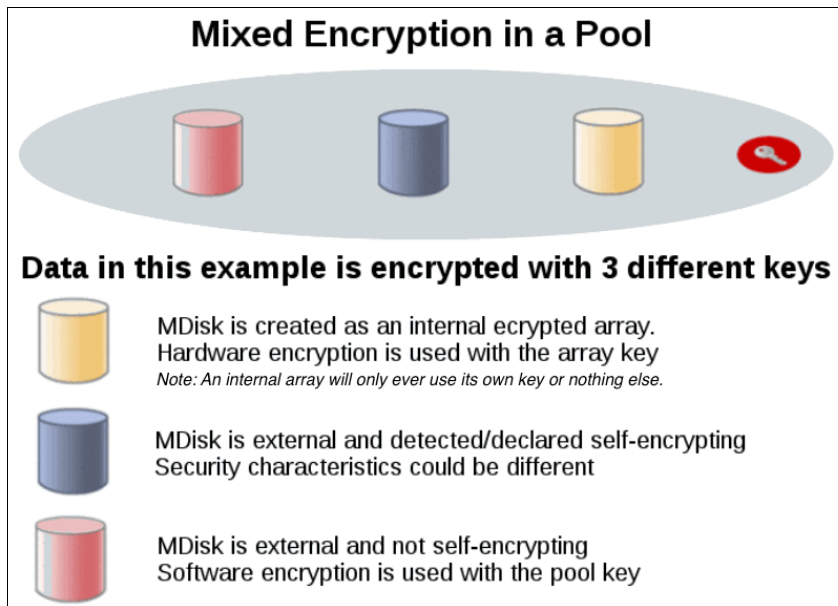


Figure 12-71 Mix and match encryption in a pool

## 12.9.2 Encrypted child pools

For more information about how to open the Create Child Pool window, see Chapter 5, “Storage pools” on page 199. If the parent pool is encrypted, every child pool also must be encrypted. The GUI enforces this requirement by automatically selecting **Encryption Enabled** in the Create Child Pool window and preventing changes to this setting, as shown in Figure 12-72.

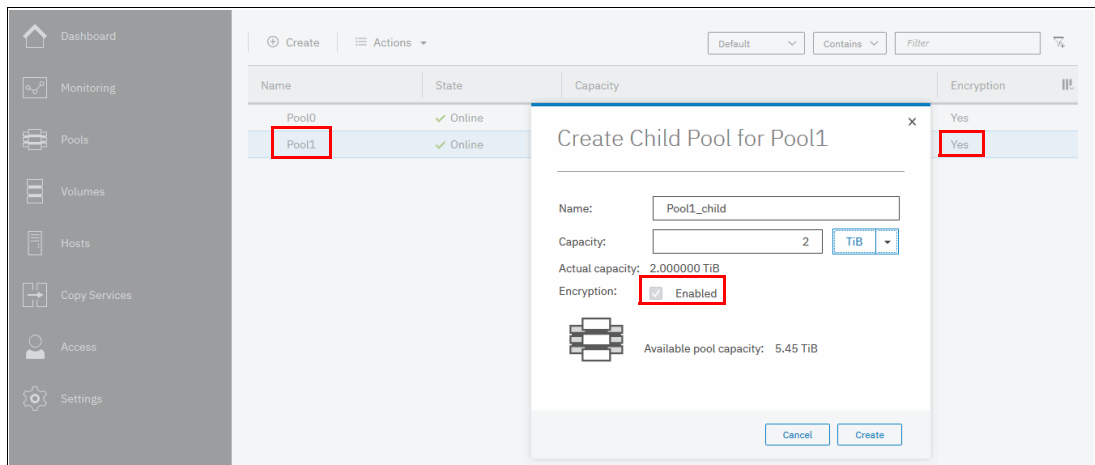


Figure 12-72 Create a child pool of an encrypted parent pool

However, if you want to create encrypted child pools from an unencrypted storage pool that contains a mix of internal arrays and external MDisks, the following restrictions apply:

- ▶ The parent pool must not contain any unencrypted internal arrays. If any unencrypted internal array is in the unencrypted pool, when you try to create a child pool and select the option to set as encrypted, it is created as unencrypted.
- ▶ All IBM SAN Volume Controller nodes in the system must support software encryption and have the encryption license activated.

**Note:** An encrypted child pool that is created from an unencrypted parent storage pool reports as unencrypted if the parent pool contains any unencrypted internal arrays. Remove these arrays to ensure that the child pool is fully encrypted.

If you modify the Pools view, you see the encryption status of child pools, as shown in Figure 12-73. The example shows an encrypted child pool with non-encrypted parent pool.

Name	State	Capacity	Encryption
Pool0	Online	0 bytes / 5.45 TiB (0%)	Yes
Pool1	Online	0 bytes / 5.45 TiB (0%)	Yes
Pool1_child	Online	0 bytes / 2.00 TiB (0%)	Yes
Pool2	Online	0 bytes / 5.45 TiB (0%)	No
Pool2_child	Online	0 bytes / 2.00 TiB (0%)	Yes

Figure 12-73 Child pool encryption state

### 12.9.3 Encrypted arrays

For more information about how to add internal storage to a pool, see Chapter 5, “Storage pools” on page 199. After encryption is enabled, all newly built arrays are hardware encrypted by default. In this case, the GUI does not allow you to create an unencrypted array. To create an unencrypted array, the command-line interface (CLI) must be used. Example 12-1 shows how to create an unencrypted array by using the CLI.

*Example 12-1 Creating an unencrypted array using CLI with IBM SAN Volume Controller*

```
IBM_SAN:ITS0-SVC:superuser>svctask mkarray -drive 6:4 -level raid1 -sparegoal 0
-strip 256 -encrypt no Pool2
MDisk, id [2], successfully created
IBM_SAN:ITS0-SVC:superuser>
```

**Note:** It is not possible to add unencrypted arrays to an encrypted pool.

You can customize the MDisks by Pools view to show array encryption status. Select **Pools** → **MDisk by Pools**, and then, click **Actions** → **Customize Columns** → **Encryption**. You also can right-click the table header to customize columns and select **Encryption**, as shown in Figure 12-74.

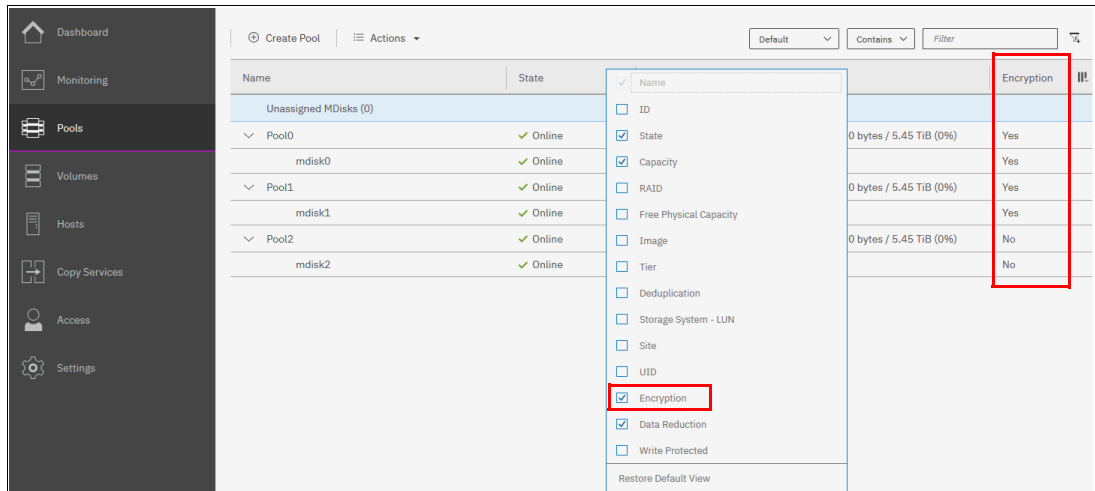


Figure 12-74 Array encryption state

You can also check the encryption state of an array by reviewing its drives in **Pools** → **Internal Storage** view. The internal drives that are associated with an encrypted array are assigned an encrypted property that can be seen, as shown in Figure 12-75.

Drive ID ↑	Capacity	Use	Status	MDisk Name	Slot ID	Encrypted
0	5.46 TiB	Member	✓ Online	mdisk0	11	✓
1	5.46 TiB	Spare	✓ Online		1	
2	5.46 TiB	Member	✓ Online	mdisk0	9	✓
3	5.46 TiB	Member	✓ Online	mdisk1	2	✓
4	5.46 TiB	Member	✓ Online	mdisk2	6	
5	5.46 TiB	Member	✓ Online	mdisk1	10	✓
6	5.46 TiB	Member	✓ Online	mdisk2	5	

Figure 12-75 Drive encryption state

## 12.9.4 Encrypted MDisks

For more information about how to add external storage to a pool, see Chapter 5, “Storage pools” on page 199. Each MDisk that belongs to external storage that is added to an encrypted pool or child pool is automatically encrypted by using the pool or child pool key, unless the MDisk is detected or declared as self-encrypting.

The user interface gives no method to see which extents contain encrypted data and which do not. However, if a volume is created in a correctly configured encrypted pool, all data that is written to this volume is encrypted.

You can use the MDisk by Pools view to show the object encryption state by selecting **Pools** → **MDisk by Pools**. Figure 12-76 shows an example in which self-encrypting MDisk is in an unencrypted pool.

The screenshot shows a table with columns: Name, State, Capacity, and Encryption. Pool0 contains mdisk0 (Yes encryption). Pool1 contains mdisk2 (No encryption) and mdisk1 (Yes encryption). The 'No' for mdisk2 is highlighted with a red box.

Name	State	Capacity	Encryption
Unassigned MDisks (0)			
Pool0	Online	0 bytes / 5.45 TiB (0%)	Yes
mdisk0	Online	5.46 TiB	Yes
Pool1	Online	0 bytes / 10.91 TiB (0%)	No
mdisk2	Online	5.46 TiB	No
mdisk1	Online	5.46 TiB	Yes

Figure 12-76 MDisk encryption state

When working with MDisk encryption, take extra care when configuring MDisk and pools.

If the MDisk was used earlier for storage of unencrypted data, the extents can contain stale unencrypted data. This issue occurs because file deletion only marks disk space as free. The data is *not* removed from the storage. Therefore, if the MDisk is not self-encrypting and was a part of an unencrypted pool and later was moved to an encrypted pool, it contains stale data from its previous life.

Another mistake that can occur is to misconfigure an external MDisk as self-encrypting, while in reality it is not self-encrypting. In that case, the data that is written to this MDisk is not encrypted by IBM SAN Volume Controller because IBM SAN Volume Controller expects that the storage system that is hosting the MDisk encrypts the data. At the same time, the MDisk does not encrypt the data because it is not self-encrypting; therefore, the system ends up with unencrypted data on an extent in an encrypted storage pool.

However, all data that is written to any MDisk that is a part of correctly configured encrypted storage pool is going to be encrypted.

### Self-encrypting MDisk

When adding external storage to a pool, be exceptionally diligent when declaring the MDisk as self-encrypting. Correctly declaring an MDisk as self-encrypting avoids waste of resources, such as CPU time. However, when used incorrectly, it might lead to unencrypted data at-rest.



To declare an MDisk as self-encrypting, select **Externally encrypted** when adding external storage in the Assign Storage view, as shown in Figure 12-77.

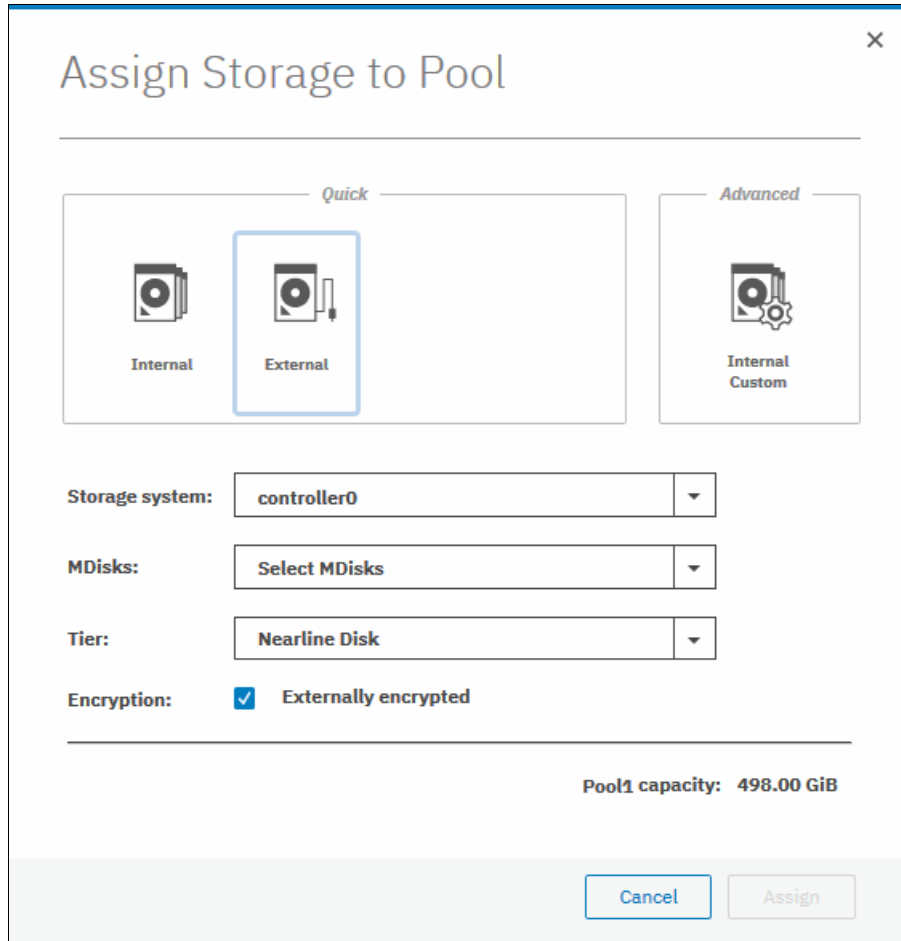


Figure 12-77 Declaring MDisk as externally encrypted

IBM Spectrum Virtualize products can detect that an MDisk is self-encrypting by using the SCSI Inquiry page C2. MDisks that are provided by other IBM Spectrum Virtualize products report this page correctly. For these MDisks, the **Externally encrypted** option that is shown in Figure 12-77 is not selected. However, when added, they are still considered as self-encrypting.

**Note:** You can override the external encryption setting of a detected MDisk as self-encrypting and configure it as unencrypted by running the CLI command `chmdisk -encrypt no`. However, only do so if you plan to decrypt the data on the backend or if the backend uses inadequate data encryption.

To check whether an MDisk was detected or declared as self-encrypting, select **Pools** → **MDisk by Pools** and verify the information in the Encryption column, as shown in Figure 12-78.

Name	State	Capacity	Encryption	!!!
Unassigned MDisks (1)				
Pool0	✓ Online	24.00 GiB / 99.00 GiB (24%)	No	
mdisk0	✓ Online	100.00 GiB	No	
Pool1	✓ Online	0 bytes / 498.00 GiB (0%)	Yes	
mdisk1	✓ Online	200.00 GiB	Yes	
mdisk2	✓ Online	300.00 GiB	No	

Figure 12-78 MDisk self-encryption state

The value that is shown in the Encryption column shows the property of objects in respective rows. That means that in the configuration that is shown in Figure 12-78, Pool1 is encrypted, so every volume created from this pool is encrypted. However, that pool is formed by two MDisks, out of which one is self-encrypting and one is not. Therefore, a value of No next to mdisk2 does not imply that encryption of Pool1 is in any way compromised. It indicates only that encryption of the data that is placed on mdisk2 is done by using software encryption. Data that is placed on mdisk1 is encrypted by the back-end storage that is providing these MDisks.

**Note:** You can change the self-encrypting attribute of an MDisk that is unmanaged or member of an unencrypted pool. However, you cannot change the self-encrypting attribute of an MDisk after it is added to an encrypted pool.

## 12.9.5 Encrypted volumes

For more information about how to create and manage volumes, see Chapter 6, “Volumes” on page 255. The encryption status of a volume depends on the pool encryption status. Volumes that are created in an encrypted pool are automatically encrypted.

You can modify Volumes view to show if the volume is encrypted. Select **Volumes** → **Volumes**. Then, click **Actions** → **Customize Columns** → **Encryption** to customize the view to show volumes encryption status, as shown in Figure 12-79.

Name	State	Synchronized	Pool	UID	Encryption	!!!
Volume000	✓ Online (formatting)		Pool0	6005076801B807F934000000000000...	Yes	
Volume001	✓ Online (formatting)		Pool0	6005076801B807F934000000000000...	Yes	
Volume002	✓ Online (formatting)		Pool0	6005076801B807F934000000000000...	Yes	
Volume003	✓ Online		Pool1	6005076801B807F934000000000000...	No	
Volume004	✓ Online		Pool1	6005076801B807F934000000000000...	No	
Volume005	✓ Online		Pool1	6005076801B807F934000000000000...	No	
Volume006	✓ Online		Pool0	6005076801B807F934000000000000...	No	
Volume007	✓ Online		Pool0	6005076801B807F934000000000000...	No	
Volume008	✓ Online		Pool0	6005076801B807F934000000000000...	Yes	
Volume009	✓ Online		Pool0	6005076801B807F934000000000000...	Yes	
Volume010	✓ Online		Pool1	6005076801B807F934000000000000...	No	
Volume011	✓ Online		Pool1	6005076801B807F934000000000000...	No	

Figure 12-79 Volume view customization

A volume is reported as encrypted only if all the volume copies are encrypted, as shown in Figure 12-80.

Name	State	Synchronized	Pool	Encryption
Volume003	Online		Pool0	Yes
Copy 0*	Online	Yes	Pool0	Yes
Copy 1	Online	Yes	Pool0	Yes
Volume004	Online		Pool1	No
Copy 0*	Online	Yes	Pool1	No
Copy 1	Online	Yes	Pool0	Yes

Figure 12-80 Volume encryption status depending on volume copies encryption

When creating volumes, make sure to select encrypted pools to create encrypted volumes, as shown in Figure 12-81.

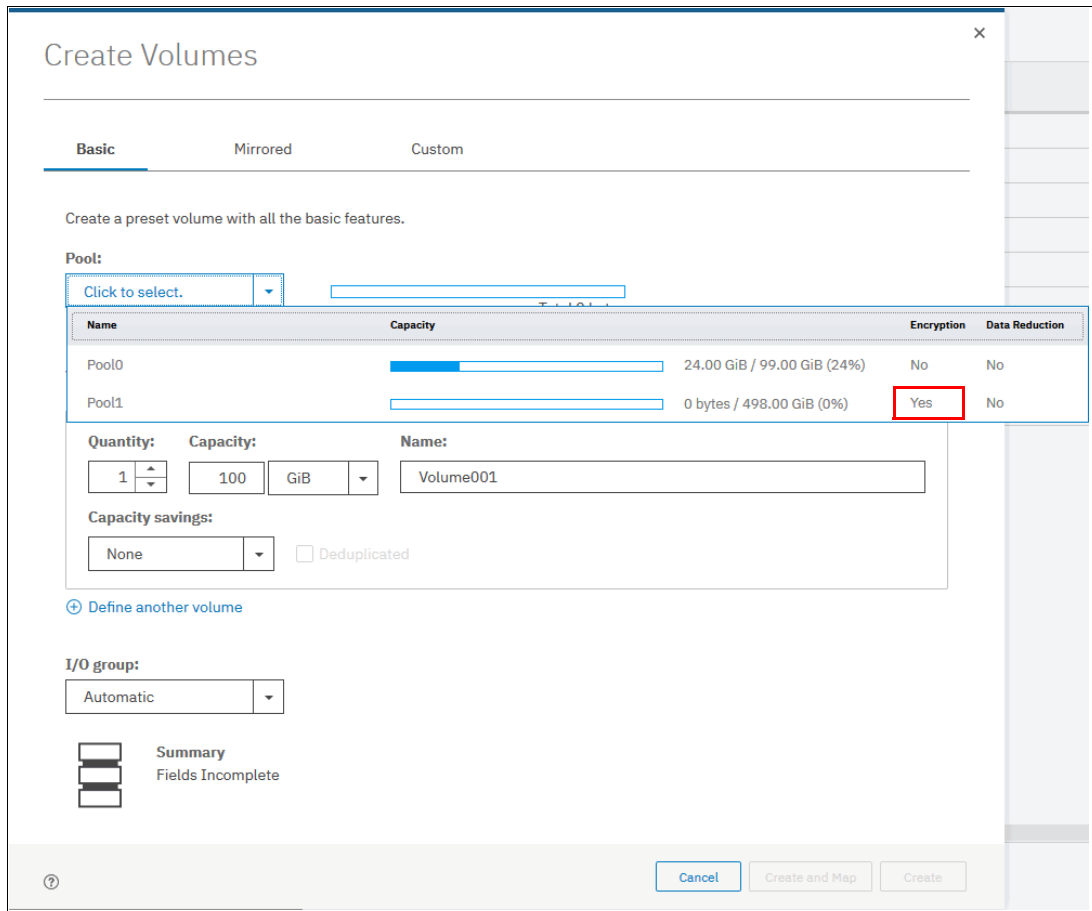


Figure 12-81 Create an encrypted volume by selecting an encrypted pool

You cannot change an unencrypted volume to an encrypted version of itself dynamically. However, this conversion is possible by using one of the following migration options:

- ▶ Migrate a volume to an encrypted pool or child pool.
- ▶ Mirror a volume to an encrypted pool or child pool and delete the unencrypted copy.

For more information about these methods, see Chapter 12, “Encryption” on page 685.

## 12.9.6 Restrictions

The following restrictions apply to encryption:

- ▶ Image mode volumes cannot be in encrypted pools.
- ▶ You cannot add external non self-encrypting MDisks to encrypted pools unless all nodes in the system support encryption.

## 12.10 Rekeying an encryption-enabled system

Changing the master access key is a security requirement. *Rekeying* is the process of replacing current master access key with a newly generated one. The rekey operation works whether encrypted objects exist. The rekeying operation requires access to a valid copy of the original master access key on an encryption key provider that you plan to rekey. Use the rekey operation according to the schedule defined in your organization's security policy and whenever you suspect that the key might have been compromised.

If you have both USB and key server enabled, rekeying is done separately for each of the providers.

**Important:** Before you create a master access key, ensure that all nodes are online and that the current master access key is accessible.

No method is available to directly change data encryption keys. If you must change the data encryption key that is used to encrypt data, the only available method is to migrate that data to a new encrypted object (for example, an encrypted child pool). Because the data encryption keys are defined per encrypted object, such migration forces a change of the key that is used to encrypt that data.

### 12.10.1 Rekeying by using a key server

Ensure that all the configured key servers can be reached by the system and that service IPs are configured on all your nodes.

To rekey the master access key kept on the key server provider, complete the following steps:

1. Select **Settings** → **Security** → **Encryption**. Ensure that Encryption Keys shows that all configured SKLM servers are reported as `Accessible`, as shown in Figure 12-82 on page 747. Click **Key Servers** to expand the section.

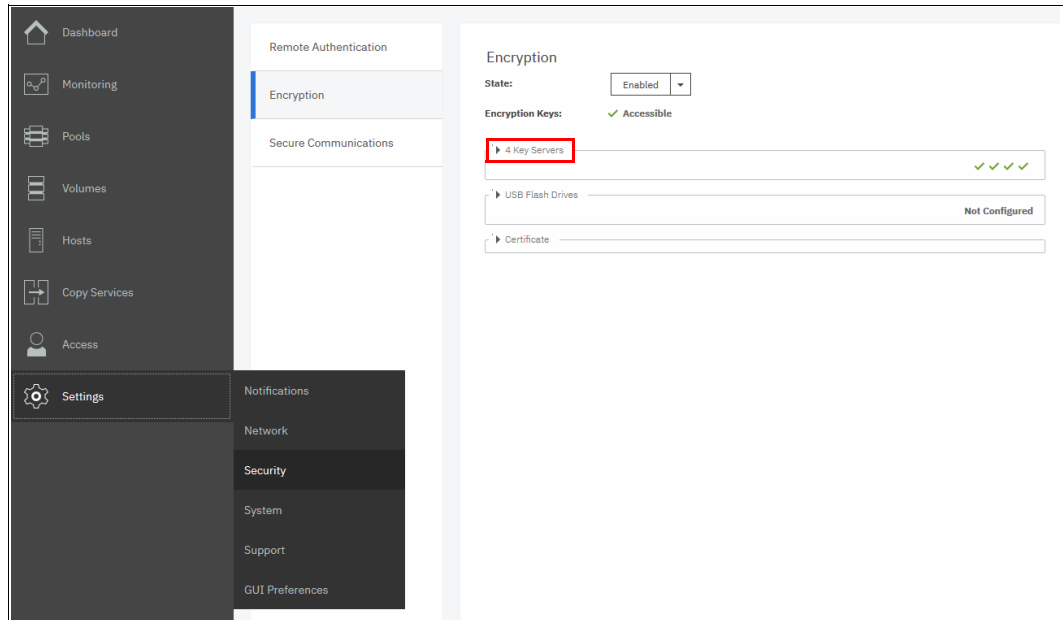


Figure 12-82 Locate Key Servers section on Encryption window

2. Click **Rekey**, as shown in Figure 12-83.

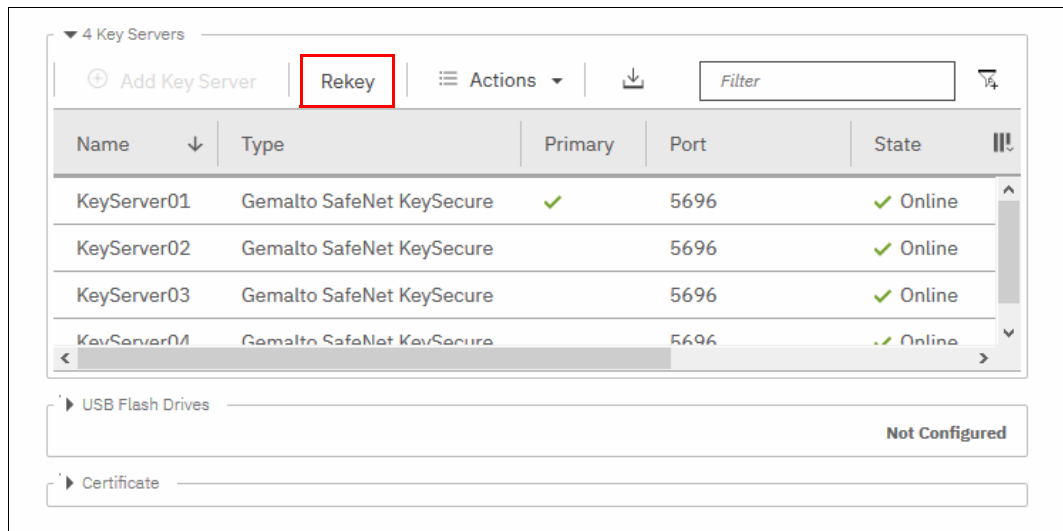


Figure 12-83 Start rekey on SKLM key server

3. Click **Yes** in the next window to confirm the rekey operation, as shown in Figure 12-84.

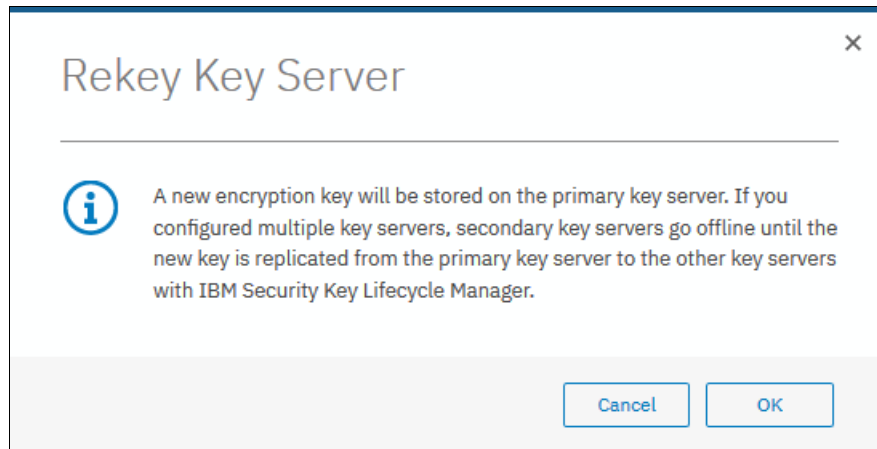


Figure 12-84 Confirm key server rekey operation

**Note:** The rekey operation is performed on only the primary key server that is configured in the system. If more key servers are configured apart from the primary key, they do not hold the updated encryption key until they obtain it from the primary key server. To restore encryption key provider redundancy after a rekey operation, replicate the encryption key from the primary key server to the secondary key servers.

You receive a message confirming that the rekey operation was successful, as shown in Figure 12-85.

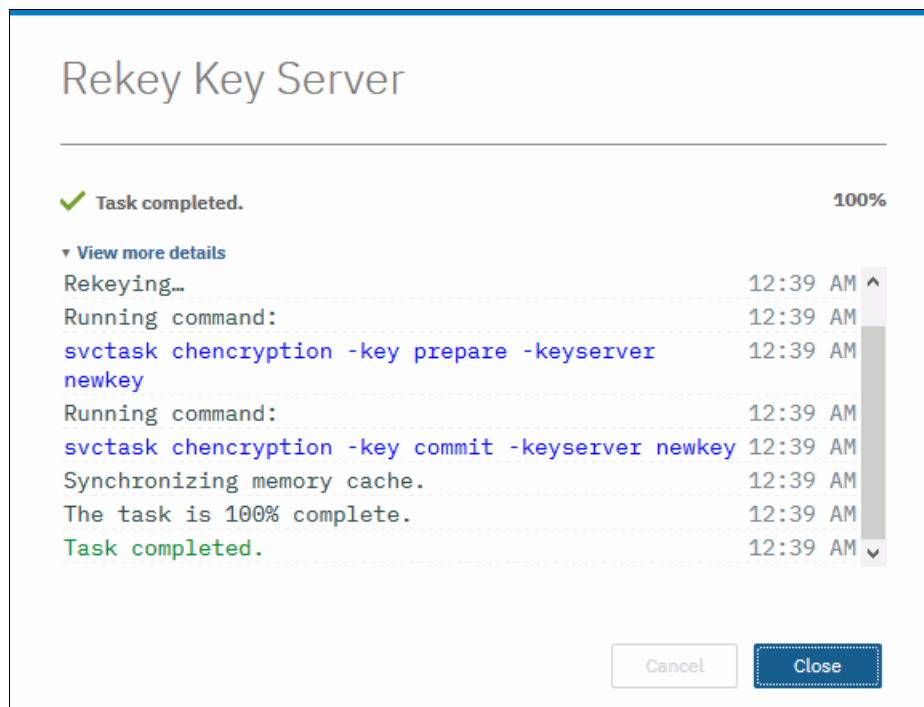


Figure 12-85 Successful key server rekey operation

## 12.10.2 Rekeying by using USB flash drives

During the rekey process, new keys are generated and copied to the USB flash drives. These keys are then used instead of the current keys. The rekey operation fails if at least one of the USB flash drives does not contain the current key. To rekey the system, you need at least three USB flash drives to store the master access key copies.

After the rekey operation is complete, update all other copies of the encryption key, including copies stored on other media. Take the same precautions to securely store all copies of the new encryption key as when you were enabling encryption for the first time.

To rekey the master access key on USB flash drives, complete the following steps:

1. Select **Settings** → **Security** → **Encryption**. Click **USB Flash Drives** to expand the section, as shown in Figure 12-86.

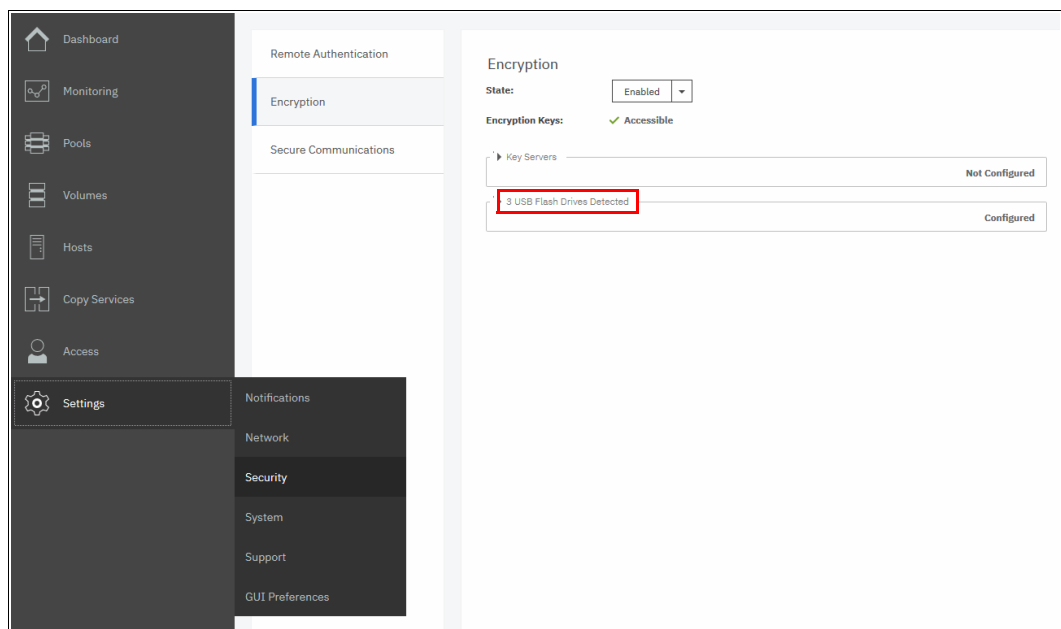


Figure 12-86 Locate USB Flash Drive section in the Encryption view

2. Verify that all USB drives are plugged into the system, detected, and show as Validated, as shown in Figure 12-87. Click **Rekey**. You need at least three USB flash drives, with at least one reported as Validated to process with rekey.

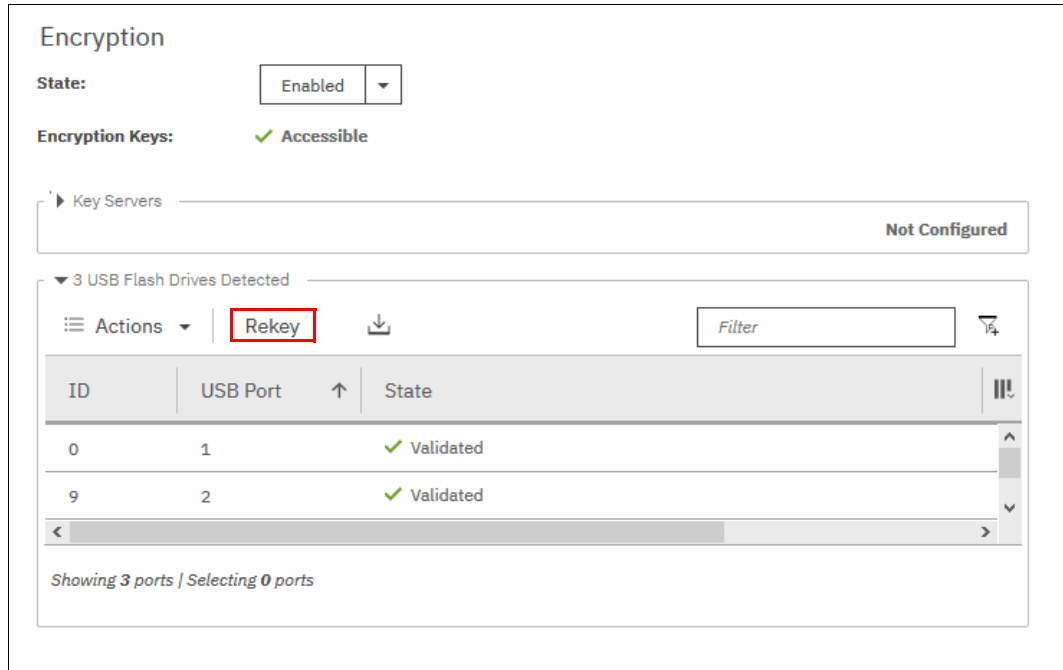


Figure 12-87 Start rekey on USB flash drives provider

3. If the system detects a validated USB flash drive and at least three available USB flash drives, new encryption keys are automatically copied on the USB flash drives, as shown in Figure 12-88 on page 751. Click **Commit** to finalize the rekey operation.



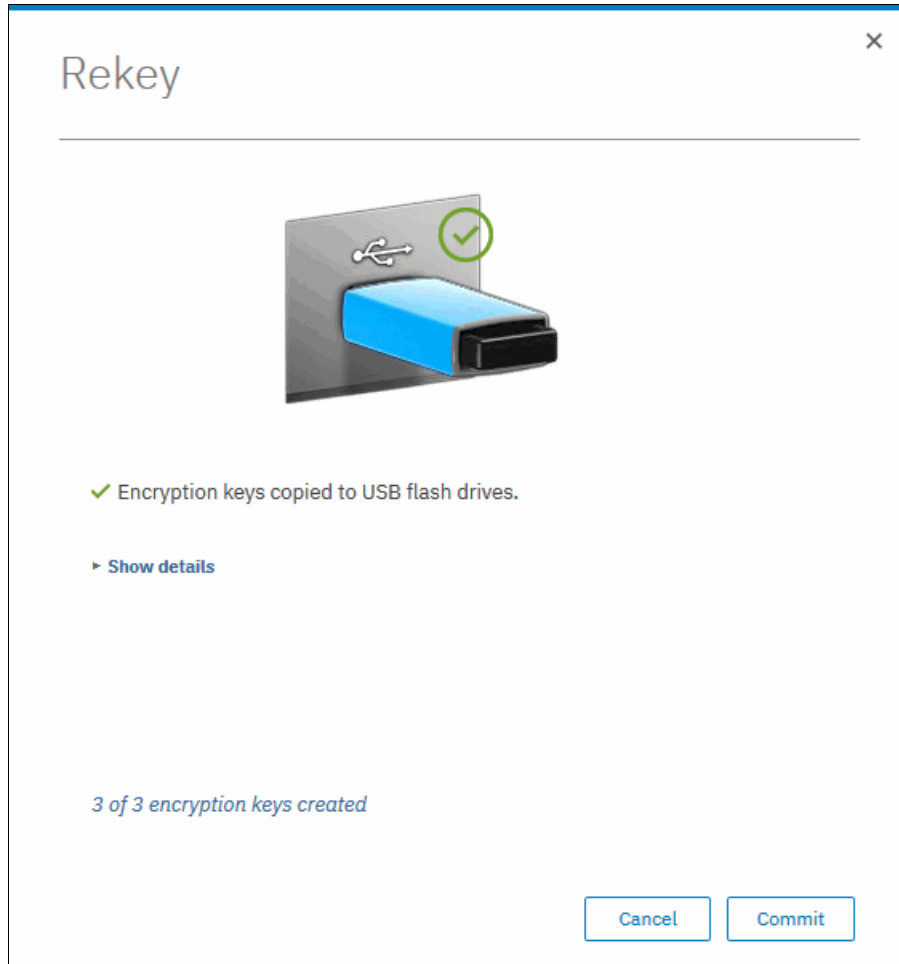


Figure 12-88 Writing new keys to USB flash drives

4. You receive a message confirming the rekey operation was successful, as shown in Figure 12-89. Click **Close**.

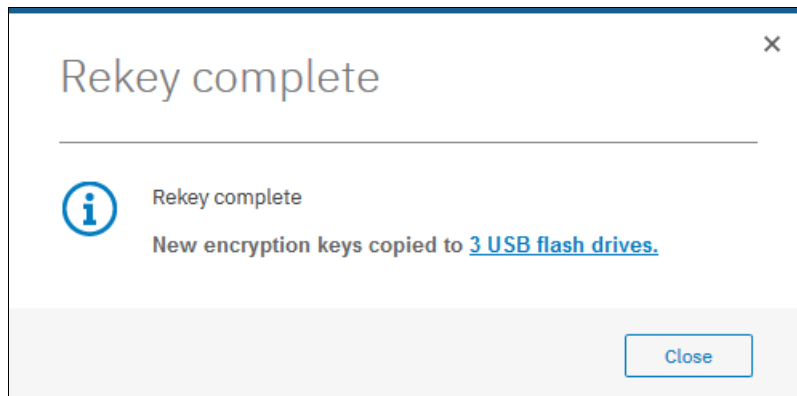


Figure 12-89 Successful rekey operation using USB flash drives

## 12.11 Disabling encryption

You are prevented from disabling encryption if any encrypted objects are defined apart from self-encrypting MDisks. You can disable encryption in the same way whether you use USB flash drives, key server, or both providers.

To disable encryption, complete the following steps:

1. Select **Settings** → **Security** → **Encryption** and click **Enabled**. If no encrypted objects exist, a menu is displayed. Click **Disabled** to disable encryption on the system. Figure 12-90 shows an example for a system with both encryption key providers configured.

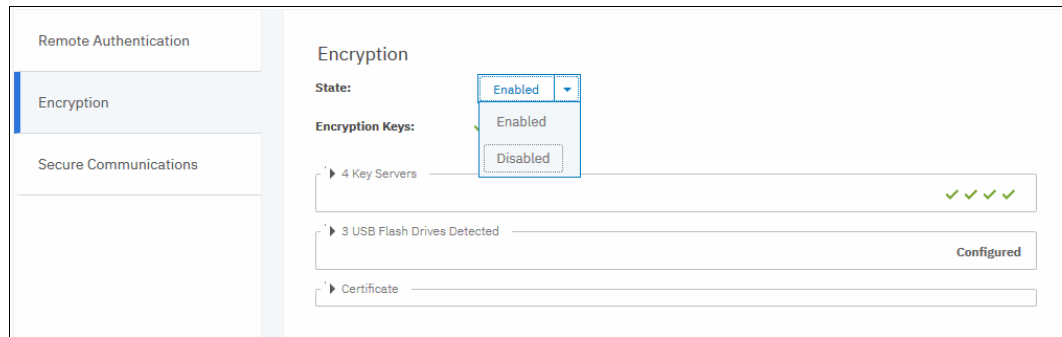


Figure 12-90 Disabling encryption on a system with both providers

2. You receive a message confirming that encryption was disabled. Figure 12-91 shows the message when a key server is used.

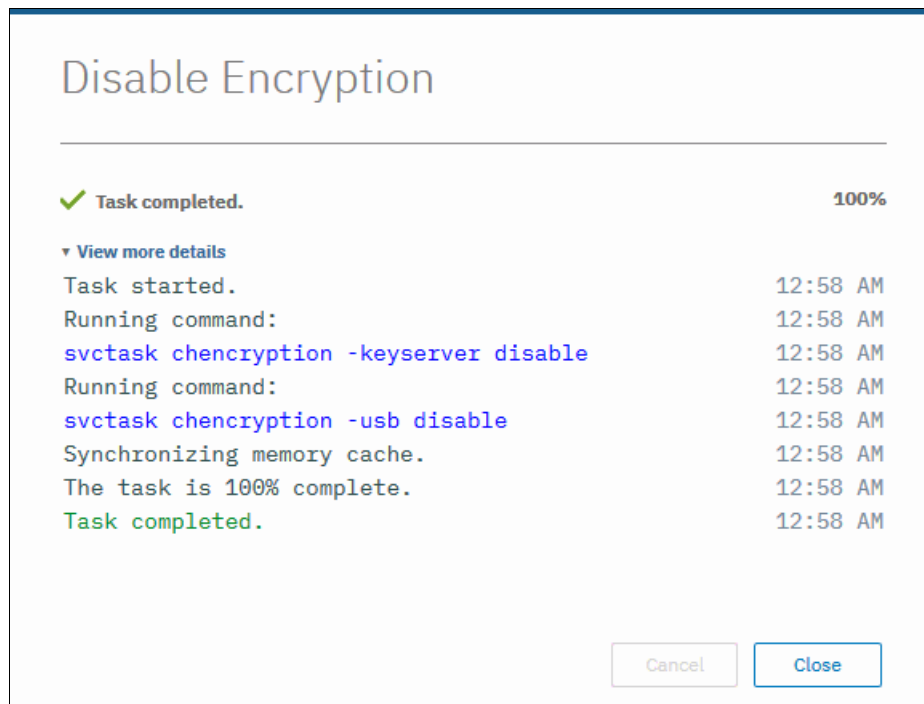



Figure 12-91 Encryption disabled



# Reliability, availability, and serviceability, and monitoring and troubleshooting

Many ways are available to manage, monitor, and troubleshoot IBM Spectrum Virtualize.

This chapter introduces useful, common procedures to maintain IBM Spectrum Virtualize. It includes the following topics:

- ▶ 13.1, “Reliability, availability, and serviceability” on page 754
- ▶ 13.2, “Shutting down an IBM SAN Volume Controller cluster” on page 759
- ▶ 13.3, “Removing or adding a node to or from the system” on page 762
- ▶ 13.4, “Configuration backup” on page 765
- ▶ 13.5, “Software update” on page 769
- ▶ 13.6, “Health checker feature” on page 781
- ▶ 13.7, “Troubleshooting and fix procedures” on page 783
- ▶ 13.8, “Monitoring” on page 790
- ▶ 13.9, “Audit log” on page 806
- ▶ 13.10, “Collecting support information by using the GUI, CLI, and USB” on page 809
- ▶ 13.11, “Service Assistant Tool” on page 816
- ▶ 13.12, “IBM Storage Insights Monitoring” on page 819

## 13.1 Reliability, availability, and serviceability

Reliability, availability, and serviceability (RAS) are important concepts in the design of the IBM Spectrum Virtualize system. Hardware features, software features, design considerations, and operational guidelines all contribute to make the system reliable.

Fault tolerance and high levels of availability are achieved by the following methods:

- ▶ The Redundant Array of Independent Disks (RAID) capabilities of the underlying disks
- ▶ IBM SAN Volume Controller nodes clustering by using a *Compass* architecture
- ▶ Auto-restart of hung nodes
- ▶ Integrated battery back up units (BBUs) to provide memory protection if there is a site power failure
- ▶ Host system failover capabilities by using N\_Port ID Virtualization (NPIV)
- ▶ Deploying advanced multi-site configurations such as HyperSwap and Stretched Cluster
- ▶ Hot Spare Node option to provide complete node redundancy and failover

High levels of serviceability are available through the following methods:

- ▶ Cluster error logging
- ▶ Asynchronous error notification
- ▶ Automatic dump capabilities to capture software detected issues
- ▶ Concurrent diagnostic procedures
- ▶ Directed Maintenance Procedures (DMPs) with guided online replacement procedures
- ▶ Concurrent log analysis and memory dump data recovery tools
- ▶ Concurrent maintenance of all IBM SAN Volume Controller components
- ▶ Concurrent upgrade of IBM Spectrum Virtualize software and firmware
- ▶ Concurrent addition or deletion of nodes in the clustered system
- ▶ Automatic software version leveling when replacing a node
- ▶ Detailed status and error conditions that are displayed by light-emitting diode (LED) indicators
- ▶ Error and event notification through Simple Network Management Protocol (SNMP), syslog, and email
- ▶ Optional Remote Support Assistant
- ▶ IBM Storage Insights

The heart of an IBM Spectrum Virtualize system is one or more pairs of *nodes*. The nodes share the read and write data workload from the attached hosts and to the disk arrays. This section examines the RAS features of an IBM SAN Volume Controller system, monitoring, and troubleshooting.

### 13.1.1 IBM SAN Volume Controller nodes

IBM SAN Volume Controller nodes work as a redundant clustered system. Each IBM SAN Volume Controller node is an individual server within the clustered system on which the IBM Spectrum Virtualize software runs.

IBM SAN Volume Controller nodes are always installed in pairs, forming an I/O group. A minimum of one pair and a maximum of four pairs of nodes constitute a clustered IBM SAN Volume Controller system. Many of the components that make up IBM SAN Volume Controller nodes include LEDs that indicate the status and activity of that component.

Figure 13-1 shows the ports, and indicator lights on the IBM SAN Volume Controller node model 2145-SV1.

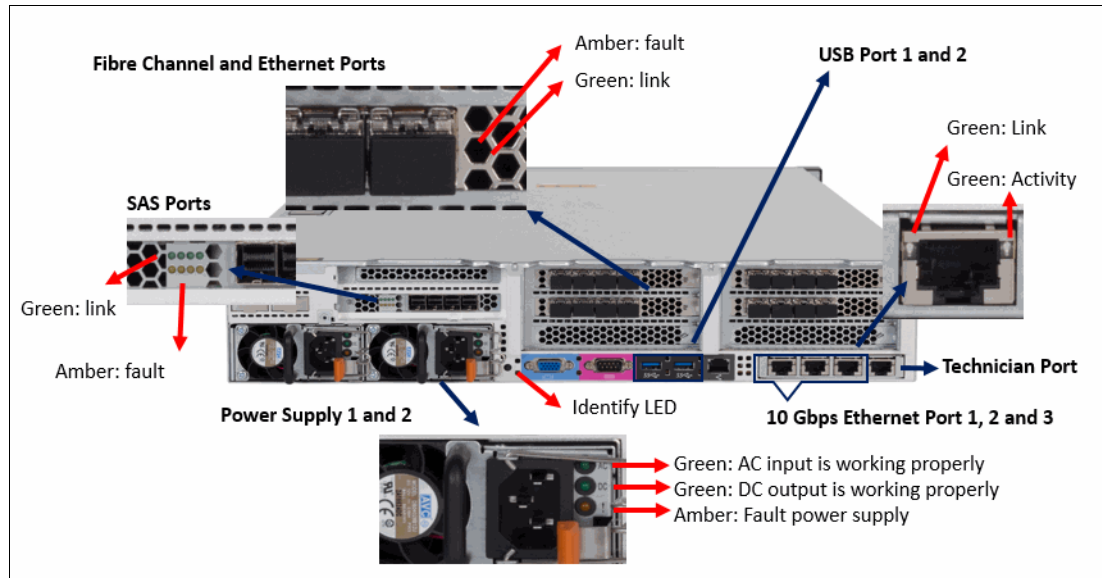


Figure 13-1 2145-SV1 Ports and Components

Figure 13-2 shows the ports, and indicator lights on the IBM SAN Volume Controller node model 2145-SV2/SA2.

