**Technical white paper**

# HP Big Data Reference Architecture: Hortonworks Data Platform reference architecture implementation

## HP Big Data Reference Architecture with Hortonworks Data Platform for Apache Hadoop

# Table of contents

**Hortonworks®**

# Executive summary

HP Big Data Reference Architecture (BDRA) is a modern architecture for the deployment of big data solutions. It is designed to improve access to big data, rapidly deploy big data solutions, and provide the flexibility needed to optimize the infrastructure in response to ever-changing requirements in a Hadoop ecosystem.

HP BDRA challenges the conventional wisdom used in current big data solutions, where compute resources are typically co-located with storage resources on a single server node. Instead, HP BDRA utilizes a much more flexible approach that leverages an asymmetric cluster featuring de-coupled tiers of workload-optimized compute and storage nodes in conjunction with modern networking. The result is a big data platform that makes it easier to deploy, use, and scale data management applications.

As customers start to retain the dropped calls, mouse-clicks and other information that was previously discarded, they are adopting the use of single data repositories (also known as data lakes) to capture and store big data for later analysis. Thus, there is a growing need for a scalable, modern architecture for the consolidation, storage, access, and processing of big data. HP BDRA meets this need by offering an extremely flexible platform for the deployment of data lakes, while also providing access to a broad range of applications. Moreover, HP BDRA's tiered architecture allows organizations to grow specific elements of the platform without the chore of having to redistribute data.

Clearly, big data solutions are evolving from a simple model where each application was deployed on a dedicated cluster of identical nodes. By integrating the significant changes that have occurred in fabrics, storage, container-based resource management, and workload-optimized servers, HP has created the next-generation, data center-friendly, big data cluster.

This white paper describes the components and capabilities of the reference architecture, highlights recognizable benefits, and provides guidance on selecting the appropriate configuration for building a Hadoop cluster based on Hortonworks Data Platform (HDP) to meet particular business needs. In addition to simplifying the procurement process, the paper also provides guidelines for setting up HDP once the system has been deployed.

**Target audience:** This paper is intended for decision makers, system and solution architects, system administrators and experienced users that are interested in reducing design time for or simplifying the purchase of a big data architecture based on solution components provided by HP and Hortonworks. An intermediate knowledge of Apache Hadoop and scale-out infrastructure is recommended.

Testing performed in May – October 2014 is described.

# Introduction

HP BDRA is a reference architecture that is based on hundreds of man-hours of research and testing by HP engineers. To allow customers to deploy solutions based on this architecture, HP offers detailed Bills of Materials (BOMs) based on a proof of concept (refer to Appendix A – Bill of Materials  and Appendix B – Alternate parts for storage nodes). Recommended solution components – HP Software, HP Moonshot System components, HP ProLiant servers, and HP Networking switches – and their respective configurations have been carefully tested with a variety of I/O-, CPU-, network-, and memory-bound workloads. The resulting architecture provides optimal price/performance for Hortonworks Data Platform (HDP).

To facilitate installation, HP has developed a broad range of Intellectual Property (IP) that allows solutions to be implemented by HP or, jointly, with partners and customers. Additional assistance is available from HP Technical Services. For more information, refer to Appendix C – HP value added services and support.
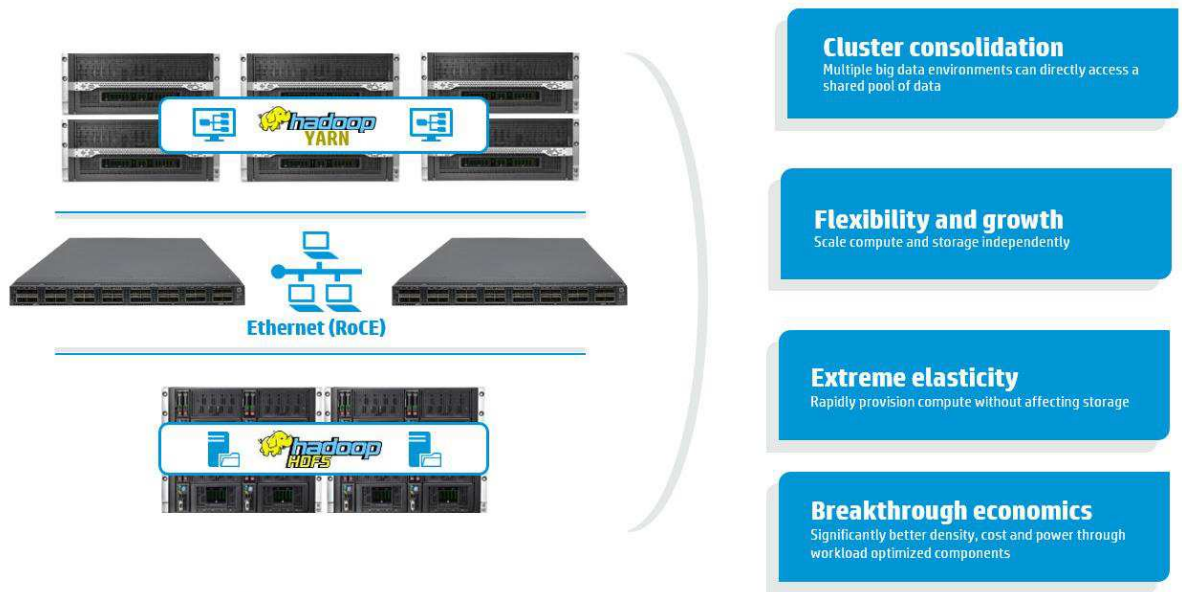
# HP BDRA overview