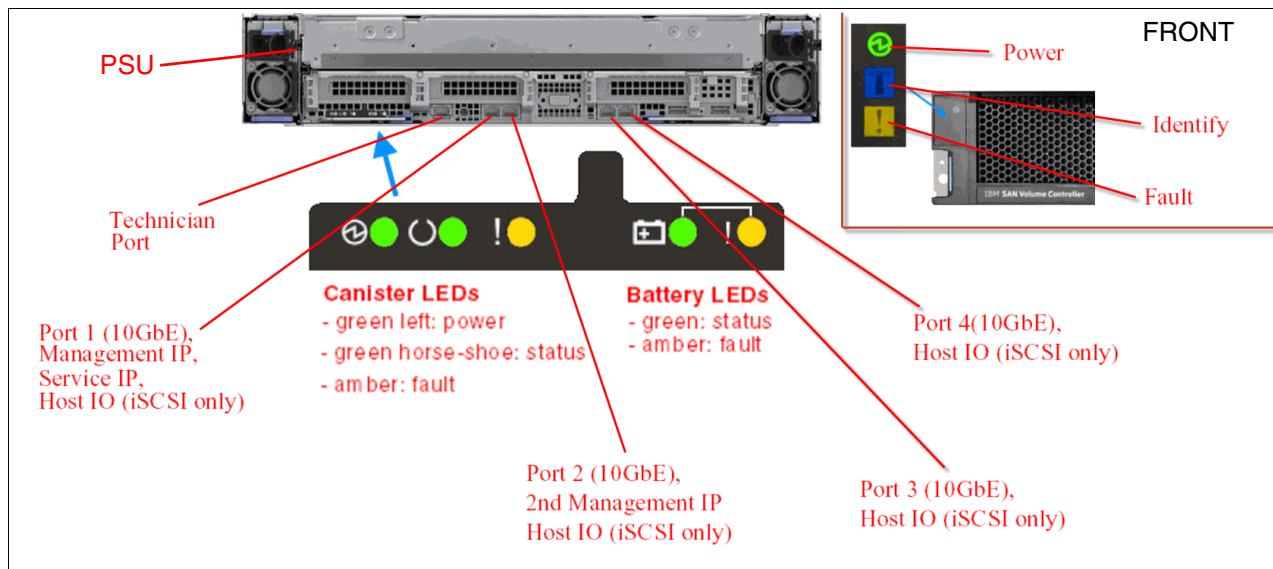


Figure 13-2 Rear ports and indicators of IBM SAN Volume Controller model 2145-SV2

**Note:** Unlike previous IBM SAN Volume Controller models, SV2/SA2 does not support direct serial-attached Small Computer System Interface (SCSI) (SAS) attached drive enclosures. Hence, these models do not ship with any SAS ports.

## Host interface cards

In this section, we focus on the new SV2/SA2 models. Each canister has three host interface card (HIC) slots. Nodes in the same I/O group must have the same HIC configuration. Available HIC's are listed in Table 13-1.

Table 13-1 Available host interface cards

Protocol	Supported number of cards	Ports	Allowed slots
16 Gb Fibre Channel (FC)/Fibre Channel over NVMe	0-3	4	1, 2, 3
32 Gb FC/Fibre Channel over NVMe	0-3	4	1, 2, 3
25 Gb iSCSI/RDMA over Converged Ethernet (RoCE)/Non-Volatile Memory Express (NVMe) over Ethernet	0-3	2	1, 2, 3
25 Gb iSCSI/Internet Wide-area RDMA Protocol (iWARP)/NVMe over Ethernet	0-3	2	1, 2, 3

**Note:** A Minimum of two host bus adapters (HBAs) of any type must be ordered to hit bandwidth requirements. Also, Fibre Channel over Ethernet (FCoE) is not supported.

Table 13-2 lists the meaning of port LEDs for an FC configuration.

Table 13-2 Fibre Channel link LED statuses

Port LED	Color	Meaning
Link status	Green	Link is up, connection established.
Speed	Amber	Link is not up or there is a speed fault.

## USB ports

Two active USB connectors are available in the horizontal position. They have no numbers and no indicators are associated with them. These ports can be used for initial cluster setup, encryption key backup, and node status or log collection.

## Ethernet Ports and their LED status

Five 10-Gigabit Ethernet (GbE) ports are on each canister. However, not all ports are equal, and their function are listed in Table 13-3. The location of the technician port on a node canister is shown in Figure 13-2 on page 755.

Table 13-3 Onboard Ethernet Ports and their function

Onboard Ethernet port	Speed	Function
1	10 GbE	Management IP, Service IP, and Host I/O (Internet Small Computer Systems Interface (iSCSI) only)

Onboard Ethernet port	Speed	Function
2	10 GbE	Secondary Management IP, and Host I/O (iSCSI only)
3	10 GbE	Host I/O (iSCSI only)
4	10 GbE	Host I/O (iSCSI only)
T	1 GbE	Technician Port: Dynamic Host Configuration Protocol (DHCP)/Domain Name System (DNS) for direct attach service management

Ethernet LED Statuses are listed in Table 13-4.

Table 13-4 Ethernet LED statuses

LED	Color	Meaning
Link state	Green	It is on when there is an Ethernet link.
Activity	Amber	It is flashing when there is activity on the link.

### Node front status LEDs

There are three LEDs in a row on the front of the node that indicates the status and the functions of the node (see Table 13-5).

Table 13-5 Node Front LEDs

Position	Color	Name	State	Meaning
Top	Green	Power	On	The CPU is active (not in standby) in the node (canister).
			Slow flash	Canister is in X86 standby mode
			Off	No power to the canister or it is running on battery.
Middle	Blue	Identify	On	Node is being identified. Fault LED on rear of canister will also be flashing.
			Off	Default
Bottom	Amber	Fault	On	Error. Node is in service state or there is a problem with the node starting. Also illuminates when battery is charging and node is waiting to come up.
			Off	No Faults. Candidate/Active State

### Node rear status LEDs

There are three LEDs in a row on the rear of the node that indicates the status and the functions of the node (see Table 13-6 on page 758).

Table 13-6 Rear Node LEDs

Position	Color	Name	State	Meaning
Left	Green	Power	On	The node is started and active. If the fault LED is off, the node is an active member of a cluster or candidate. If the fault LED is also on, node is in service state or in error, preventing the software to start.
			Flashing (2 Hz)	Node is started and in standby mode.
			Flashing (4 Hz)	Node is running power-on self-test (POST).
			Off	No power to the node or it is running on battery.
Middle	Green	Status	On	The node is a member of a cluster.
			Flashing (2 Hz)	The node is a candidate for or in a service state.
			Flashing (4 Hz)	The node is performing a fire hose dump. Never unplug the node now.
			Off	No power to the node or node is in standby mode.
Right	Amber	Fault	On	The node is in a service state, or in error, for example, a POST error that is preventing the software from starting.
			Flashing (2 Hz)	Node is being identified.
			Off	Node is either in candidate or active state.

## Battery LEDs

Immediately to the right of the canister LEDs, with a short gap between them, are the Battery LEDs that provide the status of the battery (see Table 13-7).

Table 13-7 Battery LEDs

Position	Color	Name	State	Meaning
Left	Green	Status	On	Indicates that the battery is fully charged and thus has sufficient charge to complete two fire hose dumps.
			Flashing (2 Hz)	Indicates that the battery has sufficient charge to complete a single fire hose dump only.
			Flashing (4 Hz)	Indicates that the battery is charging and has insufficient charge to complete a single fire hose dump.
			Off	Indicates that the battery is not available for use (for example, missing or containing a fault)
Right	Amber	Fault	On	Indicates that a battery has a fault or condition occurred. The node enters service state.
			Off	Indicates that there are no known battery faults or conditions. An exception to this would be where a battery has insufficient charge to complete a single fire hose dump. Refer to 'status' LED.



## 13.1.2 Power

IBM SAN Volume Controller nodes accommodate two power supply units (PSUs) for normal operation. For this reason, it is highly advised to supply AC power to each PSU from different power distribution units (PDUs).

A fully charged battery can perform two *fire hose dumps*. It supports the power outage for 5 seconds before starting safety procedures. A fully charged battery can perform two fire hose dumps. A *fire hose dump* is a process where a node stores cache and system data to an internal drive if a power failure occurs.

Figure 13-3 shows two PSUs that are present in the new IBM SAN Volume Controller SV2/SA2 node. The controller PSU on a SV2/SA2 model has one LED that can be green or amber, depending on the status of the PSU. If the LED is off, that means there is no AC power to the entire enclosure.

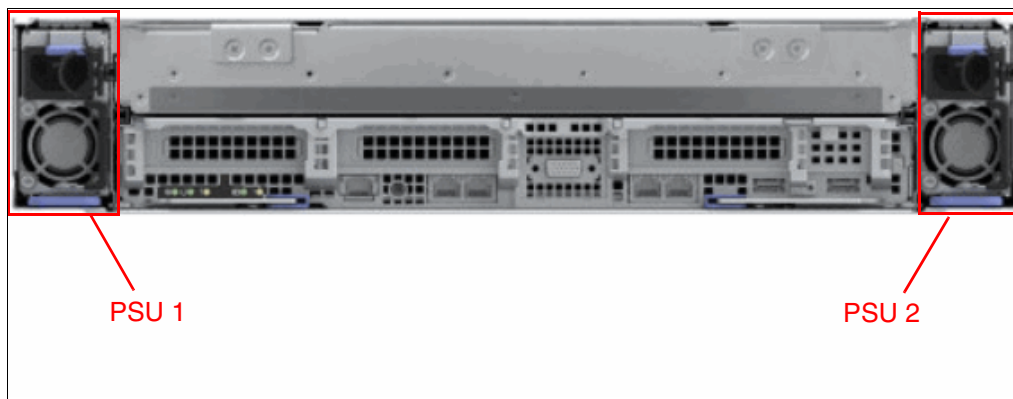


Figure 13-3 IBM SAN Volume Controller PSU 1 and 2

Power supplies in nodes are hot-swappable and replaceable without needing to shut down a node or cluster.

If the power is interrupted in one node for less than 5 seconds, the node does not perform a fire hose dump, and continues operation from the battery. This feature is useful for a case of, for example, maintenance of UPS systems in the data center or replugging the power to a different power source or PDU unit.

## 13.2 Shutting down an IBM SAN Volume Controller cluster

You can safely shut down an IBM SAN Volume Controller cluster by using the GUI or CLI.

**Important:** Never shut down your IBM SAN Volume Controller cluster by powering off the PSUs, removing both PSUs, or removing both power cables from the nodes. Doing so can lead to inconsistency or loss of the data that is staged in the cache.

Before shutting down the cluster, stop all hosts that allocated volumes from the device. This step can be skipped for hosts that have volumes that are also provisioned with mirroring (host-based mirror) from different storage devices. However, doing so incurs errors that are related to lost storage paths or disks on the host error log.

You can shut down a single node or shut down the entire cluster. When you shut down only one node, all activities remain active. When you shut down the entire cluster, you must power on the nodes locally to start the system again.

If all input power to the IBM SAN Volume Controller clustered system is removed for more than a few minutes (for example, if the machine room power is shut down for maintenance), it is important that you shut down the IBM SAN Volume Controller system before you remove the power.

Shutting down the system while it is still connected to the main power ensures that the internal node batteries are still fully charged when the power is restored.

If you remove the main power while the system is still running, the internal batteries detect the loss of power and start the node shutdown process. This shutdown can take several minutes to complete. Although the internal batteries have sufficient power to perform the shutdown, you drain the nodes batteries unnecessarily.

When power is restored, the nodes start. However, if the nodes batteries have insufficient charge to survive another power failure so that the node can perform another clean shutdown, the node enters service mode. You do *not* want the batteries to run out of power in the middle of the node's shutdown.

It can take approximately 3 hours to charge the batteries sufficiently for a node to come online.

**Important:** When a node shuts down because of a power loss, the node dumps the cache to an internal flash drive so that the cached data can be retrieved when the system starts again.

The IBM SAN Volume Controller internal batteries are designed to survive at least two power failures in a short period. After that period, the nodes do not come online until the batteries have sufficient power to survive another immediate power failure.

During maintenance activities, if the batteries detect power and then detect a loss of power multiple times (the nodes start and shut down more than once in a short time), you might discover that you unknowingly drained the batteries. In this case, you must wait until they are charged sufficiently before the nodes start again.

To shut down your IBM SAN Volume Controller system, complete the following steps:

1. From the **Monitoring** → **System** pane, click **System Actions**, as shown in Figure 13-4 on page 761. Click **Power Off System**.

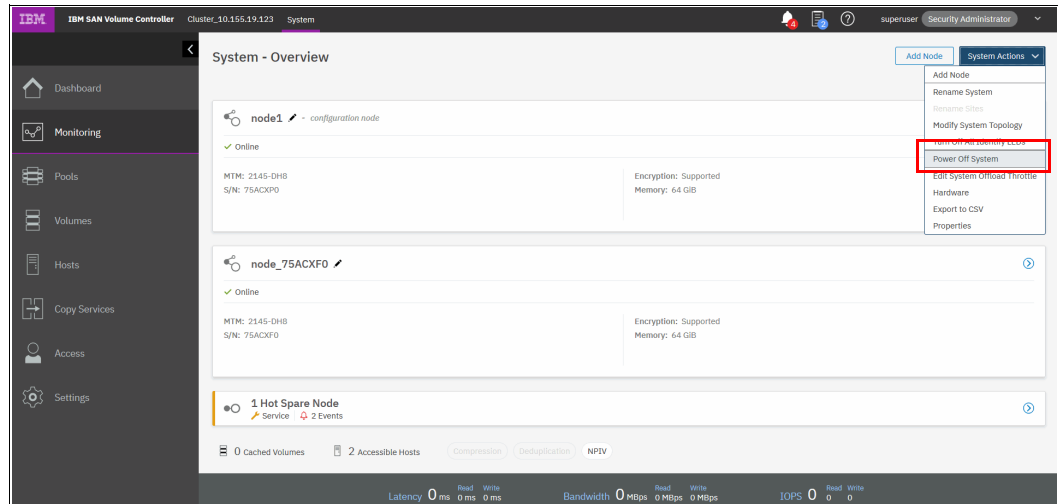


Figure 13-4 Action pane to power off the system

A confirmation window opens, as shown in Figure 13-5.

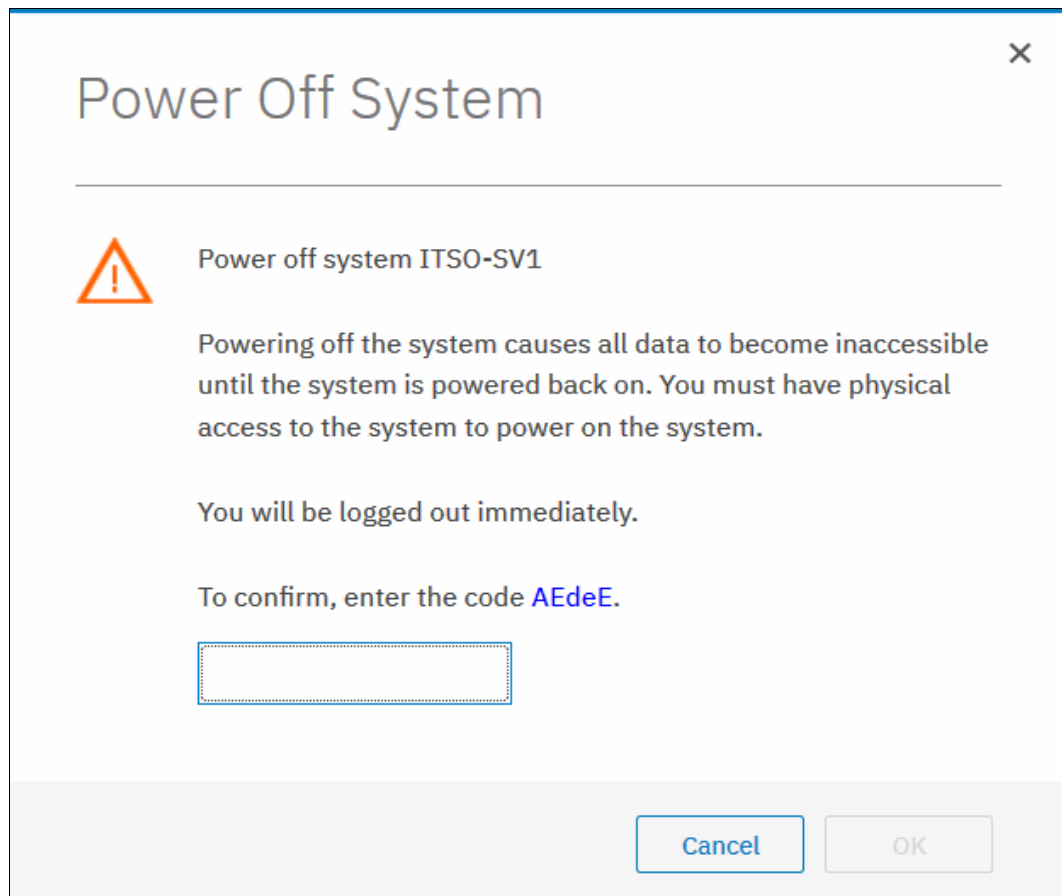


Figure 13-5 Confirmation window to confirm the shutdown of the system

2. Before you continue, ensure that you stopped all FlashCopy mappings, Remote Copy (RC) relationships, data migration operations, and forced deletions.
3. Enter the generated confirmation code and click **OK** to begin the shutdown process.

**Attention:** Pay special attention when encryption is enabled on some storage pools. You must insert a USB drive with the stored encryption keys or you must ensure that your IBM SAN Volume Controller can communicate with the IBM Security Key Lifecycle Manager (SKLM) server or clone servers to retrieve the encryptions keys. Otherwise, the data is unreadable after restart.

## 13.3 Removing or adding a node to or from the system

There are situations in which IBM Support might ask you to remove a node from the system briefly. One typical use case for this circumstance is if a node becomes stuck in a code upgrade, you can remove the node from the cluster briefly to commit the upgrade and complete (or cancel) the procedure depending on how many nodes upgraded so far. However, this procedure should only be done under the direction of IBM Support.

The easiest way is by using a CLI command. Run the `svcinfo lsnode` command to display all nodes and their ID and status, as shown in Figure 13-6. You can make sure that each IOgroup has two nodes online (or that if you remove a node, that one node remains in the IOgroup to continue serving I/O).

```
IBM_Storwize:fab2shared1-cl:superuser>svcinfo lsnode
id name   UFS_serial_number WWNN          status IO_group_id IO_group_name config_node UFS_unique_id hardware iscsi_name          iscsi_alias panel_name enclosure
e_id canister_id enclosure serial_number site_id site_name
1 node1   7822PBR      500507680B0080FA online 0      io_grp0    no          600          iqn.1986-03.com.ibm:2145.fab2shared1-cl.node1 01-1      1
3 node2   7822PBR      500507680B0080FB online 0      io_grp0    yes         600          iqn.1986-03.com.ibm:2145.fab2shared1-cl.node2 01-2      1
4 node3   7822PBR      500507680B0080E8 online 1      io_grp1    no          600          iqn.1986-03.com.ibm:2145.fab2shared1-cl.node3 02-1      2
5 node4   7822PFG      500507680B0080E9 online 1      io_grp1    no          600          iqn.1986-03.com.ibm:2145.fab2shared1-cl.node4 02-2      2
IBM_Storwize:fab2shared1-cl:superuser>C
IBM_Storwize:fab2shared1-cl:superuser>
```

Figure 13-6 *lsnode* output

In this example, we remove node 1 from the cluster. Run the `svctask rmnode 1` command, as shown in Example 13-1.

### Example 13-1 *rmnode*

```
ssh superuser@9.71.47.7
Password:
IBM_Storwize:fab2shared1-cl:superuser>svctask rmnode 1
```

A node can also be removed by using the GUI. Select **Monitoring** and then, **System**. You can then select the relevant node that you want to remove, which brings up the Node Details Window. Select **Node Actions** → **Remove**, as shown in Figure 13-7 on page 763, which brings up a confirmation window before removing the node from the cluster.

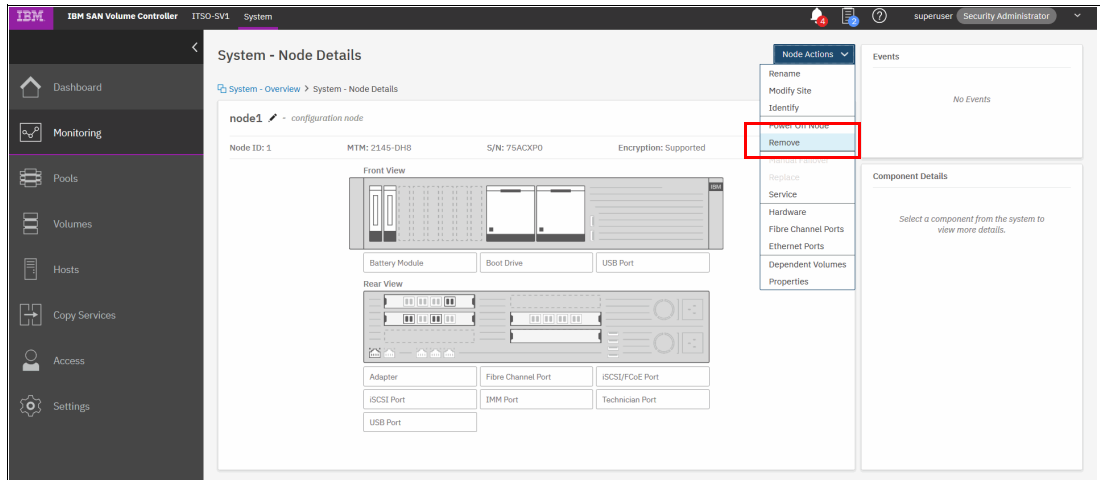


Figure 13-7 Removing a node by using the GUI

After you remove the node, if you rerun the `svcinfg 1snode` command, you will see that it disappeared from the cluster, as shown in Figure 13-8. The Service Assistant Tool (SAT) and GUI also reflect that there are now only three nodes in the cluster.

```

IBM_storwize:fab2shared1-cl:superuser>svcinfg 1snode
id name UPS_serial_number WNN status IO_group_id IO_group_name config_node UPS_unique_id hardware iscsi_name iscsi_alias panel_name enclosure
0 id canister_id enclosure_serial_number site_id site_name
3 node2 7822PBR 500507680B00808 online 0 io_grp0 yes 600 iqn.1986-03.com.ibm:2145.fab2shared1-cl.node2 01-2 1
4 node3 7822PBR 500507680B00808 online 1 io_grp1 no 600 iqn.1986-03.com.ibm:2145.fab2shared1-cl.node3 02-1 2
5 node4 7822PFG 500507680B00808 online 1 io_grp1 no 600 iqn.1986-03.com.ibm:2145.fab2shared1-cl.node4 02-2 2
IBM_storwize:fab2shared1-cl:superuser>

```

Figure 13-8 Isnode output after removing a node

**Note:** By default, the cache is flushed before the node is deleted to prevent data loss if a failure occurs on the other node in the I/O group. This process incurs a delay after you remove a node to when it comes back up as candidate status.

After a brief wait period, check the SAT, which shows the node you removed in service or candidate status, as shown in Figure 13-9.

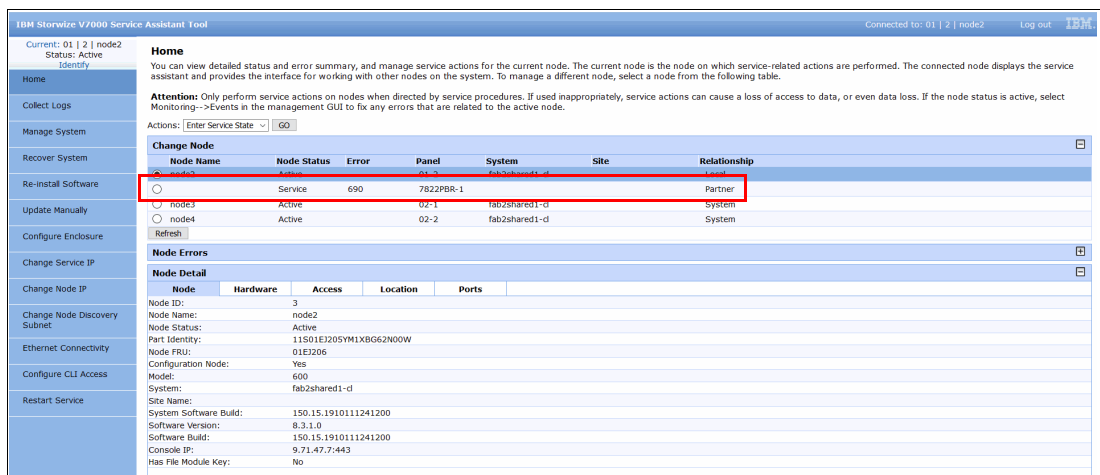


Figure 13-9 Service Assistant Tool post-node removal

Select the radio button of the node that is in service and then select **Exit Service State** from the Actions drop-down menu. Click **GO** and a confirmation pane opens, as shown in Figure 13-10.

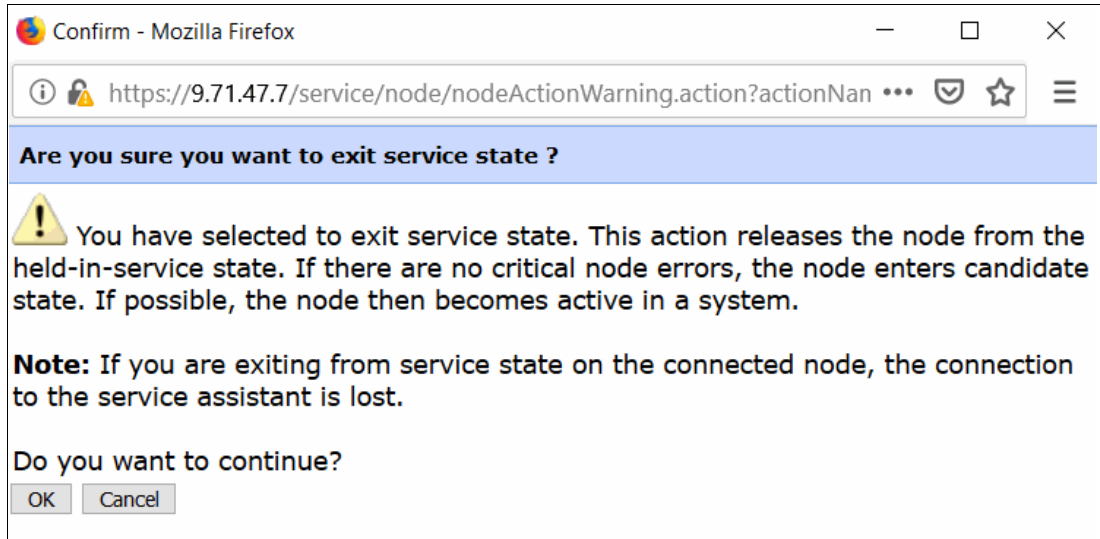


Figure 13-10 Exit service state

A confirmation window opens in which it is confirmed that the node exited the service state. Click **OK** or close the window and click the **Refresh** button under the list of the nodes. The new node now shows as in Candidate status, as shown in Figure 13-11.

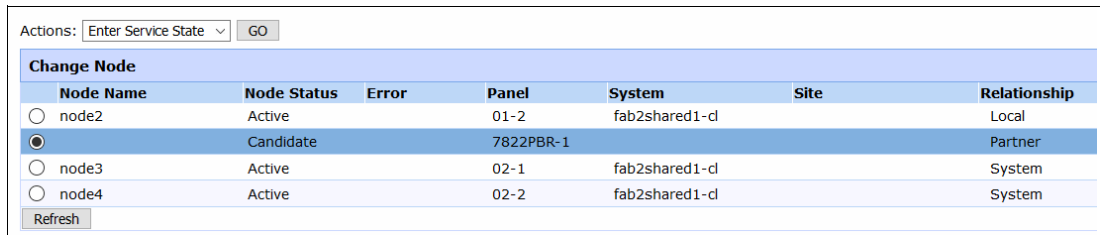


Figure 13-11 Node in candidate status

The node should be automatically added again to the system. If not, make a note of the numbers in the Panel column, and return to your CLI session. Run the **addnode** command and specify the panel ID to add the node back into the cluster, as shown in Example 13-2.

*Example 13-2 Addnode*

```
ssh superuser@9.71.47.7
Password:
IBM_Storwize:fab2shared1-cl:superuser>svctask addnode -panelname 7822PBR-1 -iogrp
io_grp0
```

Now, check by using the CLI command **svcinfo lsnode** again or check the SAT to ensure that the node was added back, as shown in Figure 13-12 on page 765.

Actions: Enter Service State GO

Node Name	Node Status	Error	Panel	System	Site	Relationship
<input checked="" type="radio"/> node2	Active		01-2	fab2shared1-d		Local
<input type="radio"/> node1	Active		01-1	fab2shared1-d		Partner
<input type="radio"/> node3	Active		02-1	fab2shared1-d		System
<input type="radio"/> node4	Active		02-2	fab2shared1-d		System

Refresh

Figure 13-12 The node is back in the cluster

## 13.4 Configuration backup

You can download and save the configuration backup file by using the IBM Spectrum Virtualize GUI or command-line interface (CLI). On an *ad hoc* basis, manually perform this procedure because you can save the file directly to your workstation. The CLI option requires you to log in to the system and download the dumped file by using a specific Secure Copy Protocol (SCP). The CLI option is a best practice for an automated backup of the configuration.

**Important:** Generally, perform a daily backup of the IBM Spectrum Virtualize configuration backup file. The best practice is to automate this task. Always perform an extra backup before any critical maintenance task, such as an update of the IBM Spectrum Virtualize software version.

The backup file is updated by the cluster every day. Saving it after you make any changes to your system configuration is also important. It contains the configuration data of arrays, pools, volumes, and so on. The backup does *not* contain any data from the volumes.

The following prerequisites must be met to successfully perform the configuration backup:

- ▶ All nodes must be online.
- ▶ No independent operations that change the configuration can be running in parallel.
- ▶ No object name can begin with an underscore.

**Important:** An *ad hoc* backup of configuration can be done only from the CLI by running the **svconfig backup** command. Then, the output of the command can be downloaded from the GUI.

### 13.4.1 Backing up by using the CLI

You can use the CLI to trigger a configuration backup manually or by using a regular automated process. The **svconfig backup** command generates a new backup file. Triggering a backup by using the GUI is not possible. However, you can choose to save the automated 1AM cron backup if no configuration changes were made,

Example 13-3 shows output of the **svconfig backup** command.

*Example 13-3 Saving the configuration by using the CLI*

```
IBM_2145:ITS0-SV1:superuser>svconfig backup
.....
.....
.....
CMMVC6155I SVCONFIG processing completed successfully
IBM_2145:ITS0-SV1:superuser>
```

The **svcconfig backup** command generates three files that provide information about the backup process and cluster configuration. These files are dumped into the /tmp directory on the configuration node. Use the **lsdumps** command to list them (see Example 13-4).

*Example 13-4 Listing the backup files by using the CLI*

```
IBM_2145:ITS0-SV1:superuser>lsdumps |grep backup
87 svc.config.backup.log_CAY0009
88 svc.config.backup.sh_CAY0009
89 svc.config.backup.xml_CAY0009
IBM_2145:ITS0-SV1:superuser>
```

Table 13-8 lists the three files that are created by the backup process.

*Table 13-8 Files that are created by the backup process*

File name	Description
svc.config.backup.xml	This file contains your cluster configuration data.
svc.config.backup.sh	This file contains the names of the commands that were run to create the backup of the cluster.
svc.config.backup.log	This file contains details about the backup, including any error information that might be reported.

Save the current backup to a secure and safe location. The files can be downloaded by using **scp** (UNIX) or **pscp** (Windows), as shown in Example 13-5. Replace the IP address with the cluster IP address of your IBM SAN Volume Controller and specify a local folder on your workstation. In this example, we are saving to C:\SVCbackups.

*Example 13-5 Saving the config backup files to your workstation*

```
C:\putty>
pscp -unsafe superuser@10.41.160.201:/dumps/svc.config.backup.* c:\SVCbackups
Using keyboard-interactive authentication.
Password:
svc.config.backup.log_CAY | 33 kB | 33.6 kB/s | ETA: 00:00:00 | 100%
svc.config.backup.sh_CAYO | 13 kB | 13.9 kB/s | ETA: 00:00:00 | 100%
svc.config.backup.xml_CAY | 312 kB | 62.5 kB/s | ETA: 00:00:00 | 100%
```

```
C:\>dir SVCbackups
Volume in drive C has no label.
Volume Serial Number is 0608-239A
```

```
Directory of C:\SVCbackups

17.10.2018 09:20 <DIR> .
17.10.2018 09:20 <DIR> ..
17.10.2018 09:20          34.415 svc.config.backup.log_CAY0009
17.10.2018 09:20          14.219 svc.config.backup.sh_CAY0009
17.10.2018 09:20          319.820 svc.config.backup.xml_CAY0009
                3 File(s)          368.454 bytes
                2 Dir(s) 78.243.868.672 bytes free

C:\>
```



By using the **-unsafe** option, you can use a wildcard for downloading all the `svc.config.backup` files by using a single command.

**Tip:** If you encounter the Fatal: Received unexpected end-of-file from server error when the `pscp` command is run, consider upgrading your version of PuTTY.

### 13.4.2 Saving the backup by using the GUI

Although it is not possible to generate an ad hoc backup, you can save the backup files by using the GUI. To do so, complete the following steps:

1. Select **Settings** → **Support** → **Support Package**.
2. Click the **Manual Upload Instructions** twistie to expand it.
3. Click **Download Support Package** and then, **Download Existing Package**, as shown in Figure 13-13.

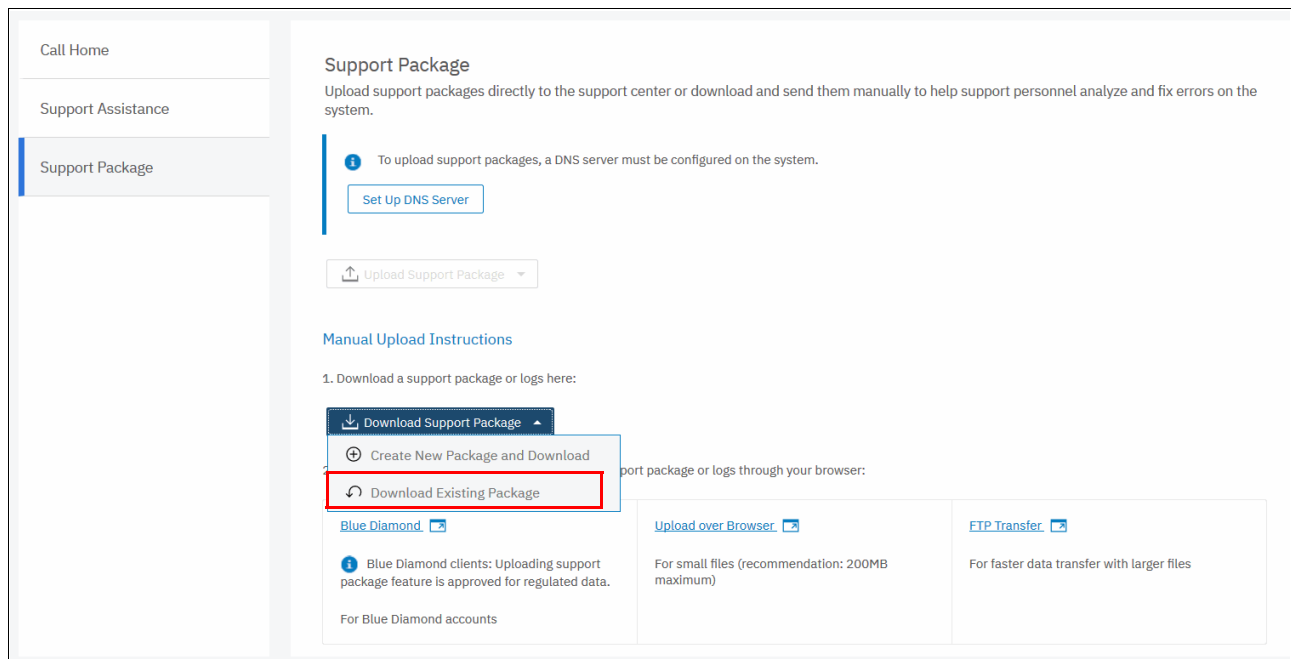


Figure 13-13 Download Support Package

The Support Package selection window opens, as shown in Figure 13-14.

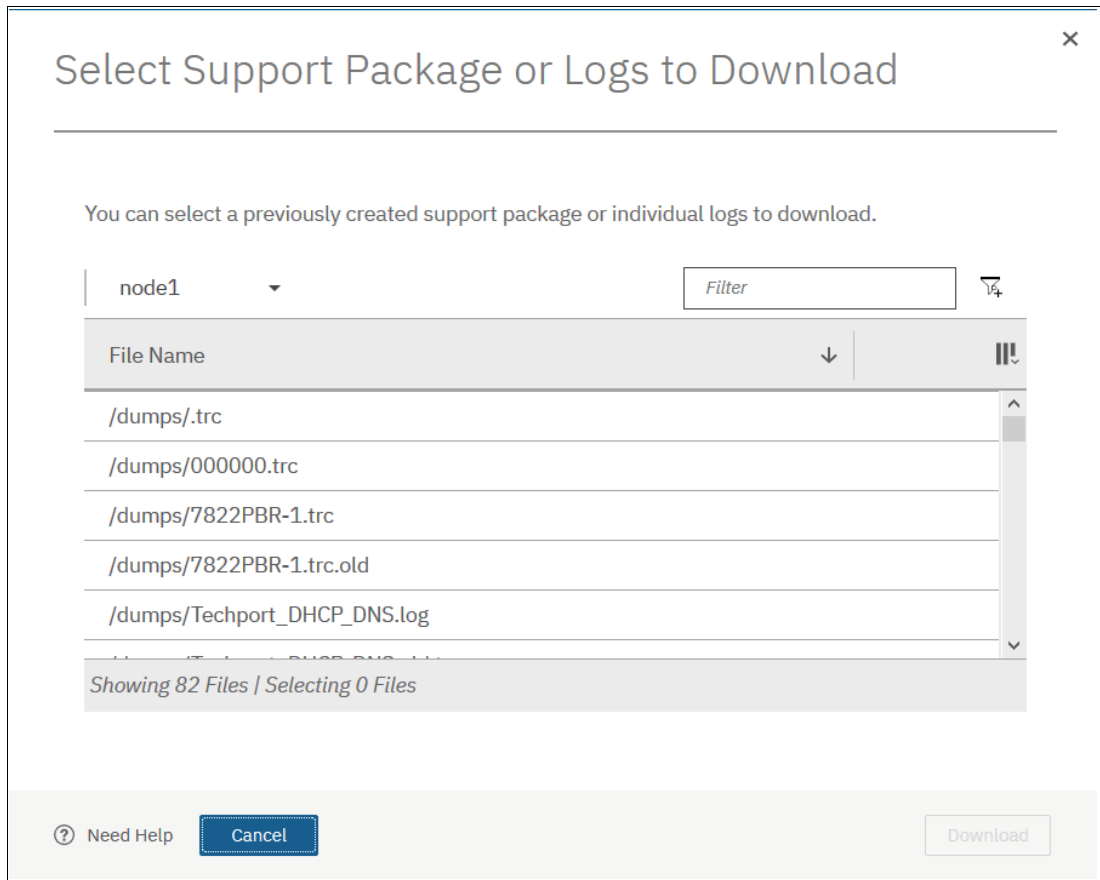


Figure 13-14 Select Support Package

4. Filter the view by clicking in the **Filter** box, entering backup, and pressing **Enter**, as shown in Figure 13-15.

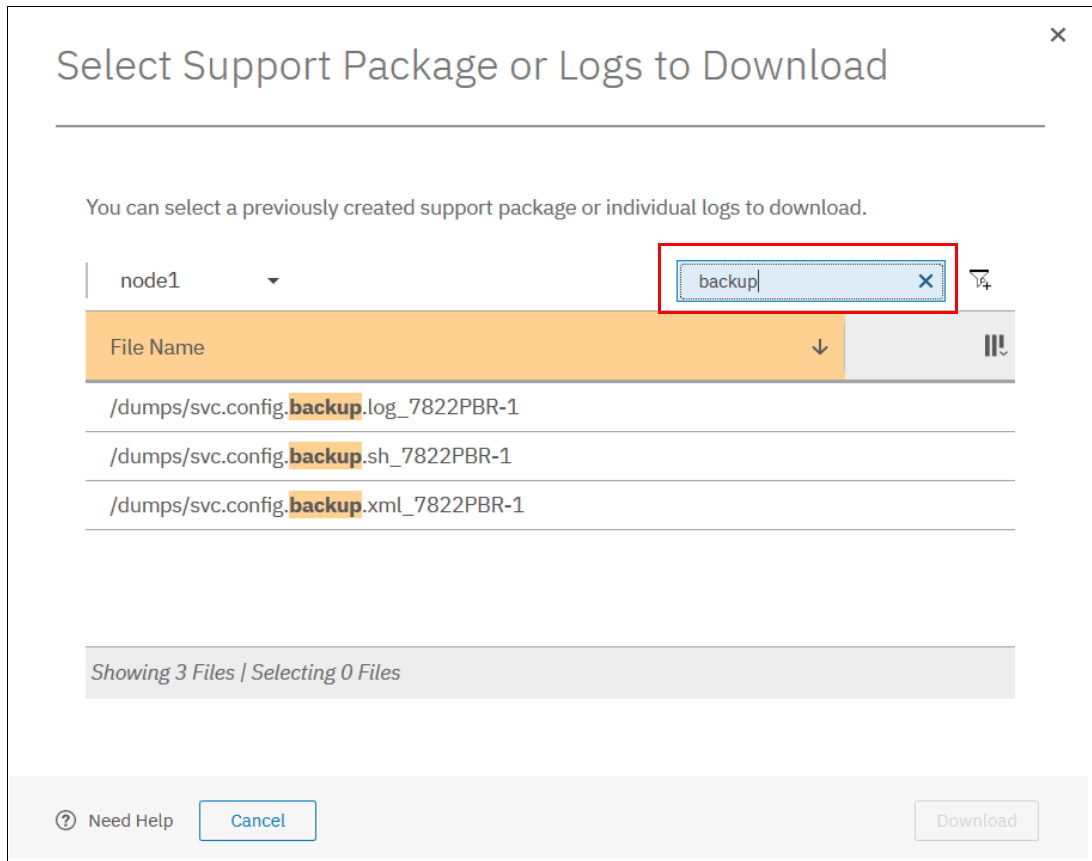


Figure 13-15 Filtering specific files for download

Select all of the files to include in the compressed file; then, click **Download**. Depending on your browser preferences, you might be prompted where to save the file or it downloads to your defined download directory.

**Note:** You must select the configuration node from the node drop-down menu in the upper left corner of the pane. This is because the backup files are stored on the configuration node.

## 13.5 Software update

This section describes the operations to update IBM Spectrum Virtualize V8.3.1.

The format for the software update package name ends in four positive integers that are separated by dots. For example, a software update package might have the following name:

IBM2145\_INSTALL\_8.3.1.0.

### 13.5.1 Precautions before the update

This section describes the precautions that you should take before you attempt an update.

**Important:** Before you attempt any IBM Spectrum Virtualize code update, read and understand the concurrent compatibility and code cross-reference matrix. For more information, see the [SAN Volume Controller Code Page](#) and click **Latest IBM Spectrum Virtualize code**.

During the update, each node in your clustered system is automatically shut down and restarted by the update process. Because each node in an I/O group provides an alternative path to volumes, use the Subsystem Device Driver (SDD) to make sure that all I/O paths between all hosts and storage area networks (SANs) work.

If you do not perform this check, certain hosts might lose connectivity to their volumes and experience I/O errors.

### 13.5.2 IBM Spectrum Virtualize upgrade test utility

The software upgrade test utility is a software instruction utility that checks for known issues that can cause problems during an IBM Spectrum Virtualize software update. More information about the utility is available at [this web page](#).

Download the software upgrade test utility from [this web page](#). This procedure ensures that you receive the current version of this utility. You can use the `svcupgradetest` utility to check for known issues that might cause problems during a software update.

The software upgrade test utility can be downloaded in advance of the update process. Alternately, it can be downloaded and run directly during the software update, as guided by the update wizard.

You can run the utility multiple times on the same system to perform a readiness check-in preparation for a software update. Run this utility for a final time immediately before you apply the IBM Spectrum Virtualize update to ensure that there were no new releases of the utility since it was originally downloaded.

The installation and use of this utility is nondisruptive, and does not require a restart of any nodes. Therefore, there is no interruption to host I/O. The utility is installed on only the current configuration node.

System administrators must continue to check whether the version of code that they plan to install is the latest version. You can obtain the current information from the [SAN Volume Controller Support Page](#).

This utility is intended to supplement rather than duplicate the existing tests that are performed by the IBM Spectrum Virtualize update procedure (for example, checking for unfixed errors in the error log).

A concurrent software update of all components is supported through the standard Ethernet management interfaces. However, during the update process, most of the configuration tasks are restricted.

### 13.5.3 Updating IBM Spectrum Virtualize V8.3.1

To update the IBM Spectrum Virtualize software, complete the following steps:

1. Open a supported web browser and go to your cluster IP address. A login window opens (see Figure 13-16).

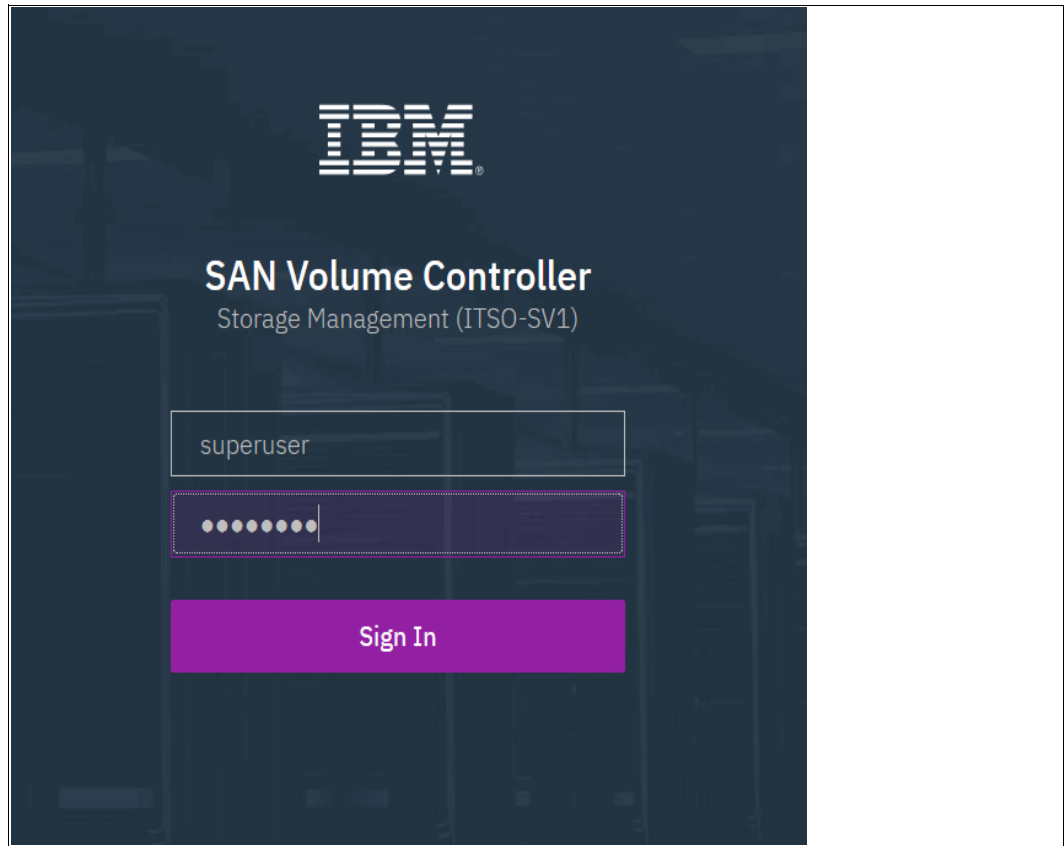


Figure 13-16 IBM SAN Volume Controller GUI login window

2. Log in with superuser rights. The IBM SAN Volume Controller management home window opens. Click **Settings** and click **System** (see Figure 13-17).