As companies grow their big data implementations they often find themselves deploying multiple clusters to support their needs. This could be to support different big data environments (Hadoop, NoSQLs, MPP DBMS's, etc.) with optimal hardware, to support rigid workload partitioning for departmental requirements or simply as a byproduct of multi-generational hardware. This often leads to data duplication and movement of large amounts of data between systems to accomplish an organization's business requirements. We find some customers searching for a way to recapture some of the traditional benefits of a converged infrastructure such as the ability to more easily share data between different applications running on different platforms, the ability to scale compute and storage separately and the ability to rapidly provision new servers without repartitioning data to achieve optimal performance.

To address these needs, HP engineers challenged the conventional wisdom that compute should always be collocated with data. While this approach works very well we believe that the power of taking the work to the data in Hadoop is that we are parallelizing the work into a pure shared nothing model where tasks run against specific slices of the data without the need for coordination such as distributed lock management or distributed cache management as you might have in older parallel database designs. Actually collocating the work with the data within the sheet metal of a computer chassis is irrelevant next to the power of the function shipping paradigm that Hadoop offers. In fact in a modern server we often have more network

bandwidth available to ship data off the server than we have bandwidth to disk. Given this and the dramatic increase in Ethernet networking bandwidth we decided to deploy the HP Big Data Reference Architecture as an asymmetric cluster with some nodes dedicated to compute and others dedicated to Hadoop Distributed File System (HDFS). Hadoop still works the same way – tasks still have complete ownership of their portion of the data (functions are still being shipped to the data) but the compute is executed on a node optimized for compute work and the file system is executed on a node that is optimized for file system. Of particular interest is that we have found that this approach can actually perform better than a traditional Hadoop cluster for several reasons. For more information on this architecture and the benefits that it provides see the overview master document at http://h20195.www2.hp.com/V2/GetDocument.aspx?docname=4AA5-6141ENW.

It is very important to understand that the HP Big Data Reference Architecture is a highly optimized configuration built using unique servers offered by HP – the HP ProLiant SL4540 for the high density storage layer, and HP Moonshot for the high density computational layer. It is also the result of a great deal of testing and optimization done by our HP engineers which resulted in the right set of software, drivers, firmware and hardware to yield extremely high density and performance. Simply deploying Hadoop onto a collection of traditional servers in an asymmetric fashion will not yield the kind of benefits that we see with our reference architecture. In order to simplify the build for customers, HP will provide the exact bill of materials in this document to allow a customer to purchase their cluster. Then, our Technical Services Consulting organization will provide a service that will install our prebuilt operating system images and verify all firmware and versions are correctly installed then run a suite of tests that verify that the configuration is performing optimally. Once this has been done, the customer is free to do a fairly standard Hortonworks installation on the cluster using the recommended guidelines in this document.

**Figure 1.** HP BDRA, changing the economics of work distribution in big data



The HP BDRA design is anchored by the following HP technologies:

- **Storage nodes** – HP ProLiant SL4540 Gen8 2 Node Servers (SL4540) make up the HDFS storage layer, providing a single repository for big data.
- **Compute nodes** – HP Moonshot System cartridges deliver a scalable, high-density layer for compute tasks and provide a framework for workload-optimization.

High-speed networking separates compute nodes and storage nodes, creating an asymmetric architecture that allows each tier to be scaled individually; there is no commitment to a particular CPU/storage ratio. Since big data is no longer co-located with storage, Hadoop does need to achieve node locality. However, rack locality works in exactly the same way as in a traditional converged infrastructure; that is, as long as you scale within a rack, overall scalability is not affected.

With compute and storage de-coupled, you can again enjoy many of the advantages of a traditional converged system. For example, you can scale compute and storage independently, simply by adding compute nodes or storage nodes. Testing carried out by HP indicates that most workloads respond almost linearly to additional compute resources.

**Hadoop YARN**

YARN is a key feature of the latest generation of Hadoop and of HP BDRA. It decouples MapReduce's resource management and scheduling capabilities from the data processing components, allowing Hadoop to support more varied processing approaches and a broader array of applications.