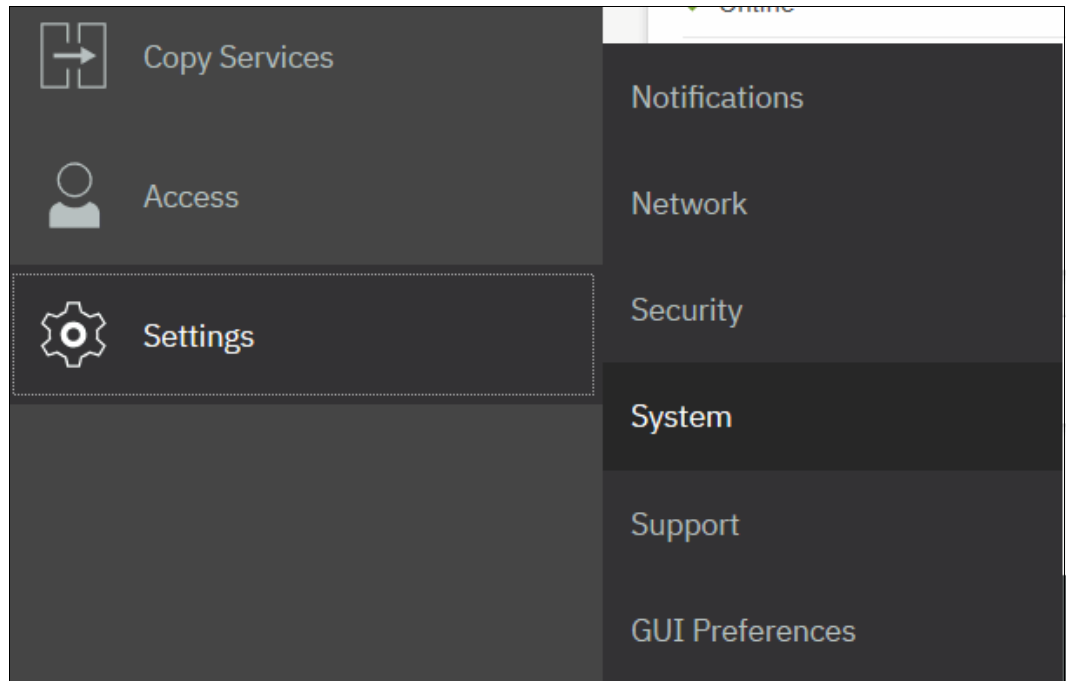


Figure 13-17 Settings window

3. In the **System** menu, click **Update System**. The Update System window opens (see Figure 13-18).

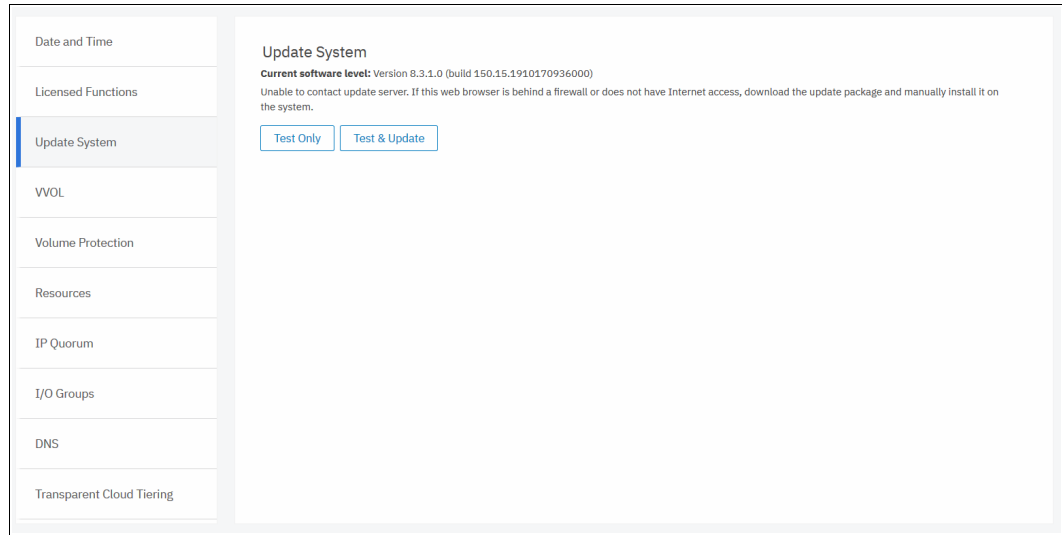


Figure 13-18 Update System window

4. From this window, you can select to run the update test utility and continue with the code update, or run only the test utility. For this example, we click **Test and Update**.

**My Notifications:** Use the My Notifications tool to receive notifications of new and updated support information to better maintain your system environment. This feature is especially useful in an environment where a direct internet connection is not possible.

Go to the [My Notifications site](#) (an IBM account is required) and add your system to the notifications list to be advised of support information, and to download the current code to your workstation for later upload.

5. Because you downloaded both files from the [SAN Volume Controller Support Page](#), you can click each folder, browse to the location where you saved the files, and upload them to the IBM SAN Volume Controller cluster. If the files are correct, the GUI detects and updates the target code level, as shown in Figure 13-19.

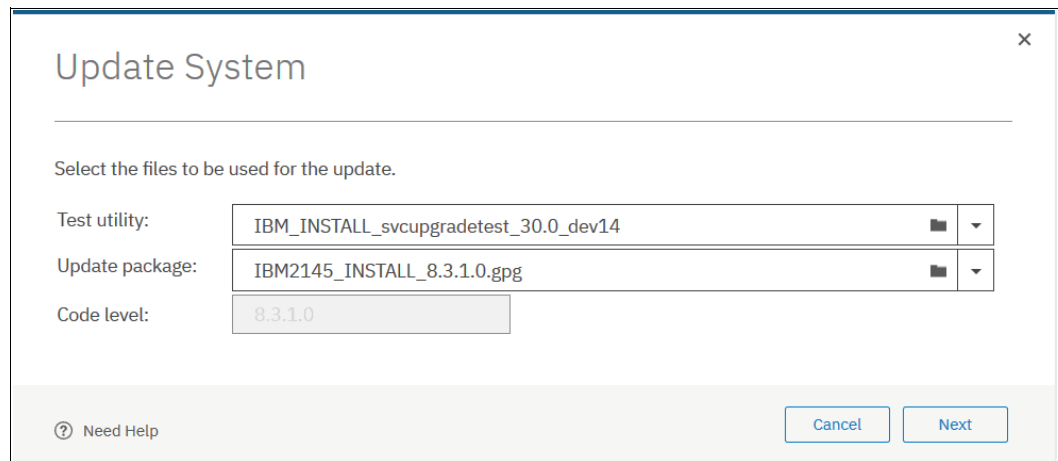


Figure 13-19 Upload option for both the Test utility and the Update package

6. Select the type of update that you want to perform, as shown in Figure 13-20. Select **Automatic update** unless IBM Support suggested a **Service Assistant Manual update**. The manual update might be preferable in cases where misbehaving host multipathing is known to cause loss of access. Click **Finish** to begin the update package upload process.

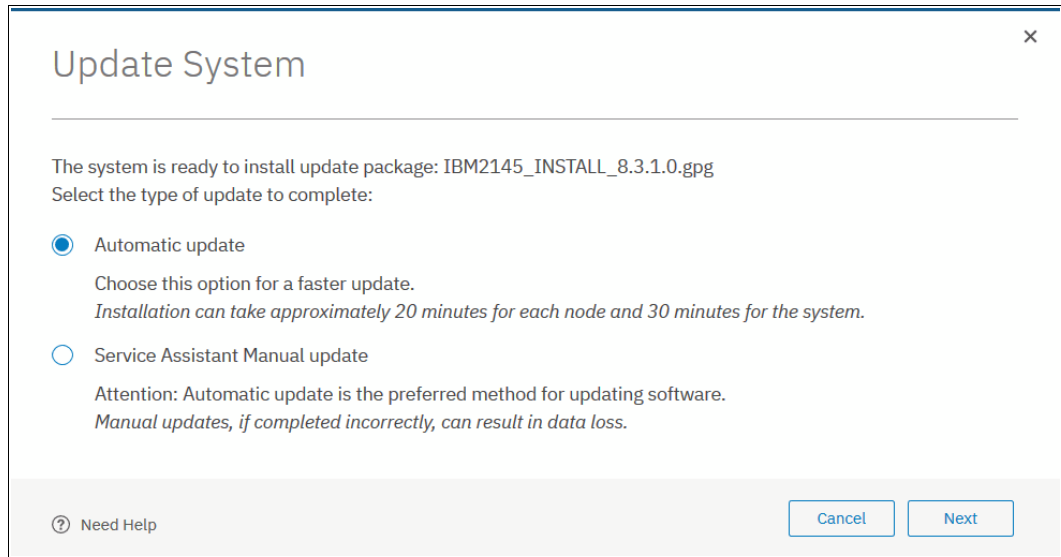


Figure 13-20 Software update type selection



- When updating from Version 8.1 or a later level, an extra window opens where you can choose a fully automated update, such as one that pauses when half the nodes have completed an update or after each node update, as shown in Figure 13-21. Click **Finish**.

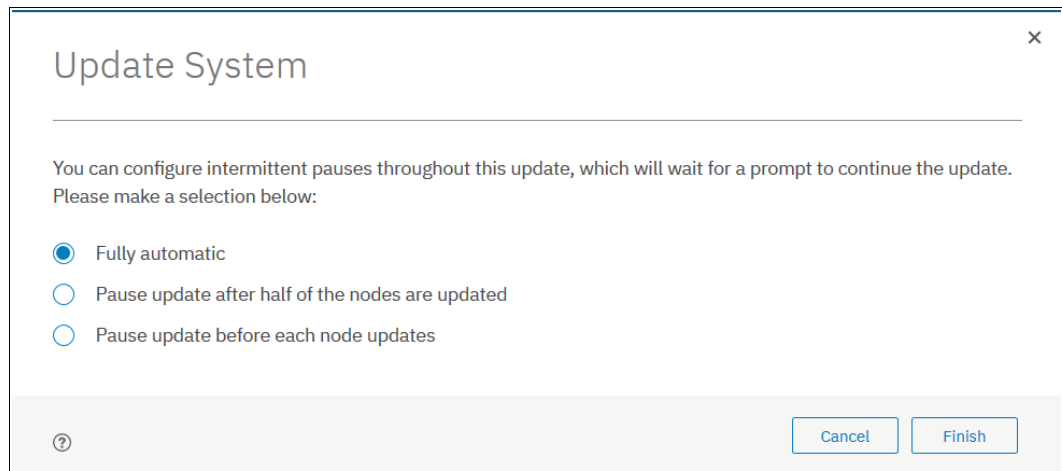


Figure 13-21 New Version 8.1 update pause options

- After the update packages are uploaded, the update test utility looks for any known issues that might affect a concurrent update of your system. The GUI helps identify any detected issues, as shown in Figure 13-22.

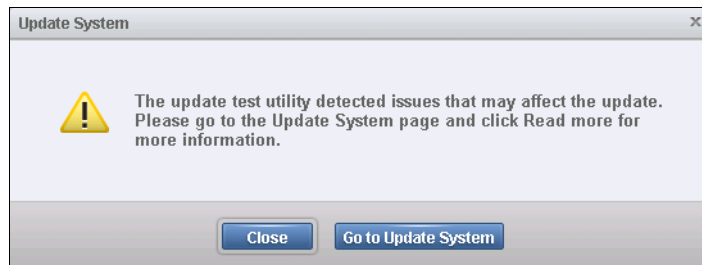


Figure 13-22 Issue detected

- Click **Go to Update System** to return to the Update System window. Then, click **Read more** (see Figure 13-23).

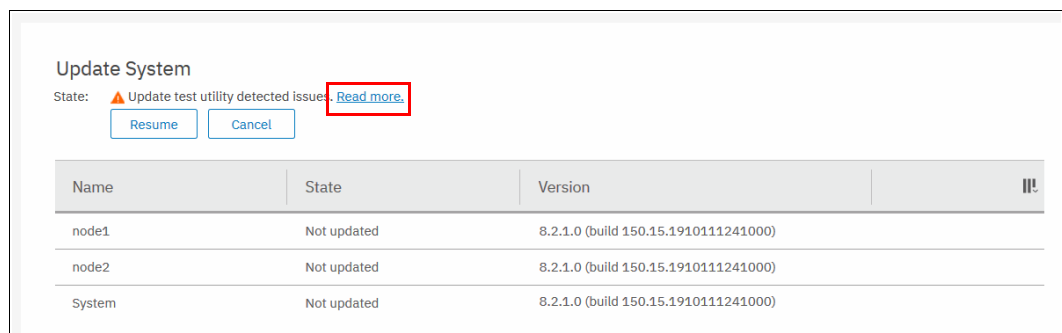


Figure 13-23 Issues that are detected by the update test utility

The results pane opens and shows you the issues that were detected (see Figure 13-24). In our example, the warning is that email notification (Call Home) is not enabled. Although this is not a recommended condition, it does not prevent the system update from running. Therefore, we can click **Close** and proceed with the update. However, you might need to contact IBM Support to assist with resolving more serious issues before continuing.

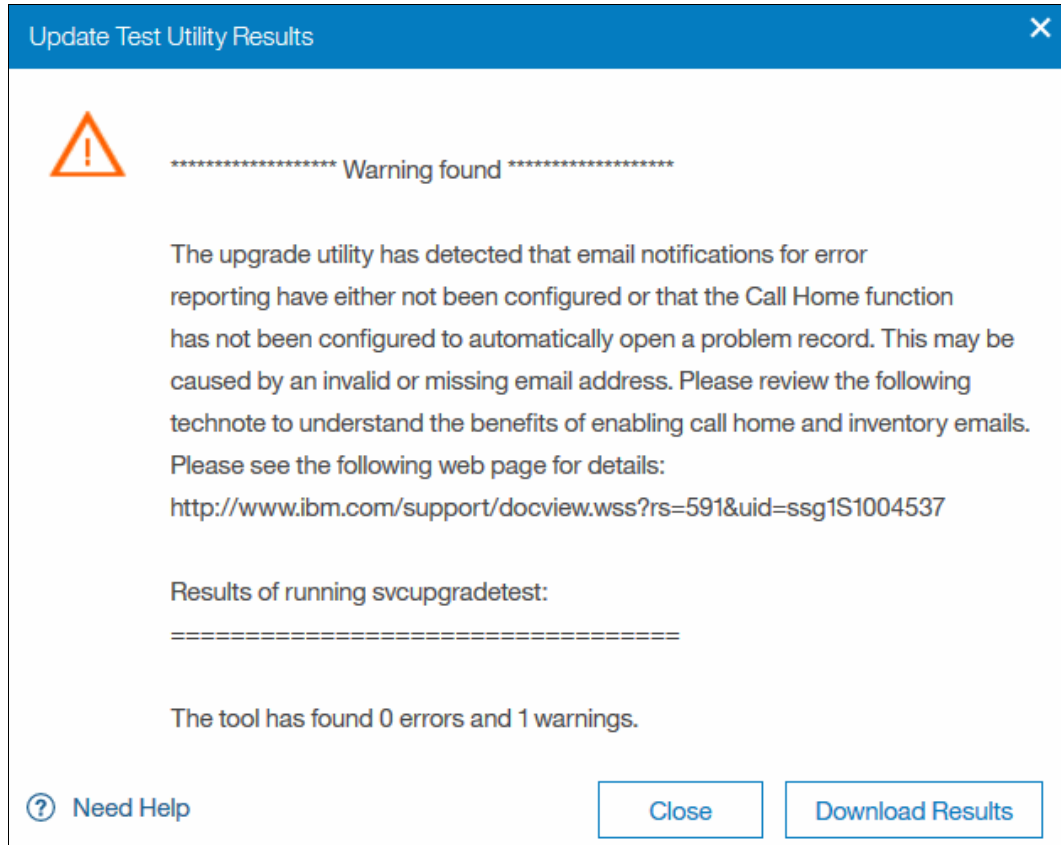


Figure 13-24 Description of the warning from the test utility

10. Click **Resume** in the Update System window and the update proceeds, as shown in Figure 13-25.

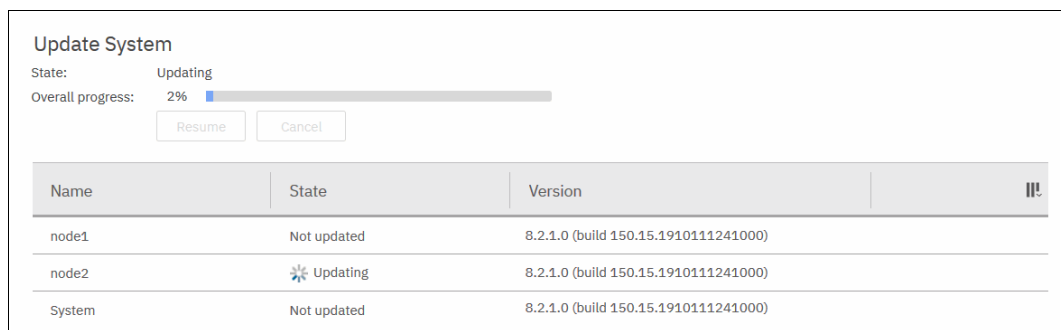


Figure 13-25 Resuming the system update

11. Because of the utility detecting issues, another warning is issued to ensure that you investigated them and are certain you want to proceed, as shown in Figure 13-26. When you are ready to proceed, click **Yes**.

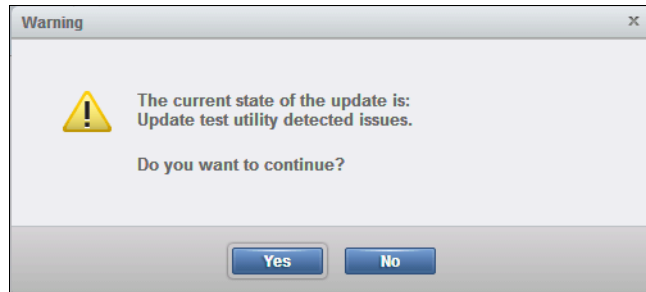


Figure 13-26 Warning before you can continue

12. The system begins updating the IBM Spectrum Virtualize software by taking one node offline and installing the new code. This process takes approximately 20 minutes. After the node returns from the update, it is listed as complete, as shown in Figure 13-27.

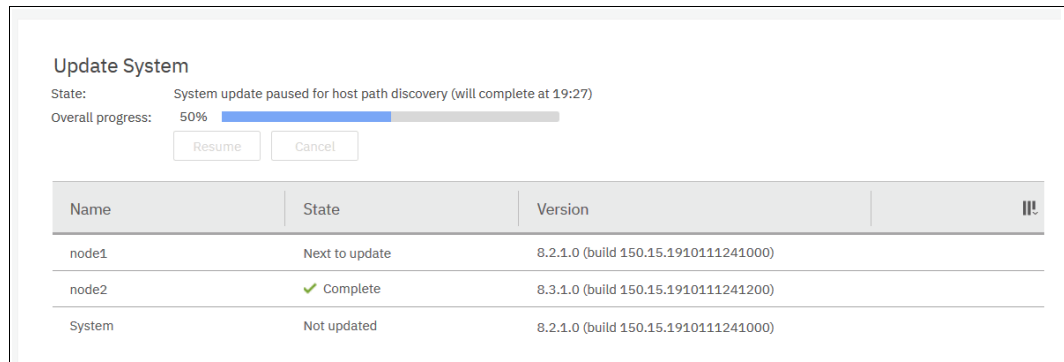


Figure 13-27 Update process starts

13. After a 30-minute pause to ensure that multipathing recovered on all attached hosts, a node failover occurs and you temporarily lose connection to the GUI. A warning window opens, in which you are prompted to refresh the current session, as shown in Figure 13-28.

**Tip:** If you are updating from Version 7.8 or later code, the 30-minute wait period can be adjusted by running the `applysoftware` command with the `-delay (mins)` parameter to begin the update instead of using the GUI.

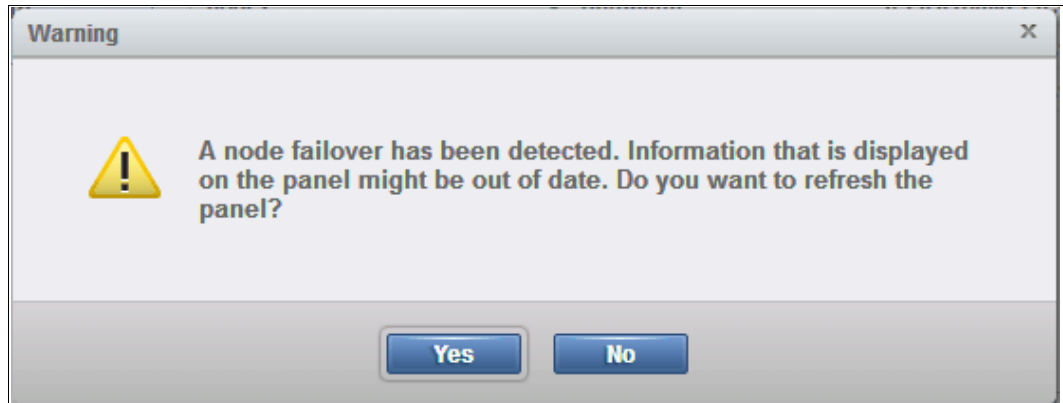


Figure 13-28 Node failover

You now see the Version 8.3 GUI and the status of the second node updating, as shown in Figure 13-29.

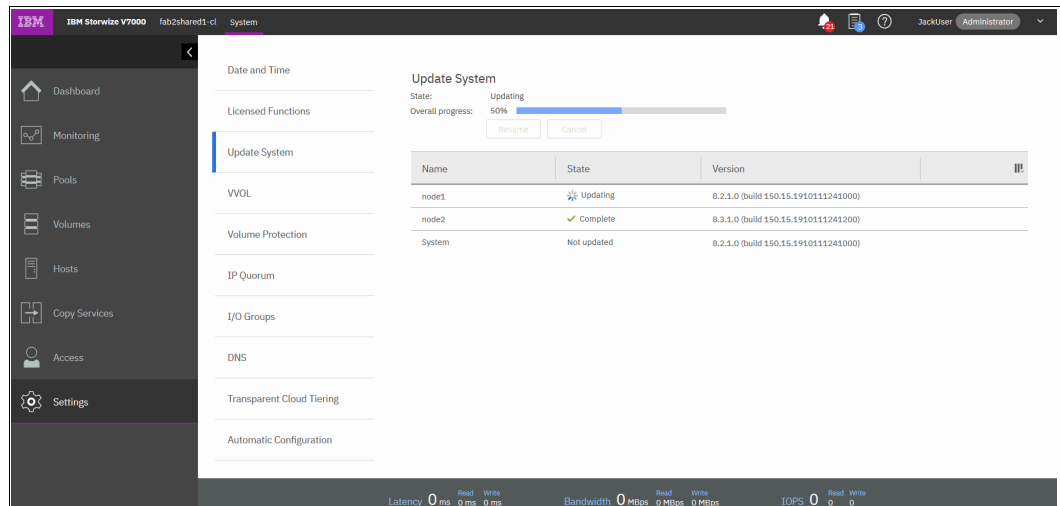


Figure 13-29 New GUI after node failover

After the second node completes, the update is committed to the system, as shown in Figure 13-30.

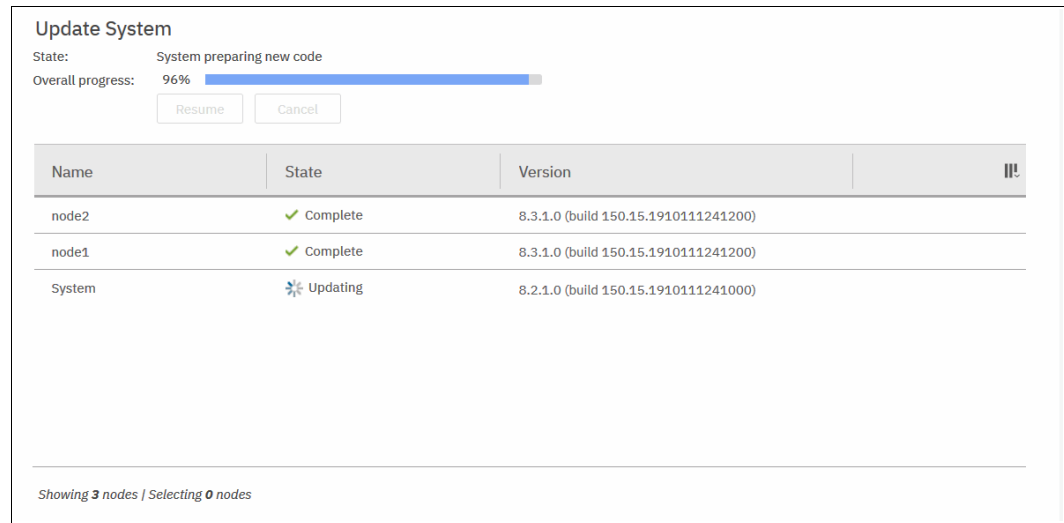


Figure 13-30 Updating the system level

- The update process completes when all nodes and the system are committed. The final status indicates the new level of code that is installed in the system.

**Note:** If your nodes have more than 64 GB of memory before updating to Version 8.1, each node posts an 841 error after the update completes. Because Version 8.1 allocates memory differently, the memory must be accepted by running the fix procedure for the event or running `svctask chnodehw <id>` for each node. For more information, see the [IBM Knowledge Center](#).

### 13.5.4 Updating IBM Spectrum Virtualize with a hot spare node

IBM Spectrum Virtualize V8.1 introduced a new optional feature of Hot Spare Node. This feature allows for IBM Spectrum Virtualize software updates to minimize any performance impact and removes any redundancy risk during the update. It does so by automatically swapping in a hot spare node after 1 minute to replace temporarily the node currently updating. After the original node is updated, the hot spare node becomes a spare again, and is ready for the next node to update. The original node rejoins the cluster.

To use this feature, the spare node must have the same amount of memory and a matching FC port configuration as the other nodes in the cluster. Up to four hot spare nodes can be added to a cluster and must be zoned as part of the IBM SAN Volume Controller cluster.

Figure 13-31 on page 780 shows how a spare node appears in the GUI. This node was added as a spare but has a different configuration to the two nodes in the IOgroup, so it is held in service state until this is fixed.

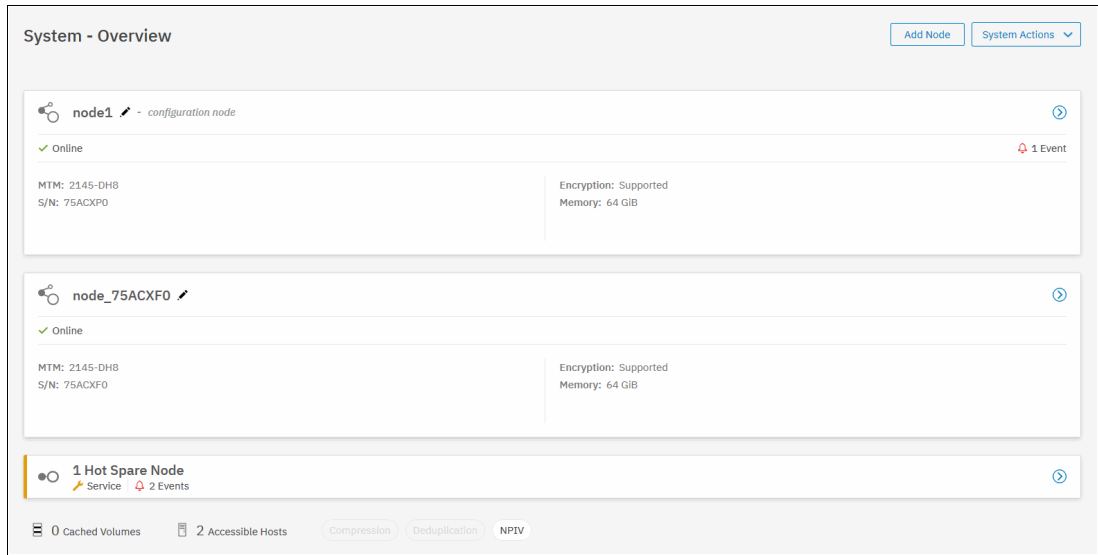


Figure 13-31 Hot Spare Node on the System Overview Pane

### 13.5.5 Updating the IBM SAN Volume Controller system manually

This example assumes that you have an 8-node cluster of the IBM SAN Volume Controller cluster, as listed in Table 13-9.

Table 13-9 The iogrp cluster

iogrp (0)	iogrp (1)	iogrp (2)	iogrp (3)
node 1 (config node)	node 3	node 5	node 7
node 2	node 4	node 6	node 8

After uploading the update utility test and software update package to the cluster by running the **pscp** command and running the utility test, complete the following steps:

1. Remove node 2 from cluster, which is the partner node of the configuration node in iogrp 0, by using the cluster GUI or CLI.
2. Log in to the service GUI to verify that the removed node is in the candidate status.
3. Select the candidate node and click **Update Manually** from the left pane.
4. Browse and locate the code that you downloaded and saved to your PC.
5. Upload the code and click **Update**.  
When the update is completed, a message caption opens that indicates that the software update completed. The node then restarts, and appears again in the service GUI after approximately 20 - 25 minutes in candidate status.
6. Select the node and verify that it is updated to the new code.
7. Add the node back by using the cluster GUI or CLI.
8. Select **node 3** from iogrp1.
9. Repeat steps 1 - 7 to remove node 3, update it manually, verify the code, and add it back to the cluster.
10. Proceed to node 5 in iogrp 2.

11. Repeat steps 1 - 7 to remove node 5, update it manually, verify the code, and add it back to the cluster.
12. Move on to node 7 in iogrp 3.
13. Repeat steps 1 - 7 to remove node 5, update it manually, verify the code, and add it back to the cluster.

**Note:** At this point, the update is 50% complete. You now have one node from each iogrp that is updated with the new code manually. Always leave the configuration node for last during a manual IBM Spectrum Control Software update.

14. Select **node 4** from iogrp 1.
15. Repeat steps 1 - 7 to remove node 4, update it manually, verify the code, and add it back to the cluster.
16. Select **node 6** from iogrp 2.
17. Repeat steps 1 - 7 to remove node 6, update it manually, verify the code, and add it back to the cluster.
18. Select **node 8** in iogrp 3.
19. Repeat steps 1 - 7 to remove node 8, update it manually, verify the code, and add it back to the cluster.
20. Select and remove node 1, which is the configuration node in iogrp 0.

**Note:** A partner node becomes the configuration node because the original config node is removed from the cluster, which keeps the cluster manageable.

The removed configuration node becomes a candidate, and you do not have to apply the code update manually. Add the node back to the cluster. It automatically updates itself and then adds itself back to the cluster with the new code.

21. After all the nodes are updated, you must confirm the update to complete the process. The confirmation restarts each node in order, which takes about 30 minutes to complete.

The update is complete.

## 13.6 Health checker feature

The IBM Spectrum Control health checker feature runs in IBM Cloud. Based on the weekly call home inventory reporting, it proactively creates recommendations. These recommendations are provided at IBM Call Home Web, which you can access by selecting **Support** → **My support** → **Call Home Web** (see Figure 13-32 on page 782).

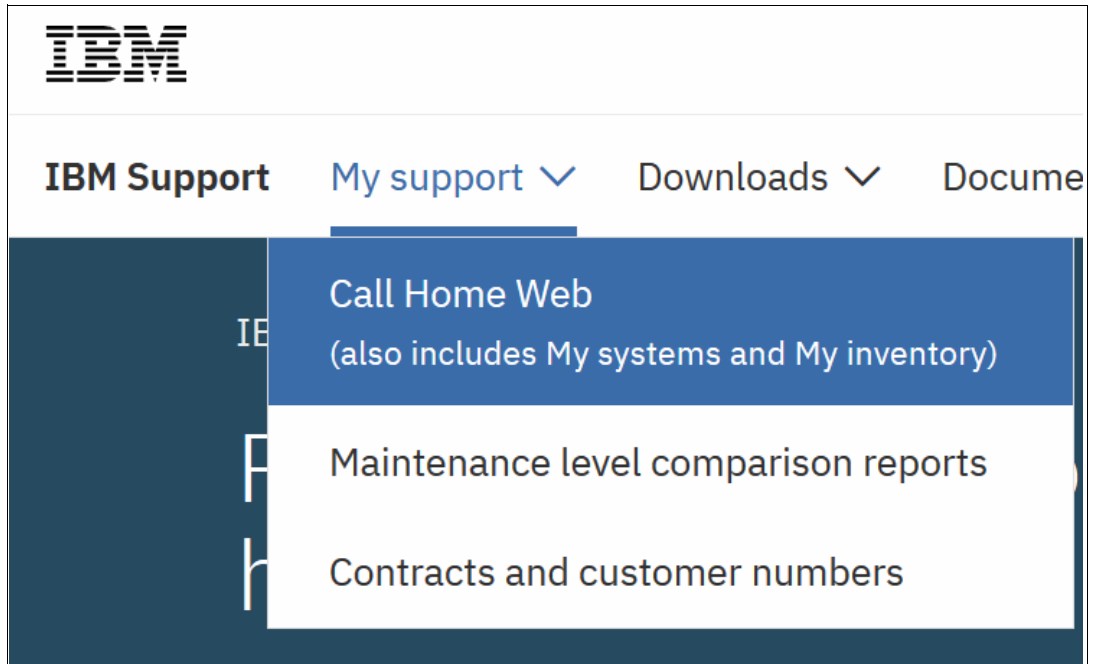


Figure 13-32 IBM Call Home Web

For a video guide about how to set up and use IBM Call Home Web, see [this YouTube video](#).

Another feature is the *Critical Fix Notification* function, which enables IBM to warn IBM Spectrum Virtualize users that a critical issue exists in the level of code that they use. The system notifies users when they log on to the GUI by using a web browser that is connected to the internet.

Consider the following information about this function:

- ▶ It warns users only about critical fixes; it does not warn them that they are running a previous version of the software.
- ▶ It works only if the browser also can access the internet. The IBM Storwize V7000 and IBM SAN Volume Controller systems do not need to be connected to the internet.
- ▶ The function cannot be disabled. Whenever a warning is displayed, it must be acknowledged (with the option to not warn the user again for that issue).

The decision about what a *critical* fix is subjective and requires judgment, which is exercised by the development team. As a result, clients might still encounter bugs in code that were not deemed critical. They should continue to review information about new code levels to determine whether they should update even without a critical fix notification.

**Important:** Inventory notification must be enabled and operational for these features to work. It is best practice to enable Call Home and Inventory reporting on your IBM Spectrum Virtualize clusters.



## 13.7 Troubleshooting and fix procedures

The management GUI of IBM Spectrum Virtualize is a browser-based GUI for configuring and managing all aspects of your system. It provides extensive facilities to help troubleshoot and correct problems. This section explains how to use effectively its features to avoid service disruption of your IBM SAN Volume Controller.

Figure 13-33 shows the Monitoring menu for System information, viewing Easy Tier Reports, viewing Events, or seeing real-time Performance statistics.

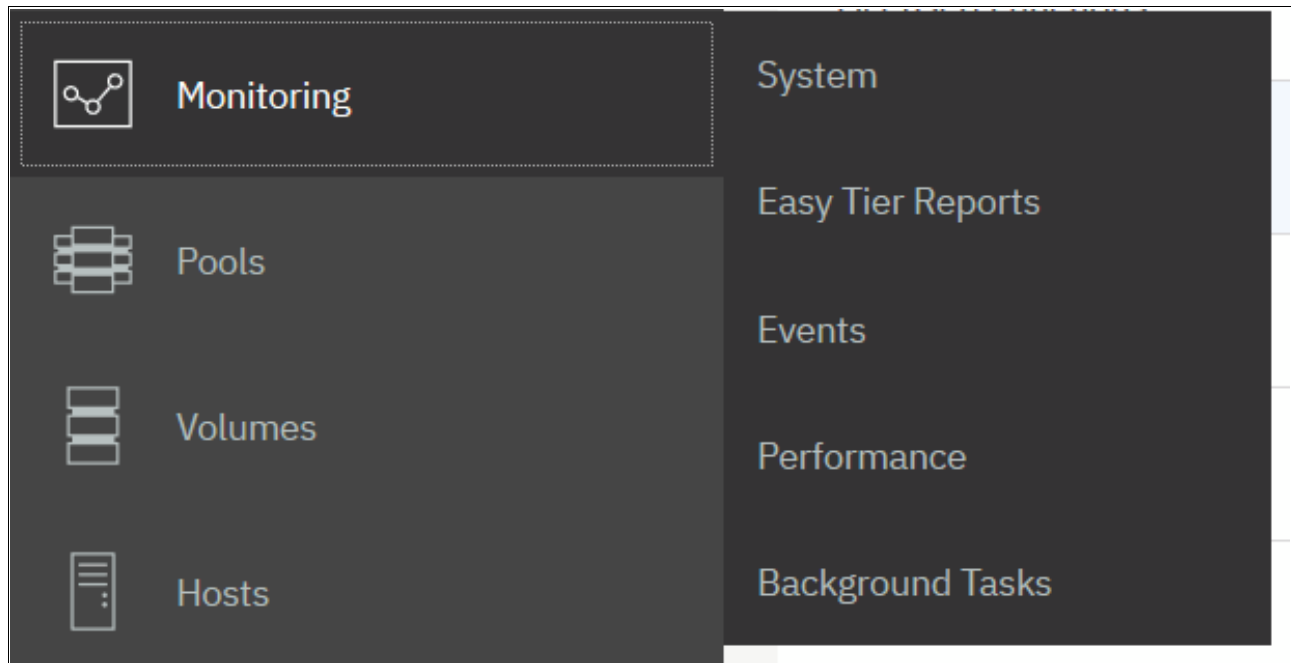


Figure 13-33 Monitoring options

Use the management GUI to manage and service your system. Click **Monitoring** → **Events** to list events that should be addressed and maintenance procedures that walk you through the process of correcting problems.

Information in the Events window can be filtered in the following ways:

► Recommended Actions

Shows only the alerts that require attention. Alerts are listed in priority order and should be resolved sequentially by using the available fix procedures. For each problem that is selected, you can perform the following tasks:

- Run a fix procedure
- View the properties

► Unfixed Alerts

Displays only the alerts that are not fixed. For each entry that is selected, you can perform the following tasks:

- Run a fix procedure
- Mark an event as fixed
- Filter the entries to show them by specific minutes, hours, or dates
- Reset the date filter
- View the properties

► Unfixed Messages and Alerts

Displays only the alerts and messages that are not fixed. For each entry that is selected, you can perform the following tasks:

- Run a fix procedure
- Mark an event as fixed
- Filter the entries to show them by specific minutes, hours, or dates
- Reset the date filter
- View the properties

► Show All

Displays all event types whether they are fixed or unfixed. For each entry that is selected, you can perform the following tasks:

- Run a fix procedure
- Mark an event as fixed
- Filter the entries to show them by specific minutes, hours, or dates
- Reset the date filter
- View the properties

Some events require a specific number of occurrences in 25 hours before they are shown as unfixed. If they do not reach this threshold in 25 hours, they are flagged as *expired*. Monitoring events are below the coalesce threshold, and are usually transient.

**Important:** The management GUI is the primary tool that is used to operate and service your system. Real-time monitoring should be established by using SNMP traps, email notifications, or syslog messaging in an automatic manner.

### 13.7.1 Managing event log

Regularly check the status of the system by using the management GUI. If you suspect a problem, first use the management GUI to diagnose and resolve the problem.

Use the views that are available in the management GUI to verify the status of the system, the hardware devices, the physical storage, and the available volumes by completing the following steps:

1. Click **Monitoring** → **Events** to see all problems that exist on the system (see Figure 13-34).

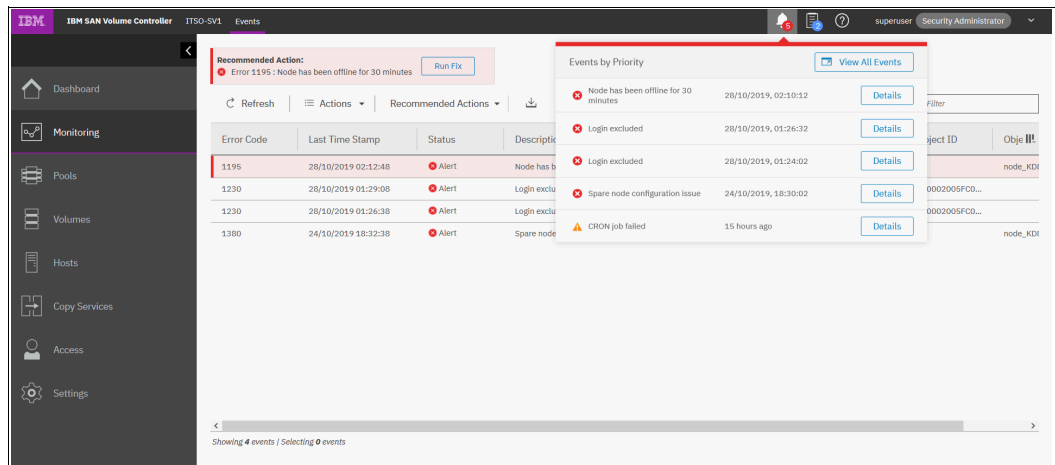


Figure 13-34 Messages in the event log

2. Select **Recommended Actions** to show the most important events to be resolved (see Figure 13-35). The Recommended Actions tab shows the highest priority maintenance procedure that must be run. Use the troubleshooting wizard so that IBM SAN Volume Controller can determine the proper order of maintenance procedures.

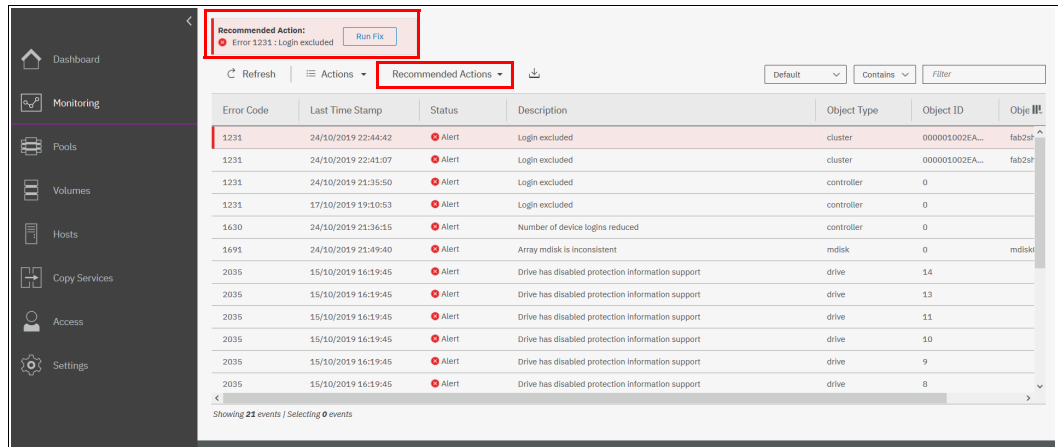


Figure 13-35 Recommended Actions

In this example, a login excluded (service error code 1231) is used. At any time and from any GUI window, you can directly go to this menu by clicking the **Status Alerts** icon at the top of the GUI (see Figure 13-36).

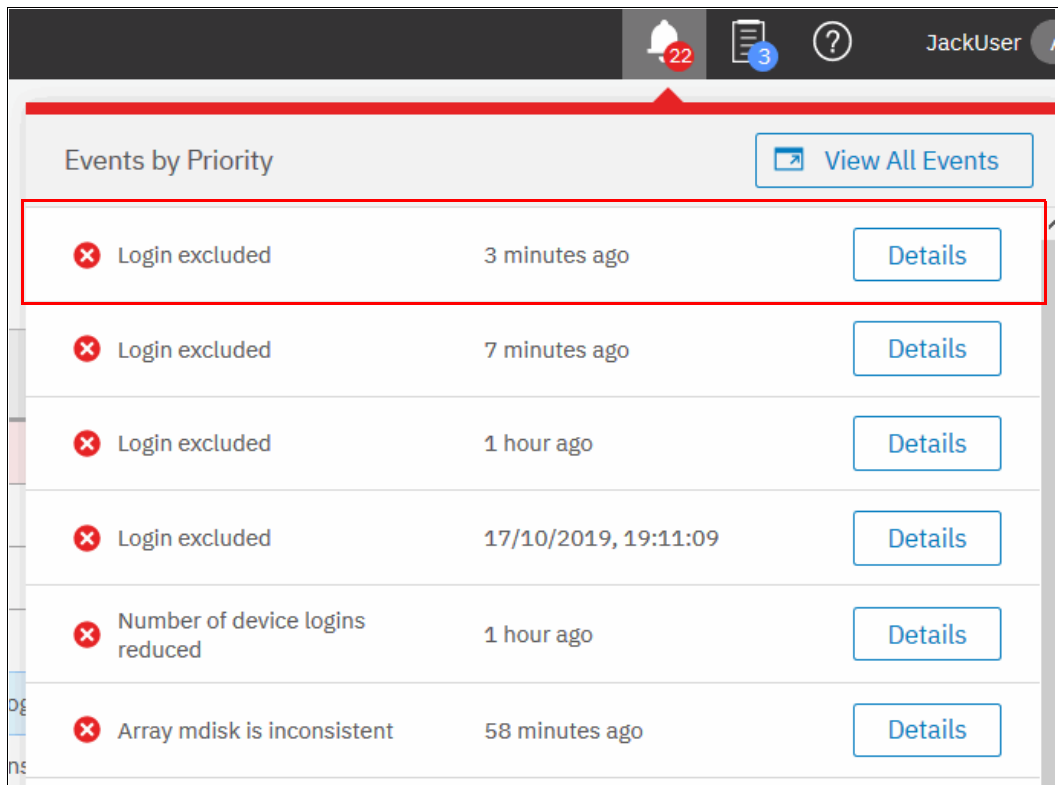


Figure 13-36 Status alerts

## 13.7.2 Running a fix procedure

If an error code exists for the alert, you should run a fix procedure to help resolve the problem. Fix procedures analyze the system and provide more information about the problem. They suggest actions to take and walk you through the actions that automatically manage the system where necessary while ensuring availability. Finally, they verify that the problem is resolved.

If an error is reported, always use the fix procedures from the management GUI to resolve the problem. Always use the fix procedures for both software configuration problems and hardware failures. The fix procedures analyze the system to ensure that the required changes do not cause volumes to become inaccessible to the hosts. The fix procedures automatically perform configuration changes that are required to return the system to its optimum state.

The fix procedure shows information that is relevant to the problem, and provides various options to correct the problem. Where possible, the fix procedure runs the commands that are required to reconfigure the system.

**Note:** After Version 7.4, you are no longer required to run the fix procedure for a failed internal enclosure drive. Hot plugging of a replacement drive automatically triggers the validation processes.

The fix procedure also checks that any other problem does not result in volume access being lost. For example, if a PSU in a node enclosure must be replaced, the fix procedure checks and warns you whether the integrated battery in the other PSU is not sufficiently charged to protect the system.

**Hint:** Always use the Run Fix function, which resolves the most serious issues first. Often, other alerts are corrected automatically because they were the result of a more serious issue.

The following example demonstrates how to clear the error that is related to an excluded FC login, likely because of errors along the link:

1. From the GUI menu on the left, select **Monitoring** → **Events**, and list only the recommended actions by using the **Actions** menu (see Figure 13-37). Click **Run Fix**.

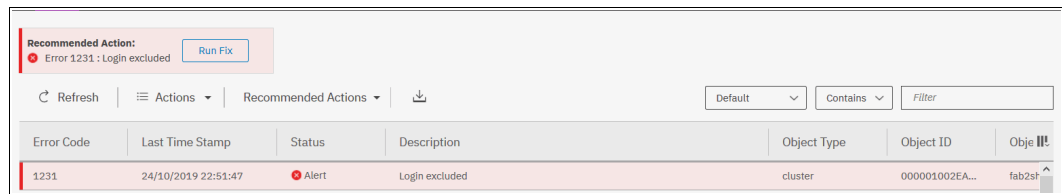


Figure 13-37 Initiate Run Fix procedure from the management GUI

- The next window shows you which login has the problem. For more information, click **Next** (see Figure 13-38).

### Login excluded

---

Local port or login excluded

*A Fibre Channel login exclusion(Logins are only excluded due to transport errors when significant numbers of errors occur over a predetermined period) has been detected on the node below.*

Machine Type and Model	Node Identifier	Node Name	Enclosure Identifier	Enclosure Serial Number	Panel Name	Canister Position In Enclosure
2076-624	1	node1	1	7822PBR	01-1	Left

Use following details on the affected Fibre Channel I/O port to resolve transport errors:

*Fibre Channel login related to the exclusion*

I/O Port ID	Local WWPN	Local I/O Port Status	Local I/O Port Speed	Remote WWNN	Remote WWPN	Remote NPort ID
4	500507680B2480FA	active	8Gb	500507680B0080E8	500507680B2380E8	0x00012200

The corresponding physical Fibre Channel port ID is 4.

Click **Next** to continue, or click **Cancel** to exit.

Figure 13-38 Determination of login problem

- In the next window (see Figure 13-39), we are shown an error count a history of login excluded errors against the port, which allows us to see whether this has a history of being a problematic link. In this case, this is a single error, and we checked the connectivity, which seems to be healthy. Therefore, we can select **Mark this event as fixed** or **Mark all the 1231 events as fixed...** Then, click **Next**.

### Login excluded

---

Event history

The following table shows the number of unfixed 1231 events currently in the log for each port.

*History of similar unfixed fibre channel exclusion events*

Fibre Channel I/O Port ID	Physical Port Type	First Timestamp	Last Timestamp	Count
3	fc	17/10/19 18:29:32	17/10/19 19:10:53	1
4	fc	15/10/19 10:05:09	24/10/19 22:56:27	1

The following table shows the unfixed events that occurred within the last 30 days:

*Unfixed events of the same type in the last 30 days*

Fibre Channel I/O Port ID	Physical Port Type	First Timestamp	Last Timestamp	Count
3	fc	17/10/19 18:29:32	17/10/19 19:10:53	1
4	fc	15/10/19 10:05:09	24/10/19 22:56:27	1

Select how you wish to proceed :

If you suspect that the problem still exists, select **Troubleshoot hardware**. This option launches a series of fix procedures to help determine the cause of the problem. If you have recently replaced some hardware on this system or suspect that the underlying problem might no longer exist, you can select **Mark all events as fixed**. If the underlying problem still exists, the event will be logged again.

Troubleshoot hardware  
 Mark all the 1231 events listed in the preceding tables as fixed  
 Mark this event as fixed

Figure 13-39 Configuration rules

- The next window (see Figure 13-40) shows a confirmation that the event was marked as fixed. **Close** to exit.

### Login excluded

---

All the events have been marked as fixed

The events on the node 1 have been marked as fixed.

Click **Close** to exit.

Figure 13-40 Correctly finished fix procedure

- The event is marked as fixed, and you can safely finish the fix procedure. Click **Close** and the event is removed from the list of events (see Figure 13-41).

Error Code	Last Time Stamp	Status	Description	Object Type	Object ID	Object Name
1630	24/10/2019 21:36:15	Alert	Number of device logins reduced	controller	0	
1691	24/10/2019 23:04:43	Alert	Array mdisk is inconsistent	mdisk	0	mdisk
2035	15/10/2019 16:19:45	Alert	Drive has disabled protection information support	drive	14	
2035	15/10/2019 16:19:45	Alert	Drive has disabled protection information support	drive	13	

Figure 13-41 Pane showing there are no more excluded local logins

### 13.7.3 Resolving alerts in a timely manner

To minimize any impact to your host systems, always perform the recommended actions as quickly as possible after a problem is reported. Your system is resilient to most single hardware failures. However, if it operates for any period with a hardware failure, the possibility increases that a second hardware failure can result in volume data that is unavailable. If several unfixed alerts exist, fixing any one alert might become more difficult because of the effects of the others.

### 13.7.4 Event log details

Multiple views of the events and recommended actions are available. The GUI works like a typical Windows menu, so the event log grid is manipulated by using the row that contains the column headings (see Figure 13-42). When you click the column icon at the right end of the table heading, a menu for the column choices opens.

Error Code	Last Time Stamp	Status	Description	Object Type	Object ID	Object Name
1195	28/10/2019 02:12:48	Alert	Node has been offline for 30 minutes			
1230	28/10/2019 01:29:08	Alert	Login excluded			
1230	28/10/2019 01:26:38	Alert	Login excluded			
1380	24/10/2019 18:32:38	Alert	Spare node configuration issue			

Figure 13-42 Grid options of the event log

Select or remove columns as needed. You can then also extend or shrink the width of the column to fit your screen resolution and size. This method is the way to manipulate it for most grids in the management GUI of IBM Spectrum Virtualize, not just the events pane.

Every field of the event log is available as a column in the event log grid. Several fields are useful when you work with IBM Support. The preferred method in this case is to use the Show All filter, with events sorted by time stamp. All fields have the sequence number, event count, and the fixed state. Using Restore Default View sets the grid back to the defaults.

You might want to see more information about each critical event. Some details are not shown in the main grid. To access properties and the sense data of a specific event, double-click the specific event anywhere in its row.

The properties window opens (see Figure 13-43) with all the relevant sense data. This data includes the first and last time of an event occurrence, worldwide port name (WWPN), and worldwide node name (WWNN), enabled or disabled automatic fix, and so on.

**Properties and Sense Data for Event 010042**

Error Code: 1627 [Run Fix](#)

**Insufficient redundancy in disk controller connectivity**

First Time Stamp	Last Time Stamp	Fixed Time Stamp	Event Count
10/18/2018 12:59:07 PM	10/18/2018 12:59:07 PM		2

Event ID	010042
Event ID Text	Only a single port on a disk controller is accessible from every node in the cluster
Sequence Number	215
Object Type	controller
Object ID	6
Object Name	xxcontroller6
Secondary Object ID	
Secondary Object Type	
Copy ID	
Reporting Node ID	2
Reporting Node Name	node2
Root Sequence Number	
Error Code	1627
Error Code Text	Insufficient redundancy in disk controller connectivity
-	----
-	----

Previous Next Close

Figure 13-43 Event sense data and properties

For more information about troubleshooting options, see [IBM Knowledge Center](#).

## 13.8 Monitoring

An important step is to correct any issues that are reported by your IBM SAN Volume Controller as soon as possible. Configure your system to send automatic notifications to a standard Call Home server or to new event Cloud Call Home server when a new event is reported. To avoid the need to monitor the management GUI for new events, select the type of event for which you want to be notified; for example, restrict notifications to only events that require action. The following event notification mechanisms exist:

<b>Call Home</b>	An event notification can be sent to one or more email addresses. This mechanism notifies individuals of problems. Individuals can receive notifications wherever they have email access, including mobile devices.
<b>Cloud Call Home</b>	Cloud services for Call Home is the optimal transmission method for error data because it ensures that notifications are delivered directly to the IBM support center.
<b>SNMP</b>	An SNMP traps report can be sent to a data center management system, such as IBM Systems Director, which consolidates SNMP reports from multiple systems. With this mechanism, you can monitor your data center from a single workstation.
<b>Syslog</b>	A syslog report can be sent to a data center management system that consolidates syslog reports from multiple systems. With this option, you can monitor your data center from a single location.

If your system is under warranty or if you have a hardware maintenance agreement, configure your IBM SAN Volume Controller cluster to send email events directly to IBM if an issue that requires hardware replacement is detected. This mechanism is known as *Call Home*. When this event is received, IBM automatically opens a problem ticket and, if suitable, contacts you to help resolve the reported problem.

**Important:** If you set up Call Home to IBM, ensure that the contact details that you configure are correct and kept current. Personnel changes can cause delays in IBM making contact.

Cloud Call Home is designed to work with new service teams, improves connectivity, and ultimately should improve customer support. The initial setup of Cloud Call Home is explained in Chapter 3, “Initial configuration” on page 81.

**Note:** If the customer does not want to open their firewall, Cloud Call Home does not work. The customer can disable Cloud Call Home and Call Home is used instead.

### 13.8.1 The Call Home function and email notification

The Call Home function of IBM Spectrum Virtualize sends an email to a specific IBM Support center. Therefore, the configuration is similar to sending emails to a specific person or system owner. Complete the following steps to configure Call Home and email notifications:

1. Prepare your contact information that you want to use for the Call Home function and verify the accuracy of the data. From the GUI menu, click **Settings** → **Support** (see Figure 13-44 on page 791).



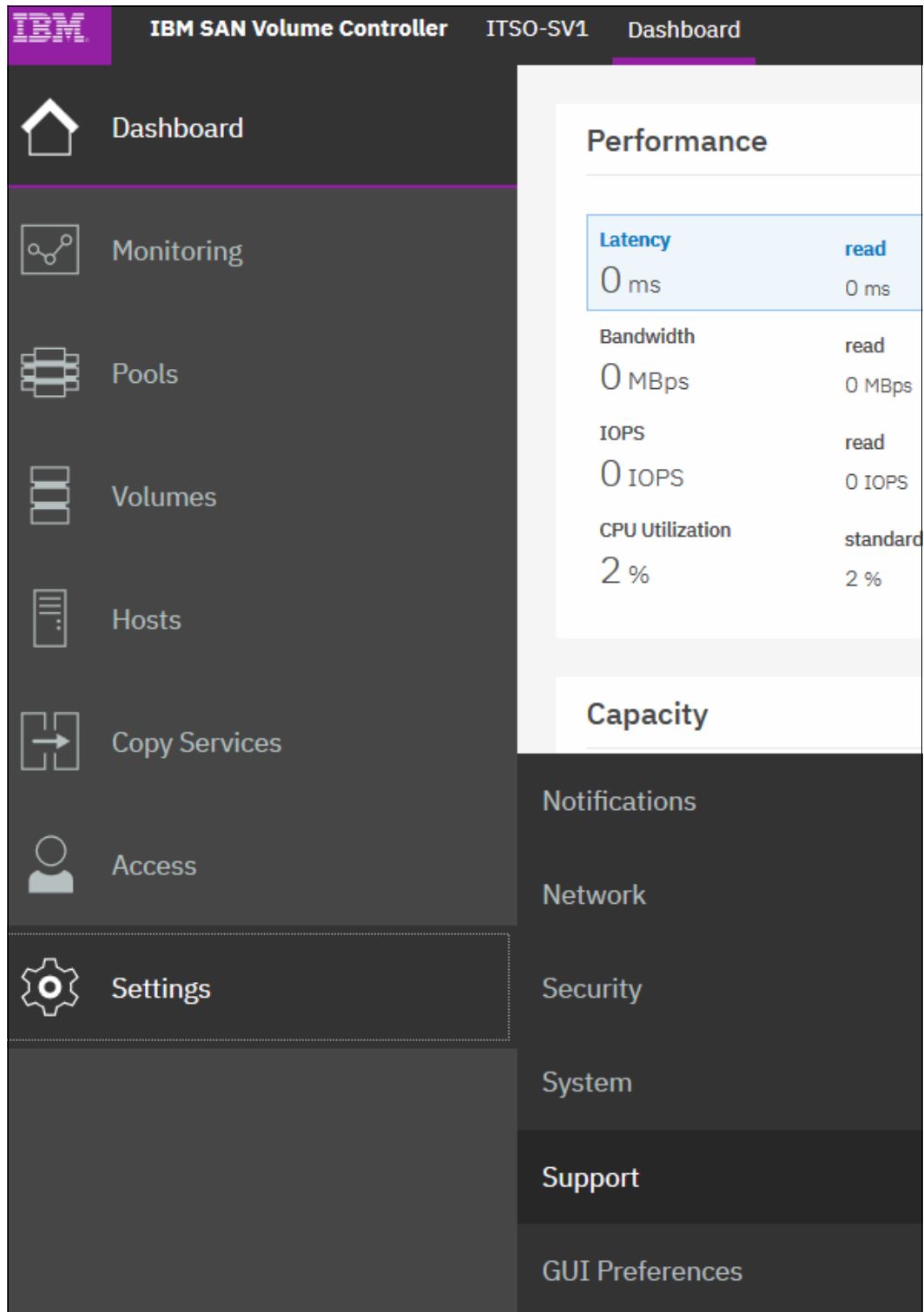


Figure 13-44 Support menu

- Click **Call Home** and then click **Enable Notifications** (see Figure 13-45 on page 792). For the correct functions of email notifications, ask your network administrator if Simple Mail Transfer Protocol (SMTP) is enabled on the management network and is not blocked by firewalls. Also, ensure that the destination @de.ibm.com is not blacklisted.

Be sure to test the accessibility to the SMTP server by running the `telnet` command (port 25 for a non-secured connection, port 465 for Secure Sockets Layer (SSL) -encrypted communication) by using any server in the same network segment.

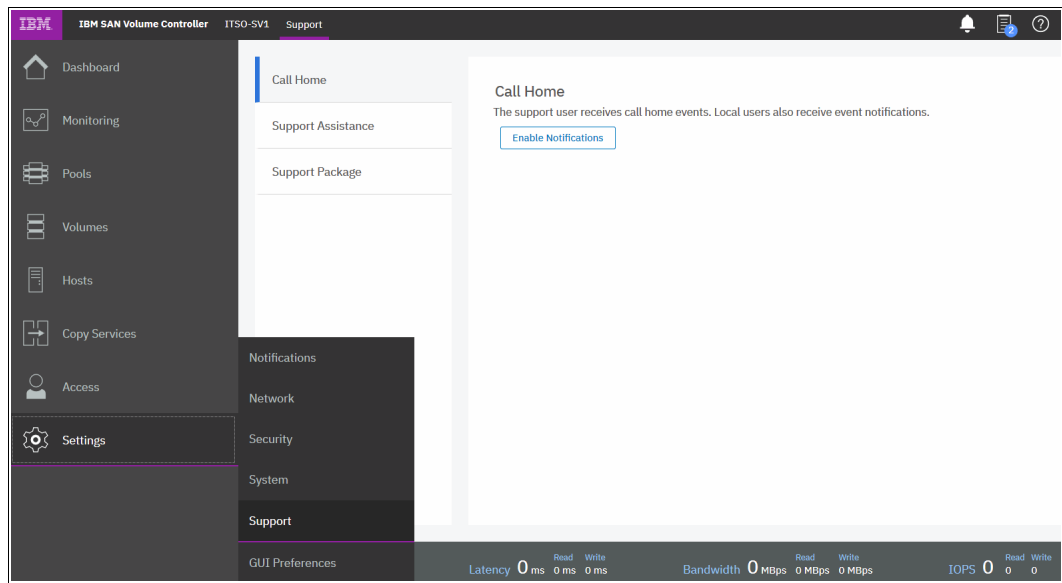


Figure 13-45 Configuration of Call Home function

Figure 13-46 shows the option to enable Cloud Call Home.

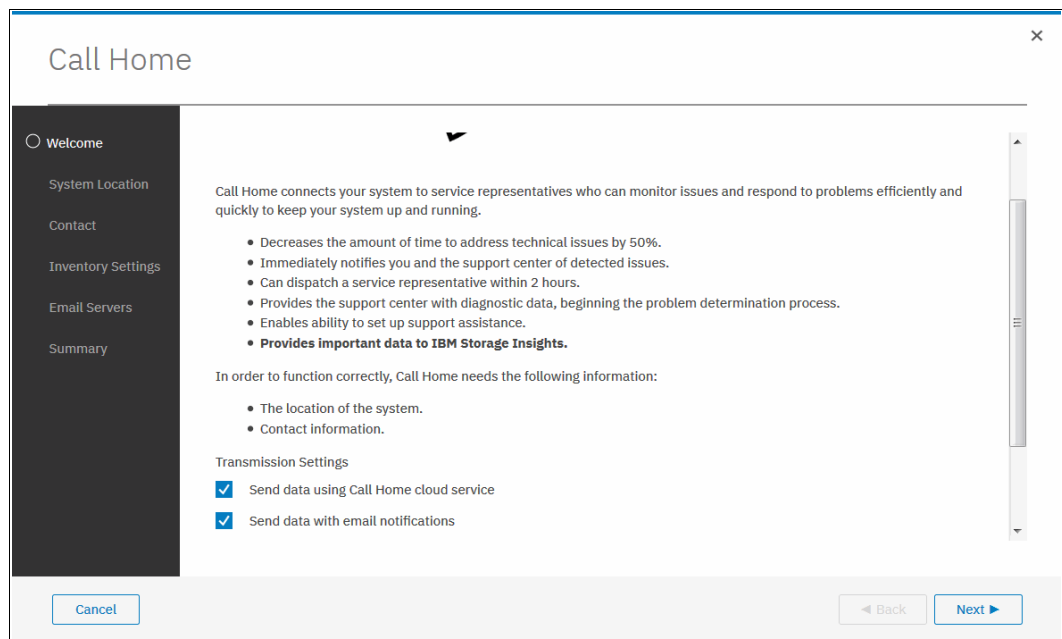


Figure 13-46 Cloud Call Home service

3. After clicking **Next** in the welcome window, enter the information about the location of the system (see Figure 13-47) and contact information of the IBM SAN Volume Controller administrator (see Figure 13-48) to be contacted by IBM Support. *Always* keep this information current.

The screenshot shows a window titled "Call Home" with a sidebar on the left containing navigation options: Welcome (checked), System Location (selected), Contact, Inventory Settings, Email Servers, and Summary. The main content area is titled "System Location" and includes the instruction: "Service parts should be shipped to the same physical location as the system." Below this are several input fields: "Company name" (IBM ITSO), "System address" (Ridder Park Dr), "City" (San Jose), "State or province" (CA), "Postal code" (95131), and "Country or region" (United States). At the bottom, there are "Cancel", "Back", and "Next" buttons.

Figure 13-47 Location of the device

Figure 13-48 shows the contact information of the owner.

The screenshot shows the same "Call Home" window, but now the "Contact" option in the sidebar is selected. The main content area features an information icon and a note: "Enter business-to-business contact information. To comply with privacy regulations, personal contact information for individuals with your organization is not recommended." Below this are input fields for "Name" (System Administrator), "Email" (name@company.com), "Phone (primary)" (+123456789), and "Phone (alternate)". At the bottom, there are "Cancel", "Back", and "Apply and Next" buttons.

Figure 13-48 Contact information

The next window allows you to enable Inventory Reporting and Configuration Reporting, as shown in Figure 13-49.

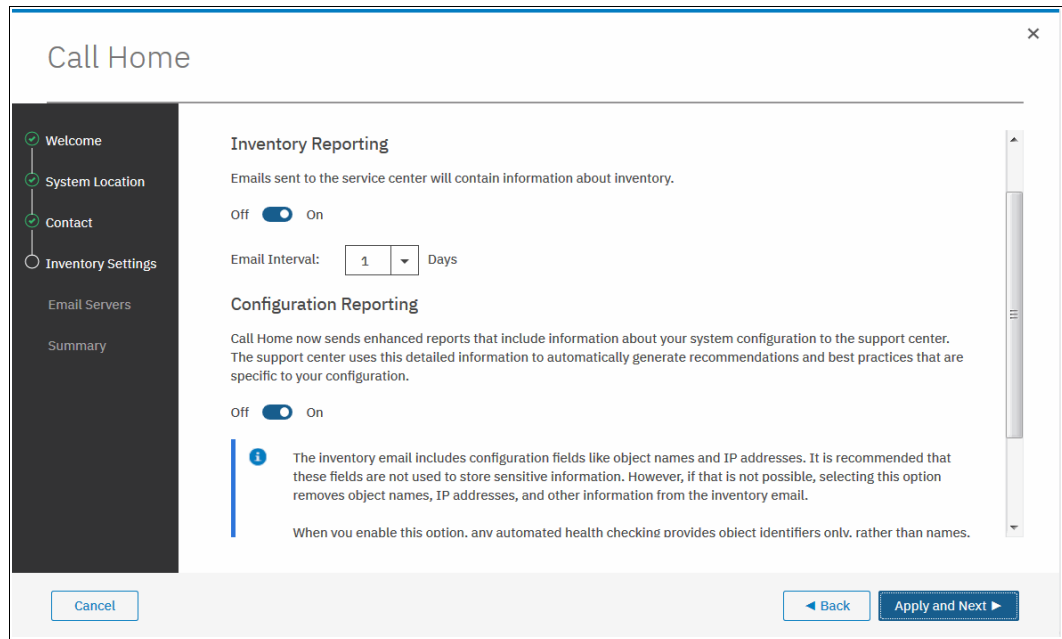


Figure 13-49 Inventory Reporting and Configuration Reporting

4. Configure the IP address of your company SMTP server, as shown in Figure 13-50. When the correct SMTP server is provided, you can test the connectivity by pinging its IP address. You can configure more SMTP servers by clicking the **Plus** sign (+) at the end of the entry line. When you are done, click **Apply and Next**.

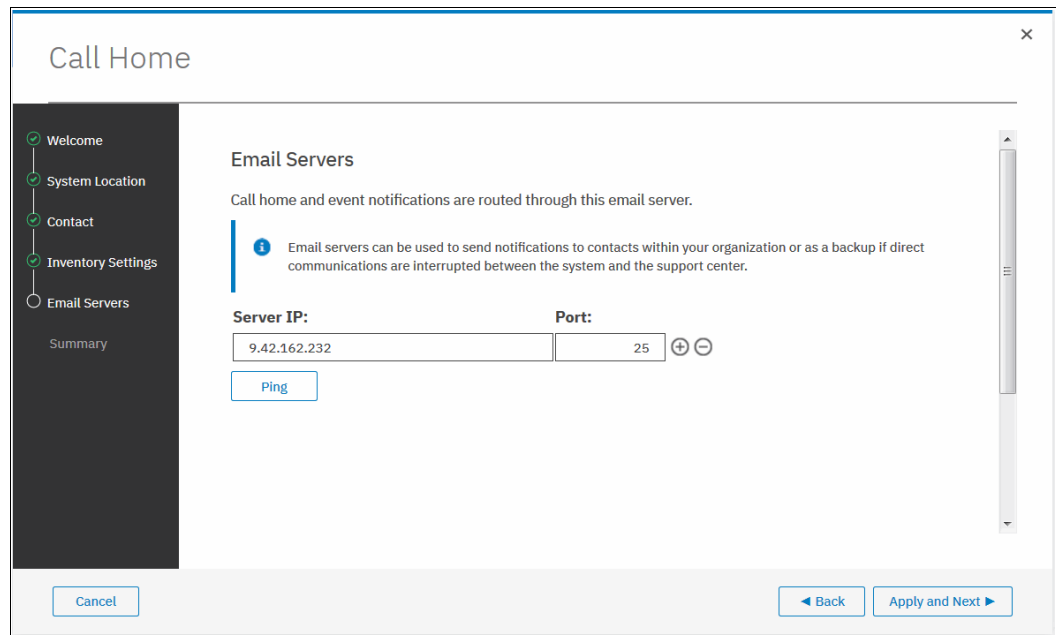


Figure 13-50 Configure email server IP settings

5. A summary window opens. Verify all of the information, and then click **Finish**. You are then returned to the Call Home window where you can verify email addresses of IBM Support (callhome0@de.ibm.com) and optionally add local users who also need to receive notifications. For more information, see Figure 13-51.

**Note:** The default support email address callhome0@de.ibm.com is predefined by the system to receive Error Events and Inventory. Do not change these settings.

You can modify or add local users by using Edit mode after the initial configuration is saved.

The Inventory Reporting function is enabled by default for Call Home. Rather than reporting a problem, an email is sent to IBM that describes your system hardware and critical configuration information. Object names and other information, such as IP addresses, are not included. By default, the inventory email is sent weekly, allowing an IBM Cloud service to analyze and inform you whether the hardware or software that you are using requires an update because of any known issue, as detailed in 13.6, “Health checker feature” on page 781.

The screenshot shows the 'Call Home' configuration page. On the left is a navigation menu with 'Call Home' selected. The main content area has a title 'Call Home' and a subtitle 'The support user receives call home events. Local users also receive event notifications.' Below this are 'Edit' and 'Disable Notifications' buttons. The 'Transmission Settings' section has two checked checkboxes: 'Send data using Call Home cloud service' and 'Send data with email notifications'. The 'Call Home with cloud services' section shows 'Connection: Active' with a green checkmark and 'Last Connection: Failure at'. The 'Call Home with email notifications' section displays 'IP Address: 9.42.162.232', 'Server Port: 25', and 'Status: Untried'. The 'Support Center Email' section shows 'Email Address: callhome0@de.ibm.com' and two checked checkboxes for 'Error Events' and 'Inventory', along with a 'Test' button.

Figure 13-51 Call Home settings, email recipients, and alert types

6. After completing the configuration wizard, test the email function. To do so, enter Edit mode, as shown in Figure 13-52 on page 796. In the same window, you can define more email recipients or alter any contact and location details as needed.

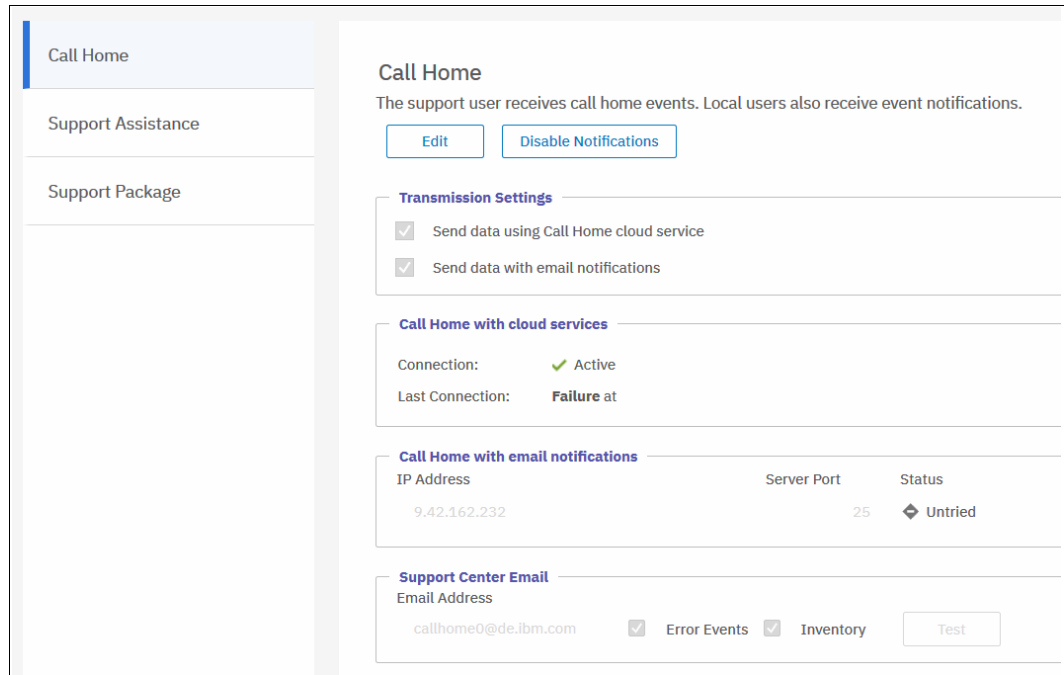


Figure 13-52 Entering Edit mode

We strongly suggest that you keep the sending inventory option enabled to IBM Support. However, it might not be of interest to local users, although inventory content can serve as a basis for inventory and asset management.

7. In Edit mode, you can change any of the previously configured settings. After you are finished editing these parameters, adding more recipients, or just testing the connection, save the configuration to make the changes take effect (see Figure 13-53).

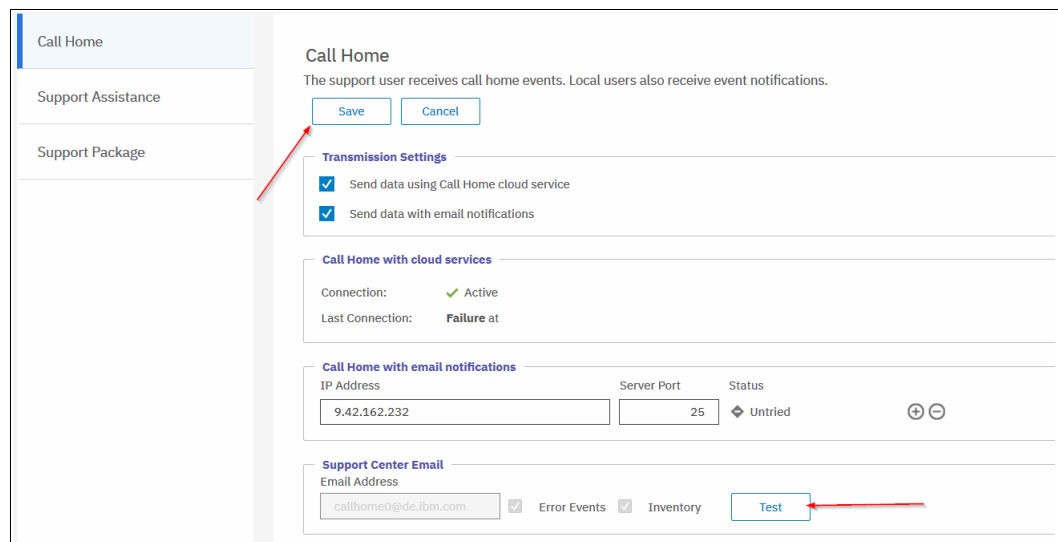


Figure 13-53 Saving modified configuration

**Note:** The **Test** button appears for new email users after first saving and then editing again.

## 13.8.2 Disabling and enabling notifications

At any time, you can temporarily or permanently disable email notifications, as shown in Figure 13-54. This is a best practice when performing activities in your environment that might generate errors on your IBM Spectrum Virtualize cluster, such as SAN reconfiguration or replacement activities. After the planned activities, remember to re-enable the email notification function. The same results can be achieved by running the `svctask stopmail` and `svctask startmail` commands.

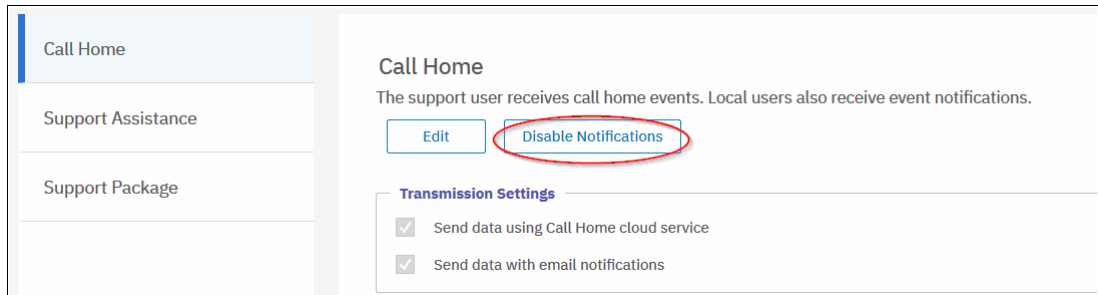


Figure 13-54 Disabling or enabling email notifications

## 13.8.3 Remote Support Assistance

Introduced with Version 8.1, Remote Support Assistance, allows IBM Support to connect remotely to the IBM SAN Volume Controller through a secure tunnel to perform analysis, log collection, and software updates. The tunnel can be enabled *ad hoc* by the client, or the client can enable a permanent connection if wanted.

**Note:** Clients who purchased Enterprise Class Support (ECS) are entitled to IBM Support by using Remote Support Assistance to connect and diagnose problems quickly. However, IBM Support might choose to use this feature on non-ECS systems at their discretion. Therefore, configure and test the connection on all systems.

If you are enabling Remote Support Assistance, ensure that the following prerequisites are met:

- ▶ Cloud Call Home or a valid email server are configured (Cloud Call Home is used as the primary method to transfer the token when you start a session, with email as back up).
- ▶ Ensure that a valid service IP address is configured on each node on the IBM Spectrum Virtualize cluster.
- ▶ If your IBM SAN Volume Controller is behind a firewall or if you want to route traffic from multiple storage systems to the same place, you must configure a Remote Support Proxy server. Before you configure Remote Support Assistance, the proxy server must be installed and configured separately. During the setup for support assistance, specify the IP address and the port number for the proxy server in the Remote Support Centers window.
- ▶ If you do not have firewall restrictions and the IBM SAN Volume Controller nodes are directly connected to the internet, request your network administrator to allow connections to 129.33.206.139 and 204.146.30.139 on port 22.
- ▶ Both uploading support packages and downloading software require direct connections to the internet. A DNS server must be defined on your IBM SAN Volume Controller for both of these functions to work.

- ▶ To ensure that support packages are uploaded correctly, configure the firewall to allow connections to the following IP addresses on port 443: 129.42.56.189, 129.42.54.189, and 129.42.60.189.
- ▶ To ensure that software is downloaded correctly, configure the firewall to allow connections to the following IP addresses on port 22: 170.225.15.105,170.225.15.104, 170.225.15.107, 129.35.224.105, 129.35.224.104, and 129.35.224.107.

Figure 13-55 shows a window that opens after you update your IBM Spectrum Virtualize software to Version 8.1 or above. It prompts you to configure your IBM SAN Volume Controller for remote support. You can choose to not enable it, open a tunnel when needed, or to open a permanent tunnel to IBM.

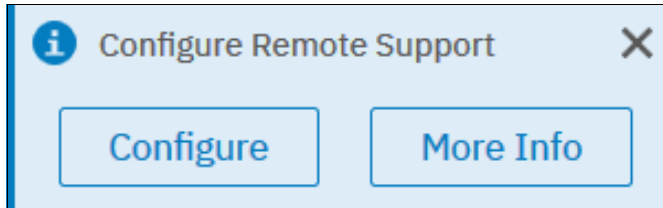


Figure 13-55 Prompt to configure Remote Support Assistance

You can choose to configure IBM SAN Volume Controller, learn some more about the feature, or close the window by clicking the X. Figure 13-56 shows how you can find the Setup Remote Support Assistance if you closed the window.

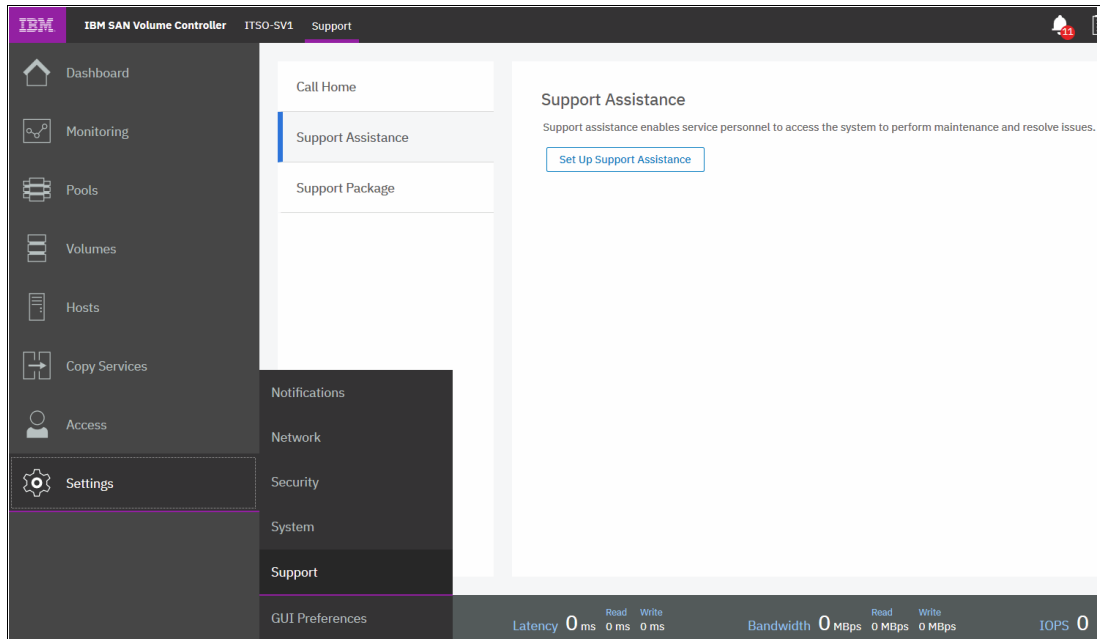


Figure 13-56 Remote Support Assistance menu



Choosing to set up support assistance opens a wizard to guide you through the configuration. Complete the following steps:

1. Figure 13-57 shows the first wizard window. Select **I want support personnel to work on-site only** or enable remote assistance by selecting **I want support personnel to access my system both on-site and remotely**. Click **Next**.

**Note:** Selecting **I want support personnel to work on-site only** does not entitle you to expect IBM Support to be onsite for all issues. Most maintenance contracts are for customer-replaceable unit (CRU) support, where IBM diagnoses your problem and send a replacement component for you to replace, if required. If you prefer to have IBM perform replacement tasks for you, contact your local sales person to investigate an upgrade to your current maintenance contract.

Set Up Support Assistance

Remote support enables technical experts to instantly access your system through the Internet to analyze and resolve problems rapidly and efficiently. Support personnel can do the following tasks during remote support sessions:

- Address configuration or performance issues without waiting for on-site support personnel.
- Determine whether more detailed logs are required and start any necessary uploads of these logs to support.
- Initiate a software download, assist you with software installation, and verify updates have resolved issues.
- Reduce your effort to resolve issues and ensure minimum interruption to operations.

How do you want to set up support assistance?

I want support personnel to work on-site only

I want support personnel to access my system both on-site and remotely

? Cancel < Back Next

Figure 13-57 Remote Support wizard enable or disable

2. The next window, which is shown in Figure 13-58, lists the IBM Support center's IP addresses and Secure Shell (SSH) port that must be open in your firewall. You can also define a Remote Support Assistance Proxy if you have multiple Storwize V7000 or IBM SAN Volume Controller systems in the data center so that the firewall configuration is required for only the proxy server rather than every storage system. Because we do not have a proxy server, we leave the field blank and click **Next**.

**Set Up Support Assistance**

**Support Centers**

Support centers respond to manual and automatic service requests from the system. The following support centers are configured on the system:

Name	IP Address	Port
default_support_center0	129.33.206.139	22
default_support_center1	204.146.30.139	22

**Remote Support Proxy (Optional)**

**i** A proxy is required for network configurations using a firewall, or for systems without direct connections to the network.

Name  IP  Port  **+**

?

Figure 13-58 Remote Support wizard proxy setup

- The next window opens and prompts you to open a tunnel to IBM permanently so that IBM can connect to your IBM SAN Volume Controller. Your options are **At Any Time** or **On Permission Only**, as shown in Figure 13-59. **On Permission Only** requires a storage administrator to log on to the GUI and enable the tunnel when required. Click **Finish**.

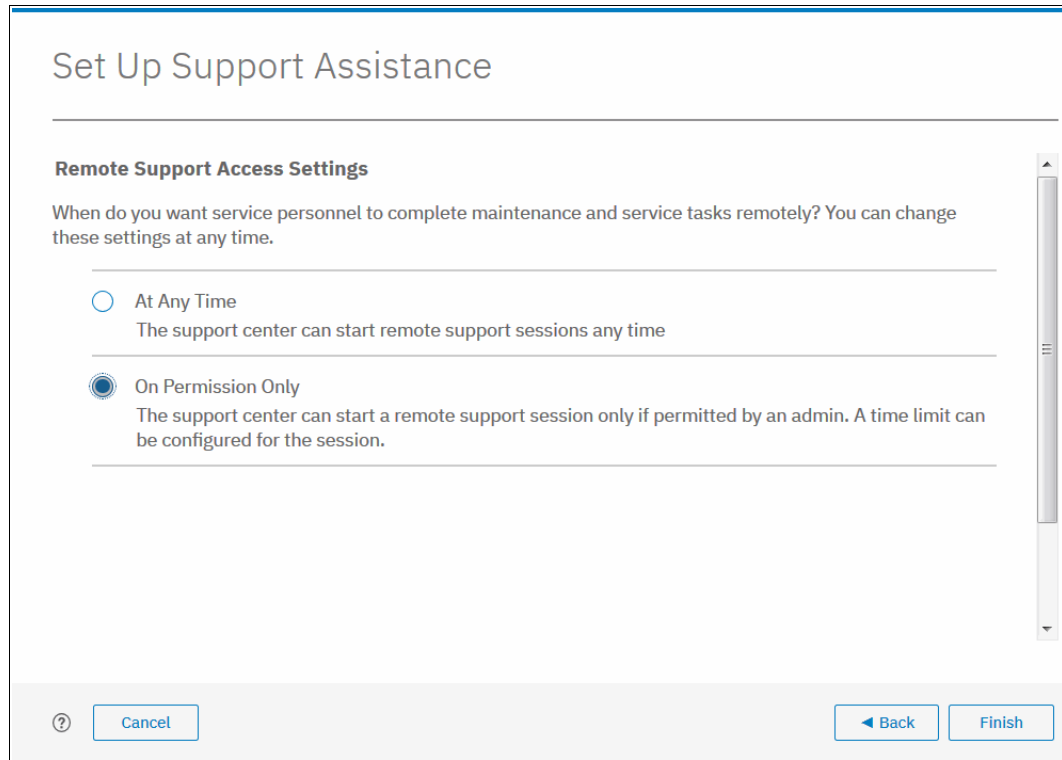


Figure 13-59 Remote Support wizard access choice

- After completing the remote support setup, you can view the status of any remote connection, start a session, test the connection to IBM, and reconfigure the setup. In Figure 13-60, we successfully tested the connection. Click **Start New Session** to open a tunnel for IBM Support to connect through.

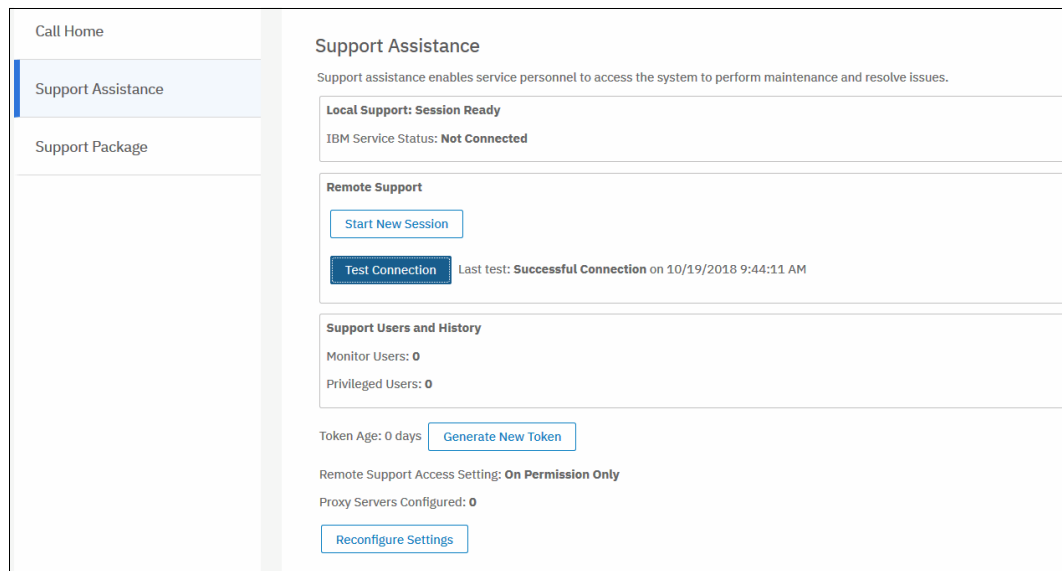


Figure 13-60 Remote support status and session management

5. A window opens and prompts you to set a timeout value for when to close the tunnel to if there is no activity for a period. As shown in Figure 13-61, the connection is established and waits for IBM Support to connect.

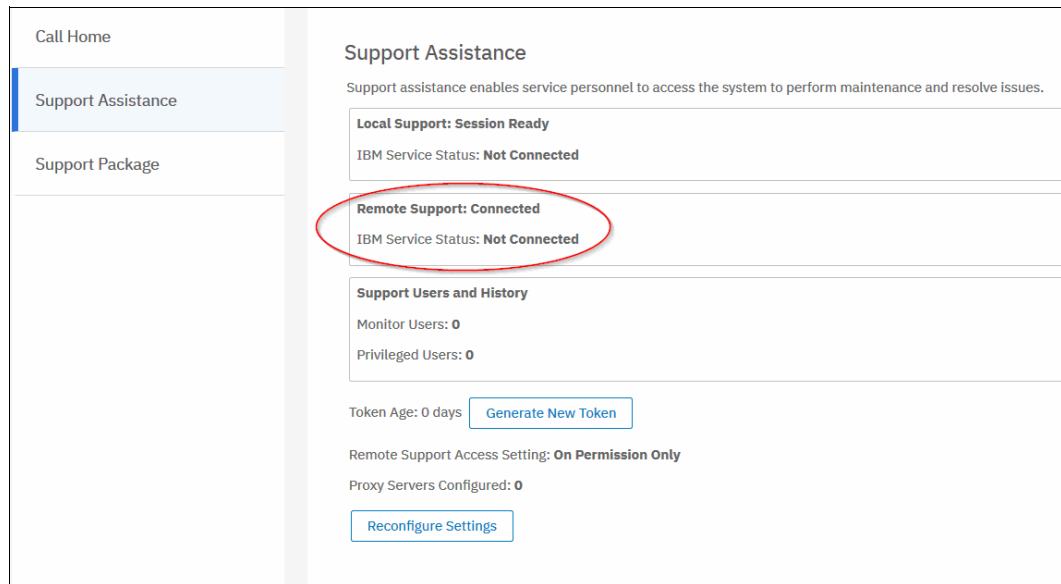


Figure 13-61 Remote Assistance tunnel connected

## 13.8.4 SNMP configuration

SNMP is a standard protocol for managing networks and exchanging messages. The system can send SNMP messages that notify personnel about an event. You can use an SNMP manager to view the SNMP messages that are sent by the IBM SAN Volume Controller.

You can configure an SNMP server to receive various informational, error, or warning notifications by entering the following information (see Figure 13-62 on page 803):

► **IP Address**

The address for the SNMP server.

► **Server Port**

The remote port number for the SNMP server. The remote port number must be a value of 1 - 65535, where the default is port 162 for SNMP.

► **Community**

The SNMP community is the name of the group to which devices and management stations that run SNMP belong. Typically, the default of `public` is used.

► **Event Notifications**

Consider the following points about event notifications:

- Click **Error** if you want the user to receive messages about problems, such as hardware failures, that require prompt action.

**Important:** Select **Recommended Actions** to run the fix procedures on these notifications.

- Click **Warning** if you want the user to receive messages about problems and unexpected conditions. Investigate the cause immediately to determine any corrective action, such as a space-efficient volume running out of space.

**Important:** Select **Recommended Actions** to run the fix procedures on these notifications.

- Click **Info** if you want the user to receive messages about expected events. No action is required for these events.

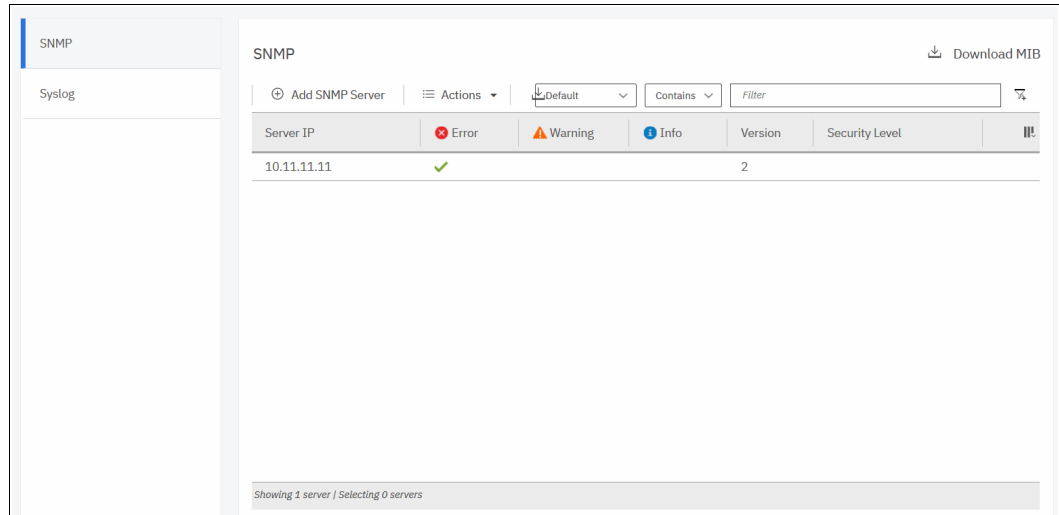


Figure 13-62 SNMP configuration

To add an SNMP server, click **Actions** → **Add** and complete the **Add SNMP Server** window, as shown in Figure 13-63. To remove an SNMP server, click the line with the server you want to remove, and click **Actions** → **Remove**.

The screenshot shows the 'Add SNMP Server' dialog box. It has a title bar with a close button (X). The dialog contains several input fields and checkboxes. The fields are: 'Server IP\*' with the value '9.78.41.9', 'Community\*' with the value 'public', 'Port\*' with the value '162', 'Engine ID' with the placeholder 'Enter ID', 'Security Name' with the placeholder 'Enter Username', 'Authentication Protocol' with a dropdown menu showing 'Select an option', 'Authentication Passphrase (8 characters min.)' with the placeholder 'Enter Passphrase', 'Privacy Protocol' with a dropdown menu showing 'Select an option', and 'Privacy Passphrase (8 characters min.)' with the placeholder 'Enter Passphrase'. There are also checkboxes for 'Events\*': 'Error' (checked), 'Warning' (unchecked), and 'Info' (unchecked). At the bottom left is a 'Need Help' link, and at the bottom right are 'Cancel' and 'Add' buttons.

Figure 13-63 Add SNMP Server

**Note:** It is optional to define the following properties:

- ▶ Engine ID: Indicates the unique identifier (UID) in hexadecimal that identifies the SNMP server
- ▶ Security Name: Indicates which security controls are configured for the SNMP server. Supported security controls are: none, authentication, or authentication and privacy.
- ▶ Authentication Protocol: Indicates the authentication protocol that is used to verify the system to the SNMP server.
- ▶ Privacy Protocol: Indicates the encryption protocol that is used to encrypt data between the system and the SNMP server.
- ▶ Privacy Passphrase: Indicates the user-defined passphrase that is used to verify encryption between the system and SNMP server.

## 13.8.5 Syslog notifications

The syslog protocol is a standard protocol for forwarding log messages from a sender to a receiver on an IP network. The IP network can be Internet Protocol Version 4 (IPv4) or Internet Protocol Version 6 (IPv6). The system can send syslog messages that notify personnel about an event.

You can configure a syslog server to receive log messages from various systems and store them in a central repository by entering the following information into the window that is shown in Figure 13-64

IP Address or Domain	Port	Protocol	Facility	Message Format	Error	Warning	Info	Audit Log	Authentication
9.80.45.1	514	UDP	Level 0	Concise	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>

Figure 13-64 Syslog configuration

► IP Address

The IP address for the syslog server.

► Facility

The facility determines the format for the syslog messages. You can use the facility to determine the source of the message.

► Message Format

The message format depends on the facility. The system can transmit syslog messages in the following formats:

- The concise message format provides standard details about the event.
- The expanded format provides more details about the event.

► Event Notifications

Consider the following points about event notifications:

- Click **Error** if you want the user to receive messages about problems, such as hardware failures, which must be resolved immediately.

**Important:** Browse to **Recommended Actions** to run the fix procedures on these notifications.

- Click **Warning** if you want the user to receive messages about problems and unexpected conditions. Investigate the cause immediately to determine whether any corrective action is necessary.

**Important:** Browse to **Recommended Actions** to run the fix procedures on these notifications.

- Click **Info** if you want the user to receive messages about expected events. No action is required for these events.

To remove a syslog server, click the **Minus** sign (-).





- ▶ dumpintervallog
- ▶ svcserVICetak dumperrlog
- ▶ svcserVICetak finderr

Figure 13-65 shows the access to the audit log. Click **Audit Log** in the left menu to see which configuration CLI commands ran on the IBM SAN Volume Controller system.

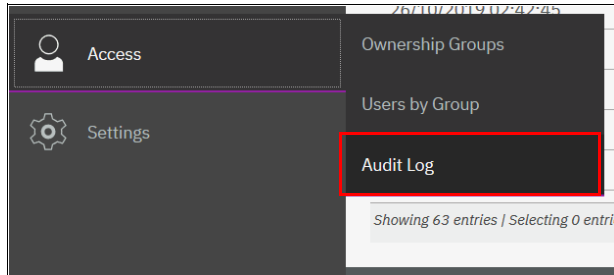


Figure 13-65 Audit Log from Access menu

Figure 13-66 shows an example of the audit log shows an example of the audit log after a volume is created and mapped to a host.

Date and Time	User Name	Command	Object ID
25/10/2019 17:54:43	JackUser	svctask mkvdiskhostmap -force -gui -host 0 -scsi 1 35	
25/10/2019 17:54:25	JackUser	svctask mkvdisk -gui -iogrp io_grp1 -mdiskgrp 0 -name Redbook...	35
25/10/2019 17:53:46	superuser	svctask startfcmap -gui -prep 1	
25/10/2019 17:52:00	superuser	svctask chfcmap -cleanrate 50 -copyrate 100 -gui 1	
25/10/2019 17:51:53	superuser	svctask chenclosure -gui -managed yes 2	
25/10/2019 17:44:51	superuser	svctask addcontrolenclosure -gui -iogrp 1 -sernum 7822PFG	
25/10/2019 02:00:20	superuser	satask cfiles -prefix /dumps/svc.config.cron.*_7822PBR-2 -sour...	
25/10/2019 02:00:02	superuser	svctask detectmdisk	
25/10/2019 00:23:49	JackUser	svctask mksnmpserver -community public -error on -gui -info off ...	1
25/10/2019 00:11:04	JackUser	svctask chsra -gui -idletimeout 60 -remotesupport enable	
25/10/2019 00:09:25	JackUser	svctask chsra -enable -gui	
24/10/2019 23:53:31	JackUser	svctask startemail -gui	
24/10/2019 23:53:17	JackUser	svctask mkemailserver -gui -ip 9.71.47.10 -port 25	0
24/10/2019 23:53:14	JackUser	svctask mkemailuser -address callhome1@de.ibm.com -error on...	

Figure 13-66 Audit log

Changing the view of the Audit Log grid is also possible by right-clicking column headings (see Figure 13-67). The grid layout and sorting is under the user’s control. Therefore, you can view everything in the audit log, sort different columns, and reset the default grid preferences.