New to YARN is a concept of labels for groupings of compute nodes. Jobs submitted through YARN can now be flagged to perform their work on a particular set of nodes when the appropriate label name is included with the job. Thus, you can now create groups of compute resources that are designed, built, or optimized for particular types of computational work, allowing for jobs to be passed out to groups of hardware that are more geared to the type of computation work. In addition, the use of labels allows you to isolate compute resources that may be required to perform a high-priority job, for example, ensuring that sufficient resources are available at any given time.

# Benefits of the HP BDRA solution

While the most obvious benefits of the HP BDRA solution center on density and price/performance, other benefits include:

- **Elasticity** – HP BDRA has been designed for flexibility.
  - Compute nodes can be allocated very flexibly without redistributing data; for example, based on time-of-day or even a single job.
  - You are no longer committed to yesterday's CPU/storage ratios; now, you have much more flexibility in design and cost.
  - You only need to grow your system in areas that are needed.
  - YARN-compliant workloads access big data directly via HDFS; other workloads can access the same data via appropriate connectors
- **Consolidation** – HP BDRA is based on HDFS, which has enough performance and can scale to large enough capacities to be the single source for big data within any organization. You can consolidate the various pools of data currently being used in your big data projects into a single, central repository.
- **Workload-optimization** – There is no one go-to software for big data; instead there is a federation of data management tools. After selecting the appropriate tool for your requirements, you can then run your job using the compute node that is best optimized for the workload, such as a low-power Moonshot cartridge or a compute-intense cartridge.
- **Enhanced capacity management**
  - Compute nodes can be provisioned on the fly.
  - Storage nodes are a smaller subset of the cluster and, as such, are less costly to overprovision.
  - Managing a single data repository rather than multiple different clusters reduces overall management costs.
- **Faster time-to-solution**

  Processing big data typically requires the use of multiple data management tools. If these tools are deployed on their own Hadoop clusters with their own – often fragmented – copies of the data, time-to-solution can be lengthy. With HP BDRA, data is unfragmented, consolidated in a single data lake; and tools access the same data via YARN or a connector. Thus, there is more time spent on analysis, less on shipping data; time-to-solution is typically faster.
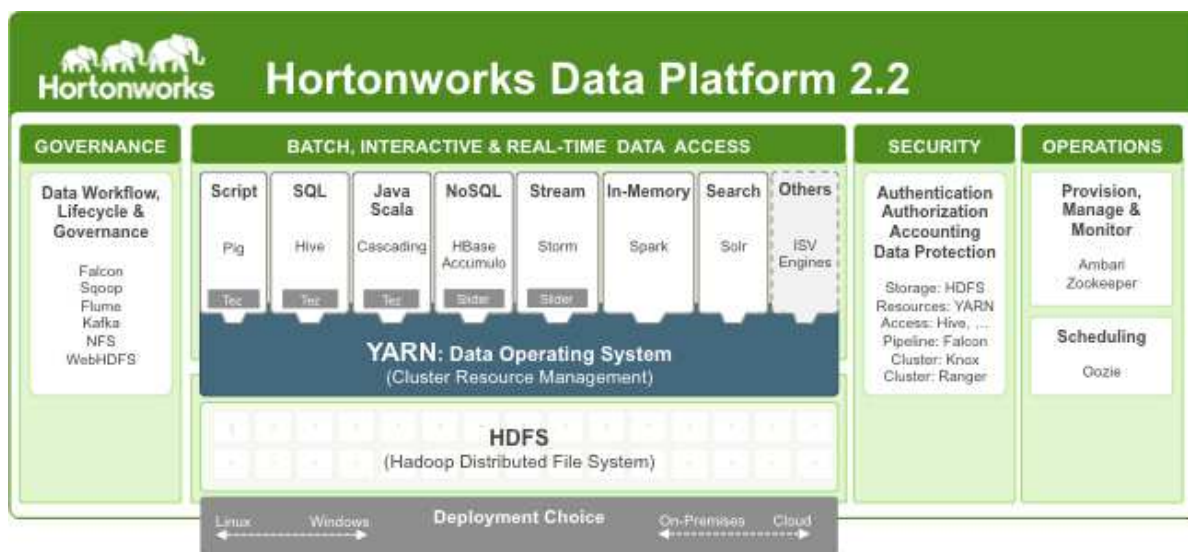
# Hortonworks Data Platform overview

Apache Hadoop is an open-source project administered by the Apache Software Foundation. Hadoop's contributors work for some of the world's biggest technology companies, creating a diverse, motivated community that has produced an innovative platform for consolidating, combining, and understanding big data.