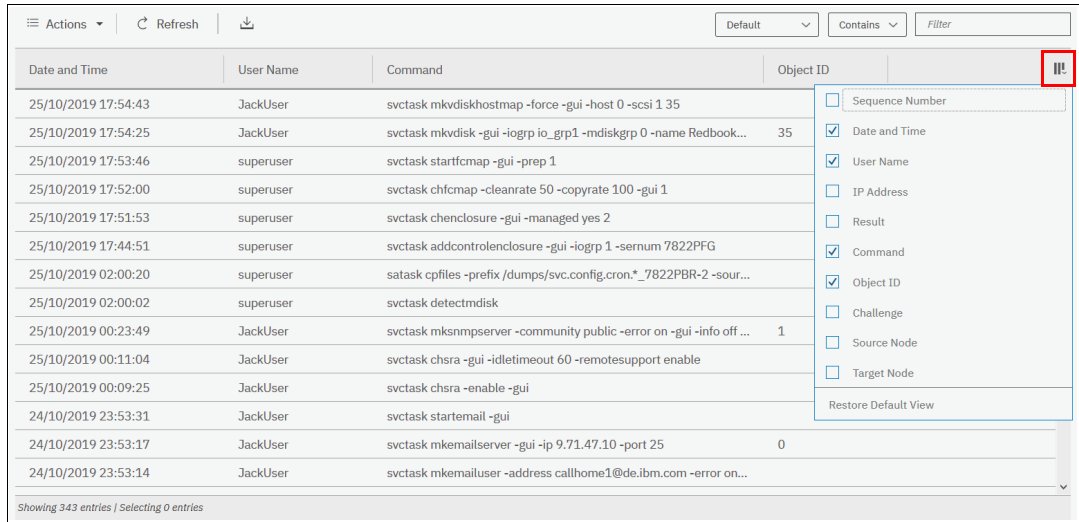


Figure 13-67 Right-click audit log column headings

## 13.10 Collecting support information by using the GUI, CLI, and USB

At times, if you have a problem and call the IBM Support Center, they most likely ask you to provide a support package. You can collect and upload this package by clicking **Settings** → **Support**.

### 13.10.1 Collecting information by using the GUI

To collect information by using the GUI, complete the following steps:

1. Click **Settings** → **Support**, and then click the **Support Package** tab (see Figure 13-68).
2. Click **Upload Support Package**.

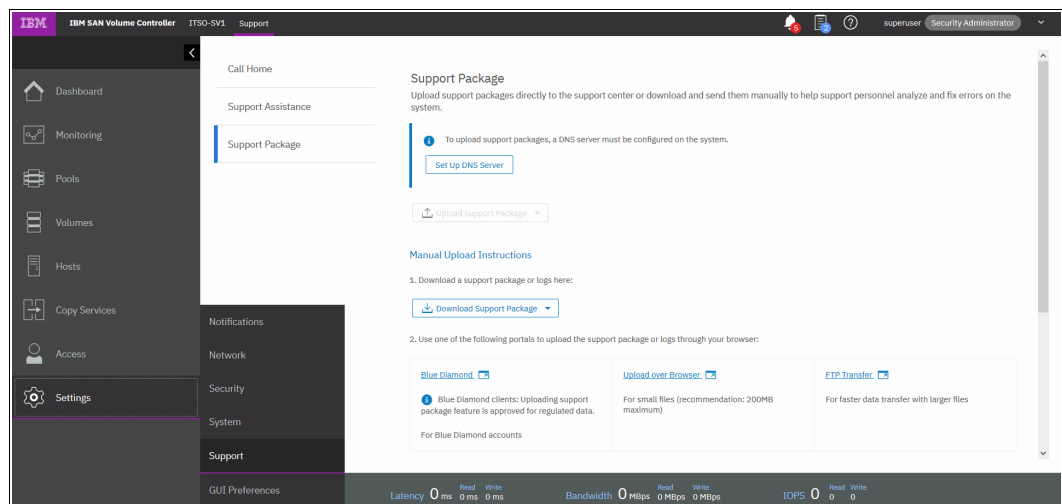


Figure 13-68 Support Package option

Assuming that the problem that was encountered was an unexpected node restart that logged a 2030 error, collect the default logs plus the most recent statesave from each node to capture the most relevant data for support.

**Note:** When a node unexpectedly restarts, it first dumps its current statesave information before it restarts to recover from an error condition. This statesave is critical for support to analyze what happened. Collecting a snap type 4 creates new statesaves at the time of the collection, which is not useful for understanding the restart event.

3. The **Upload Support Package** window provides four options for data collection. If you are contacted by IBM Support because your system is calling home or you manually open a call with IBM Support, you are given a Problem Management Record (PMR) number. Enter that PMR number into the **PMR** field and select the snap type, often referred to as an *option 1, 2, 3, 4 snap*, as requested by IBM Support (see Figure 13-69). In our case, we enter our PMR number, select **Snap Type 3** (option 3) because this automatically collects the statesave that is created at the time the node restarted, and click **Upload**.

**Tip:** To open a service request online, see the IBM Support Service requests and PMRs [web page](#).

Upload Support Package

PMR Number: [Don't have PMR?](#)

ppppp,bbb,ccc

Select the type of new support package to generate and upload to the IBM support center:

- Snap Type 1: Standard logs  
Contains the most recent logs for the system, including the event and audit logs.
- Snap Type 2: Standard logs plus one existing statesave  
Contains all the standard logs plus one existing statesave from any of the nodes in the system.
- Snap Type 3: Standard logs plus most recent statesave from each node  
Contains all the standard logs plus each node's most recent statesave.

? Need Help   Cancel   Upload

Figure 13-69 Upload Support Package window

- The procedure to create the snap on an IBM SAN Volume Controller system, including the latest statesave from each node, begins. This process might take a few minutes (see Figure 13-70).

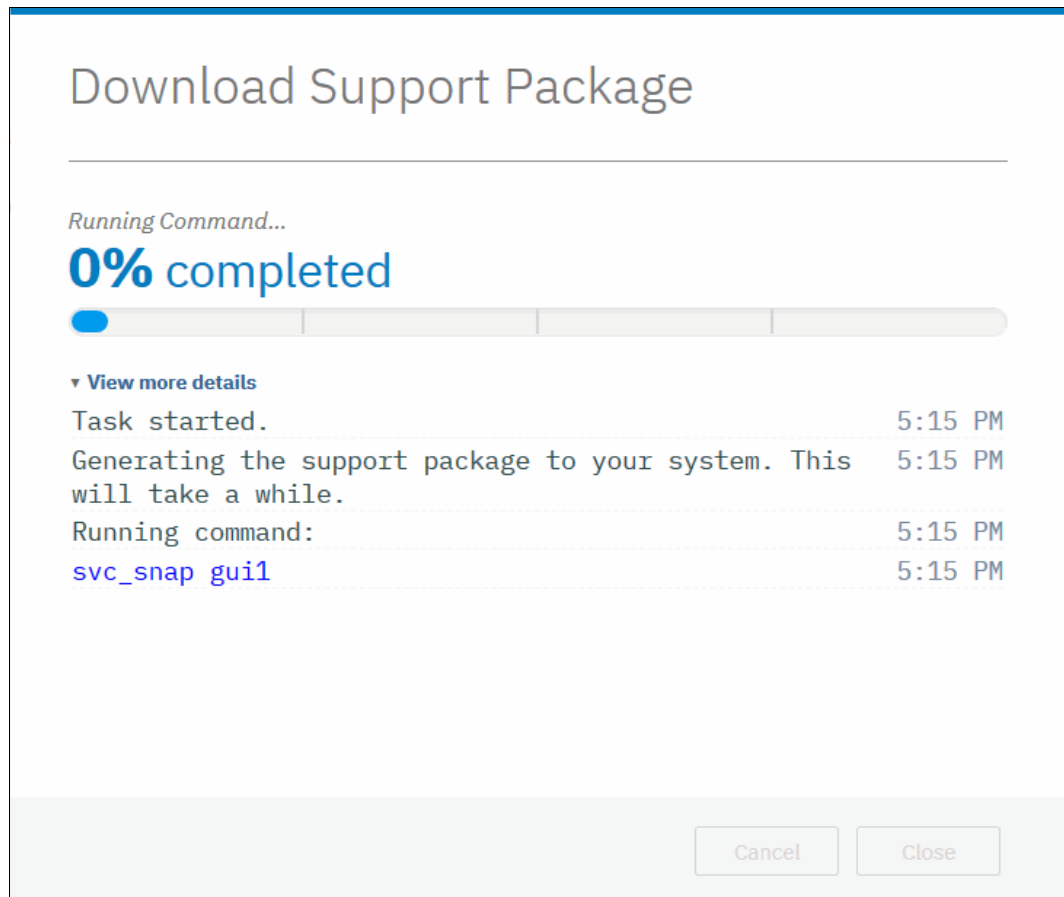


Figure 13-70 Task detail window

The time it takes to generate the SNAP and the size of the file that is generated depends mainly on two factors: the SNAP option you selected, and the size of your system.

An option 1 SNAP takes much less time than an option 4 SNAP; for example, because nothing new needs to be gathered for an option 1 SNAP, whereas an option 4 SNAP requires the system to collect new statesaves from each node. In an 8-node cluster, this process can take quite some time, so you should always collect the SNAP option that IBM Support recommends to you.

The approximate file sizes for each SNAP option are listed in Table 13-10.

Table 13-10 Types of SNAP

Option	Description	Approximate Size (1 I/O group, 30 volumes)	Approximate Size (4 I/O groups, 250 Volumes)
1	Standard logs	10 MB	340 MB
2	Standard logs plus one existing statesave	50 MB	520 MB

Option	Description	Approximate Size (1 I/O group, 30 volumes)	Approximate Size (4 I/O groups, 250 Volumes)
3	Standard logs plus most recent statesave from each node	90 MB	790 MB
4	Standard logs plus new statesaves	90 MB	790 MB

### 13.10.2 Collecting logs by using the CLI

The CLI can be used to collect and upload a support package as requested by IBM Support by completing the following steps:

1. Log in to the CLI and run the **svc\_snap** command that matches the type of snap that is requested by IBM Support:

- Standard logs (type 1):

```
svc_snap upload pmr=ppppp,bbb,ccc gui1
```

- Standard logs plus one existing statesave (type 2):

```
svc_snap upload pmr=ppppp,bbb,ccc gui2
```

- Standard logs plus most recent statesave from each node (type 3):

```
svc_snap upload pmr=ppppp,bbb,ccc gui3
```

- Standard logs plus new statesaves:

```
svc_livedump -nodes all -yes
svc_snap upload pmr=ppppp,bbb,ccc gui3
```

2. We collect the type 3 (option 3) information, which is automatically uploaded to the PMR number that is provided by IBM Support, as shown in Example 13-8.

*Example 13-8 The svc\_snap command*

---

```
ssh superuser@10.18.228.64
Password:
IBM_2145:ITS0 DH8_B:superuser>>svc_snap upload pmr=04923,215,616 gui3
```

---

3. If you do not want to upload automatically the snap to IBM, do not specify the **upload pmr=ppppp,bbb,ccc** part of the commands. When the snap creation completes, it creates a file that is named with the following format:

```
/dumps/snap.<panel_id>.YYMMDD.hhmmss.tgz
```

It takes a few minutes for the snap file to complete (longer if it includes statesaves).

- The generated file can then be retrieved from the GUI by clicking **Settings** → **Support**, and selecting **Manual Upload Instructions** → **Download Support Package**, and then, clicking **Download Existing Package**, as shown in Figure 13-71.

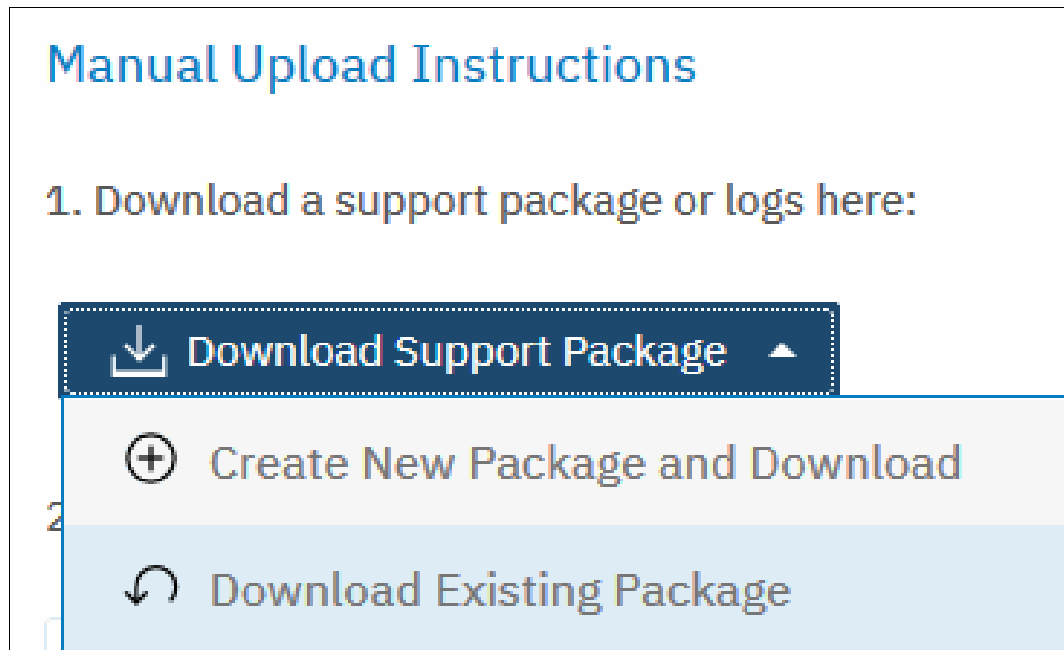


Figure 13-71 Downloaded Existing Package

- Click **Filter** and enter `snap` to see a list of snap files, as shown in Figure 13-72 on page 814. Locate the exact name of the snap that was generated by the `svc_snap` command that ran earlier, select that file, and then click **Download**.

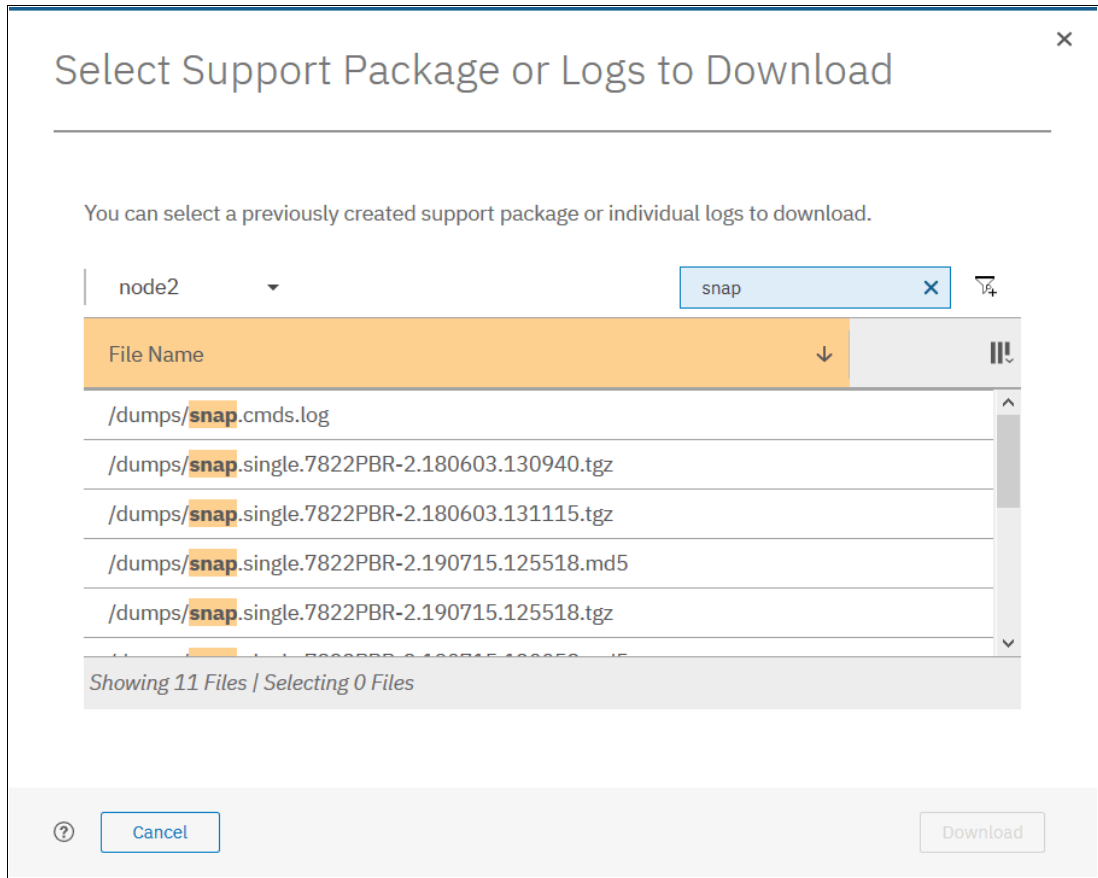


Figure 13-72 Filtering on snap to download

6. Save the file to a folder of your choice on your workstation.

### 13.10.3 Collecting logs by using USB

As a backup, if there is no connectivity to CLI and GUI (which is rare), it is possible to get a SNAP from a node from the USB ports on the rear.

**Note:** This procedure collects a single SNAP from the node canister, not a cluster SNAP. It is useful for determining the state of the node canister.

When a USB flash drive is plugged into a node canister, the canister code searches for a text file that is named `satask.txt` in the root directory. If the code finds the file, it attempts to run a command that is specified in the file. When the command completes, a file that is called `satask_result.html` is written to the root directory of the USB flash drive. If this file does not exist, it is created. If it exists, the data is inserted at the start of the file. The file contains the details and results of the command that was run and the status and the configuration information from the node canister. The status and configuration information matches the detail that is shown on the service assistant home page panels.

Complete the following steps to collect a SNAP:

1. Ensure that your USB drive is formatted with a FAT32 file system on its first partition.
2. Create a text (.txt) file on the USB drive in the root directory called `satask.txt` (case-sensitive).



3. In the `satask.txt` file, write `satask snap` and save it.
4. Put the USB into one of the USB ports on the rear of the canister, and wait for a short time. The fault LED on the node canister flashes when the USB service action is being completed. When the fault LED stops flashing, it is safe to remove the USB flash drive.
5. Unplug the USB drive from the system and plug it into your workstation. If the procedure was successful, the `satask.txt` file was deleted and you have a `satask_result.html` file and a single SNAP from the node canister. This SNAP can then be uploaded to the support center, as shown in 13.10.4, “Uploading files to the Support Center” on page 815.

**Note:** If a problem occurs with the procedure, the html file is still generated and includes the reasons why the procedure failed.

For more information about the commands that can be run by using USB, see [IBM Knowledge Center](#).

### 13.10.4 Uploading files to the Support Center

If you choose to not have the system upload the support package automatically, it can still be uploaded for analysis from the Enhanced Customer Data Repository (ECuRep). Any uploads should be associated with a specific PMR. The PMR is also known as a *service request* and is a mandatory requirement when uploading.

To upload the information, complete the following steps:

1. Using a web browser, go to the [ECuRep web page](#) (see Figure 13-73).

Figure 13-73 ECuRep details

2. Complete the required fields:
  - **PMR number** (mandatory) that is provided by IBM Support for your specific case. This number should be in the format of `ppppp,bbb,ccc`, for example, `04923,215,616`, using a comma (,) as a separator.
  - **Upload is for** (mandatory).
  - Select **Hardware** from the menu.

3. When the form is complete, click **Continue** to open the input window (see Figure 13-74).

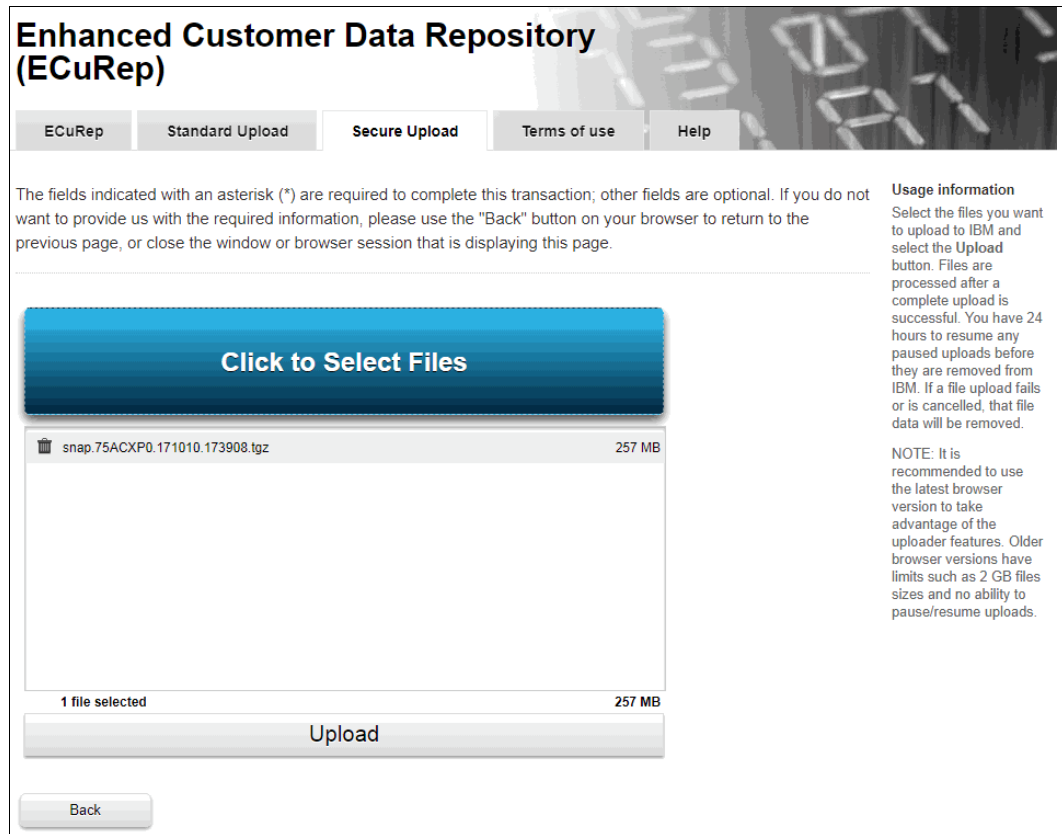


Figure 13-74 ECuRep File upload

4. Select one or more files, click **Upload** to continue, and follow the directions.

## 13.11 Service Assistant Tool

The SAT is a web-based GUI that is used to service individual node canisters, primarily when a node encounters a fault and is in a service state. A node is not an active part of a clustered system while it is in service state.

Typically, an IBM Spectrum Virtualize cluster is initially configured with the following IP addresses:

- ▶ One service IP address for each IBM SAN Volume Controller node.
- ▶ One cluster management IP address, which is set when the cluster is created.

The SAT is available even when the management GUI is not accessible. The following information is available and tasks can be accomplished by using the Service Assistance Tool:

- ▶ Status information about the connections and the IBM SAN Volume Controller nodes
- ▶ Basic configuration information, such as configuring IP addresses
- ▶ Service tasks, such as restarting the Common Information Model Object Manager (CIMOM) and updating the WWNN
- ▶ Details about node error codes

- Details about the hardware, such as IP address and Media Access Control (MAC) addresses

The SAT GUI is available by using a service assistant IP address that is configured on each IBM SAN Volume Controller node. It can also be accessed through the cluster IP addresses by appending /service to the cluster management IP.

It is also possible to access the SAT GUI of the config node if you enter the URL of the service IP of the config node into any web browser and click **Service Assistant Tool**, as shown in Figure 13-75.

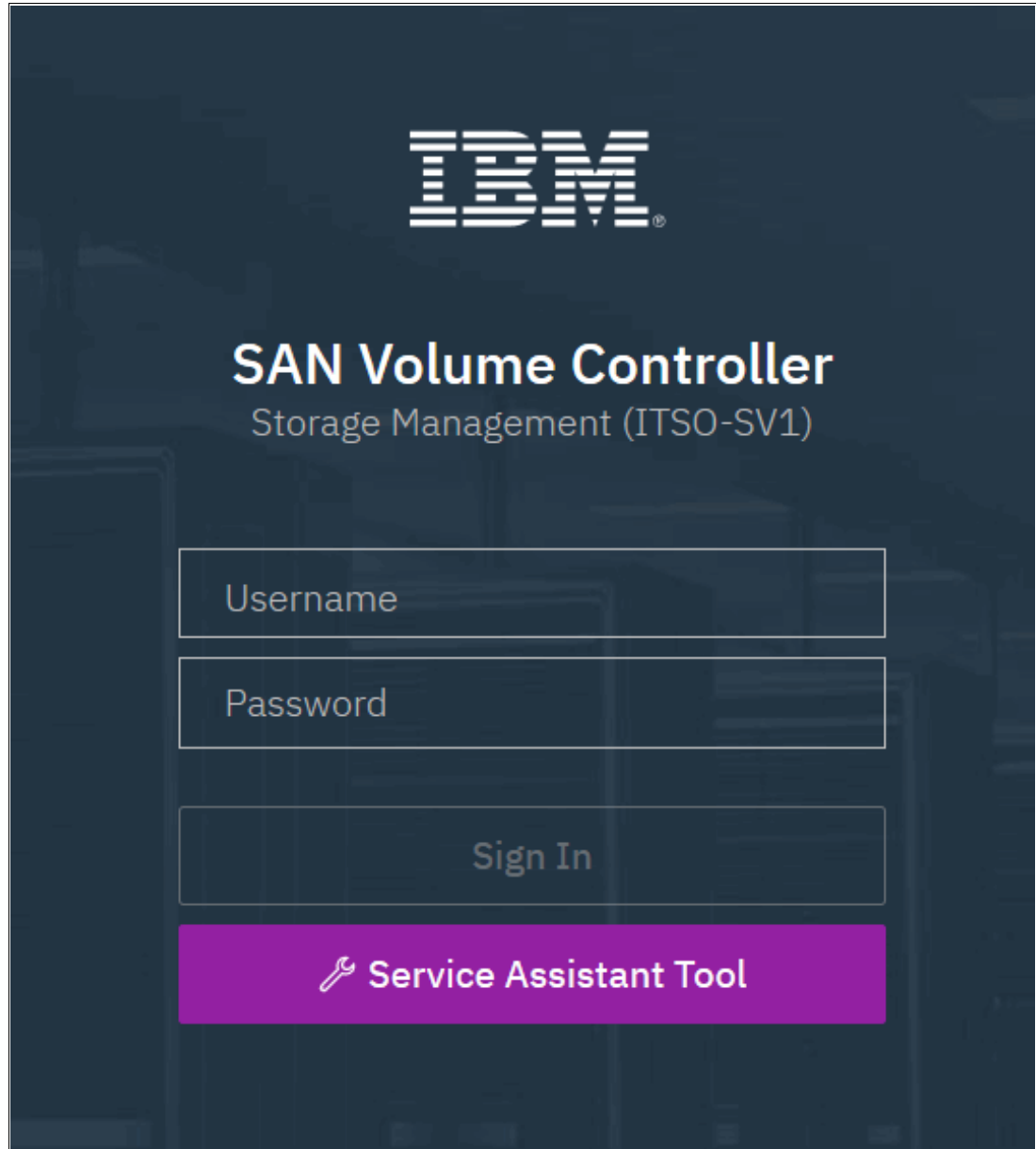


Figure 13-75 Service Assistant Tool

If the clustered system is down, the only method of communicating with the nodes is through the SAT IP address directly. Each node can have a single service IP address on Ethernet port 1 and should be configured on all nodes of the cluster, including any hot spare nodes.

To open the SAT GUI, enter one of the following URLs into any web browser, and then click **Service Assistant Tool**:

- ▶ `http(s)://<cluster IP address of your cluster>/service`
- ▶ `http(s)://<service IP address of a node>/service`
- ▶ `http(s)://<service IP address of config node>`

To access the SAT, complete the following steps:

1. When you are accessing SAT by using `cluster IP address/service`, the configuration node canister SAT GUI login window opens. Enter the Superuser Password, as shown in Figure 13-76.

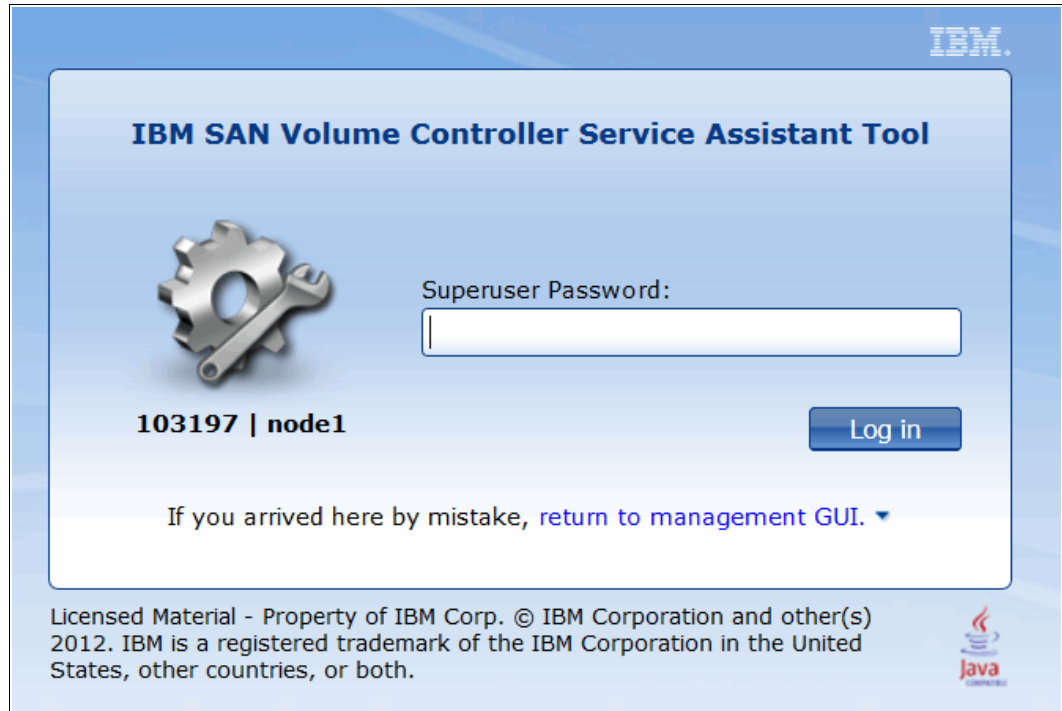


Figure 13-76 Service Assistant Tool Login GUI

- After you are logged in, you see the Service Assistant Home window, as shown in Figure 13-77. The SAT can view the status and run service actions on other nodes, in addition to the node that the user is logged in to.

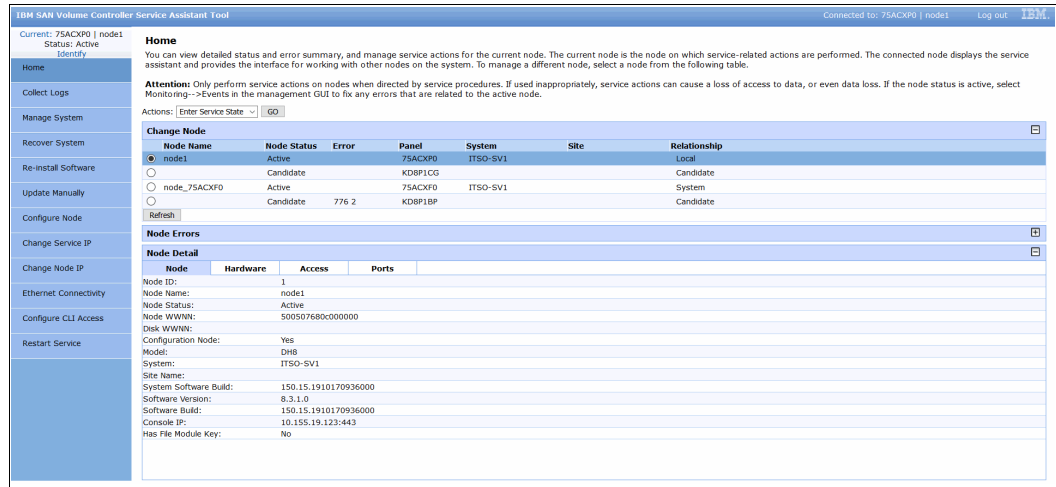


Figure 13-77 Service Assistant Tool GUI

- The current selected IBM SAN Volume Controller node is shown in the upper left corner of the GUI. In Figure 13-77, this is node ID 2. Select the node that you want in the Change Node section of the window. You see the details in the upper left change to reflect the selected node.

**Note:** The SAT GUI provides access to service procedures and shows the status of the nodes. It is advised that these procedures should be carried out only if you are directed to do so by IBM Support.

For more information about how to use the SAT, see [this website](#).

## 13.12 IBM Storage Insights Monitoring

IBM Storage Insights allows you to optimize your storage infrastructure by using this cloud-based storage management and support platform with predictive analytics. For more information about how to set up IBM Storage Insights and the features it offers, see Chapter 2, “Planning” on page 53.

The monitoring capabilities that IBM Storage Insights provides are useful for tasks, such as capacity planning, workload optimization, and managing support tickets for ongoing issues.

After you add your systems to IBM Storage Insights, you see the Dashboard, in which you can select a system that you want to see the overview for, as shown in Figure 13-78 on page 820.

**Note:** The system used to demonstrate IBM Storage Insights is a FlashSystem 9150, but because you can add a large range of systems to IBM Storage Insights, the processes are standard across all. This gives you a single pane of glass to manage your devices.

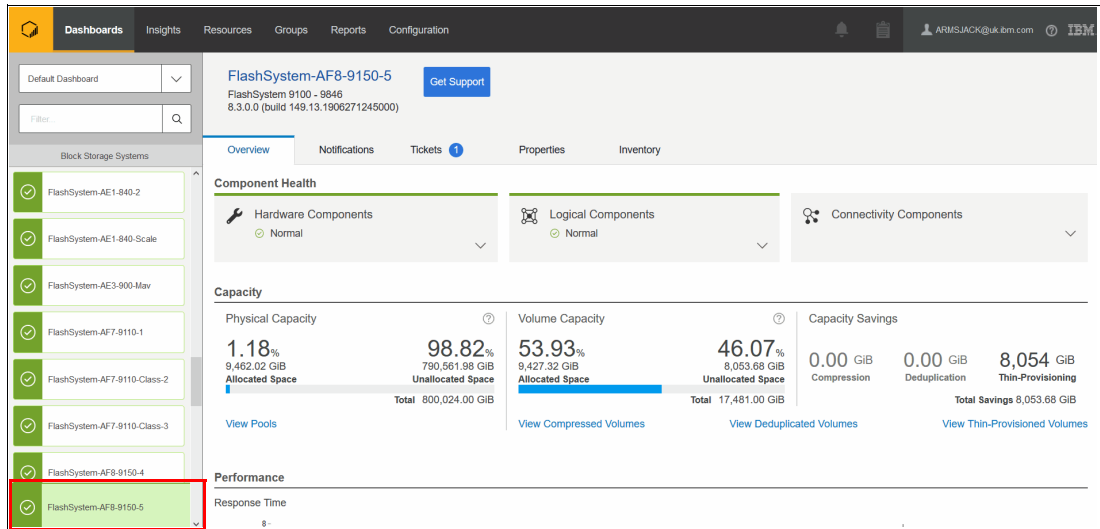


Figure 13-78 IBM Storage Insights System overview

Component health can be seen at the top in the most prominent part of the pane. Here, if there is a problem with one of the Hardware, Logical or Connectivity components, errors are displayed, as shown in Figure 13-79.

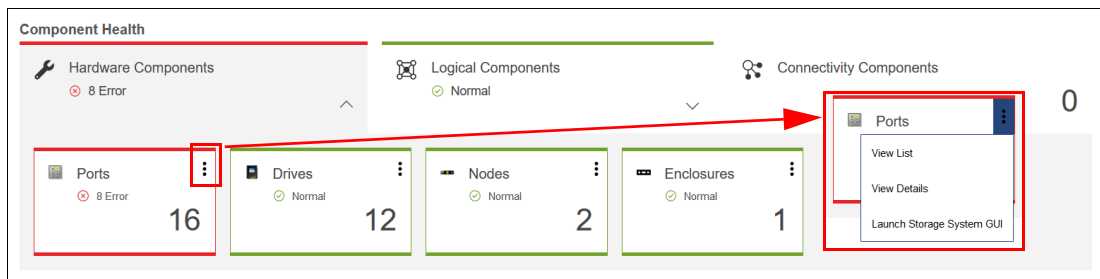


Figure 13-79 Component Overview

The errors can be expanded to obtain more details by selecting the three dots in the upper right corner of the Component that has an error and then selecting **View Details**. The relevant part of the more detailed System View opens, depending on which component has the error, as shown in Figure 13-80.

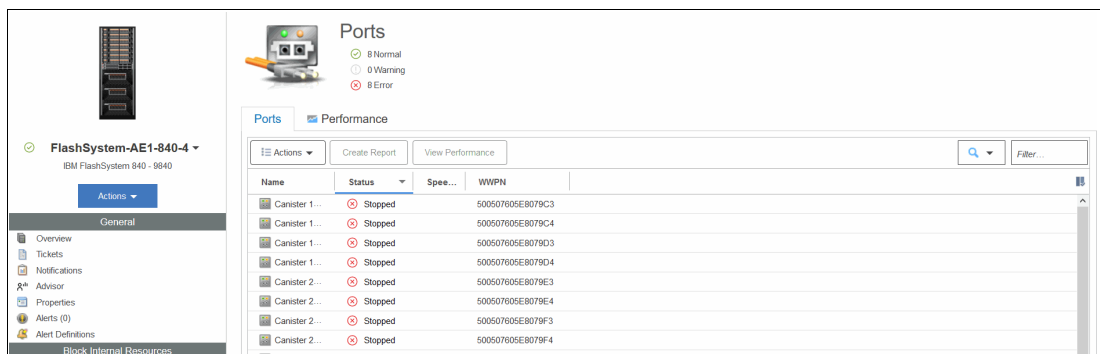


Figure 13-80 Ports in error

From here, it is now obvious which components have the problem and exactly what is wrong with them, which enables you to log a support ticket with IBM if necessary.

## 13.12.1 Capacity Monitoring

You can see key statistics, such as Physical Capacity, Volume Capacity, and Capacity Savings, depending on what is configured after you select a system as shown in Figure 13-81.

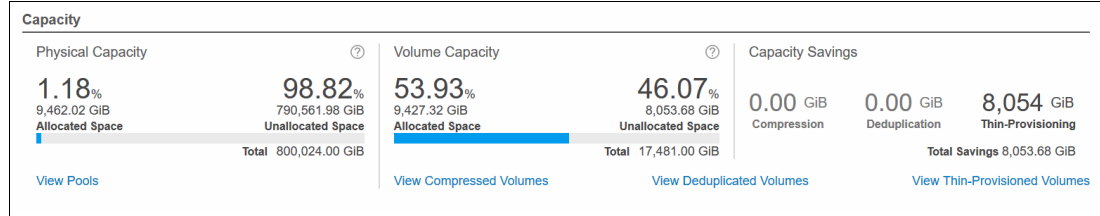


Figure 13-81 Capacity Area of the IBM Storage Insights System Overview

This menu allows the user to View Pools, View Compress Volumes, View Deduplicated Volumes, and View Thin-Provisioned Volumes. Clicking these options takes the user to the detailed System View for the selected option. From there, Capacity can be selected to get a historical view of how system capacity changed over time, as shown in Figure 13-82. At any time, the user can select the wanted timescale, resources, and metrics to be displayed on the graph by clicking any of the options that are around the graph.

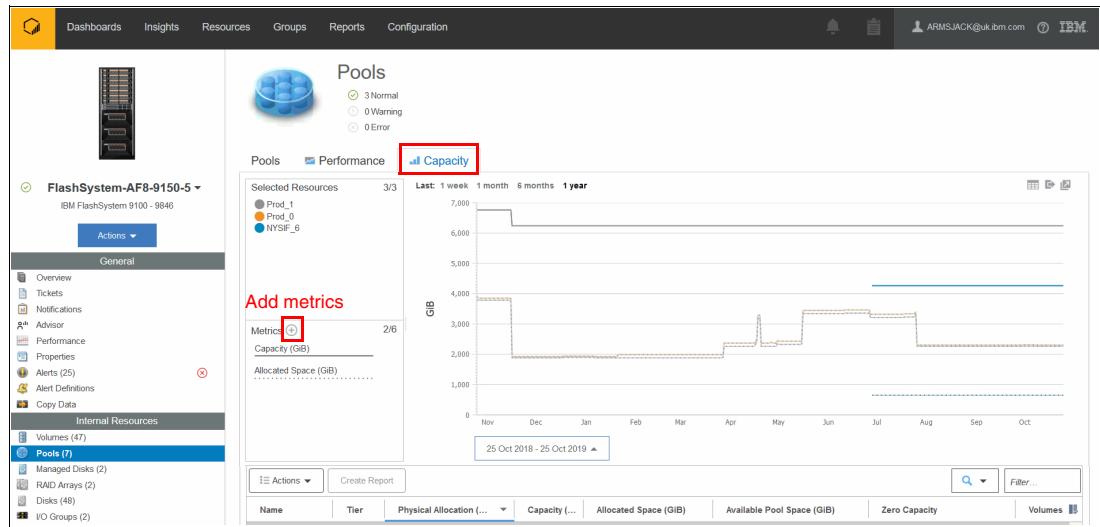


Figure 13-82 IBM Storage Insights Capacity View

If you scroll down, a list view of the selected option can be found below the graph. In this example, we selected **View Pools** so the configured pools are shown with the relevant key capacity metrics, as shown in Figure 13-83. Double-clicking a pool in the table display the pool's properties.

Name	Tier	Physical Allocation (...)	Capacity (...)	Allocated Space (GiB)	Available Pool Space (GiB)	Zero Capacity	Volumes
Prod_0	Tier 1	41 %	6,239.00	2,304.79	3,889.00	None	<a href="#">232</a>
Prod_1	Tier 0	40 %	6,239.00	2,262.50	3,731.00	None	<a href="#">232</a>
NYSIF_6		15 %	4,284.00	650.00	3,610.00	None	<a href="#">6</a>

Showing 3 items | Selected 0 items | 25 Sep 2019 00:00:00 – 25 Oct 2019 21:51:48 | Refreshed a few moments ago

Figure 13-83 Pools List View

## 13.12.2 Performance monitoring

From the system overview, you can scroll down and receive the three key performance statistics for your system, as shown in Figure 13-84. For the performance overview, these stats are aggregated across the whole system, and you cannot drill down by Pool, Volume, and so on.

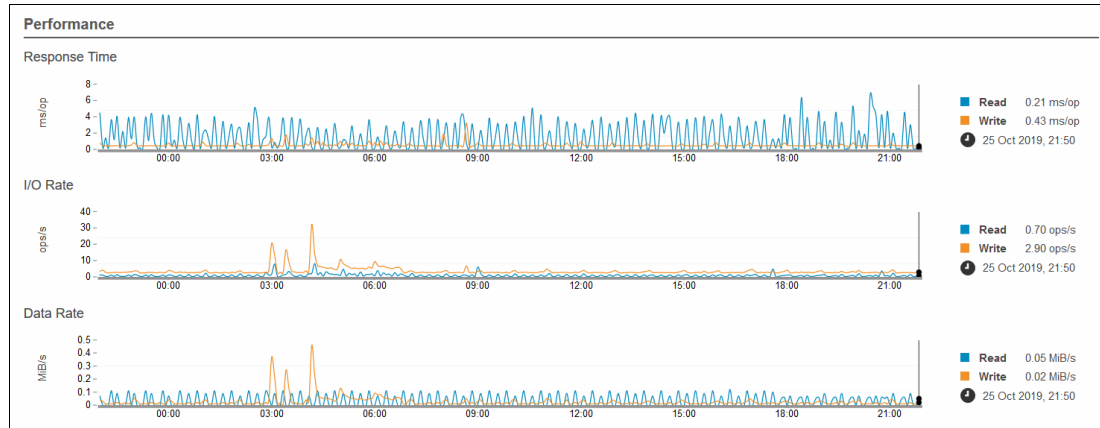


Figure 13-84 System overview Performance

For more information about performance statistics, we must enter the system view again as described in the 13.12.1, “Capacity Monitoring” on page 821 example.

For this performance example, we again select **View Pools**, and then select **Performance** from the System View pane, as shown in Figure 13-85.

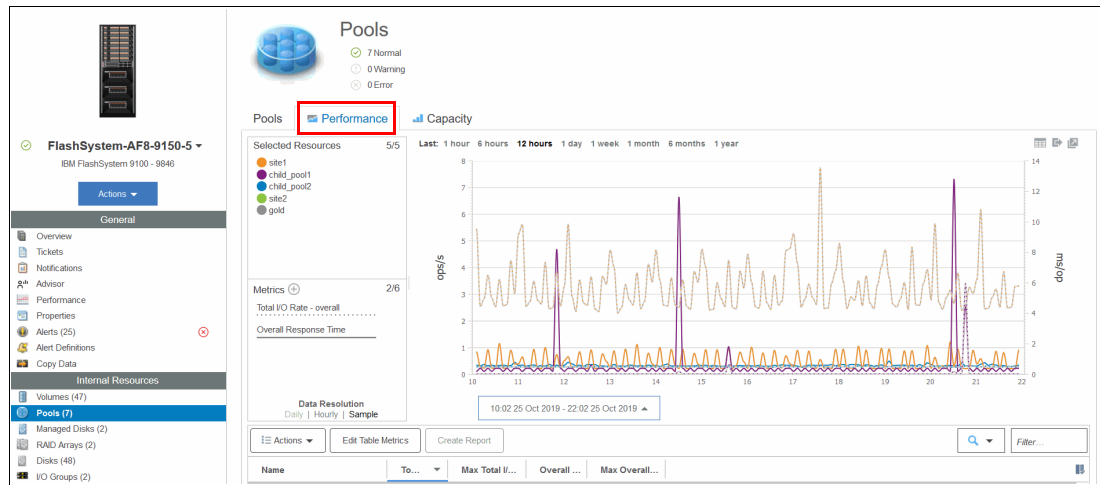


Figure 13-85 IBM Storage Insights Performance View®



It is then possible to customize what can be seen on the graph by selecting the wanted metrics and resources. In Figure 13-86, the Overall Response Time for one Pool over a 12-hour period is displayed.

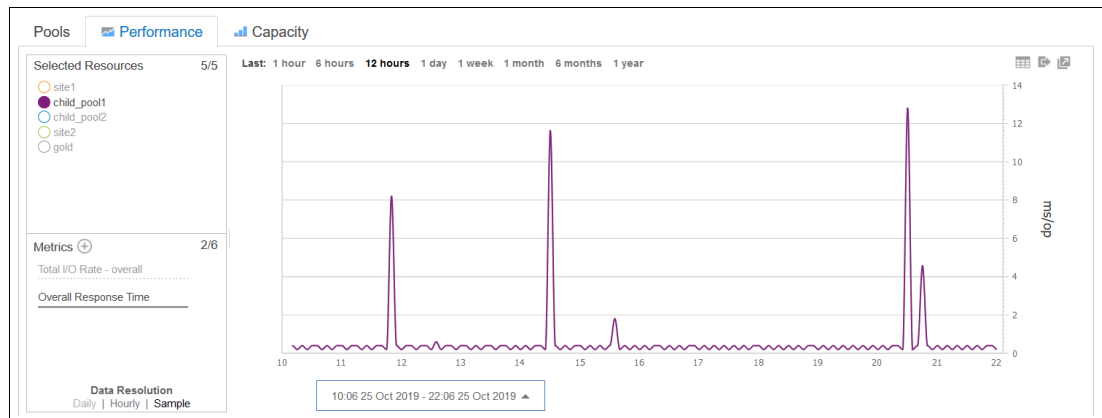


Figure 13-86 Filtered performance graph

Scrolling down below the graph, the Performance List View is visible, as shown in Figure 13-87. Wanted metrics can be selected by clicking the filter button on the right of the column headers. If you select a row, the graph changes to be filtered for that selection only. Multiple rows can be selected by holding down **SHIFT** or **CTRL**.

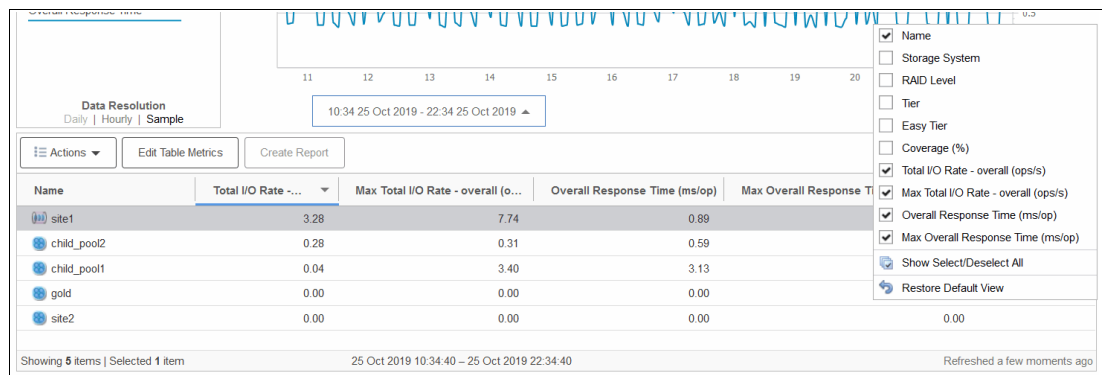


Figure 13-87 Performance List View

### 13.12.3 Logging Support Tickets by using IBM Storage Insights

IBM Storage Insights allows the user to log support tickets that greatly complements the enhanced monitoring opportunities that the software brings. When an issue is detected, and the user wants to engage IBM Support, complete the following steps:

1. Select the system to open the System Overview and click the **Get Support** button, as shown in Figure 13-88.

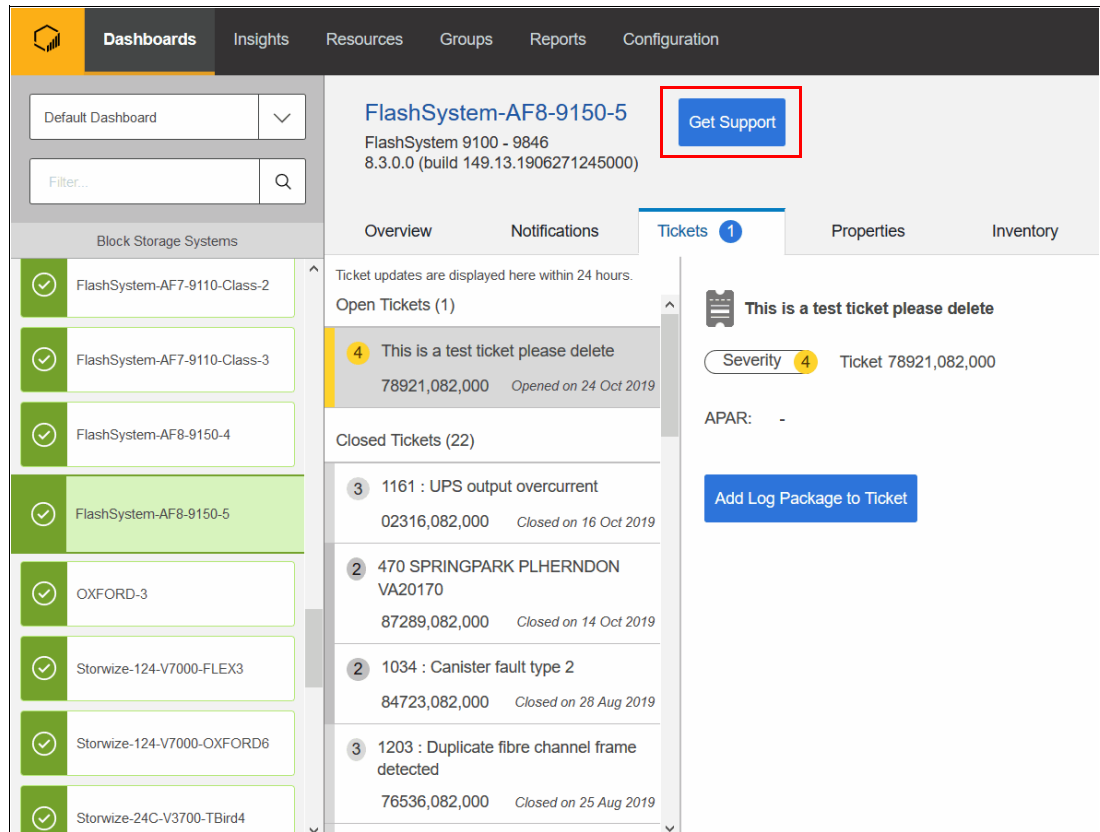


Figure 13-88 Get Support Button

This brings up a window where the user can opt to create a ticket, or update a ticket, as shown in Figure 13-89.

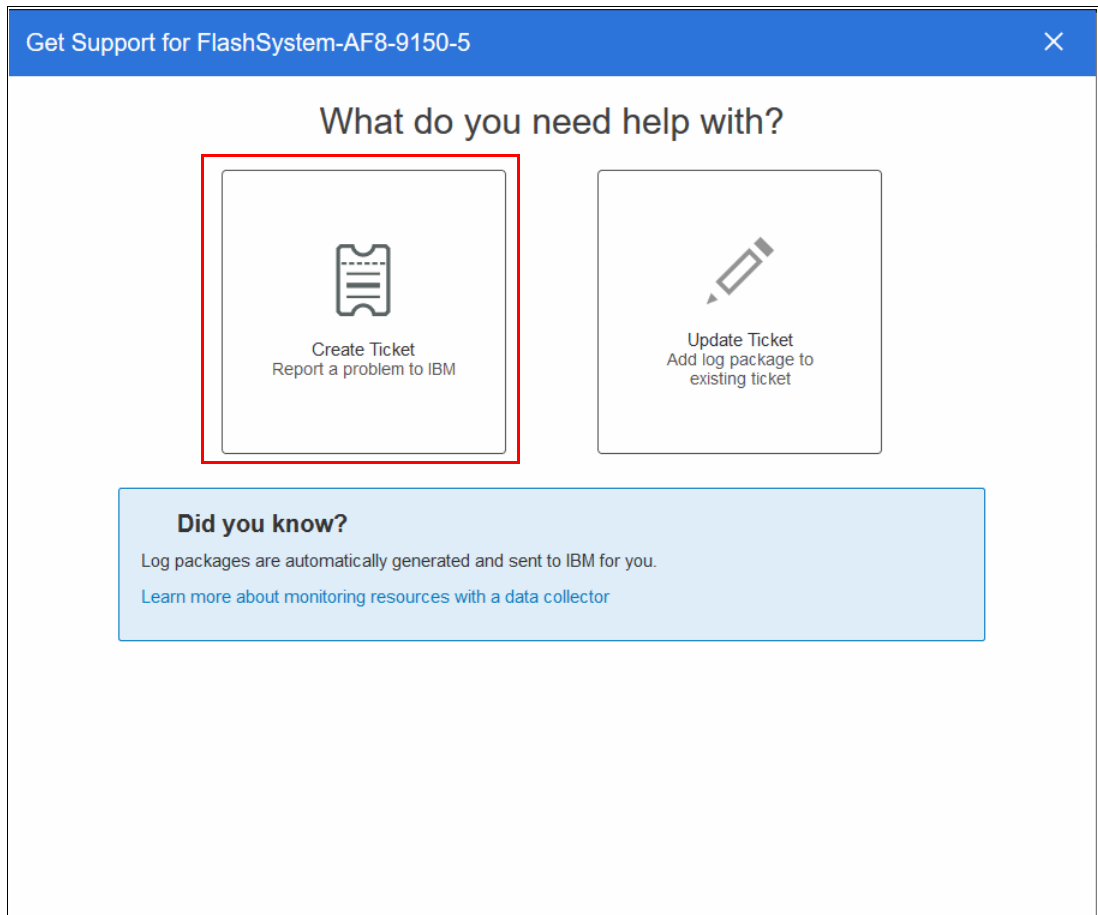


Figure 13-89 Get Support Window

2. Select **Create Ticket** and the ticket creation wizard opens, as shown in Figure 13-90.

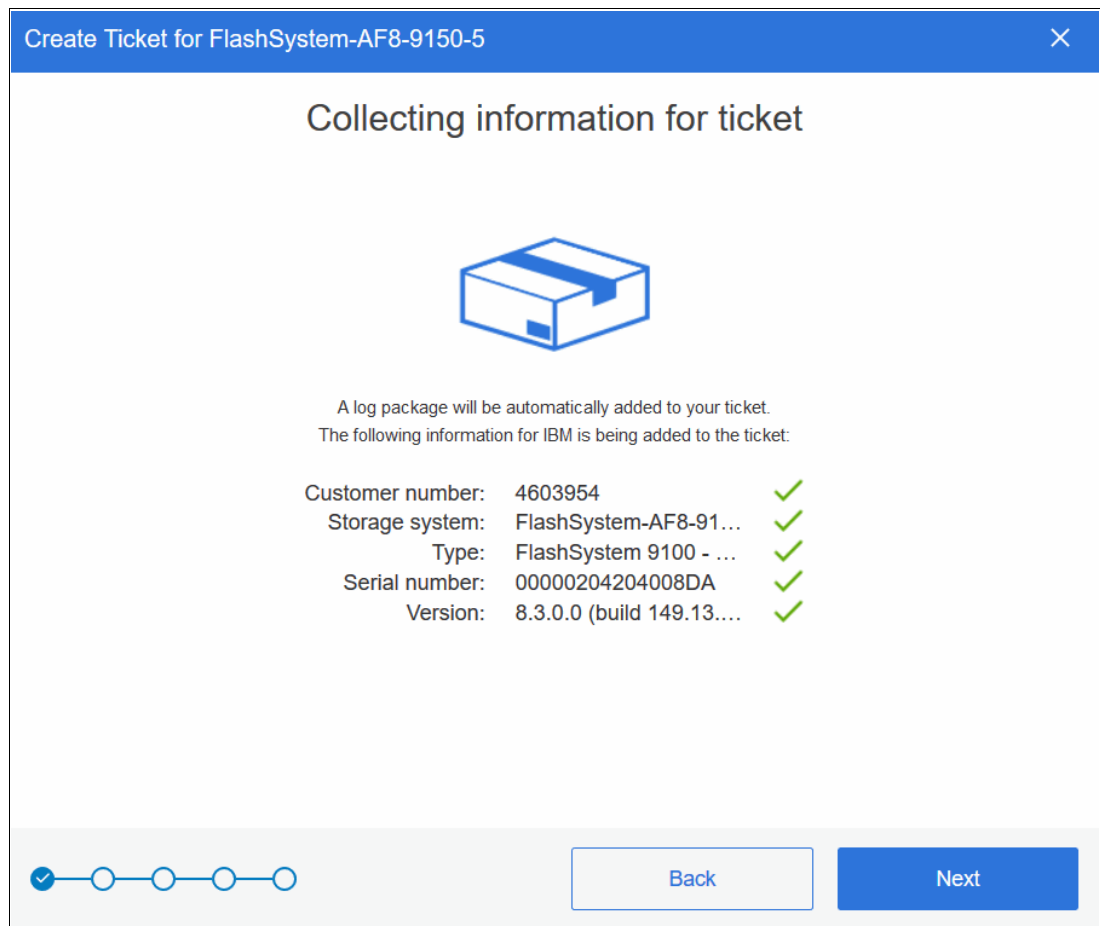


Figure 13-90 Create Ticket Wizard

3. Details of the system are automatically populated, including the customer number. Select **Next** (see Figure 13-91).

The screenshot shows a web form titled "Add a note or attachment" within a window labeled "Create Ticket for FlashSystem-AF8-9150-5". The form contains two text input fields. The first field contains the text "Ports are offline" and has a character count of 55. Below it is a hint: "Hint: Include what happened and the error code, if any." The second field contains the text "I have seen on storage insights that I have storage ports offline. Please assist." and has a character count of 2919. Below it is a hint: "Hint: Include the time the problem or error occurred, the affected resources, and details of any maintenance or other activities that occurred before the problem." There are two options for attaching files: a blue "Browse" button under the heading "Attach Image or File:" and a dashed blue box with an upward arrow icon and the text "Drag file here". At the bottom left is a progress indicator with five circles, the first two of which are checked. At the bottom right are "Back" and "Next" buttons.

Figure 13-91 Add note or attachment

4. The user can enter relevant details about your problem into the ticket, as shown in Figure 13-91 on page 827. It is also possible to attach images or files to the ticket, such as PuTTY logs and screen captures. After you are done, select **Next**.

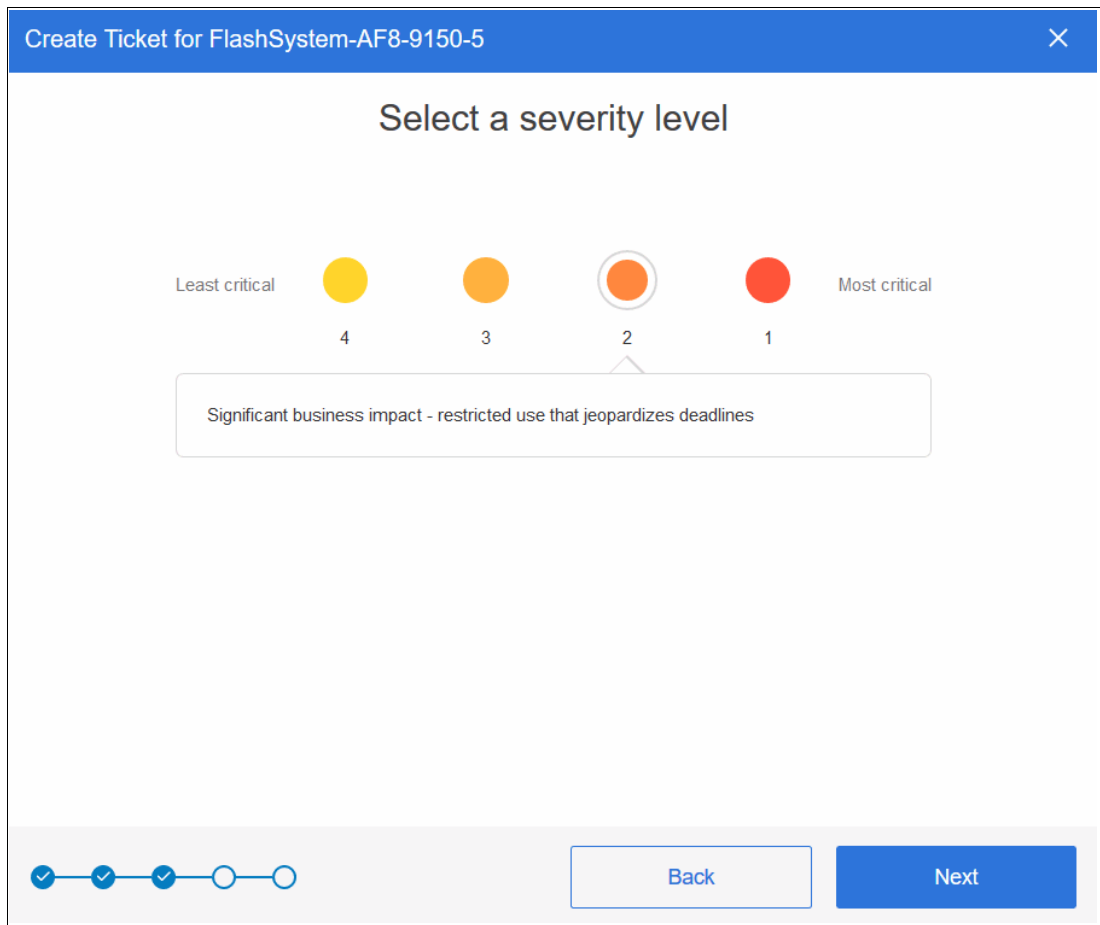


Figure 13-92 Select Severity Level

- In the third pane, you can select severity for the ticket; examples of what severity you should select are displayed, as shown in Figure 13-92 on page 828. Because in our example storage ports are offline with no effect, we select severity 2 because we have lost redundancy.

**Create Ticket for FlashSystem-AF8-9150-5** [X]

### Review the ticket

**Problem summary:** Ports are offline

**Description:** I have seen on storage insights that I have storage ports offline. ...

**Severity level:** 2 Significant business impact - restricted use that jeopardizes ...

**Type of problem:** Hardware

**Contact name:**

**Contact email:**

**Contact phone:**

**Customer number:**  

**Storage system:** FlashSystem-AF8-9150-5

**Type:** FlashSystem 9100 - 9846

**Serial number:** 00000204204008DA

**Version:** 8.3.0.0 (build 149.13.1906271245000)

**Log package:** 

[Progress: 5 steps, all checked]

Figure 13-93 Review Ticket

- The fourth pane allows the user to choose whether this is a hardware or a software problem. Select the relevant option (for this example, offline ports likely are caused by a physical layer hardware problem). After you are done, click **Next**.

The final page allows the user to review the details of the ticket that will be logged with IBM as shown in Figure 13-93. Contact details must be entered so that IBM Support can respond to the correct person. The user also must choose which type of logs should be attached to the ticket. For more information about the types of SNAP, see Table 13-10 on page 811.

7. After you are done, select **Create Ticket**. A confirmation window opens, as shown in Figure 13-94, and IBM Storage Insights automatically uploads the SNAP to the ticket when it is collected, which requires no other intervention from the user.

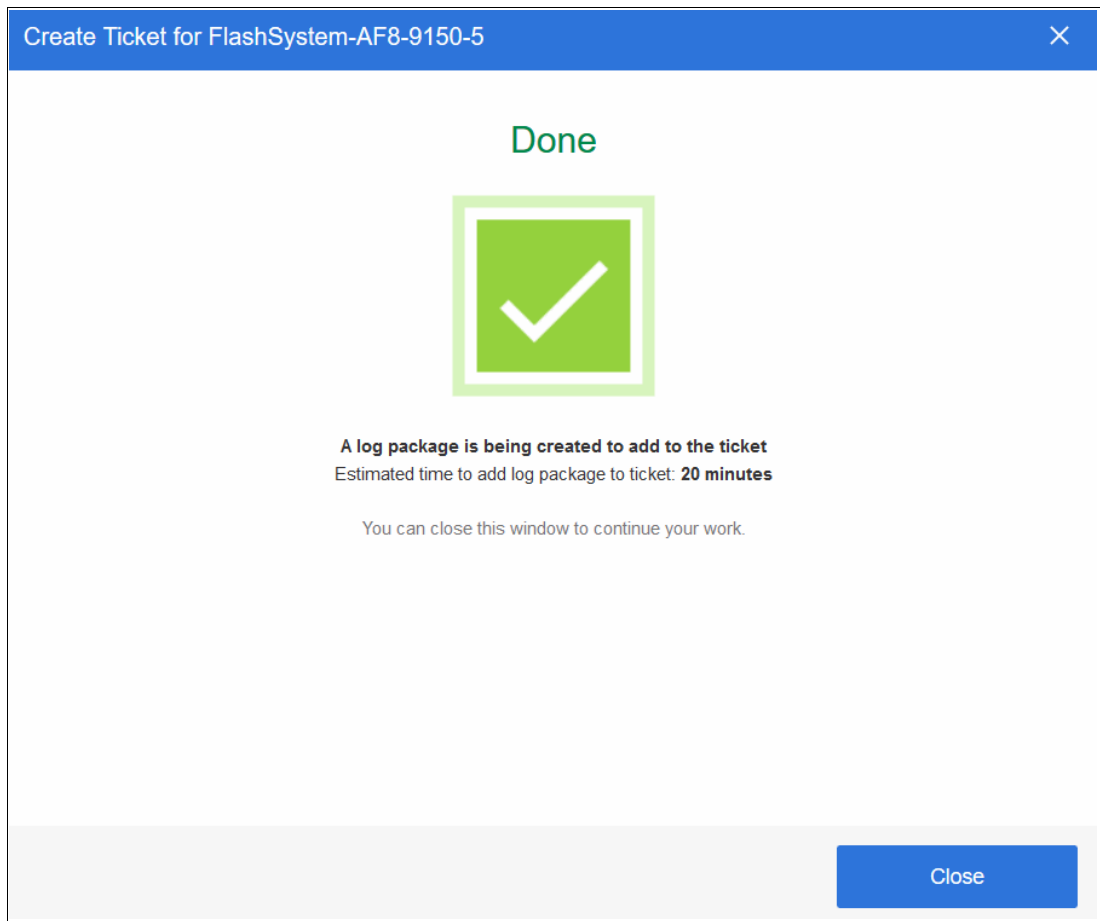


Figure 13-94 Ticket Creation confirmation



## 13.12.4 Managing support tickets by using IBM Storage Insights and uploading logs

IBM Storage Insights allows the user to track support tickets and upload logs to them. From the System Overview pane, select **Tickets**, as shown in Figure 13-95.

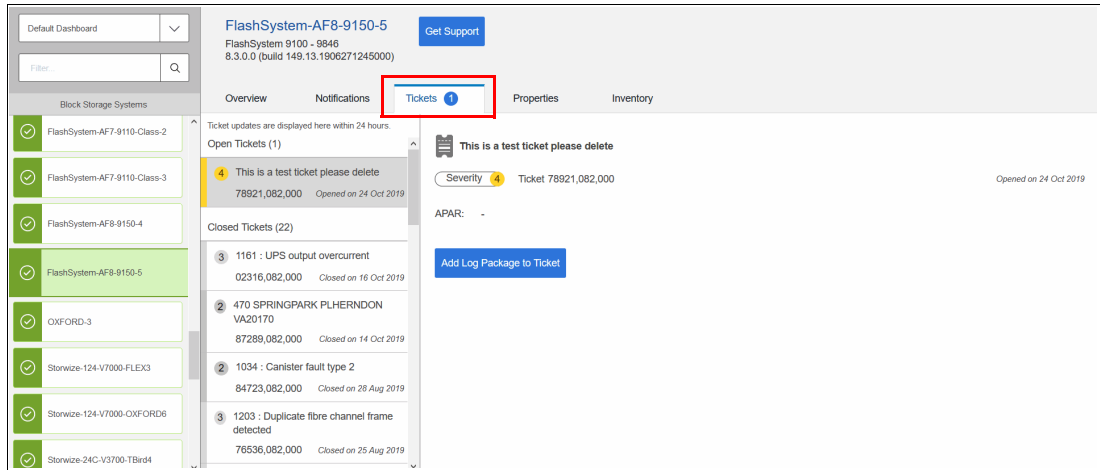


Figure 13-95 View Tickets

This pane allows the user to see a large history of support tickets that were logged through IBM Storage Insights for the system. Tickets that are not open are listed under Closed Tickets; open tickets are listed under Open Tickets.

The user selects **Add Log Package to Ticket** to quickly add logs to a ticket without the need to browse to the system GUI or IBM ECuRep. A pop-up window opens that guides the user through the process, as shown in Figure 13-96. The user can select which type of log package they want and add a note to the ticket with the logs.

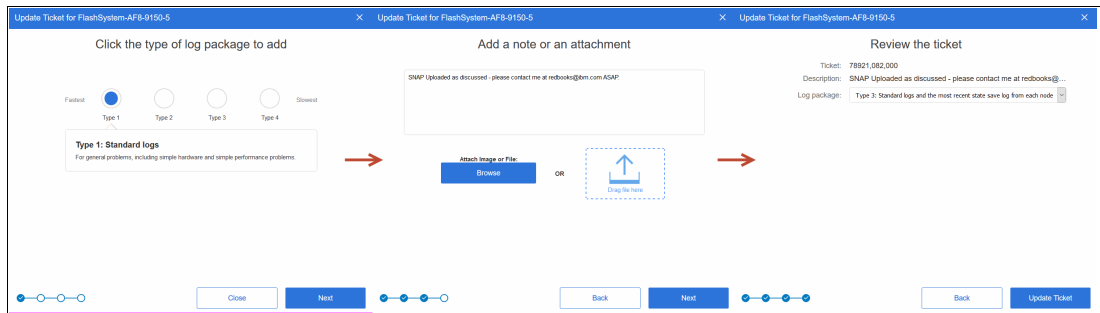


Figure 13-96 Adding a Log Package to the ticket

After clicking **Update Ticket**, a confirmation window opens, as shown in Figure 13-97, and the user can exit the wizard. IBM Storage Insights runs in the background to gather the logs and upload them to the ticket, which requires no other intervention from the user.

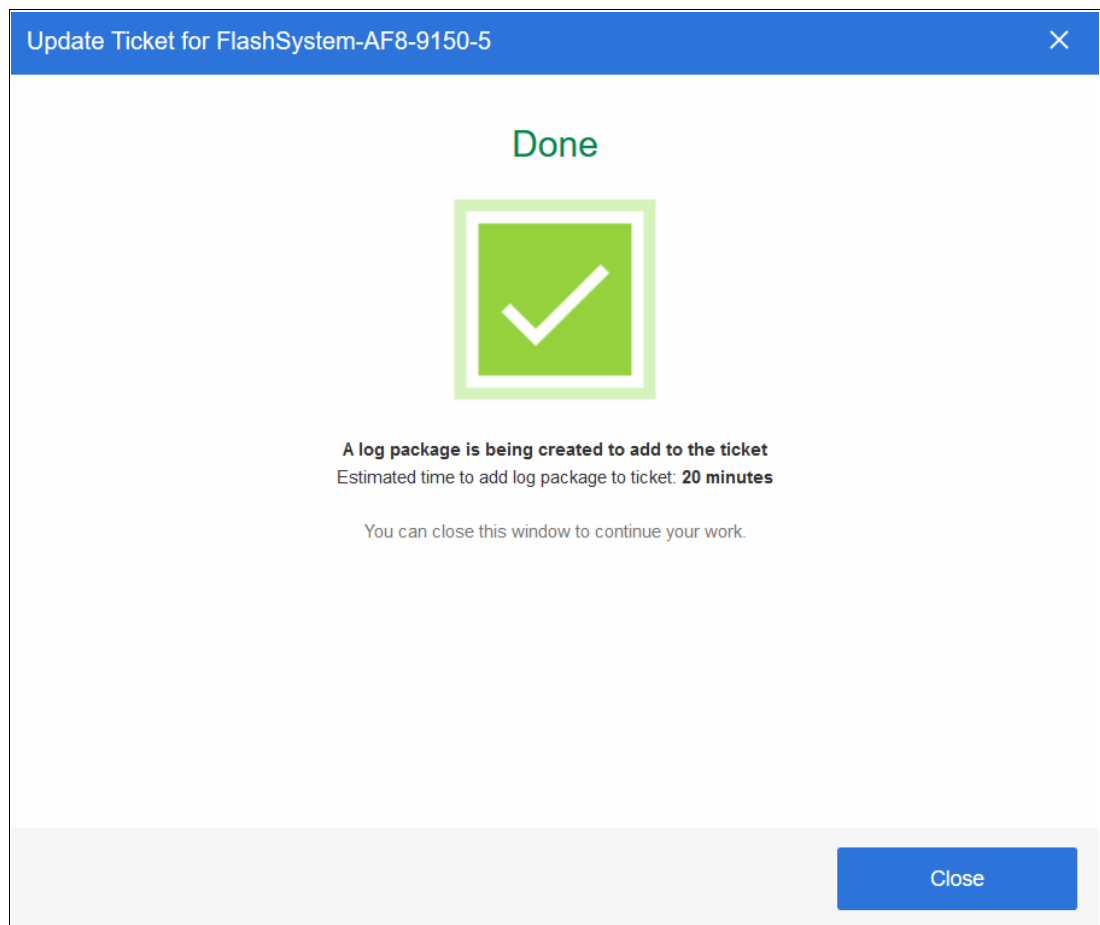


Figure 13-97 Confirmation of log upload



# Performance data and statistics gathering

This appendix provides a brief overview of the performance analysis capabilities of the IBM SAN Volume Controller and IBM Spectrum Virtualize V8.3. It also describes a method that you can use to collect and process IBM Spectrum Virtualize performance statistics.

It is beyond the intended scope of this book to provide an in-depth understanding of performance statistics or to explain how to interpret them. For more information about the performance of the SAN Volume Controller, see *IBM System Storage SAN Volume Controller, IBM Storwize V7000, and IBM FlashSystem 7200 Best Practices and Performance Guidelines*, SG24-7521.

This appendix includes the following topics:

- ▶ “IBM SAN Volume Controller performance overview” on page 834
- ▶ “Performance monitoring” on page 836

# IBM SAN Volume Controller performance overview

Storage virtualization with IBM Spectrum Virtualize provides many administrative benefits. In addition, it can provide an increase in performance for some workloads. The caching capability of the IBM Spectrum Virtualize and its ability to stripe volumes across multiple disk arrays can provide a performance improvement over what can otherwise be achieved when midrange disk subsystems are used.

To ensure that the performance levels of your system are maintained, monitor performance periodically to provide visibility into potential problems that exist or are developing so that they can be addressed in a timely manner.

## Performance considerations

When you are designing the IBM Spectrum Virtualize infrastructure or maintaining an infrastructure, you must consider many factors in terms of their potential effect on performance. These factors include dissimilar workloads that are competing for the same resources, overloaded resources, insufficient available resources, poor performing resources, and similar performance constraints.

Remember the following high-level rules when you are designing your storage area network (SAN) and IBM Spectrum Virtualize layout:

- ▶ Host-to-SAN Volume Controller inter-switch link (ISL) oversubscription.

This area is the most significant input/output (I/O) load across ISLs. The recommendation is to maintain a maximum of 7-to-1 oversubscription. A higher ratio is possible, but it tends to lead to I/O bottlenecks. This suggestion also assumes a core-edge design, where the hosts are on the edges and the IBM SAN Volume Controller is the core.

- ▶ Storage-to-SAN Volume Controller ISL oversubscription.

This area is the second most significant I/O load across ISLs. The maximum oversubscription is 7-to-1. A higher ratio is not supported. Again, this suggestion assumes a multiple-switch SAN fabric design.

- ▶ Node-to-node ISL oversubscription.

This area is the least significant load of the three possible oversubscription bottlenecks. In standard setups, this load can be ignored. Although this load is not entirely negligible, it does not contribute significantly to the ISL load. However, node-to-node ISL oversubscription is referenced here in relation to the split-cluster capability that was made available since Version 6.3 (Stretched Cluster and HyperSwap).

When the system is running in this manner, the number of ISL links becomes more important. As with the storage-to-SAN Volume Controller ISL oversubscription, this load also requires a maximum of 7-to-1 oversubscription. Exercise caution and careful planning when you determine the number of ISLs to implement. If you need assistance, contact your IBM representative and request technical assistance.

- ▶ ISL trunking and port channeling.

For the best performance and availability, use ISL trunking or port channeling. Independent ISL links can easily become overloaded and turn into performance bottlenecks. Bonded or trunked ISLs automatically share load and provide better redundancy in a failure.

- ▶ Number of paths per host multipath device.

The maximum supported number of paths per multipath device that is visible on the host is eight. Although the IBM Subsystem Device Driver Path Control Module (SDDPCM), related products, and most vendor multipathing software can support more paths, the IBM SAN Volume Controller expects a maximum of eight paths. In general, you see only an effect on performance from more paths than eight. Although the IBM Spectrum Virtualize system can work with more than eight paths, this design is technically unsupported.

- ▶ Do not intermix dissimilar array types or sizes.

Although the IBM Spectrum Virtualize supports an intermix of differing storage within storage pools, it is best to always use the same array model, which is Redundant Array of Independent Disks (RAID) mode, RAID size (RAID 5 6+P+S does not mix well with RAID 6 14+2), and drive speeds.

Rules and guidelines are no substitution for monitoring performance. Monitoring performance can provide a validation that design expectations are met, and identify opportunities for improvement.

## IBM Spectrum Virtualize performance perspectives

The IBM Spectrum Virtualize software was developed by the IBM Research Group. It is designed to run on IBM SAN Volume Controller, IBM Storwize products, and commodity hardware (mass-produced Intel-based processors with mass-produced expansion cards). It is also designed to provide distributed cache and a scalable cluster architecture.

Currently, the IBM SAN Volume Controller cluster is scalable up to eight nodes and these nodes can be swapped for newer hardware while online. This capability provides a great investment value because the nodes are relatively inexpensive and a node swap can be done online. This capability provides an instant performance boost with no license changes. The IBM SAN Volume Controller node model 2145-SV1, which has 64 GB of cache and can be upgraded to up to 256 GB per node, provides an extra benefit on top of the typical refresh cycle.

For more information about replacing nodes nondisruptively, see [IBM Knowledge Center](#).

For more information about setting up Fibre Channel (FC) port masking when upgrading from nodes 2145-CF8, 2145-CG8, or 2145-DH8 to 2145-SV1, see [this web page](#).

The performance is near linear when nodes are added into the cluster until performance eventually becomes limited by the attached components. Although virtualization provides significant flexibility in terms of the components that are used, it does not diminish the necessity of designing the system around the components so that it can deliver the level of performance that you want.

The key item for planning is your SAN layout. Switch vendors have slightly different planning requirements, but the goal is that you always want to maximize the bandwidth that is available to the IBM SAN Volume Controller ports. The IBM SAN Volume Controller is one of the few devices that can drive ports to their limits on average, so it is imperative that you put significant thought into planning the SAN layout.

Essentially, performance improvements are gained by spreading the workload across a greater number of back-end resources and by more caching. These capabilities are provided by the IBM SAN Volume Controller cluster. However, the performance of individual resources eventually becomes the limiting factor.

## Performance monitoring

This section highlights several performance monitoring techniques.

### Collecting performance statistics

IBM Spectrum Virtualize is constantly collecting performance statistics. The default frequency by which files are created is 15-minute intervals. The collection interval can be changed by running the **startstats** command.

The statistics files for volumes, managed disks (MDisks), nodes, and drives are saved at the end of the sampling interval. A maximum of 16 files (each) are stored before they are overlaid in a rotating log fashion. This design then provides statistics for the most recent 240-minute period if the default 15-minute sampling interval is used. IBM Spectrum Virtualize supports user-defined sampling intervals of 1 - 60 minutes.

For each type of object (volumes, MDisks, nodes, and drives), a separate file with statistic data is created at the end of each sampling period and stored in `/dumps/iostats`.

Run the **startstats** command to start the collection of statistics, as shown in Example A-1.

*Example A-1 The startstats command*

---

```
IBM_2145:Cluster_10.155.19.123:superuser>startstats -interval 2
```

---

This command starts statistics collection and gathers data at 2-minute intervals.

To verify the statistics status and collection interval, display the system properties, as shown in Example A-2.

*Example A-2 Statistics collection status and frequency*

---

```
IBM_2145:Cluster_10.155.19.123:superuser>lssystem
id 000002032140083A
name Cluster_10.155.19.123
statistics_status on
statistics_frequency 2
-- The output has been shortened for easier reading. --
```

---

Starting with V8.1, it is not possible to stop statistics collection by running the **stopstats** command.

**Collection intervals:** Although more frequent collection intervals provide a more detailed view of what happens within IBM Spectrum Virtualize and IBM SAN Volume Controller, they shorten the amount of time that the historical data is available on the IBM Spectrum Virtualize system. For example, rather than a 240-minute period of data with the default 15-minute interval, if you adjust to 2-minute intervals, you have a 32-minute period instead.

Statistics are collected per node. The sampling of the internal performance counters is coordinated across the cluster so that when a sample is taken, all nodes sample their internal counters at the same time. It is important to collect all files from all nodes for a complete analysis. Tools, such as IBM Spectrum Control, perform this intensive data collection for you.

## Statistics file naming

The statistics files that are generated are written to the `/dumps/iostats/` directory. The file name is in the following formats:

- ▶ `Nm_stats_<node_frontpanel_id>_<date>_<time>` for MDisks statistics
- ▶ `Nv_stats_<node_frontpanel_id>_<date>_<time>` for Volumes statistics
- ▶ `Nn_stats_<node_frontpanel_id>_<date>_<time>` for node statistics
- ▶ `Nd_stats_<node_frontpanel_id>_<date>_<time>` for drives statistics

The `node_frontpanel_id` is the pane name of the node on which the statistics were collected. The date is in the form `<yymmdd>` and the time is in the form `<hhmmss>`. The following example shows an MDisk statistics file name:

```
Nm_stats_Cluster_75ACXF0_180113_121848
```

Example A-3 shows typical MDisk, volume, node, and disk drive statistics file names.

### Example A-3 File names of per node statistics

---

```
IBM_2145:ITS0-SV1:superuser>lsdumps -prefix /dumps/iostats
id filename
0  Nv_stats_75ACXF0_191027_121848
1  Nd_stats_75ACXF0_191027_121848
2  Nm_stats_75ACXF0_191027_121848
3  Nn_stats_75ACXF0_191027_121848
4  Nm_stats_75ACXF0_191027_121948
5  Nd_stats_75ACXF0_191027_121948
6  Nn_stats_75ACXF0_191027_121948
7  Nv_stats_75ACXF0_191027_121948
...
60 Nn_stats_75ACXF0_191027_123348
61 Nm_stats_75ACXF0_191027_123348
62 Nv_stats_75ACXF0_191027_123348
63 Nd_stats_75ACXF0_191027_123348
```

---

**Tip:** The performance statistics files can be copied from the IBM SAN Volume Controller nodes to a local drive on your workstation by using `pscp.exe` (included with PuTTY) from an MS-DOS command prompt, as shown in the following example:

```
C:\Program Files\PuTTY>pscp -unsafe -load 75ACXF0
superuser@192.168.100.100:/dumps/iostats/* c:\statsfiles
```

Use the `-load` parameter to specify the session that is defined in PuTTY.

Specify the `-unsafe` parameter when you use wildcards.

You can download PuTTY from [this web page](#).

## Real-time performance monitoring

IBM SAN Volume Controller supports real-time performance monitoring. Real-time performance statistics provide short-term status information for the IBM SAN Volume Controller. The statistics are shown as graphs in the management GUI, or can be viewed from the CLI.

With system-level statistics, you can quickly view the CPU usage and the bandwidth of volumes, interfaces, and MDisks. Each graph displays the current bandwidth in megabytes per second (MBps) or input/output operations per second (IOPS), and a view of bandwidth over time.

Each node collects various performance statistics, mostly at 5-second intervals, and the statistics that are available from the config node in a clustered environment. This information can help you determine the performance effect of a specific node. As with system statistics, node statistics help you to evaluate whether the node is operating within normal performance metrics.

Real-time performance monitoring gathers the following system-level performance statistics:

- ▶ CPU usage
- ▶ Port usage and I/O rates
- ▶ Volume and MDisk I/O rates
- ▶ Bandwidth
- ▶ Latency

Real-time statistics are not a configurable option and cannot be disabled.

### Real-time performance monitoring with the CLI

The `lsnodestats` and `lssystemstats` commands are available for monitoring the statistics through the CLI.

The `lsnodestats` command provides performance statistics for the nodes that are part of a clustered system, as shown in Example A-4. This output is truncated and shows only part of the available statistics. You can also specify a node name in the command to limit the output for a specific node.

*Example A-4 The `lsnodestats` command output*

---

```
IBM_2145:Cluster_10.155.19.123:superuser>lsnodestats
node_id node_name      stat_name          stat_current  stat_peak  stat_peak_time
1       node1        compression_cpu_pc 0              0          191027230859
1       node1        cpu_pc             1              1          191027230859
1       node1        fc_mb              0              0          191027230859
1       node1        fc_io              276            305        191027230414
1       node1        sas_mb             0              0          191027230859
1       node1        sas_io            0              0          191027230859
1       node1        iscsi_mb          0              0          191027230859
1       node1        iscsi_io          0              0          191027230859
1       node1        write_cache_pc    0              0          191027230859
1       node1        total_cache_pc    0              0          191027230859
1       node1        vdisk_mb          0              0          191027230859
1       node1        vdisk_io          0              0          191027230859
1       node1        vdisk_ms          0              0          191027230859
1       node1        mdisk_mb          0              0          191027230859
1       node1        mdisk_io          0              0          191027230859
1       node1        mdisk_ms          0              0          191027230859
```



1	node1	drive_mb	0	0	191027230859
1	node1	drive_io	0	0	191027230859
...					
2	node_75ACXF0	compression_cpu_pc	0	0	191027230902
2	node_75ACXF0	cpu_pc	1	1	191027230902
2	node_75ACXF0	fc_mb	0	0	191027230902
2	node_75ACXF0	fc_io	285	294	191027230547
2	node_75ACXF0	sas_mb	0	0	191027230902
2	node_75ACXF0	sas_io	0	0	191027230902
2	node_75ACXF0	iscsi_mb	0	0	191027230902
2	node_75ACXF0	iscsi_io	0	0	191027230902
2	node_75ACXF0	write_cache_pc	0	0	191027230902
2	node_75ACXF0	total_cache_pc	0	0	191027230902
2	node_75ACXF0	vdisk_mb	0	0	191027230902
2	node_75ACXF0	vdisk_io	0	0	191027230902

Example A-4 on page 838 shows statistics for the two nodes members of Cluster\_10.155.19.123. For each node, the following columns are displayed:

- ▶ **stat\_name**: The name of the statistic field
- ▶ **stat\_current**: The current value of the statistic field
- ▶ **stat\_peak**: The peak value of the statistic field in the last 5 minutes
- ▶ **stat\_peak\_time**: The time that the peak occurred

The **l1nodestats** command can also be used with a node name or ID as an argument. For example, you run the **l1nodestats node1** command to display the statistics of node with name node1 only.

The **l1systemstats** command lists the same set of statistics that is listed by running the **l1nodestats** command, but representing all nodes in the cluster. The values for these statistics are calculated from the node statistics values in the following way:

- ▶ **Bandwidth**: Sum of bandwidth of all nodes.
- ▶ **Latency**: Average latency for the cluster, which is calculated by using data from the whole cluster, not an average of the single node values.
- ▶ **IOPS**: Total IOPS of all nodes.
- ▶ **CPU percentage**: Average CPU percentage of all nodes.

Example A-5 shows the resulting output of the **l1systemstats** command.

*Example A-5 The l1systemstats command output*

```

IIBM_2145:Cluster_10.155.19.123:superuser>l1systemstats
stat_name      stat_current  stat_peak  stat_peak_time
compression_cpu_pc  0           0          191027231453
cpu_pc         1           1          191027231453
fc_mb          0           0          191027231453
fc_io          595         595        191027231453
sas_mb         0           0          191027231453
sas_io         0           0          191027231453
iscsi_mb       0           0          191027231453
iscsi_io       0           0          191027231453
write_cache_pc 0           0          191027231453
total_cache_pc 0           0          191027231453
vdisk_mb       0           0          191027231453
vdisk_io       0           0          191027231453

```

vdisk_ms	0	0	191027231453
mdisk_mb	0	0	191027231453
mdisk_io	0	0	191027231453
mdisk_ms	0	0	191027231453
drive_mb	0	0	191027231453
drive_io	0	0	191027231453
drive_ms	0	0	191027231453
vdisk_r_mb	0	0	191027231453
vdisk_r_io	0	0	191027231453
vdisk_r_ms	0	0	191027231453
vdisk_w_mb	0	0	191027231453
vdisk_w_io	0	0	191027231453
vdisk_w_ms	0	0	191027231453
mdisk_r_mb	0	0	191027231453
mdisk_r_io	0	0	191027231453
mdisk_r_ms	0	0	191027231453
mdisk_w_mb	0	0	191027231453
mdisk_w_io	0	0	191027231453
mdisk_w_ms	0	0	191027231453
drive_r_mb	0	0	191027231453
drive_r_io	0	0	191027231453
drive_r_ms	0	0	191027231453
drive_w_mb	0	0	191027231453
drive_w_io	0	0	191027231453
drive_w_ms	0	0	191027231453
iplink_mb	0	0	191027231453
iplink_io	0	0	191027231453
iplink_comp_mb	0	0	191027231453
cloud_up_mb	0	0	191027231453
cloud_up_ms	0	0	191027231453
cloud_down_mb	0	0	191027231453
cloud_down_ms	0	0	191027231453
iser_mb	0	0	191027231453
iser_io	0	0	191027231453

Table A-1 lists the different counters that are presented by the **lssystemstats** and **lsnodestats** commands.

*Table A-1 Counters in lssystemstats and lsnodestats*

<b>Value</b>	<b>Description</b>
compression_cpu_pc	Displays the percentage of allocated CPU capacity that is used for compression.
cpu_pc	Displays the percentage of allocated CPU capacity that is used for the system.
fc_mb	Displays the total number of MBps for FC traffic on the system. This value includes host I/O and any bandwidth that is used for communication within the system.
fc_io	Displays the total IOPS for FC traffic on the system. This value includes host I/O and any bandwidth that is used for communication within the system.

<b>Value</b>	<b>Description</b>
sas_mb	Displays the total number of MBps for serial-attached Small Computer System Interface (SCSI) (SAS) traffic on the system. This value includes host I/O and bandwidth that is used for background RAID activity.
sas_io	Displays the total IOPS for SAS traffic on the system. This value includes host I/O and bandwidth that is used for background RAID activity.
iscsi_mb	Displays the total number of MBps for internet Small Computer Systems Interface (iSCSI) traffic on the system.
iscsi_io	Displays the total IOPS for iSCSI traffic on the system.
write_cache_pc	Displays the percentage of the write cache usage for the node.
total_cache_pc	Displays the total percentage for both the write and read cache usage for the node.
vdisk_mb	Displays the average number of MBps for read and write operations to volumes during the sample period.
vdisk_io	Displays the average number of IOPS for read and write operations to volumes during the sample period.
vdisk_ms	Displays the average amount of time in milliseconds that the system takes to respond to read and write requests to volumes over the sample period.
mdisk_mb	Displays the average number of MBps for read and write operations to MDisks during the sample period.
mdisk_io	Displays the average number of IOPS for read and write operations to MDisks during the sample period.
mdisk_ms	Displays the average amount of time in milliseconds that the system takes to respond to read and write requests to MDisks over the sample period.
drive_mb	Displays the average number of MBps for read and write operations to drives during the sample period.
drive_io	Displays the average number of IOPS for read and write operations to drives during the sample period.
drive_ms	Displays the average amount of time in milliseconds that the system takes to respond to read and write requests to drives over the sample period.
vdisk_w_mb	Displays the average number of MBps for read and write operations to volumes during the sample period.
vdisk_w_io	Displays the average number of IOPS for write operations to volumes during the sample period.
vdisk_w_ms	Displays the average amount of time in milliseconds that the system takes to respond to write requests to volumes over the sample period.
mdisk_w_mb	Displays the average number of MBps for write operations to MDisks during the sample period.
mdisk_w_io	Displays the average number of IOPS for write operations to MDisks during the sample period.

<b>Value</b>	<b>Description</b>
mdisk_w_ms	Displays the average amount of time in milliseconds that the system takes to respond to write requests to MDisks over the sample period.
drive_w_mb	Displays the average number of MBps for write operations to drives during the sample period.
drive_w_io	Displays the average number of IOPS for write operations to drives during the sample period.
drive_w_ms	Displays the average amount of time in milliseconds that the system takes to respond write requests to drives over the sample period.
vdisk_r_mb	Displays the average number of MBps for read operations to volumes during the sample period.
vdisk_r_io	Displays the average number of IOPS for read operations to volumes during the sample period.
vdisk_r_ms	Displays the average amount of time in milliseconds that the system takes to respond to read requests to volumes over the sample period.
mdisk_r_mb	Displays the average number of MBps for read operations to MDisks during the sample period.
mdisk_r_io	Displays the average number of IOPS for read operations to MDisks during the sample period.
mdisk_r_ms	Displays the average amount of time in milliseconds that the system takes to respond to read requests to MDisks over the sample period.
drive_r_mb	Displays the average number of MBps for read operations to drives during the sample period.
drive_r_io	Displays the average number of IOPS for read operations to drives during the sample period.
drive_r_ms	Displays the average amount of time in milliseconds that the system takes to respond to read requests to drives over the sample period.
iplink_mb	The total number of MBps for IP replication traffic on the system. This value does not include iSCSI host I/O operations.
iplink_comp_mb	Displays the average number of compressed MBps over the IP replication link during the sample period.
iplink_io	The total IOPS for IP partnership traffic on the system. This value does not include iSCSI host I/O operations.
cloud_up_mb	Displays the average number of Mbps for upload operations to a cloud account during the sample period.
cloud_up_ms	Displays the average amount of time (in milliseconds) it takes for the system to respond to upload requests to a cloud account during the sample period.
cloud_down_mb	Displays the average number of Mbps for download operations to a cloud account during the sample period.
cloud_down_ms	Displays the average amount of time (in milliseconds) that it takes for the system to respond to download requests to a cloud account during the sample period.

Value	Description
iser_mb	Displays the total number of MBps for iSCSI Extensions for RDMA (iSER) traffic on the system.
iser_io	Displays the total IOPS for iSER traffic on the system.

## Real-time performance statistics monitoring by using the GUI

The IBM Spectrum Virtualize dashboard gives you performance at a glance by displaying some information about the system. You can see the entire cluster (the system) performance by selecting the information between Bandwidth, Response Time, IOPS, or CPU usage. You can also display a Node Comparison by selecting the same information as for the cluster, and then switching the button, as shown in Figure A-1 and Figure A-2.

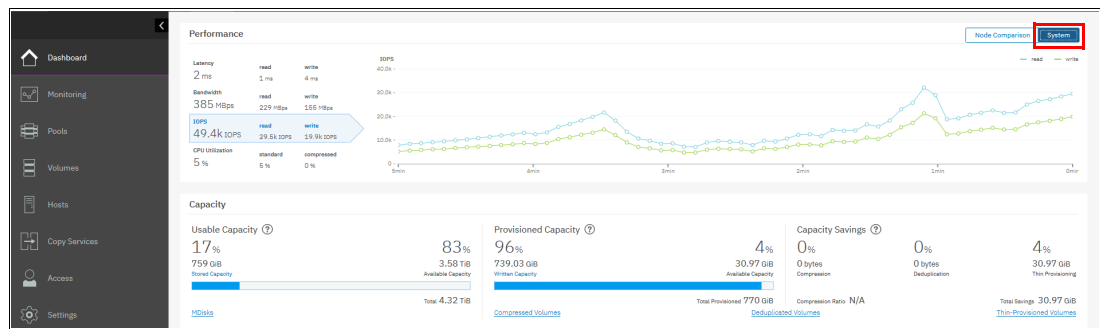


Figure A-1 IBM Spectrum Virtualize Dashboard displaying System performance overview

Figure A-2 shows the display after switching the button.

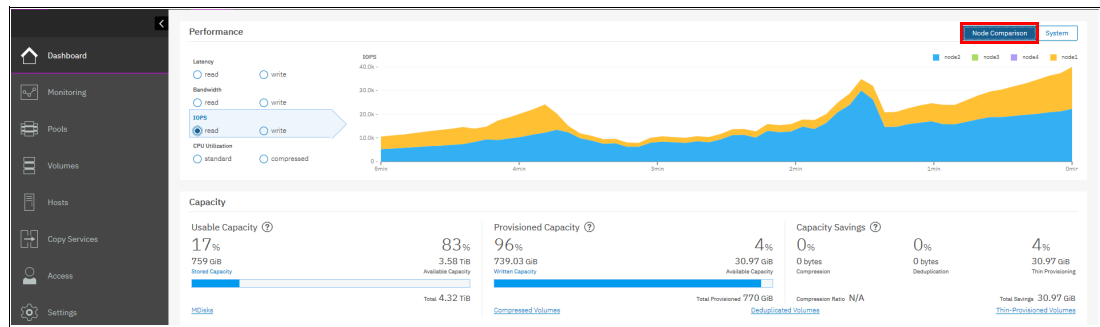


Figure A-2 IBM Spectrum Virtualize Dashboard displaying Nodes performance overview

You can also use real-time statistics to monitor CPU usage, volume, interface, and MDisk bandwidth of your system and nodes. Each graph represents 5 minutes of collected statistics and provides a means of assessing the overall performance of your system.

The real-time statistics are available from the IBM Spectrum Virtualize GUI. Click **Monitoring** → **Performance** (as shown in Figure A-3) to open the **Performance Monitoring** window.