While today's enterprises collect and generate more data than ever before, the focus of conventional relational and data warehouse products is on OLAP and OLTP, which utilize structured data. Hadoop, however, is designed to solve a different business problem — the fast, reliable analysis of both structured and complex data. As a result, many organizations are now deploying Hadoop alongside existing IT systems, thus combining semi-structured data and new data sets in powerful new ways.

HDP is a platform for multi-workload data processing; it can utilize a range of processing methods — from batch through interactive and real-time — all supported by solutions for governance, integration, security, and operations. HDP integrates with and augments solutions like HP BDRA, allowing you to maximize the value of big data, and it allows you to deploy Hadoop wherever you want — from the cloud, on-premises, across both Linux® and Microsoft® Windows®.

Figure 2 provides an overview of HDP.

**Figure 2.** HDP provides a blueprint for Enterprise Hadoop



HDP enables Enterprise Hadoop, a full suite of essential Hadoop capabilities in the following functional areas: data management, data access, data governance and integration, security, and operations.

Key highlights of HDP 2.2 include the following:

• Batch and interactive SQL queries via Apache Hive and Apache Tez, along with a cost-based optimizer powered by Apache Calcite

• High-performance ETL via Pig and Tez

• Stream processing via Apache Storm and Apache Kafka

• YARN labels

• Search via Apache Solr

• Streamlined cluster operations via Apache Ambari

• Data lifecycle management via Apache Falcon

• Perimeter security via Apache Knox

• Centralized security administration for HDFS, Hive, HBase, Storm, and Knox via Apache Ranger

For more information on HDP, refer to hortonworks.com/hdp.

# Solution components

Figure 3 provides a basic conceptual diagram of HP BDRA.

**Figure 3.** HP BDRA concept



For full BOM listings of products selected for the proof of concept, refer to Appendix A – Bill of Materials .

## Minimum configuration

Figure 4 shows a minimum HP BDRA configuration, with 90 worker nodes and six storage nodes housed in a single 42U rack.

**Figure 4.** Base HP BDRA configuration, with detailed lists of components



**Management Node**
**1x** HP ProLiant DL360p Gen8
2x E5-2650 v2 CPU, 8 cores each
128GB Memory 8x16GB 2Rx4 PC3-14900R-13
7.2TB – 8x HP 900GB 6G SAS 10K HDD
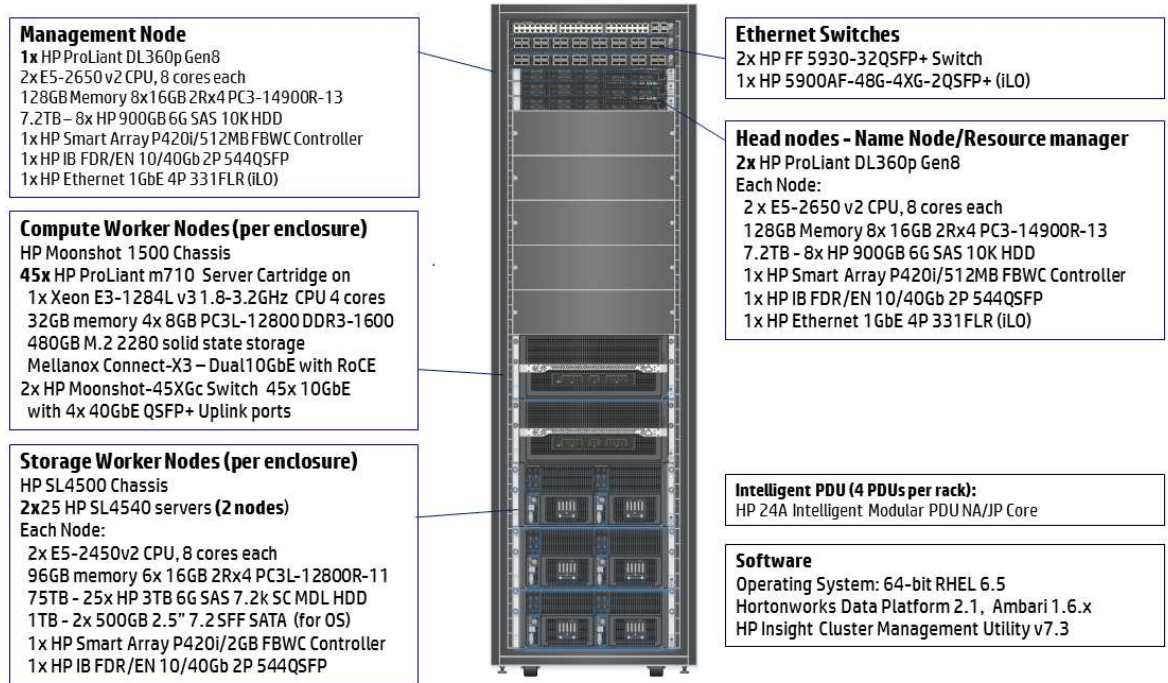1x HP Smart Array P420i/512MB FBWC Controller
1x HP IB FDR/EN 10/40Gb 2P 544QSFP
1x HP Ethernet 1GbE 4P 331FLR (iLO)

**Compute Worker Nodes (per enclosure)**
HP Moonshot 1500 Chassis
**45x** HP ProLiant m710  Server Cartridge on
  1x Xeon E3-1284L v3 1.8-3.2GHz  CPU 4 cores
  32GB memory 4x 8GB PC3L-12800 DDR3-1600
  480GB M.2 2280 solid state storage
  Mellanox Connect-X3 – Dual 10GbE with RoCE
2x HP Moonshot-45XGc Switch  45x 10GbE
  with 4x 40GbE QSFP+ Uplink ports

**Storage Worker Nodes (per enclosure)**
HP SL4500 Chassis
**2x25** HP SL4540 servers **(2 nodes)**
Each Node:
  2x E5-2450v2 CPU, 8 cores each
  96GB memory 6x 16GB 2Rx4 PC3L-12800R-11
  75TB - 25x HP 3TB 6G SAS 7.2k SC MDL HDD
  1TB - 2x 500GB 2.5" 7.2 SFF SATA (for OS)
  1x HP Smart Array P420i/2GB FBWC Controller
  1x HP IB FDR/EN 10/40Gb 2P 544QSFP

**Ethernet Switches**
2x HP FF 5930-32QSFP+ Switch
1x HP 5900AF-48G-4XG-2QSFP+ (iLO)

**Head nodes - Name Node/Resource manager**
**2x** HP ProLiant DL360p Gen8
Each Node:
  2 x E5-2650 v2 CPU, 8 cores each
  128GB Memory 8x 16GB 2Rx4 PC3-14900R-13
  7.2TB - 8x HP 900GB 6G SAS 10K HDD
  1x HP Smart Array P420i/512MB FBWC Controller
  1x HP IB FDR/EN 10/40Gb 2P 544QSFP
  1x HP Ethernet 1GbE 4P 331FLR (iLO)

**Intelligent PDU (4 PDUs per rack):**
HP 24A Intelligent Modular PDU NA/JP Core

**Software**
Operating System: 64-bit RHEL 6.5
Hortonworks Data Platform 2.1,  Ambari 1.6.x
HP Insight Cluster Management Utility v7.3

The following nodes are used in the base HP BDRA configuration:

- **Compute nodes** – The base HP BDRA configuration features two HP Moonshot 1500 Chassis containing a total of 90 m710 cartridge servers.
- **Storage nodes** – There are three SL4500 Chassis, each with two SL4540 servers, for a total of six storage nodes. Each node is configured with 25 disks; each typically runs DataNode.
- **Management/head nodes** – Three HP ProLiant DL360p Gen8 servers are configured as management/head nodes with the following functionality:
  - Management node, with Ambari, HP Insight Cluster Management Utility (Insight CMU), Hadoop ZooKeeper, and JournalNode
  - Head node, with active ResourceManager/standby NameNode, ZooKeeper, and JournalNode
  - Head node, with active NameNode/standby ResourceManager, ZooKeeper, and JournalNode

**Best Practice**
HP recommends starting with a minimum of two full Moonshot 1500 chassis, with 45 compute nodes on each. To provide high availability, you should start with a minimum of three SL4500 chassis, with a total of six storage nodes in order to provide the redundancy that comes with the default replication factor of three.

**Power and cooling**
When planning large clusters, it is important to properly manage power redundancy and distribution. To ensure servers and racks have adequate power redundancy, HP recommends that each chassis (Moonshot 1500 or SL4500) and each HP ProLiant DL360p Gen8 server should have a backup power supply and each rack should have at least two Power Distribution Units (PDUs).

There is additional cost associated with procuring redundant power supplies; however, the need for redundancy is less critical in larger clusters where the inherent redundancy within HDP ensures there would be less impact in the event of a failure.

---

**Best Practice**
For each chassis and HP ProLiant DL360p Gen8 server, HP recommends connecting each of the device's power supplies to a different PDU. Furthermore, PDUs should each be connected to a separate data center power line to protect the infrastructure from a line failure.

Distributing server power supply connections evenly to the PDUs while also distributing PDU connections evenly to data center power lines ensures an even power distribution in the data center and avoids overloading any single power line. When designing a cluster, check the maximum power and cooling that the data center can supply to each rack and ensure that the rack does not require more power and cooling than is available.

---

# Networking

Two IRF-bonded HP 5930-32QSFP+ switches are specified in each rack for high performance and redundancy. Each provides six 40GbE uplinks that can be used to connect to the desired network or, in a multi-rack configuration, to another pair of HP 5930-32QSFP+ switches that are used for aggregation. The other 22 ports are available for Hadoop nodes.