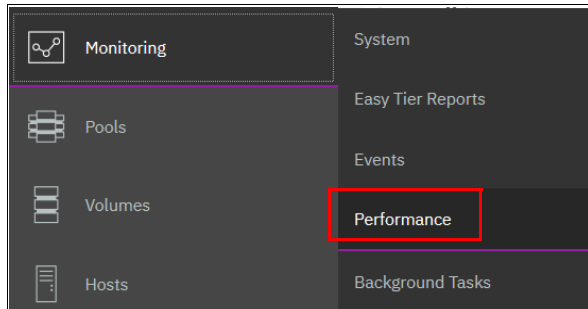


Figure A-3 Selecting Performance in the Monitoring menu

As shown in Figure A-4 on page 845, the Performance monitoring pane is divided into the following sections that provide usage views for the following resources:

- ▶ CPU Utilization: The CPU usage graph shows the current percentage of CPU usage and peaks in utilization. It can also display compression CPU usage for systems with compressed volumes.
- ▶ Volumes: Shows four metrics for the overall volume utilization graphics:
  - Read
  - Write
  - Read latency
  - Write latency
- ▶ Interfaces: The Interfaces graph displays data points for FC, iSCSI, SAS, and IP Remote Copy (RC) interfaces. You can use this information to help determine connectivity issues that might affect performance:
  - FC
  - iSCSI
  - SAS
  - IP Remote Copy
- ▶ MDisks: Shows four metrics for the overall MDisks graphics:
  - Read
  - Write
  - Read latency
  - Write latency

You can use these metrics to help determine the overall performance health of the volumes and MDisks on your system. Consistent unexpected results can indicate errors in configuration, system faults, or connectivity issues.

The system's performance also is always visible in the bottom of the IBM Spectrum Virtualize window, as shown in Figure A-4 on page 845.

**Note:** The indicated values in the graphics are averaged on a 1-second-based sample.



Figure A-4 IBM Spectrum Virtualize Performance window

You can also view performance statistics for each of the available nodes of the system, as shown in Figure A-5.

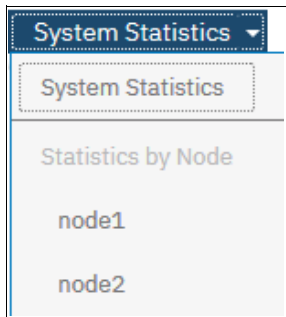


Figure A-5 View statistics per node or for the entire system

You can also change the metric between MBps or IOPS, as shown in Figure A-6.

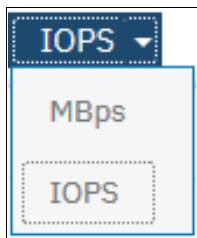


Figure A-6 View performance metrics by MBps or IOPS

On any of these views, you can select any point with your cursor to know the value and when it occurred. When you place your cursor over the timeline, it becomes a dotted line with the various values gathered, as shown in Figure A-7.

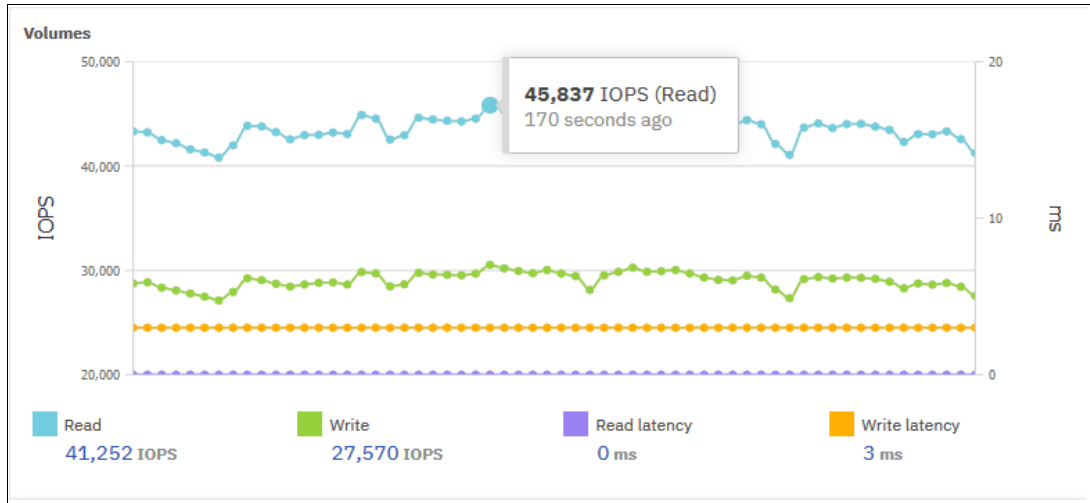


Figure A-7 Viewing performance with details

For each of the resources, various metrics are available and you can select which are shown. For example, as shown in Figure A-8, from the four available metrics for the MDisks view (Read, Write, Read latency, and Write latency), only Write and Write latency IOPS are selected.

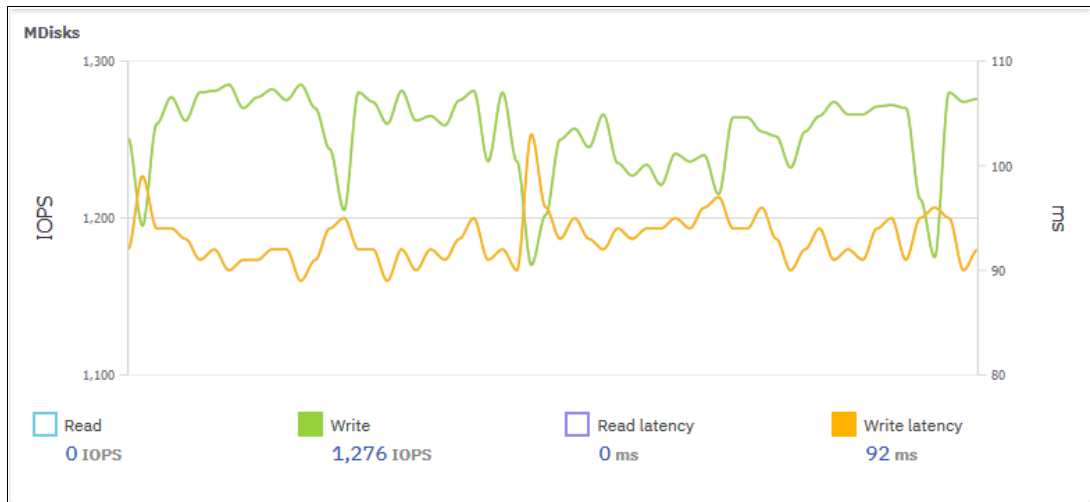


Figure A-8 Displaying performance counters

## Performance data collection and IBM Spectrum Control

Although you can obtain performance statistics in standard .xml files, the use of .xml files is a less practical and more complicated method to analyze the IBM Spectrum Virtualize performance statistics. IBM Spectrum Control is the supported IBM tool to collect and analyze SAN Volume Controller performance statistics.

IBM Spectrum Control is installed separately on a dedicated system, and is not part of the IBM Spectrum Virtualize bundle.



For more information about the use of IBM Spectrum Control to monitor your storage subsystem, see [this web page](#).

As an alternative to IBM Spectrum Control, a cloud-based tool is available that is called IBM Storage Insights, which provides a single dashboard that gives you a clear view of all your IBM block storage, showing performance and capacity information. You do not have to install this tool in your environment because it is a cloud-based solution. Only an agent is required to collect the data of the storage devices.





## CLI setup

This appendix describes the access configuration to the command-line interface (CLI) by using the local Secure Shell (SSH) authentication method.

## Setting up the CLI

The IBM Spectrum Virtualize system features a powerful CLI, which offers a few more options and flexibility as compared to the GUI. This section describes how to configure a management system by using the SSH protocol to connect to the IBM Spectrum Virtualize system for running commands by using the CLI.

For more information about the CLI, see [IBM Knowledge Center](#).

**Note:** If a task completes in the GUI, the associated CLI command is always displayed in the details, as shown throughout this book.

In the IBM Spectrum Virtualize GUI, authentication is performed by entering a user name and password. CLI uses SSH to connect from a host to the IBM Spectrum Virtualize system. A private and a public key pair or user name and password is necessary.

The use of SSH keys with a passphrase is more secure than a login with a user name and password because authenticating to a system requires the private key and the passphrase, while in the other method, only the password is required to obtain access to the system.

When SSH keys are used without a passphrase, it becomes easier to log in to a system because you provide only the private key when performing the login and you are not prompted for a password. This option is less secure than the use of SSH keys with a passphrase.

To enable CLI access with SSH keys, the following process is used:

1. A public key and a private key are generated together as a pair.
2. A public key is uploaded to the IBM Spectrum Virtualize system through the GUI.
3. A client SSH tool is configured to authenticate with the private key.
4. A secure connection is established between the client and the IBM SAN Volume Controller system.

SSH is the communication vehicle between the management workstation and the IBM Spectrum Virtualize system. The SSH client provides a secure environment from which to connect to a remote machine. It uses the principles of public and private keys for authentication.

SSH keys are generated by the SSH client software. The SSH keys include a public key, which is uploaded and maintained by the IBM SAN Volume Controller clustered system, and a private key, which is kept private on the workstation that is running the SSH client. These keys authorize specific users to access the administration and service functions on the system.

Each key pair is associated with a user-defined ID string that can consist of up to 256 characters. Up to 100 keys can be stored on the system.

New IDs and keys can be added, and unwanted IDs and keys can be deleted. To use the CLI, an SSH client must be installed on that system. To use the CLI with SSH keys, the SSH client is required, but also an SSH key pair must be generated on the client system, and the client's SSH public key must be stored on the IBM Spectrum Virtualize systems.

## Basic setup on a Windows host

The SSH client on the Windows host that is used in this book is PuTTY. A PuTTY key generator can also be used to generate the private and public key pair. The PuTTY client can be downloaded from [this web page](#).

Download the following tools:

- ▶ PuTTY SSH client: `putty.exe`
- ▶ PuTTY key generator: `puttygen.exe`

### Generating a public and private key pair

To generate a public and private key pair, complete the following steps:

1. Start the PuTTY key generator to generate the public and private key pair, as shown in Figure B-1.

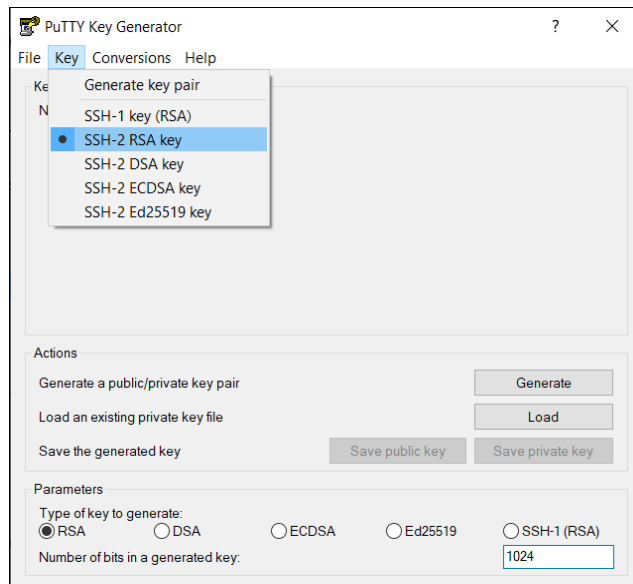


Figure B-1 PuTTY key generator

Select the following options:

- **SSH2 RSA**
- Number of bits in a generated key: **1024**

**Note:** Larger SSH keys, such as 2048 bits, are also supported.

2. Click **Generate** and move the cursor over the blank area to generate keys (see Figure B-2).

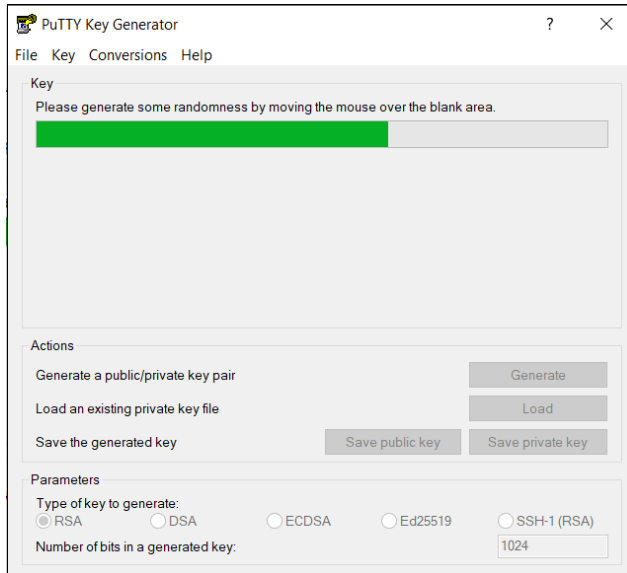


Figure B-2 Generating keys

**To generate keys:** The blank area that is indicated by the message is the large blank rectangle on the GUI inside the section of the GUI labeled **Key**. Continue to move the mouse pointer over the blank area until the progress bar reaches the far right. This action generates random characters to create a unique key pair.

3. After the keys are generated, save them for later use. Click **Save public key**, as shown in Figure B-3.

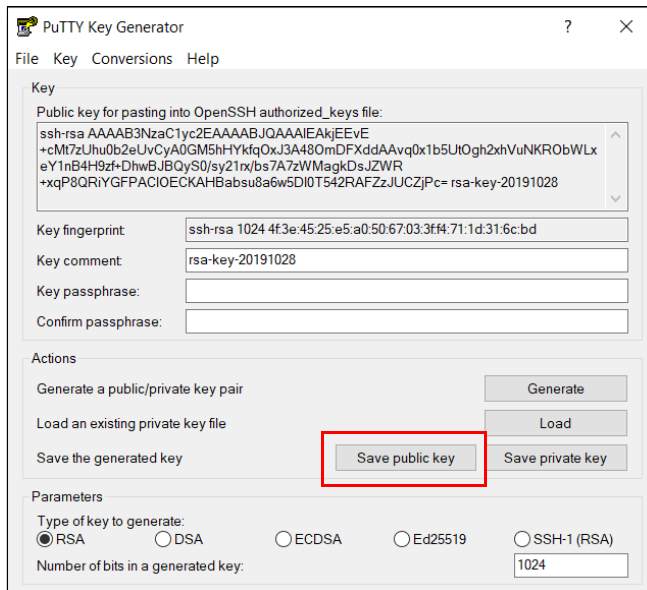


Figure B-3 Saving the public key

4. You are prompted for a name (for example, `sshkey.pub`) and a location for the public key (for example, `C:\Keys\`). Click **Save**.

Ensure that you record the name and location because the name and location of this SSH public key must be specified later.

**Public key extension:** By default, the PuTTY key generator saves the public key with no extension. Use the string `pub` for naming the public key. For example, add the extension `.pub` to the name of the file to easily differentiate the SSH public key from the SSH private key.

5. Click **Save private key**. You are prompted with a warning message (see Figure B-4). Click **Yes** to save the private key without a passphrase.

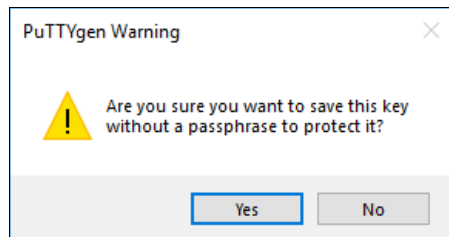


Figure B-4 Confirming the security warning

**Note:** It is possible to use a passphrase for an SSH key. This action increases security, but generates an extra step to log in with the SSH key because it requires the passphrase input.

6. When prompted, enter a name (for example, `sshkey.ppk`), select a secure place as the location, and click **Save**.

**Key generator:** The PuTTY key generator saves the PuTTY private key (PPK) with the `.ppk` extension.

7. Close the PuTTY key generator.

## Uploading the SSH public key to the IBM SAN Volume Controller

After you create your SSH key pair, upload your SSH public key onto the IBM SAN Volume Controller. Complete the following steps:

1. Open the user section in the GUI, as shown in Figure B-5.

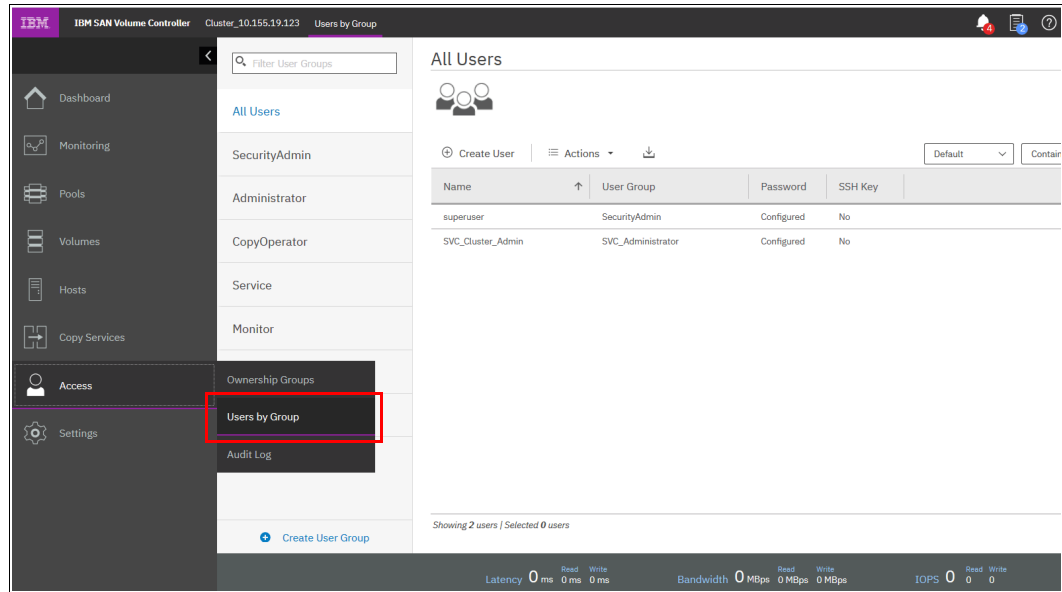


Figure B-5 Open user section

2. Right-click the user name for which you want to upload the key and click **Properties** (see Figure B-6).

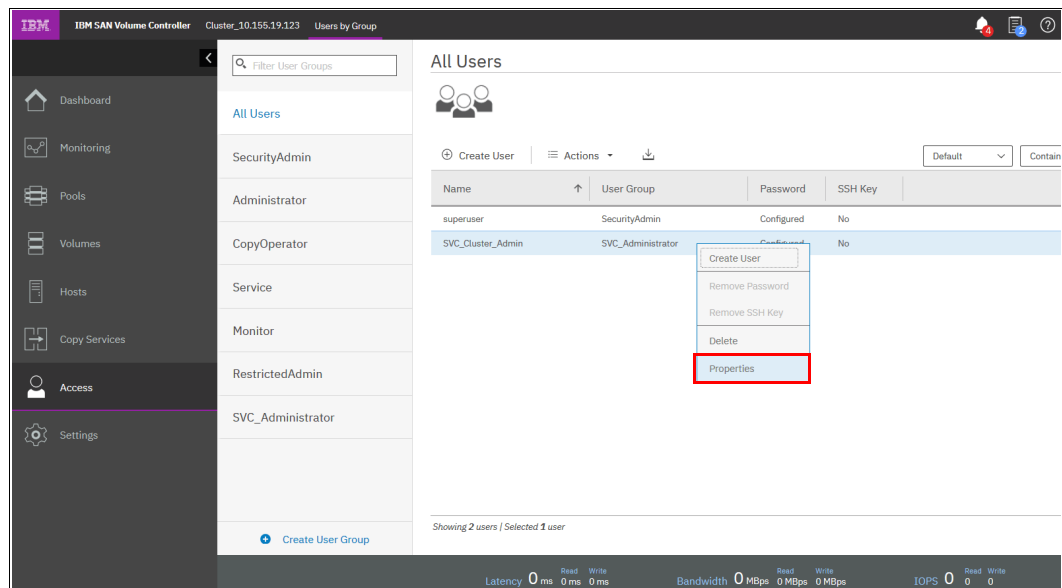


Figure B-6 User properties



- To upload the public key, click **Browse**, open the folder where you stored the public SSH key, and select the key (see Figure B-7).

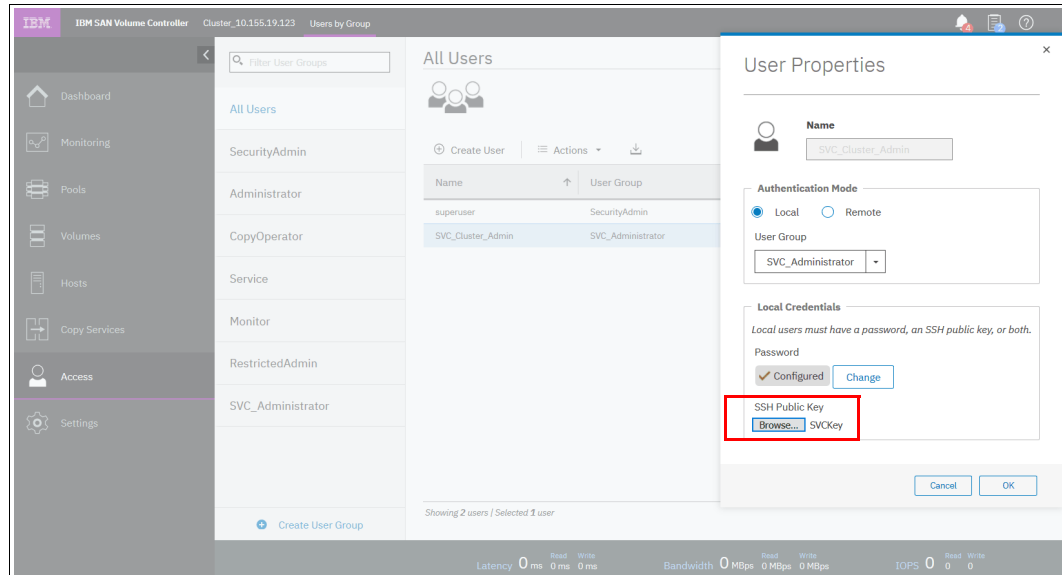


Figure B-7 Selecting the public key

- Click **OK** and the key is uploaded (see Figure 13-98).

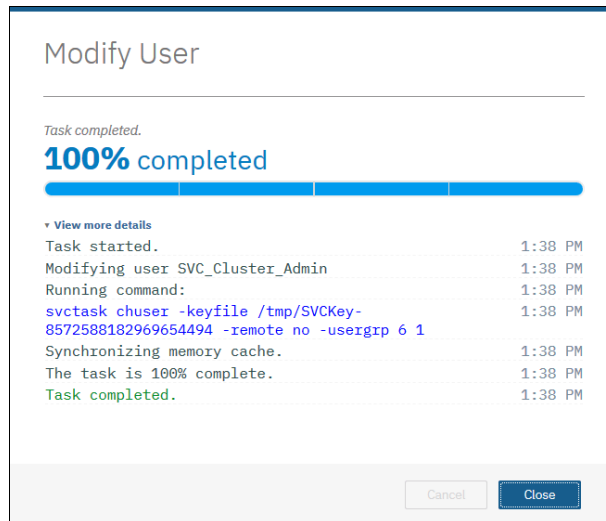


Figure 13-98 Confirmation of SSH public key upload

5. Check in the GUI to make sure that the SSH key is imported successfully, as shown in Figure B-8.

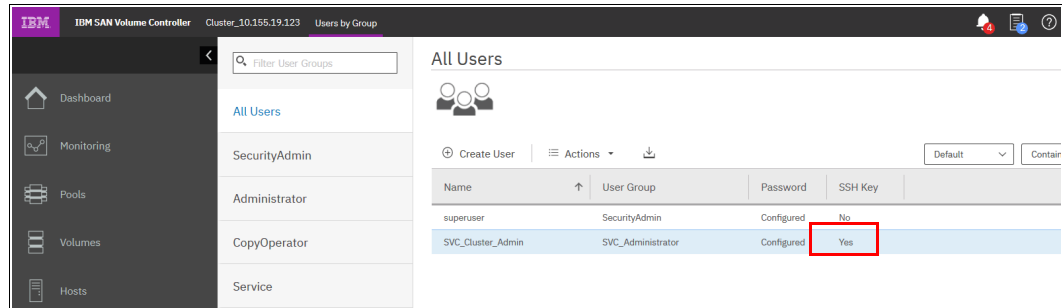


Figure B-8 Key successfully imported

## Configuring the SSH client

Before the CLI can be used, the SSH client must be configured. Complete the following steps:

1. Start PuTTY. The PuTTY Configuration window opens (see Figure B-9).

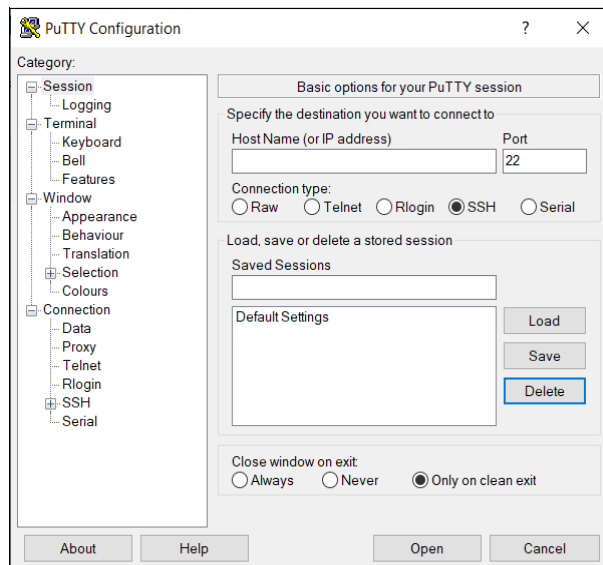


Figure B-9 PuTTY Configuration window

2. In the right pane, select **SSH** as the connection type. In the Close window on exit section, select **Only on clean exit**, which ensures that if any connection errors occur, they are shown in the user's window. These settings are shown in Figure B-9.

- In the Category pane, on the left side of the PuTTY Configuration window, click **Connection** → **Data**, as shown in Figure B-10. In the Auto-login username field, enter the IBM Spectrum Virtualize user ID that was used when uploading the public key. The admin account was used in the example of Figure B-6 on page 854.

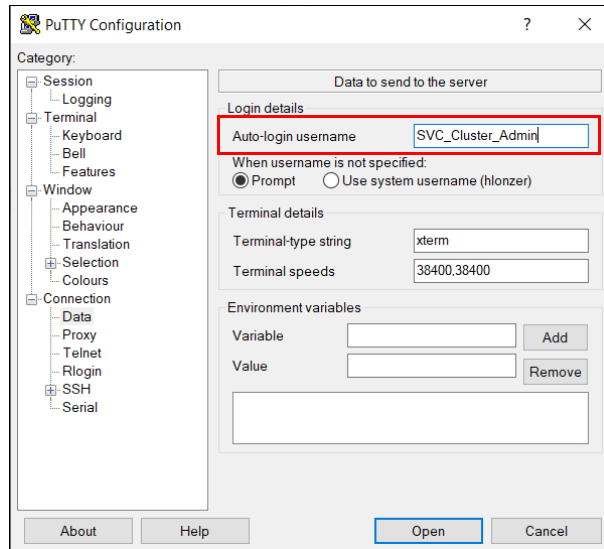


Figure B-10 PuTTY Auto-login username

- In the Category pane, on the left side of the PuTTY Configuration window (see Figure B-11), click **Connection** → **SSH** to open the PuTTY SSH Configuration window. In the SSH protocol version section, select **2**.

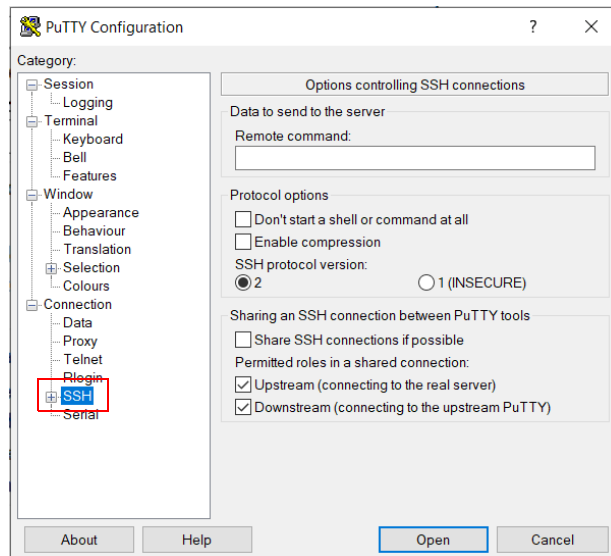


Figure B-11 SSH protocol version 2

- In the Category pane on the left, click **Connection** → **SSH** → **Auth**. More options are displayed for controlling SSH authentication.

6. In the Private key file for authentication field (see Figure B-12), browse to or enter the fully qualified directory path and file name of the SSH client private key file that was created (in this example, C:\Putty\SVCprivate.ppk).

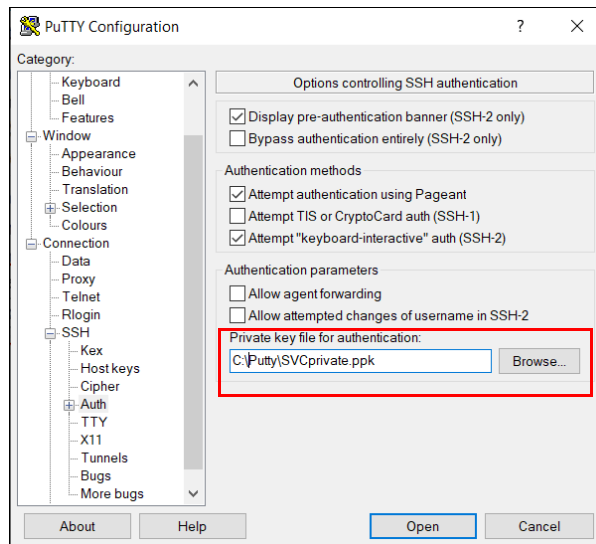


Figure B-12 SSH authentication

7. In the Category pane, click **Session** to return to the Basic options for your PuTTY session view.
8. Enter the following information in these fields in the right pane (see Figure B-13):
  - **Host Name:** Specify the host name or system IP address of the IBM Spectrum Virtualize system.
  - **Saved Sessions:** Enter a session name.
 Click **Save** to save the new session

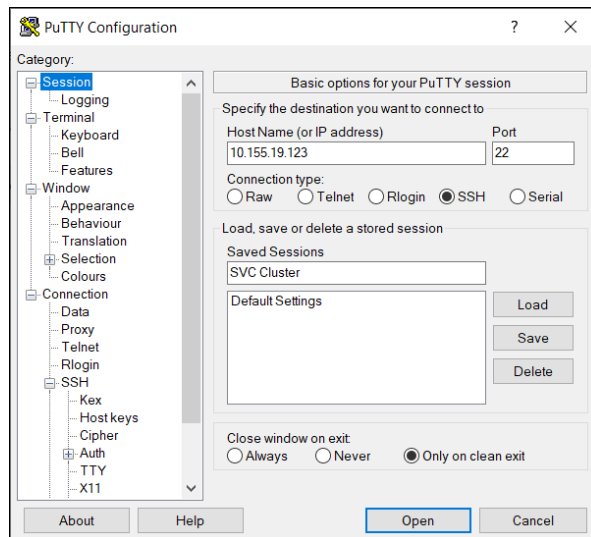


Figure B-13 Session information

9. Select the session and click **Open** to connect to the IBM Spectrum Virtualize system, as shown in Figure B-14.

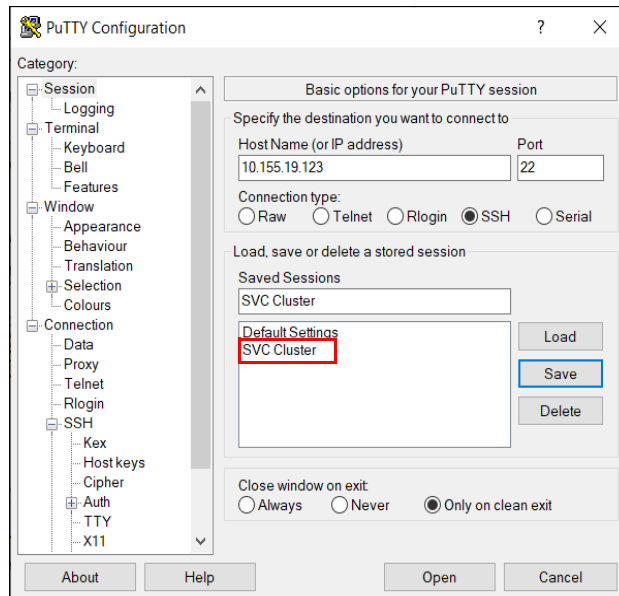


Figure B-14 Connecting to the system

10. If a PuTTY Security Alert opens, such as shown in Figure B-15, confirm it by clicking **Yes**.

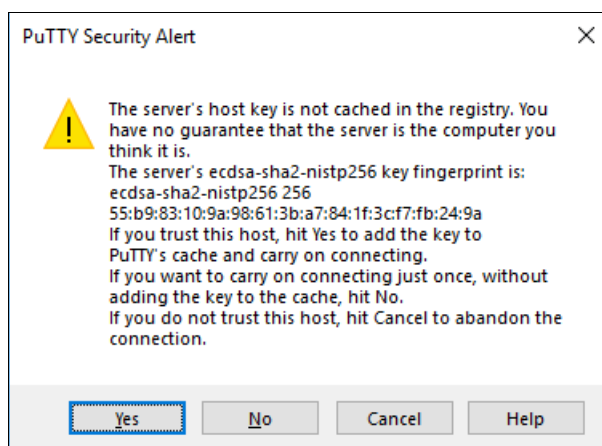


Figure B-15 Confirming the security alert

As shown in Figure B-16, PuTTY now connects to the system automatically by using the user ID that was specified earlier without prompting for password.

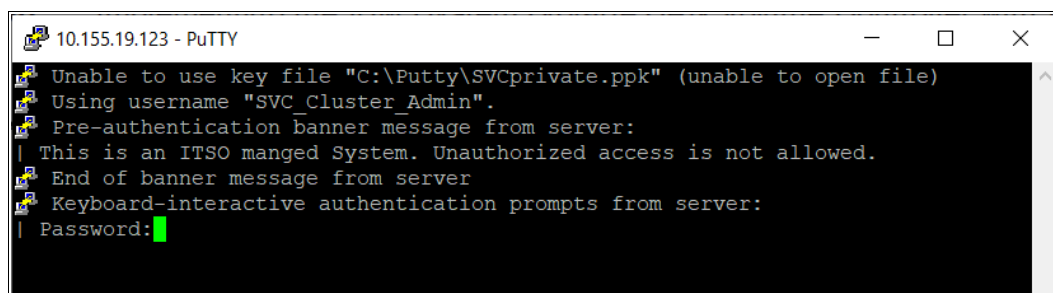


Figure B-16 PuTTY login

The CLI for IBM Spectrum Virtualize system administration is now configured.

## Basic setup on a UNIX or Linux host

OpenSSH client is the most common tool that is used on Linux or UNIX operating systems. It is installed by default on most of these types of operating systems. If it is not installed on your system, OpenSSH is available from the [OpenSSH website](#).

The OpenSSH suite consists of several tools, but the following tools are used to generate the SSH keys, transfer the SSH keys to a remote system, and establish a connection to an IBM Spectrum Virtualize device:

- ▶ **ssh**: OpenSSH SSH client
- ▶ **ssh-keygen**: Tool to generate SSH keys
- ▶ **scp**: Tool to transfer files between hosts

## Generating a public and private key pair

To generate a public and a private key to connect to IBM Spectrum Virtualize system without entering the user password, run the **ssh-keygen** tool, as shown in Example B-1.

*Example B-1 SSH keys generation with ssh-keygen*

---

```
# ssh-keygen -t rsa -b 1024
Generating public/private rsa key pair.
Enter file in which to save the key (//.ssh/id_rsa): /.ssh/sshkey
Enter passphrase (empty for no passphrase):
Enter same passphrase again:
Your identification has been saved in /.ssh/sshkey.
Your public key has been saved in /.ssh/sshkey.pub.
The key fingerprint is:
55:5e:5e:09:db:a4:11:b9:57:96:74:0c:85:ed:5b root@hostname.ibm.com
The key's randomart image is:
+--[ RSA 1024]-----+
|          .+=B0* |
|         + oB**+ |
|          . oo+o  |
|          . . . E |
|         S .  o   |
|                   |
|                   |
+-----+
#
```

---

In **ssh-keygen**, the parameter **-t** refers to the type of SSH key (RSA in Example B-1) and **-b** is the size of SSH key in bits (in Example B-1, 1024 bits was used). You also must specify the path and name for the SSH keys. The name that you provide is the name of the private key. The public key has the same name, but with the extension **.pub**. In Example B-1, the path is **/.ssh/**, the name of the private key is **sshkey**, and the name of the public key is **sshkey.pub**.

**Note:** The use of a passphrase for the SSH key is optional. If a passphrase is used, security is increased, but extra steps are required to log in by using the SSH key because the user must enter the passphrase.

## Uploading the SSH public key to the IBM SAN Volume Controller

In “Uploading the SSH public key to the IBM SAN Volume Controller” on page 854, you learned how to upload the new SSH public key to IBM Spectrum Virtualize by using the GUI.

Complete the following steps to upload the public key by using CLI:

1. On the SSH client (for example, AIX or Linux host), run **scp** to copy the public key to IBM SAN Volume Controller. The command features the following basic syntax:

```
scp <file> <user>@<hostname_or_IP_address>:<path>
```

The **/tmp** directory in the IBM Spectrum Virtualize active configuration node can be used to store the public key temporarily.

Example B-2 shows the command to copy the newly generated public key to the IBM Spectrum Virtualize system.

*Example B-2 SSH public key copy to IBM*

---

```
# scp /.ssh/sshkey.pub admin@192.168.1.100:/tmp/  
Password:  
sshkey.pub  
100% 241    0.2KB/s   00:00  
#
```

---

2. Log in to IBM SAN Volume Controller by using SSH and run the **chuser** command to import the public SSH key to a user, as shown in Example B-3.

*Example B-3 Importing an SSH public key to a user*

---

```
IBM_2145:SVC_Cluster:SVC_Cluster_Admin>chuser -keyfile /tmp/sshkey.pub admin  
IBM_2145:SVC_Cluster:SVC_Cluster_Admin>lsuser admin  
id 4  
name SVC_Cluster_Admin  
password yes  
ssh_key yes  
remote no  
usergrp_id 1  
usergrp_name Administrator  
IBM_2145:SVC_Cluster:SVC_Cluster_Admin>
```

---

When running the **lsuser** command as shown in Example B-3, it shows that a user has a configured SSH key in the field `ssh_key`.

## Connecting to an IBM Spectrum Virtualize system

Now that you uploaded the SSH key to the IBM Spectrum Virtualize system and assigned it to a user account, you can connect to the device by running the **ssh** command with the following options:

```
ssh -i <SSH_private_key> <user>@<IP_address_or_hostname>
```

Example B-4 shows the **ssh** command running from an AIX server and connecting to IBM SAN Volume Controller by using an SSH private key and no password prompt.

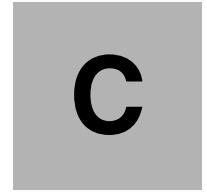
*Example B-4 Connecting to IBM SAN Volume Controller by using an SSH private key*

---

```
# ssh -i /.ssh/sshkey admin@192.168.1.100  
IBM_2145:SVC_Cluster:SVC_Cluster_Admin>
```

---





# Terminology

This appendix summarizes the IBM Spectrum Virtualize and IBM SAN Volume Controller terms that are commonly used in this book.

To see the complete set of terms that are related to the IBM SAN Volume Controller, see the Glossary that is available at [IBM Knowledge Center](#).

## Commonly used terms

This book uses the common IBM Spectrum Virtualize and IBM SAN Volume Controller terminology that are listed in this section.

### **Access mode**

One of the modes in which a logical unit (LU) in a disk controller system can operate. The three access modes are image mode, managed space mode, and unconfigured mode. See also “Image mode” on page 876, “Managed mode” on page 879, “Unconfigured mode” on page 889.

### **Activation key**

See “License key” on page 878

### **Array**

An ordered collection, or group, of physical devices (disk drive modules) that are used to define logical volumes or devices. An array is a group of drives designated to be managed with a Redundant Array of Independent Disks (RAID).

### **Asymmetric virtualization**

Asymmetric virtualization is a virtualization technique in which the virtualization engine is outside the data path and performs a metadata-style service. The metadata server contains all the mapping and locking tables, and the storage devices contain only data. See also “Symmetric virtualization” on page 888.

### **Asynchronous replication**

Asynchronous replication is a type of replication in which control is given back to the application as soon as the write operation is made to the source volume. Later, the write operation is made to the target volume. See also “Synchronous replication” on page 888.

### **Audit Log**

An unalterable record of all commands or user interactions that are issued to the system.

### **Automatic data placement mode**

Automatic data placement mode is an IBM Easy Tier operating mode in which the host activity on all the volume extents in a pool are “measured,” a migration plan is created, and then automatic extent migration is performed.

### **Auxiliary volume**

The auxiliary volume that contains a mirror of the data on the master volume. See also “Master volume” on page 879, and “Relationship” on page 884.

### **Available (usable) capacity**

See “Capacity” on page 865.

### **Back end**

See “Front end and back end” on page 873.

## Caching I/O group

The caching I/O group is the I/O group in the system that performs the cache function for a volume.

## Call home

Call home is a communication link that is established between a product and a service provider. The product can use this link to call IBM or another service provider when the product requires service. With access to the machine, service personnel can perform service tasks, such as viewing error and problem logs or initiating trace and dump retrievals.

## Canister

A canister is a single processing unit within a storage system.

## Capacity

These are the definitions IBM applies to capacity:

- ▶ Available capacity

The amount of usable capacity that is not yet used in a system, pool, array, or managed disk (MDisk).

- ▶ Capacity

The amount of data that can be contained on a storage medium.

- ▶ Data reduction

A set of techniques that can be used to reduce the amount of usable capacity that is required to store data. Examples of data reduction include data deduplication and compression.

- ▶ Data reduction savings

The total amount of usable capacity that is saved in a system, pool, or volume through the application of an algorithm such as compression or deduplication on the written data. This saved capacity is the difference between the written capacity and the used capacity.

- ▶ Effective capacity

The amount of provisioned capacity that can be created in a system or pool without running out of usable capacity given the current data reduction savings being achieved. This capacity equals the usable capacity divided by the data reduction savings percentage.

- ▶ Overhead capacity

An amount of usable capacity that is occupied by metadata in a system or pool and other data that is used for system operations.

- ▶ Over provisioned

The result of creating more provisioned capacity in a storage system or pool than there is usable capacity. Overprovisioning occurs when thin provisioning or data reduction techniques ensure that the used capacity of the provisioned volumes is less than their provisioned capacity.

- ▶ Over provisioned ratio

The ratio of provisioned capacity to usable capacity in the pool or system.

- ▶ Provisioned capacity

Total capacity of all volumes and Volume copies in a pool or system.

- ▶ **Provisioning limit - maximum provisioned capacity - over provisioning limit**  
In some storage systems, restrictions in the storage hardware or configured by the user that define a limit the maximum provisioned capacity allowed in a pool or system.
- ▶ **Raw capacity**  
The reported capacity of the drives in the system before formatting or RAID is applied.
- ▶ **Standard provisioning**  
The ability to completely use a volume's capacity for that specific volume.
- ▶ **Standard provisioned Volume**  
A volume that completely uses storage at creation.
- ▶ **Thin provisioning savings**  
The total amount of usable capacity saved in a pool, system, or volume by using usable capacity when needed as a result of write operations. The capacity saved is the difference between the provisioned capacity minus the written capacity.
- ▶ **Total capacity savings**  
The total amount of usable capacity saved in a pool, system, or volume through thin-provisioning and data reduction techniques. This capacity saved is the difference between the used usable capacity and the provisioned capacity.
- ▶ **Usable capacity**  
The amount of capacity that is provided for storing data on a system, pool, array, or MDisk after formatting and RAID techniques are applied. Usable capacity is the total of used and available capacity. For example, 50 TiB used, 50 TiB available is a usable capacity of 100 TiB.
- ▶ **Used capacity**  
The amount of usable capacity that is taken up by data or overhead capacity in a system, pool, array, or MDisk after data reduction techniques are applied.
- ▶ **Written capacity**  
The amount of usable capacity that would be used to store written data in a pool or system before data reduction is applied.
- ▶ **Written capacity limit**  
The largest amount of capacity that can be written to a drive, array, or MDisk. The limit can be reached even when usable capacity is still available.

## Capacity licensing

Capacity licensing is a licensing model that licenses features with a price-per-terabyte model. Licensed features are FlashCopy, Metro Mirror (MM), Global Mirror (GM), and virtualization. See also “FlashCopy” on page 873, “Metro Mirror” on page 879, and “Virtualized storage” on page 889.

## Capacity recycling

Capacity recycling means the amount of provisioned capacity that can be recovered without causing stress or performance degradation. This capacity identifies the amount of resources that can be reclaimed and provisioned to other objects in an environment.

## Chain

A set of enclosures that are attached to provide redundant access to the drives inside the enclosures. Each control enclosure can have one or more chains.

## **Challenge Handshake Authentication Protocol**

Challenge Handshake Authentication Protocol (CHAP) is an authentication protocol that protects against eavesdropping by encrypting the user name and password.

## **Change volume**

A volume that is used in GM that holds earlier consistent revisions of data when changes are made.

## **Channel extender**

A channel extender is a device that is used for long-distance communication that connects other storage area network (SAN) fabric components. Generally, channel extenders can involve protocol conversion to asynchronous transfer mode (ATM), IP, or another long-distance communication protocol.

## **Child pool**

Administrators can use child pools to control capacity allocation for volumes that are used for specific purposes. Rather than being created directly from MDisks, child pools are created from existing capacity that is allocated to a parent pool. As with parent pools, volumes can be created that specifically use the capacity that is allocated to the child pool. Child pools are similar to parent pools with similar properties. Child pools can be used for volume copy operation. Also, see “Parent pool” on page 881.

## **Cloud account**

An agreement with a cloud service provider (CSP) to use storage or other services at that service provider. Access to the cloud account is granted by presenting valid credentials.

## **Cloud Container**

Cloud Container is a virtual object that include all of the elements, components, or data that are common to a specific application or data.

## **Cloud Provider**

Cloud provider is the company or organization that provide off- and on-premises cloud services such as storage, server, network, and so on. IBM Spectrum Virtualize has built-in software capabilities to interact with Cloud Providers such as IBM Cloud, Amazon S3, and deployments of OpenStack Swift.

## **Cloud Tenant**

Cloud Tenant is a group or an instance that provides common access with the specific privileges to an object, software, or data source.

## **Clustered system (SAN Volume Controller)**

A clustered system, which was known as a cluster, is a group of up to eight SAN Volume Controller nodes that presents a single configuration, management, and service interface to the user.

## **Cold extent**

A cold extent is an extent of a volume that does not get any performance benefit if it is moved from a hard disk drive (HDD) to a Flash disk. A cold extent also refers to an extent that needs to be migrated onto an HDD if it is on a Flash disk drive.

## Compression

Compression is a function that removes repetitive characters, spaces, strings of characters, or binary data from the data that is being processed and replaces characters with control characters. Compression reduces the amount of storage space that is required for data.

## Compression accelerator

A compression accelerator is hardware onto which the work of compression is offloaded from the microprocessor.

## Configuration node

While the cluster is operational, a single node in the cluster is appointed to provide configuration and service functions over the network interface. This node is termed the configuration node. This configuration node manages the data that describes the clustered-system configuration and provides a focal point for configuration commands. If the configuration node fails, another node in the cluster transparently assumes that role.

## Consistency Group

A Consistency Group is a group of copy relationships between virtual volumes or data sets that are maintained with the same time reference so that all copies are consistent in time. A Consistency Group can be managed as a single entity.

## Container

A container is a software object that holds or organizes other software objects or entities.

## Contingency capacity

For thin-provisioned volumes that are configured to automatically expand, the unused real capacity that is maintained. For thin-provisioned volumes that are not configured to automatically expand, the difference between the used capacity and the new real capacity.

## Copied state

Copied is a FlashCopy state that indicates that a copy was triggered after the copy relationship was created. The Copied state indicates that the copy process is complete and the target disk has no further dependency on the source disk. The time of the last trigger event is normally displayed with this status.

## Counterpart SAN

A counterpart SAN is a non-redundant portion of a redundant SAN. A counterpart SAN provides all of the connectivity of the redundant SAN, but without the 100% redundancy. SAN Volume Controller nodes are typically connected to a *redundant SAN* that is made up of two *counterpart SANs*. A counterpart SAN is often called a *SAN fabric*.

## Cross-volume consistency

A consistency group property that ensures consistency between volumes when an application issues dependent write operations that span multiple volumes.

## Data consistency

Data consistency is a characteristic of the data at the target site where the dependent write order is maintained to ensure the recoverability of applications.

## Data deduplication

Data deduplication is a method of reducing storage needs by eliminating redundant data. Only one instance of the data is retained on storage media. Other instances of the same data are replaced with a pointer to the retained instance.

## Data encryption key

The data encryption key is used to encrypt data. It is created automatically when an encrypted object, such as an array, a pool, or a child pool, is created. It is stored in secure memory and it cannot be viewed or changed. The data encryption key is encrypted using the master access key.

## Data migration

Data migration is the movement of data from one physical location to another physical location without the disruption of application I/O operations.

## Data reduction

Data reduction is a set of techniques that can be used to reduce the amount of physical storage that is required to store data. An example of data reduction includes data deduplication and compression. See also “Data Reduction Pool” and “Capacity” on page 865.

## Data Reduction Pool

Data Reduction Pools (DRPs) are specific types of pools where more control over volume capacity is given to specific hosts (for example VMware VAAI/VASA/VVOL, Microsoft Offloaded Data Transfer (ODX)). These hosts are able to return unused space for reuse. With standard pools, the system is not aware of any unused space on host-allocated volumes. See also “Data reduction”.

## Data reduction savings

See “Capacity” on page 865.

## Deduplication

See “Data deduplication” on page 869.

## Dependent write operation

A write operation that must be applied in the correct order to maintain cross-volume consistency.