---

**Note**
IRF-bonding requires four 40GbE ports per switch, leaving six 40GbE ports on each switch for uplinks.

---

The three NICs in each of the two switch modules in a Moonshot 1500 chassis are LACP-bonded.

**iLO network**
A single HP 5900 switch is used exclusively to provide connectivity to HP Integrated Lights-Out (HP iLO) management ports, which run at or below 1GbE. The HP iLO network is used for system provisioning and maintenance.

**Cluster isolation and access configuration**
It is important to isolate the Hadoop cluster so that external network traffic does not affect the performance of the cluster. In addition, isolation allows the Hadoop cluster to be managed independently from its users, ensuring that the cluster administrator is the only person able to make changes to the cluster configuration.

Thus, HP recommends deploying ResourceManager, NameNode, and Worker nodes on their own private Hadoop cluster subnet.

---

**Key Point**
Once a Hadoop cluster is isolated, the users still need a way to access the cluster and submit jobs to it. To achieve this, HP recommends multi-homing the management node so that it can participate in both the Hadoop cluster subnet and a subnet belonging to users.

Ambari is a web application that runs on the management node, allowing users to manage and configure the Hadoop cluster and view the status of jobs without being on the same subnet – provided that the management node is multi-homed. Furthermore, this approach allows users to shell into the management node and run Apache Pig or Apache Hive command line interfaces and submit jobs to the cluster in that way.

**Staging data**

After the Hadoop cluster has been isolated on its own private network, you must determine how to access HDFS in order to ingest data. The HDFS client must be able to reach every Hadoop DataNode in the cluster in order to stream blocks of data on to the HDFS.

HP BDRA provides the following options for staging data:

- **Using the management node** – You can utilize the already multi-homed management server to stage data. HP recommends configuring this server with eight disks to provide a sufficient amount of disk capacity to provide a staging area for ingesting data into the Hadoop cluster from another subnet.
- **Using the ToR switches** – You can make use of open ports in the ToR switches. HP BDRA is designed such that if both NICs on each storage node are used, six NICs are used on each Moonshot 1500 Chassis (three ports from Switch A and three from Switch B), and two NICs are used on each management/head node, the remaining 40GbE ports on the ToR switches can be used by multi-homed systems outside the Hadoop cluster to move data into the cluster.

---

**Note**

The benefit of using dual-homed management node(s) to isolate in-cluster Hadoop traffic from the ETL traffic flowing to the cluster may be debated. This enhances security; however, it may result in ETL performance/connectivity issues, since relatively few nodes are capable of ingesting data. For example, you may wish to initiate Sqoop tasks on the compute nodes to ingest data from an external RDBMS in order to maximize the ingest rate. However, this approach requires worker nodes to be exposed to the external network to parallelize data ingestion, which is less secure.
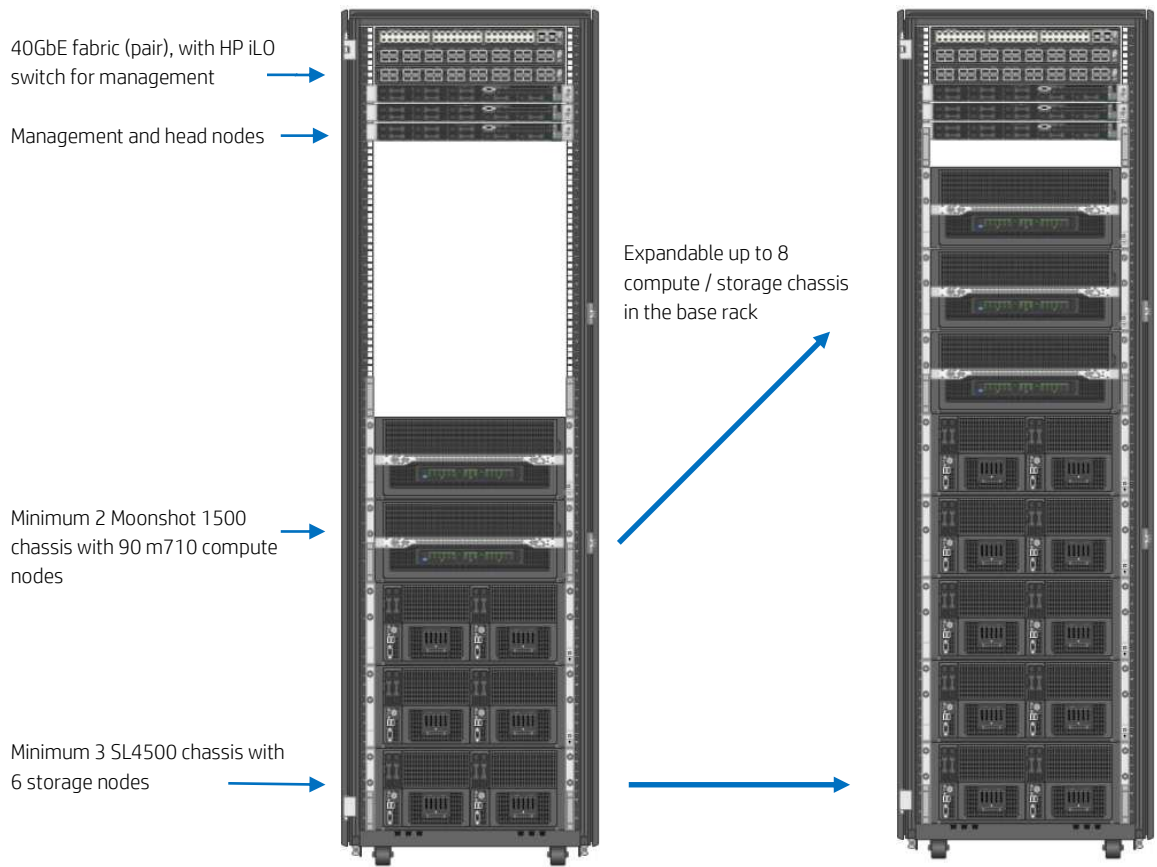
You should consider your options before committing to an optimal network design for your particular environment. HP recommends dual-homing the entire cluster, allowing every node to ingest data.

---

**Using WebHDFS** – You can leverage WebHDFS, which provides HTTP access to securely read and write data to and from HDFS. For more information on WebHDFS, refer to docs.hortonworks.com/HDPDocuments/HDP2/HDP-2.1.5/bk_system-admin-guide/content/sysadminguides_webhdfs.html.

## Expanding the base configuration

As needed, you can add compute and/or storage nodes to a base HP BDRA configuration. Figure 5 shows how the minimum HP BDRA configuration with two Moonshot 1500 chassis and three SL4500 chassis can expand to include a total of eight mixed chassis in a single rack.
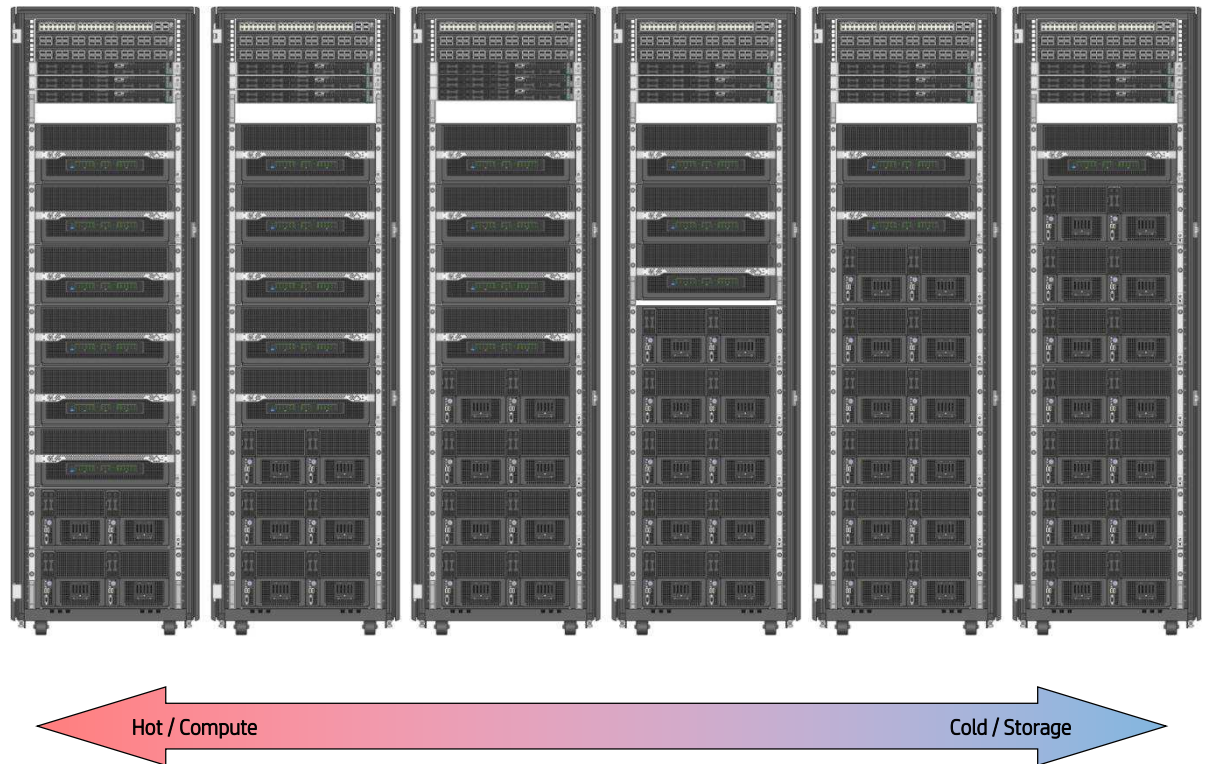
**Figure 5.** Minimum HP BDRA configuration, expanding to include a total of eight chassis