## Directed maintenance procedure

The fix procedures, which are also known as Directed Maintenance Procedures (DMPs), ensure that you fix any outstanding errors in the error log. To fix errors, from the Monitoring pane, click **Events**. The Next Recommended Action is displayed at the top of the Events window. Select **Run This Fix Procedure** and follow the instructions.

## Discovery

The automatic detection of a network topology change, for example, new and deleted nodes or links.

## Disk tier

MDisks (logical unit numbers (LUNs)) that are presented to the SAN Volume Controller cluster likely have different performance attributes because of the type of disk or RAID array on which they are installed. The MDisks can be on 15,000 revolutions per minute (RPM) Fibre Channel (FC) or serial-attached Small Computer System Interface (SCSI) (SAS) disk, Nearline (NL) SAS, or Serial Advanced Technology Attachment (SATA), or even Flash Disks. Therefore, a storage tier attribute is assigned to each MDisk and the default is *generic\_hdd*.

## Distributed Redundant Array of Independent Disks

An alternative RAID scheme where the number of drives that are used to store the array can be greater than the equivalent, typical RAID scheme. The same data stripes are distributed across a greater number of drives, which increases the opportunity for parallel I/O and hence improves performance of the array. See also “Rebuild area” on page 884.

## Drive technology

A category of a drive that pertains to the method and reliability of the data storage techniques being used on the drive. Possible values include enterprise (ENT) drive, NL drive, or solid-state drive (SSD).

## Dual Inline Memory Module

A Dual Inline Memory Module (DIMM) is small circuit board with memory-integrated circuits containing signal and power pins on both sides of the board.

- Memory Terminology

- Channel

- The memory modules are installed into matching banks, which are usually color-coded on the motherboard. These separate channels allow the memory controller access each memory module. For the Intel Cascade Lake architecture, there are 6 DIMM Memory channels per CPU and each memory channel has two DIMMs. The memory bandwidth is tied to each of these channels, and the speed of access for the memory controller is shared across the pair of DIMMs in that channel.

- Slot

- Generally the physical slot that a DIMM can fit into, however, for the purposes of clarity, slot here will be DIMM0 or DIMM1 which refers to the first or second slot within a channel on the motherboard. There are 2 slots per memory channel on the SAN Volume Controller SV2 hardware. On the motherboard, DIMM0 is the blue slot and DIMM1 is the black slot within each channel.

- Rank

- A single-rank DIMM has one set of memory chips that is accessed while writing to or reading from the memory. A dual-rank DIMM is similar to having two single-rank DIMMs on the same module, with only one rank accessible at a time. A quad-rank DIMM is, effectively, two dual-rank DIMMs on the same module. The 32G DIMMS are dual rank.

## Easy Tier

Easy Tier is a volume performance function within the SAN Volume Controller that provides automatic data placement of a volume’s extents in a multitiered storage pool. The pool normally contains a mix of Flash Disks and HDDs. Easy Tier measures host I/O activity on the volume’s extents and migrates hot extents onto the Flash Disks to ensure the maximum performance.



## **Effective capacity**

See “Capacity” on page 865.

## **Encryption key**

The encryption key, also known as the master access key, is created and stored on USB flash drives or on a key server when encryption is enabled. The master access key is used to decrypt the data encryption key.

## **Encryption key manger / server**

An internal or external system that receives and then serves existing encryption keys or certificates to a storage system.

## **Encryption recovery key**

An encryption key that allows a method to recover from an encryption deadlock situation where the normal encryption key servers are not available.

## **Encryption of data-at-rest**

Encryption of data-at-rest is the inactive encryption data that is stored physically on the storage system.

## **Enhanced Stretched Systems**

A stretched system is an extended high availability (HA) method that is supported by the SAN Volume Controller to enable I/O operations to continue after the loss of half of the system. Enhanced Stretched Systems provide the following primary benefits. In addition to the automatic failover that occurs when a site fails in a standard stretched system configuration, an Enhanced Stretched System provides a manual override that can be used to select which of two sites continues operation.

Enhanced Stretched Systems intelligently route I/O traffic between nodes and controllers to reduce the amount of I/O traffic between sites, and to minimize the effect on host application I/O latency. Enhanced Stretched Systems include an implementation of additional policing rules to ensure that the correct configuration of a standard stretched system is used.

## **Evaluation mode**

The evaluation mode is an Easy Tier operating mode in which the host activity on all the volume extents in a pool are “measured” only. No automatic extent migration is performed.

## **Event (error)**

An event is an occurrence of significance to a task or system. Events can include the completion or failure of an operation, user action, or a change in the state of a process.

## **Event code**

An event code is a value that is used to identify an event condition to a user. This value might map to one or more event IDs or to values that are presented on the service window. This value is used to report error conditions to IBM and to provide an entry point into the service guide.

## **Event ID**

An event ID is a value that is used to identify a unique error condition that was detected by the SAN Volume Controller. An event ID is used internally in the cluster to identify the error.

**Excluded condition**

The excluded condition is a status condition. It describes an MDisk that the SAN Volume Controller decided is no longer sufficiently reliable to be managed by the cluster. The user must issue a command to include the MDisk in the cluster-managed storage.

**Extent**

An extent is a fixed-size unit of data that is used to manage the mapping of data between MDisks and volumes. The size of the extent can range from 16 MB - 8 GB.

**External storage**

External storage refers to MDisks that are SCSI LUs that are presented by storage systems that are attached to and managed by the clustered system.

**Failback**

Failback is the restoration of an appliance to its initial configuration after the detection and repair of a failed network or component.

**Failover**

Failover is an automatic operation that switches to a redundant or standby system or node in a software, hardware, or network interruption. See also Failback.

**Feature activation code**

An alphanumeric code that activates a licensed function on a product. See also “License key” on page 878.

**Fibre Channel**

A technology for transmitting data between computer devices. It is especially suited for attaching computer servers to shared storage devices and for interconnecting storage controllers and drives. See also “Zoning” on page 891.

**Fibre Channel port logins**

FC port logins refer to the number of hosts that can see any one SAN Volume Controller node port. The SAN Volume Controller has a maximum limit per node port of FC logins that are allowed.

**Fibre Channel Arbitrated Loop (FC-AL)**

An implementation of the FC standards that uses a ring topology for the communication fabric; refer to American National Standards Institute (ANSI) INCITS 272-1996, (R2001). In this topology, two or more FC end points are interconnected through a looped interface.

**Fibre Channel Connection**

A FC communication protocol designed for IBM mainframe computers and peripherals.

**Fibre Channel over IP**

A network storage technology that combines the features of the Fibre Channel Protocol (FCP) and the IP to connect distributed SANs over large distances.

**Fibre Channel Protocol**

The serial SCSI command protocol used on FC networks.

## **Field-replaceable unit**

Field-replaceable units (FRUs) are individual parts that are replaced entirely when any one of the unit's components fails. They are held as spares by the IBM service organization.

## **File Transfer Protocol**

In TCP/IP, an application layer protocol that uses TCP and Telnet services to transfer bulk-data files between machines or hosts.

## **Fix procedure**

A maintenance procedure that runs within the product application and provides step-by-step guidance to resolve an error condition.

## **FlashCopy**

FlashCopy refers to a point-in-time (PiT) copy where a virtual copy of a volume is created. The target volume maintains the contents of the volume at the PiT when the copy was established. Any subsequent write operations to the source volume are not reflected on the target volume.

## **FlashCopy mapping**

A FlashCopy mapping is a continuous space on a direct-access storage volume that is occupied by or reserved for a particular data set, data space, or file.

## **FlashCopy relationship**

See FlashCopy mapping.

## **FlashCopy service**

FlashCopy service is a copy service that duplicates the contents of a source volume on a target volume. In the process, the original contents of the target volume are lost. See also "Point-in-time copy" on page 881.

## **FlashCore Module**

FlashCore Module (FCM) modules are a family of high-performance flash drives. The FCM design utilizes the Non-Volatile Memory Express (NVMe) protocol, a PCIe Gen3 interface, and high-speed NAND memory to provide high throughput and input/output operations per second (IOPS) and very low latency. FCM modules are available in different capacities. Hardware-based data compression and self-encryption are supported. The FCM modules are accessible from the front of the enclosure.

## **Flash drive**

A data storage device, which is typically removable and rewriteable that uses solid-state memory to store persistent data. See also "Flash module"

## **Flash module**

A modular hardware unit that contains flash memory, one or more flash controllers, and associated electronics.

## **Front end and back end**

The SAN Volume Controller takes MDisks to create pools of capacity from which volumes are created and presented to application servers (hosts). The volumes that are presented to the hosts are in the front end of the IBM SAN Volume Controller.

## **Full restore operation**

A copy operation where a local volume is created by reading an entire a volume snapshot from cloud storage.

## **Full snapshot**

A type of volume snapshot that contains all the volume data. When a full snapshot is created, an entire copy of the volume data is transmitted to the cloud.

## **General Parallel File System**

A high-performance shared-disk file system that can provide data access from nodes in a clustered system environment.

## **GB**

See “Gigabyte” on page 874.

## **Gigabyte**

For processor storage, real and virtual storage, and channel volume, two to the power of 30 or 1,073,741,824 bytes. For disk storage capacity and communications volume, 1,000,000,000 bytes.

## **Global Mirror**

Global Mirror (GM) is a method of asynchronous replication that maintains data consistency across multiple volumes within or across multiple systems. GM is generally used where distances between the source site and target site cause increased latency beyond what the application can accept.

## **Global Mirror with Change Volumes**

Change volumes are used to record changes to the primary and secondary volumes of a Remote Copy (RC) relationship. A FlashCopy mapping exists between a primary and its change volume, and a secondary and its change volume.

## **GPFS**

See “General Parallel File System”.

## **GPFS cluster**

A system of nodes that are defined as being available for use by GPFS file systems.

## **GPFS snapshot**

A PiT copy of a file system or file set.

## **Grain**

A grain is the unit of data that is represented by a single bit in a FlashCopy bitmap (64 kibibytes [KiB] or 256 KiB) in the SAN Volume Controller. A grain is also the unit to extend the real size of a thin-provisioned volume (32 KiB, 64 KiB, 128 KiB, or 256 KiB).

## **Graphical user interface**

A computer interface that presents a visual metaphor of a real-world scene, often of a desktop, by combining high-resolution graphics, pointing devices, menu bars and other menus, overlapping windows, icons and the object-action relationship.

**Hop**

One segment of a transmission path between adjacent nodes in a routed network.

**Host**

A physical or virtual computer system that hosts computer applications, with the host and the applications using storage.

**Host bus adapter**

A host bus adapter (HBA) is an interface card that connects a server to the SAN environment through its internal bus system, for example, PCI Express. Typically it is referred to the FC adapters.

**Host Cluster**

A configured set of physical or virtual hosts that share one or more storage volumes in order to increase scalability or availability of computer applications.

**Host ID**

A host ID is a numeric identifier that is assigned to a group of host FC ports or Internet Small Computer Systems Interface (iSCSI) host names for LUN mapping. For each host ID, SCSI IDs are mapped to volumes separately. The intent is to have a one-to-one relationship between hosts and host IDs, although this relationship cannot be policed.

**Host mapping**

Host mapping refers to the process of controlling which hosts have access to specific volumes within a cluster. Host mapping is equivalent to LUN masking.

**Host object**

A logical representation of a host within a storage system that is used to represent the host for configuration tasks.

**Host zone**

A zone that is defined in the SAN fabric in which the hosts can address the system.

**Hot extent**

A hot extent is a frequently accessed volume extent that gets a performance benefit if it is moved from an HDD onto a Flash Disk.

**Hot Spare Node**

Hot Spare Node is an online SAN Volume Controller node defined in a cluster, but not in any IO group. During a failure of any of online nodes in any IO group of cluster, it is automatically swapped with this spare node. After the recovery of an original node has finished, the spare node returns to the standby spare status.

**HyperSwap**

Pertaining to a function that provides continuous, transparent availability against storage errors and site failures, and is based on synchronous replication.

## **Image mode**

Image mode is an access mode that establishes a one-to-one mapping of extents in the storage pool (existing LUN or (image mode) MDisk) with the extents in the volume. See also “Managed mode” on page 879 and “Unconfigured mode” on page 889.

## **Image volume**

An image volume is a volume in which a direct block-for-block translation exists from the MDisk to the volume.

## **I/O group**

Each pair of SAN Volume Controller cluster nodes is known as an input/output (I/O) Group. An I/O group has a set of volumes that are associated with it that are presented to host systems. Each SAN Volume Controller node is associated with exactly one I/O group. The nodes in an I/O group provide a failover and failback function for each other.

## **Incremental restore operation**

A copy operation where a local volume is modified to match a volume snapshot by reading from cloud storage only the parts of the volume snapshot that differ from the local volume.

## **Incremental snapshot**

A type of volume snapshot where the changes to a local volume relative to the volume's previous snapshot are stored on cloud storage.

## **Input/output operations per second**

A standard computing benchmark used to determine the best configuration settings for servers.

## **Internal storage**

Internal storage refers to an array of MDisks and drives that are held in enclosures and in nodes that are part of the SAN Volume Controller cluster.

## **Internet Protocol**

A protocol that routes data through a network or interconnected networks. This protocol acts as an intermediary between the higher protocol layers and the physical network.

## **Internet Small Computer Systems Interface**

A protocol that is used by a host system to manage iSCSI targets and iSCSI discovery. iSCSI initiators use the internet Storage Name Service (iSNS) protocol to locate the appropriate storage resources.

## **iSCSI Qualified Name**

iSCSI Qualified Name (IQN) refers to special names that identify both iSCSI initiators and targets. IQN is one of the three name formats that is provided by iSCSI. The IQN format is `iqn.<yyyymm>.<reversed domain name>`. For example, the default for a SAN Volume Controller node can be in the following format:

```
iqn.1986-03.com.ibm:2145.<clustername>.<nodename>
```

## **Internet storage name service**

The iSNS protocol that is used by a host system to manage iSCSI targets and the automated iSCSI discovery, management, and configuration of iSCSI and FC devices. It was defined in Request for Comments (RFC) 4171.

## **Inter-switch link hop**

An inter-switch link (ISL) is a connection between two switches and counted as one ISL hop. The number of hops is always counted on the shortest route between two N-ports (device connections). In a SAN Volume Controller environment, the number of ISL hops is counted on the shortest route between the pair of nodes that are farthest apart. The SAN Volume Controller supports a maximum of three ISL hops.

## **Input/output group**

A collection of volumes and node relationships that present a common interface to host systems. Each pair of nodes is known as an I/O group.

## **I/O throttling rate**

The maximum rate at which an I/O transaction is accepted for a volume.

## **iSCSI**

See “Internet Small Computer Systems Interface” on page 876

## **iSCSI initiator**

An initiator functions as an iSCSI client. An initiator typically serves the same purpose to a computer as a SCSI bus adapter would, except that, instead of physically cabling SCSI devices (like hard drives and tape changers), an iSCSI initiator sends SCSI commands over an IP network.

## **iSCSI alias**

An alternative name for the iSCSI-attached host.

## **iSCSI name**

A name that identifies an iSCSI target adapter or an iSCSI initiator adapter. An iSCSI name can be an IQN or an extended-unique identifier (EUI). Typically, this identifier has the following format: iqn.datecode.reverse domain.

## **iSCSI session**

The interaction (conversation) between an iSCSI Initiator and an iSCSI Target.

## **iSCSI target**

An iSCSI target is a storage resource located on an iSCSI server.

## **Just A Bunch Of Disks**

Hard disks that haven't been configured according to the RAID system to increase fault tolerance and improve data access performance.

## **Key server**

- ▶ A server that negotiates the values that determine the characteristics of a dynamic VPN connection that is established between two endpoints.
- ▶ See “Encryption key manger / server” on page 871.

**Latency**

The time interval between the initiation of a send operation by a source task and the completion of the matching receive operation by the target task. More generally, latency is the time between a task initiating data transfer and the time that transfer is recognized as complete at the data destination.

**Least recently used**

Least recently used (LRU) pertains to an algorithm used to identify and make available the cache space that contains the data that was LRU.

**Licensed capacity**

The amount of capacity on a storage system that a user is entitled to configure.

**License key**

An alphanumeric code that activates a licensed function on a product.

**License key file**

A file that contains one or more licensed keys.

**Lightweight Directory Access Protocol**

Lightweight Directory Access Protocol (LDAP) is an open protocol that uses TCP/IP to provide access to directories that support an X.500 model. It does not incur the resource requirements of the more complex X.500 Directory Access Protocol (DAP). For example, LDAP can be used to locate people, organizations, and other resources in an Internet or intranet directory.

**Local and remote fabric interconnect**

The local fabric interconnect and the remote fabric interconnect are the SAN components that are used to connect the local and remote fabrics. Depending on the distance between the two fabrics, they can be single-mode optical fibers that are driven by long wave (LW) gigabit interface converters (GBICs) or Small Form-factor Pluggable (SFP) modules, or more sophisticated components, such as channel extenders or special SFP modules that are used to extend the distance between SAN components.

**Local fabric**

The local fabric is composed of SAN components (switches, cables, and so on) that connect the components (nodes, hosts, and switches) of the local cluster together.

**Logical drive**

See "Volume" on page 890.

**Logical unit and logical unit number**

The LU is defined by the SCSI standards as a LUN. LUN is an abbreviation for an entity that exhibits disk-like behavior, such as a volume or an MDisk.

**LUN masking**

A process where a host object can detect more LUNs than it is intended to use, and the device-driver software masks the LUNs that are not to be used by this host.



## **Machine signature**

A string of characters that identifies a system. A machine signature might be required to obtain a license key.

## **Managed disk**

An MDisk is a SCSI disk that is presented by a RAID controller and managed by the SAN Volume Controller. The MDisk is not visible to host systems on the SAN.

## **Managed disk group (storage pool)**

See “Storage pool (managed disk group)” on page 887.

## **Managed mode**

An access mode that enables virtualization functions to be performed. See also “Image mode” on page 876 and “Unconfigured mode” on page 889.

## **Management node**

A node that is used for configuring, administering, and monitoring a system.

## **Maximum replication delay**

Maximum replication delay is the number of seconds that MM or GM replication can delay a write operation to a volume.

## **Master volume**

In most cases, the volume that contains a production copy of the data and that an application accesses. See also “Auxiliary volume” on page 864, and “Relationship” on page 884.

## **Maximum replication delay**

Maximum replication delay is the number of seconds that MM or GM replication can delay a write operation to a volume.

## **MDisk**

See “Managed disk”.

## **Megabytes per second**

A unit of data transfer rate equal to 1024 \* 1024 bytes.

## **Metro Global Mirror**

Metro Global Mirror is a cascaded solution where MM synchronously copies data to the target site. This MM target is the source volume for GM that asynchronously copies data to a third site. This solution has the potential to provide disaster recovery (DR) with no data loss at GM distances when the intermediate site does not participate in the disaster that occurs at the production site.

## **Metro Mirror**

MM is a method of synchronous replication that maintains data consistency across multiple volumes within the system. MM is generally used when the write latency that is caused by the distance between the source site and target site is acceptable to application performance.

## **Mirrored volume**

A mirrored volume is a single virtual volume that has two physical volume copies. The primary physical copy is known within the SAN Volume Controller as copy 0 and the secondary copy is known within the SAN Volume Controller as copy 1.

## **Nearline SAS drive**

A drive that combines the high capacity data storage technology of a SATA drive with the benefits of a SAS interface for improved connectivity.

## **Node**

A SAN Volume Controller node is a hardware entity that provides virtualization, cache, and copy services for the cluster. The SAN Volume Controller nodes are deployed in pairs that are called I/O groups. One node in a clustered system is designated as the configuration node.

## **Node canister**

A node canister is a hardware unit that includes the node hardware, fabric and service interfaces, and SAS expansion ports. Node canisters are specifically recognized on IBM Storwize products. In SAN Volume Controller, all these components are spread within the whole system chassis, so we usually do not consider node canisters in SAN Volume Controller, but just the node as a whole.

## **Node rescue**

The process by which a node that has no valid software installed on its HDD can copy software from another node connected to the same FC fabric.

## **N\_Port ID Virtualization**

N\_Port ID Virtualization (NPIV) is a FC feature whereby multiple FC N\_Port IDs can share a single physical N\_Port.

## **Object-based access control**

Refer to “Ownership Groups”.

## **Object Storage**

Object storage is a general term that refers to the entity in which Cloud Object Storage organizes, manages, and stores with units of storage, or just *objects*.

## **Oversubscription**

Oversubscription refers to the ratio of the sum of the traffic on the initiator N-port connections to the traffic on the most heavily loaded ISLs, where more than one connection is used between these switches. Oversubscription assumes a symmetrical network, and a specific workload that is applied equally from all initiators and sent equally to all targets. A symmetrical network means that all the initiators are connected at the same level, and all the controllers are connected at the same level.

## **Over provisioned**

See “Capacity” on page 865.

## **Over provisioned ratio**

See “Capacity” on page 865.

## Ownership Groups

Ownership Groups feature provides a method of implementing multi-tenant solution on the system. Ownership groups allow the allocation of storage resources to several independent tenants with the assurance that one tenant cannot access resources associated with another tenant. Ownership groups restrict access for users in the ownership group to only those objects that are defined within that ownership group.

## Parent pool

Parent pools receive their capacity from MDisks. All MDisks in a pool are split into extents of the same size. Volumes are created from the extents that are available in the pool. You can add MDisks to a pool at any time either to increase the number of extents that are available for new volume copies or to expand existing volume copies. The system automatically balances volume extents between the MDisks to provide the best performance to the volumes. See also “Child pool” on page 867.

## Partnership

In MM or GM operations, the relationship between two clustered systems. In a clustered-system partnership, one system is defined as the local system and the other system as the remote system.

## Performance group

A collection of volumes that is assigned the same performance characteristics. See also “Performance policy”.

## Performance policy

A policy that specifies performance characteristics, for example quality of service (QoS). See also “Performance group”.

## Point-in-time copy

A PiT copy is an instantaneous copy that the FlashCopy service makes of the source volume. See also “FlashCopy service” on page 873.

## Pool

See “Storage pool (managed disk group)” on page 887.

## Pool pair

Two storage pools that are required to balance workload. Each storage pool is controlled by a separate node.

## Preferred Node

When you create a volume, you can specify a preferred node. Many of the multipathing driver implementations that the system supports use this information to direct I/O to the preferred node. The other node in the I/O group is used only if the preferred node is not accessible. If you do not specify a preferred node for a volume, the system selects the node in the I/O group that has the fewest volumes to be the preferred node. After the preferred node is chosen, it can be changed only when the volume is moved to a different I/O group. Note: The management GUI provides a wizard that moves volumes between I/O groups without disrupting host I/O operations.

## Preparing phase

Before you start the FlashCopy process, you must prepare a FlashCopy mapping. The preparing phase flushes a volume's data from cache in preparation for the FlashCopy operation.

## Primary volume

In a stand-alone MM or GM relationship, the target of write operations that are issued by the host application.

## Priority flow control

A link-level flow control mechanism, IEEE standard 802.1Qbb. Priority flow control (PFC) operates on individual priorities. Instead of pausing all traffic on a link, PFC is used to selectively pause traffic according to its class.

## Private fabric

Configure one SAN per fabric so that it is dedicated for node-to-node communication. This SAN is referred to as a private SAN.

## Provisioned capacity

See "Capacity" on page 865.

## Provisioning Group

A provisioning group is an object that represents a set of MDisks that share physical resources. Provisioning groups are used for capacity reporting and monitoring of overprovisioned storage resources.

## Public fabric

Configure one SAN per fabric so that it is dedicated for host attachment, storage system attachment, and RC operations. This SAN is referred to as a public SAN. You can configure the public SAN to allow SAN Volume Controller node-to-node communication also. You can optionally use the `-localportfcmask` parameter of the `chsystem` command to constrain the node-to-node communication to use only the private SAN.

## Qualifier

- ▶ A value that provides additional information about a class, association, indication, method, method parameter, instance, property, or reference.
- ▶ A modifier that makes a name unique.

## Quorum disk

A disk that contains a reserved area that is used exclusively for system management. The quorum disk is accessed when it is necessary to determine which half of the clustered system continues to read and write data. Quorum disks can either be MDisks or drives.

## Quorum index

The quorum index is the pointer that indicates the order that is used to resolve a tie. Nodes attempt to lock the first quorum disk (index 0), followed by the next disk (index 1), and finally the last disk (index 2). The tie is broken by the node that locks them first.

**Quota**

The amount of disk space and number of files and directories assigned as upper limits for a specified user, group of users, or file set.

**Random Access Compression Engine**

The Random Access Compression Engine (RACE) compresses data on volumes in real time with minimal effect on performance. See “Compression” on page 868 or “Real-time Compression” on page 883.

**RAID**

See “Redundant Array of Independent Disks” on page 883.

**Raw capacity**

See “Capacity” on page 865.

**Real capacity**

Real capacity is the amount of storage that is allocated to a volume copy from a storage pool. See also “Capacity” on page 865.

**Real-time Compression**

IBM Real-time Compression (RtC) is an IBM integrated software function for storage space efficiency. RACE compresses data on volumes in real time with minimal effect on performance.

**Redundant Array of Independent Disks**

RAID refers to two or more physical disk drives that are combined in an array in a certain way, which incorporates a RAID level for failure protection or better performance. The most common RAID levels are 0, 1, 5, 6, and 10. Some storage administrators refer to the RAID group as traditional RAID (TRAIID). For distributed RAID (DRAID) see “Distributed Redundant Array of Independent Disks” on page 870.

**RAID 0**

A data striping technique, which is commonly called RAID Level 0 or RAID 0 because of its similarity to common, RAID, data-mapping techniques. It includes no data protection, however, so, strictly speaking, the appellation RAID is a misnomer. RAID 0 is also known as data striping.

**RAID 1**

RAID 1 is a mirroring technique that is used on a storage array in which two or more identical copies of data are maintained on separate mirrored disks.

**RAID 10**

A collection of two or more physical drives that present to the host an image of one or more drives. In the event of a physical device failure, the data can be read or regenerated from the other drives in the RAID due to data redundancy.

**RAID 5**

RAID 5 is an array that has a data stripe, which includes a single logical parity drive. The parity check data is distributed across all the disks of the array.

## **RAID 6**

RAID 6 is a RAID level that has two logical parity drives per stripe, which are calculated with different algorithms. Therefore, this level can continue to process read and write requests to all of the array's virtual disks (VDisks) in the presence of two concurrent disk failures.

## **RAID controller**

See "Node canister" on page 880.

## **Read-intensive drives**

The read-intensive flash drives that are available on Storwize V7000 Gen2, Storwize V5000 Gen2, and IBM SAN Volume Controller 2145-DH8, SV1, and 24F enclosures are one Drive Write Per Day (DWPD) read-intensive drives.

## **Real capacity**

The amount of storage that is allocated to a volume copy from a storage pool.

## **Rebuild area**

Reserved capacity that is distributed across all drives in a redundant array of drives. If a drive in the array fails, the lost array data is systematically restored into the reserved capacity, returning redundancy to the array. The duration of the restoration process is minimized because all drive members simultaneously participate in restoring the data. See also "Distributed Redundant Array of Independent Disks" on page 870.

## **Reclaimable (or reclaimed) capacity**

Reclaimable Data is the capacity that is no longer needed. Reclaimable capacity is created when data is overwritten and the new data is stored in a new location, when data is marked as unneeded by a host using the SCSI `unmap` command, or when a volume is deleted.

## **Recovery key**

See "Encryption recovery key" on page 871.

## **Redundant storage area network**

A redundant SAN is a SAN configuration in which there is no single point of failure (SPOF). Therefore, data traffic continues no matter what component fails. Connectivity between the devices within the SAN is maintained (although possibly with degraded performance) when an error occurs. A redundant SAN design is normally achieved by splitting the SAN into two independent counterpart SANs (two SAN fabrics). In this configuration, if one path of the counterpart SAN is destroyed, the other counterpart SAN path keeps functioning. See also "Counterpart SAN" on page 868.

## **Relationship**

In MM or GM, a relationship is the association between a master volume and an auxiliary volume. These volumes also have the attributes of a primary or secondary volume. See also "Auxiliary volume" on page 864, "Master volume" on page 879, "Primary volume" on page 882, "Secondary volume".

## **Reliability, availability, and serviceability**

Reliability, availability, and serviceability (RAS) are a combination of design methodologies, system policies, and intrinsic capabilities that, when taken together, balance improved hardware availability with the costs that are required to achieve it.

Reliability is the degree to which the hardware remains free of faults. Availability is the ability of the system to continue operating despite predicted or experienced faults. Serviceability is how efficiently and nondisruptively broken hardware can be fixed.

### **Remote Copy**

- ▶ See “Global Mirror” on page 874.
- ▶ See “Metro Mirror” on page 879
- ▶ See “Master volume” on page 879

### **Remote fabric**

The remote fabric is composed of SAN components (switches, cables, and so on) that connect the components (nodes, hosts, and switches) of the remote cluster together. Significant distances can exist between the components in the local cluster and those components in the remote cluster.

### **Remote Support Server and Client**

Remote Support Client is a software toolkit that resides in the SAN Volume Controller and opens a secured tunnel to the Remote Support Server. Remote Support Server resides in the IBM network and collects key health check and troubleshooting information that is required by IBM Support personnel.

### **IBM SAN Volume Controller**

The SAN Volume Controller is an appliance that is designed for attachment to various host computer systems. The IBM Spectrum Virtualize is a software engine of SAN Volume Controller that performs block-level virtualization of disk storage.

### **SCSI initiator**

The system component that initiates communications with attached targets.

### **SCSI target**

A device that acts as a subordinate to a SCSI initiator and consists of a set of one or more LUs, each with an assigned LUN. The LUs on the SCSI target are typically I/O devices.

### **Secondary volume**

Pertinent to RC, the volume in a relationship that contains a copy of data written by the host application to the primary volume. See also “Relationship” on page 884.

### **Secure Copy Protocol**

The secure transfer of computer files between a local and a remote host or between two remote hosts, using the Secure Shell (SSH) protocol.

### **Secure Sockets Layer certificate**

Secure Sockets Layer (SSL) is the standard security technology for establishing an encrypted link between a web server and a browser. This link ensures that all data passed between the web server and browsers remain private. To be able to create an SSL connection, a web server requires an SSL Certificate.

### **IBM Security Key Lifecycle Manager**

IBM Security Key Lifecycle Manager (SKLM) centralizes, simplifies, and automates the encryption key management process to help minimize risk and reduce operational costs of encryption key management.

## **Sequential volume**

A volume that uses extents from a single MDisk.

## **Serial-attached SCSI**

SAS is a method that is used in accessing computer peripheral devices that employs a serial (one bit at a time) means of digital data transfer over thin cables. The method is specified in the American National Standard Institute standard called SAS. In the business enterprise, SAS is useful for access to mass storage devices, particularly external HDDs.

## **Service Location Protocol**

Service Location Protocol (SLP) is an Internet service discovery protocol that enables computers and other devices to find services in a local area network (LAN) without prior configuration. It was defined in the RFC 2608.

## **Simple Network Management Protocol**

A set of protocols for monitoring systems and devices in complex networks. Information about managed devices is defined and stored in a Management Information Base (MIB).

## **Small Computer System Interface**

Small Computer System Interface (SCSI) is an ANSI-standard electronic interface with which PCs can communicate with peripheral hardware, such as disk drives, tape drives, CD-ROM drives, printers, and scanners, faster and more flexibly than with previous interfaces.

## **Snapshot**

A snapshot is an image backup type that consists of a PiT view of a volume.

## **Solid-state drive**

A solid-state drive (SSD) or Flash Disk is a disk that is made from solid-state memory and therefore has no moving parts. Most SSDs use NAND-based flash memory technology. It is defined to the SAN Volume Controller as a disk tier generic\_ssd.

## **Space efficient**

See “Thin provisioning” on page 888.

## **Spare**

An extra storage component, such as a drive or tape, that is predesignated for use as a replacement for a failed component.

## **Spare goal**

The optimal number of spares that are needed to protect the drives in the array from failures. The system logs a warning event when the number of spares that protect the array drops below this number.

## **Space-efficient volume**

For more information about a space-efficient volume, see “Thin-provisioned volume” on page 888.

## **Stand-alone relationship**

In FlashCopy, MM, and GM, relationships that do not belong to a consistency group and that have a null consistency-group attribute.



## **Statesave**

Binary data collection that is used for problem determination by service support.

## **Storage area network**

A SAN is a dedicated storage network that is tailored to a specific environment, which combines servers, systems, storage products, networking products, software, and services.

## **Storage-class memory**

Storage-class memory (SCM) is a type of NAND flash that includes a power source to ensure that data won't be lost due to a system crash or power failure. SCM treats non-volatile memory as DRAM and includes it in the memory space of the server. Access to data in that space is significantly quicker than access to data in local, PCI-connected solid state drives (SSDs), direct-attached HDDs or external storage arrays. SCM read/write technology is up to 10 times faster than NAND flash drives and is more durable.

SCM drives can be installed only in drives slots 21 - 24 in an IBM Storage System. The highest capacity drive must be installed in the highest available drive slot.

## **Storage Capacity Unit**

Storage Capacity Unit (SCU) is a SAN Volume Controller license metric that measures the managed capacity in a way that the price is differentiated by the technology used to store the data.

## **Storage node**

A component of a storage system that provides internal storage or a connection to one or more external storage systems.

## **Storage pool (managed disk group)**

A storage pool is a collection of storage capacity, which is made up of MDisks, that provides the pool of storage capacity for a specific set of volumes. A storage pool can contain more than one tier of disk, which is known as a multitier storage pool and a prerequisite of Easy Tier automatic data placement. Before SAN Volume Controller V6.1, this storage pool was known as a managed disk group (MDG).

## **Stretched system**

A stretched system is an extended HA method that is supported by SAN Volume Controller to enable I/O operations to continue after the loss of half of the system. A stretched system is also sometimes referred to as a split system. One half of the system and I/O group is usually in a geographically distant location from the other, often 10 kilometers (6.2 miles) or more. A third site is required to host a storage system that provides a quorum disk.

## **Striped**

Pertaining to a volume that is created from multiple MDisks that are in the storage pool. Extents are allocated on the MDisks in the order specified.

## **Support Assistance**

A function that is used to provide support personnel remote access to the system to perform troubleshooting and maintenance tasks.

## **Symmetric virtualization**

Symmetric virtualization is a virtualization technique in which the physical storage, in the form of a RAID, is split into smaller chunks of storage known as extents. These extents are then concatenated, by using various policies, to make volumes. See also “Asymmetric virtualization” on page 864.

## **Synchronous replication**

Synchronous replication is a type of replication in which the application write operation is made to both the source volume and target volume before control is given back to the application. See also “Asynchronous replication” on page 864.

## **Syslog**

A standard for transmitting and storing log messages from many sources to a centralized location to enhance system management.

## **Tie-break**

In a case of a cluster split in 2 groups of nodes, tie-break is a role of a quorum device used to decide which group continues to operate as the system, handling all I/O requests.

## **Thin-provisioned volume**

A thin-provisioned volume is a volume that allocates storage when data is written to it.

## **Thin provisioning**

Thin provisioning refers to the ability to define storage, usually a storage pool or volume, with a “logical” capacity size that is larger than the actual physical capacity that is assigned to that pool or volume. Therefore, a thin-provisioned volume is a volume with a virtual capacity that differs from its real capacity. Before SAN Volume Controller V6.1, this thin-provisioned volume was known as *space efficient*.

## **Thin provisioning savings**

See “Capacity” on page 865.

## **Throttles**

Throttling is a mechanism to control the amount of resources that are used when the system is processing I/Os on supported objects. The system supports throttles on hosts, host clusters, volumes, copy offload operations, and storage pools. If a throttle limit is defined, the system either processes the I/O for that object, or delays the processing of the I/O to free resources for more critical I/O operations.

## **Throughput**

A measure of the amount of information transmitted over a network in a given period of time. Throughput is generally measured in bits per second (bps), kilobits per second (Kbps), or megabits per second (Mbps).

## **Transparent Cloud Tiering**

Transparent Cloud Tiering (TCT) is a separately installable feature of IBM Spectrum Scale that provides a native cloud storage tier.

## **Trial License**

A temporary entitlement to use a licensed function.

## **Total capacity savings**

See “Capacity” on page 865.

## **T10 DIF**

T10 DIF is a Data Integrity Field (DIF) extension to SCSI to enable end-to-end protection of data from host application to physical media.

## **Unconfigured mode**

An access mode in which an external storage MDisk is not configured in the system, so no operations can be performed. See also “Image mode” on page 876 and “Managed mode” on page 879.

## **Unique identifier**

A unique identifier (UID) is an identifier that is assigned to storage-system LUs when they are created. It is used to identify the LU regardless of the LUN, the status of the LU, or whether alternative paths exist to the same device. Typically, a UID is used only once.

## **VDisk**

See “Virtual disk”.

## **VDisk-to-host mapping**

See “Host mapping” on page 875.

## **Virtualization**

In the storage industry, virtualization is a concept in which a pool of storage is created that contains several storage systems. Storage systems from various vendors can be used. The pool can be split into volumes that are visible to the host systems that use them.

## **Virtualized storage**

Virtualized storage is physical storage that has virtualization techniques applied to it by a virtualization engine.

## **Virtual capacity**

The amount of storage that is available. In a thin-provisioned volume, the virtual capacity can be different from the real capacity. In a standard volume, the virtual capacity and real capacity are the same.

## **Virtual disk**

See “Volume”.

## **Virtual local area network**

Virtual local area network (VLAN) tagging separates network traffic at the layer 2 level for Ethernet transport. The system supports VLAN configuration on both Internet Protocol Version 4 (IPv4) and Internet Protocol Version 6 (IPv6) connections.

## **Virtual storage area network**

A virtual storage area network (VSAN) is a logical fabric entity defined within the SAN. It can be defined on a single physical SAN switch or across multiple physical switched or directors. In VMware terminology, the VSAN is defined as a logical layer of storage capacity built from physical disk drives attached directly into the ESXi hosts. This solution is not considered for the scope of our publication.

## **Vital product data**

Vital product data (VPD or VDP) is information that uniquely defines system, hardware, software, and microcode elements of a processing system.

## **Volume**

A volume is a SAN Volume Controller logical device that appears to host systems that are attached to the SAN as a SCSI disk. Each volume is associated with exactly one I/O group. A volume has a preferred node within the I/O group. Before SAN Volume Controller 6.1, this volume was known as a VDisk.

## **Volume copy**

A volume copy is a physical copy of the data that is stored on a volume. Mirrored volumes have two copies. Non-mirrored volumes have one copy.

## **Volume protection**

To prevent active volumes or host mappings from inadvertent deletion, the system supports a global setting that prevents these objects from being deleted if the system detects that they have recent I/O activity. When you delete a volume, the system checks to verify whether it is part of a host mapping, FlashCopy mapping, or remote-copy relationship. In these cases, the system fails to delete the volume, unless the **-force** parameter is specified. Using the **-force** parameter can lead to unintentional deletions of volumes that are still active. Active means that the system detected recent I/O activity to the volume from any host.

## **Volume snapshot**

A collection of objects on a cloud storage account that represents the data of a volume at a particular time.

## **Worldwide ID**

A worldwide ID (WWID) is a name identifier that is unique worldwide and that is represented by a 64-bit value that includes the IEEE-assigned organizationally unique identifier (OUI).

## **Worldwide name**

A worldwide name (WWN) is a 64-bit, unsigned name identifier that is unique.

## **Worldwide node name**

A worldwide node name is a unique 64-bit identifier for a host containing a FC port. See also "Worldwide port name".

## **Worldwide port name**

A worldwide port name (WWPN) is a unique 64-bit identifier associated with a FC adapter port. The WWPN is assigned in an implementation-independent and protocol-independent manner. See also "Worldwide node name".

**Write-through mode**

Write-through mode is a process in which data is written to a storage device at the same time that the data is cached.

**Written capacity**

See “Capacity” on page 865.

**Zoning**

The grouping of multiple ports to form a virtual, private, storage network. Ports that are members of a zone can communicate with each other, but are isolated from ports in other zones. See also “Fibre Channel” on page 872.



# Abbreviations and acronyms

<b>AD</b>	Active Directory	<b>ETS</b>	Enhanced Transmission Selection
<b>AES</b>	Advanced Encryption Standard	<b>EUI</b>	extended-unique identifier
<b>AI</b>	artificial intelligence	<b>FC</b>	Fibre Channel
<b>ANSI</b>	American National Standards Institute	<b>FC-NVMe</b>	NVMe over Fibre Channel
<b>API</b>	application programming interface	<b>FCIP</b>	Fibre Channel over IP
<b>ASCII</b>	American Standard Code for Information Interchange	<b>FCM</b>	FlashCore Module
<b>ATM</b>	asynchronous transfer mode	<b>FCoE</b>	Fibre Channel over Ethernet
<b>BBU</b>	battery backup unit	<b>FCP</b>	Fibre Channel Protocol
<b>CA</b>	certificate authority	<b>FICON</b>	Fibre Channel Connection
<b>CHAP</b>	Challenge Handshake Authentication Protocol	<b>FRU</b>	field-replaceable unit
<b>CIMOM</b>	Common Information Model Object Manager	<b>FTP</b>	File Transfer Protocol
<b>CLI</b>	command-line interface	<b>GA</b>	generally available
<b>COTS</b>	commercial off the shelf	<b>GbE</b>	Gigabit Ethernet
<b>CoW</b>	Copy on Write	<b>GBICs</b>	gigabit interface converters
<b>CSP</b>	cloud service provider	<b>GM</b>	Global Mirror
<b>CRU</b>	customer-replaceable unit	<b>GMCV</b>	Global Mirror with Change Volumes
<b>CSV</b>	comma-separated value	<b>GPFS</b>	General Parallel File System
<b>DAP</b>	Directory Access Protocol	<b>HA</b>	high availability
<b>DCBx</b>	Data Center Bridging Exchange	<b>HBA</b>	host bus adapter
<b>DHCP</b>	Dynamic Host Configuration Protocol	<b>HDD</b>	hard disk drive
<b>DIF</b>	Data Integrity Field	<b>HIC</b>	host interface card
<b>DIMM</b>	Dual Inline Memory Module	<b>HSM</b>	hardware security module
<b>DMP</b>	Directed Maintenance Procedure	<b>laaS</b>	infrastructure as a service
<b>DNS</b>	Domain Name System	<b>IBM</b>	International Business Machines Corporation
<b>DR</b>	disaster recovery	<b>IBM SSR</b>	IBM System Services Representative
<b>DRAID</b>	Distributed RAID	<b>IDA</b>	Information Dispersal Algorithms
<b>DRAID 6</b>	Distributed RAID 6	<b>IOPS</b>	input/output operations per second
<b>DRET</b>	Data Reduction Estimation Tool	<b>IPv4</b>	Internet Protocol Version 4
<b>DRP</b>	Data Reduction Pool	<b>IPv6</b>	Internet Protocol Version 6
<b>DSFA</b>	Data Storage Feature Activation	<b>IQN</b>	iSCSI Qualified Name
<b>DWPD</b>	Drive Write Per Day	<b>iSCSI</b>	Internet Small Computer Systems Interface
<b>ECS</b>	Enterprise Class Support	<b>iSER</b>	iSCSI Extensions for RDMA
<b>ECuRep</b>	Enhanced Customer Data Repository	<b>ISL</b>	inter-switch link
<b>EMEA</b>	Europe, Middle East, and Africa	<b>iSNS</b>	internet Storage Name Service
<b>ENT</b>	enterprise	<b>iWARP</b>	internet Wide-area RDMA Protocol
<b>ESC</b>	Enhanced Stretched Cluster	<b>JBOD</b>	Just A Bunch Of Disks
		<b>JRE</b>	Java Runtime Environment
		<b>KB</b>	kilobytes

<b>Kbps</b>	kilobits per second	<b>PaaS</b>	platform as a service
<b>KiB</b>	kibibyte	<b>PDU</b>	power distribution unit
<b>KMIP</b>	Key Management Interoperability Protocol	<b>PFC</b>	priority flow control
<b>LAN</b>	local area network	<b>PiT</b>	point-in-time
<b>LBA</b>	logical block address	<b>PMP</b>	Project Management Professional
<b>LDAP</b>	Lightweight Directory Access Protocol	<b>PMR</b>	Problem Management Report
<b>LED</b>	light-emitting diode	<b>POST</b>	power-on self-test
<b>LRU</b>	least recently used	<b>PPK</b>	PuTTY private key
<b>LSA</b>	log structured array	<b>PSU</b>	power supply unit
<b>LU</b>	logical unit	<b>QoS</b>	quality of service
<b>LUN</b>	logical unit number	<b>RACE</b>	Random Access Compression Engine
<b>LVM</b>	Logical Volume Manager	<b>RAID</b>	Redundant Array of Independent Disks
<b>LW</b>	long wave	<b>RAS</b>	reliability, availability, and serviceability
<b>LZ</b>	Lempel-Ziv	<b>RC</b>	Remote Copy
<b>MAC</b>	Media Access Control	<b>RDMA</b>	Remote Direct Memory Access
<b>Mb</b>	megabits	<b>RFC</b>	Request for Comments
<b>MBps</b>	megabytes per second	<b>RoCE</b>	RDMA over Converged Ethernet
<b>Mbps</b>	megabits per second	<b>RPM</b>	revolutions per minute
<b>MDG</b>	managed disk group	<b>RPO</b>	recovery point objective
<b>MDisk</b>	managed disk	<b>RtC</b>	Real-time Compression
<b>MIB</b>	Management Information Base	<b>RTT</b>	round-trip time
<b>MiB</b>	mebibytes	<b>SAN</b>	storage area network
<b>MLC</b>	multi-level cell	<b>SAS</b>	serial-attached SCSI
<b>MM</b>	Metro Mirror	<b>SAT</b>	Service Assistant Tool
<b>MSCS</b>	Microsoft Cluster Server or Microsoft Clustering Service	<b>SATA</b>	Serial Advanced Technology Attachment
<b>MT</b>	machine type	<b>SCM</b>	storage-class memory
<b>MTBF</b>	mean time between failures	<b>SCP</b>	Secure Copy Protocol
<b>MTM</b>	machine type and model	<b>SCSI</b>	Small Computer System Interface
<b>MTU</b>	maximum transmission unit	<b>SCU</b>	Storage Capacity Unit
<b>NAA</b>	Network Address Authority	<b>SDD</b>	Subsystem Device Driver
<b>NDVM</b>	Non-Disruptive Volume Move	<b>SDDPCM</b>	Subsystem Device Driver Path Control Module
<b>NIC</b>	network interface controller	<b>SEM</b>	Secondary Expander Module
<b>NL</b>	nearline	<b>SFP</b>	Small Form-factor Pluggable
<b>NPIV</b>	N_Port ID Virtualization	<b>SKLM</b>	Security Key Lifecycle Manager
<b>NQN</b>	NVMe Qualified Name	<b>SLA</b>	service level agreement
<b>NTP</b>	Network Time Protocol	<b>SLC</b>	single-level cell
<b>NVMe</b>	Non-Volatile Memory Express	<b>SLP</b>	Service Location Protocol
<b>NVMe-oF</b>	NVMe over Fabric	<b>SME</b>	subject matter expert
<b>OBAC</b>	object-based access control	<b>SMTF</b>	Simple Mail Transfer Protocol
<b>ODX</b>	Offloaded Data Transfer	<b>SNIA</b>	Storage Networking Industry Association
<b>OS</b>	operating system		
<b>OUI</b>	organizationally unique identifier		



<b>SNMP</b>	Simple Network Management Protocol
<b>SPOF</b>	single point of failure
<b>SSD</b>	solid-state drive
<b>SSH</b>	Secure Shell
<b>SSIC</b>	IBM System Storage Interoperation Center
<b>SSL</b>	Secure Sockets Layer
<b>STAT</b>	Storage Tier Advisor Tool
<b>T0</b>	time-zero
<b>TCT</b>	Transparent Cloud Tiering
<b>TPGS</b>	Target Port Group Support
<b>TRAIID</b>	traditional RAID
<b>UDID</b>	unit device identifier
<b>UID</b>	unique identifier
<b>VC</b>	virtual connections
<b>VDisk</b>	virtual disk
<b>VIOS</b>	Virtual I/O Server
<b>VLAN</b>	virtual local area network
<b>VSAN</b>	virtual storage area network
<b>VVOLs</b>	VMware vSphere Virtual Volumes
<b>WWID</b>	worldwide ID
<b>WWN</b>	worldwide name
<b>WWNN</b>	worldwide node name
<b>WWPN</b>	worldwide port name



# Related publications

The publications that are listed in this section are considered particularly suitable for a more detailed discussion of the topics covered in this book.

## IBM Redbooks

The following IBM Redbooks publications provide more information about the topic in this document. Note that some publications that are referenced in this list might be available in softcopy only:

- ▶ *Implementing the IBM Storwize V7000 with IBM Spectrum Virtualize V8.2.1*, SG24-7938
- ▶ *Implementing the IBM System Storage SAN Volume Controller with IBM Spectrum Virtualize V8.2.1*, SG24-7933
- ▶ *Implementing the IBM Storwize V5000 Gen2 (including the Storwize V5010, V5020, and V5030) with IBM Spectrum Virtualize V8.2.1*, SG24-8162
- ▶ *IBM FlashSystem 9100 Architecture, Performance, and Implementation*, SG24-8425
- ▶ *IBM FlashSystem 5000 Family Products*, SG24-8449
- ▶ *Introduction and Implementation of Data Reduction Pools and Deduplication*, SG24-8430
- ▶ *IBM System Storage SAN Volume Controller, IBM Storwize V7000, and IBM FlashSystem 7200 Best Practices and Performance Guidelines*, SG24-7521
- ▶ *IBM FlashSystem 9200 and 9100 Best Practices and Performance Guidelines*, SG24-8448

You can search for, view, download, or order these documents and other Redbooks, Redpapers, Web Docs, draft, and additional materials, at the following website:

[ibm.com/redbooks](https://ibm.com/redbooks)

## Help from IBM

IBM Support and downloads

[ibm.com/support](https://ibm.com/support)

IBM Global Services

[ibm.com/services](https://ibm.com/services)





# Implementing the IBM System Storage SAN Volume Controller

SG24-8465-00

ISBN 0738458902



(1.5" spine)  
1.5" x 1.998"  
789 <-> 1051 pages







SG24-8465-00

ISBN 0738458902

Printed in U.S.A.

Get connected