40GbE fabric (pair), with HP iLO switch for management

Management and head nodes

Expandable up to 8 compute / storage chassis in the base rack

Minimum 2 Moonshot 1500 chassis with 90 m710 compute nodes

Minimum 3 SL4500 chassis with 6 storage nodes

**Range of single-rack options**

As shown in Figure 6, you have a range of options for configuring a single-rack HP BDRA solution, ranging from hot (with a large number of compute nodes and minimal storage) to cold (with a large number of storage nodes and minimal compute).

**Figure 6.** HP BDRA compute/storage mix options
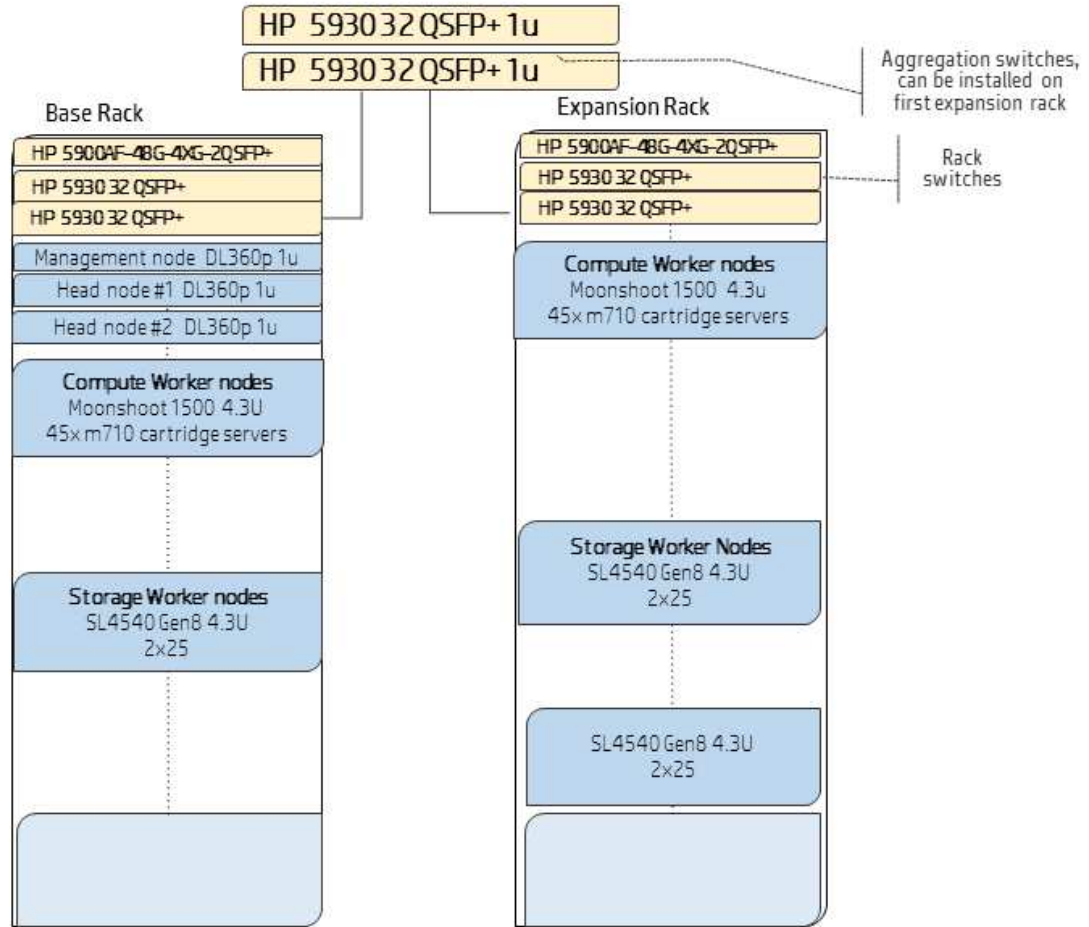


Hot / Compute        Cold / Storage

## Multi-rack configuration

The single-rack HP BDRA configuration is designed to perform well as a standalone solution but also form the basis of a much larger multi-rack solution, as shown in Figure 7. When moving from a single rack to a multi-rack solution, you simply add racks without having to change any components within the base rack.

A multi-rack solution assumes the base rack is already in place and extends its scalability. For example, the base rack already provides sufficient management services for a large scale-out system.

**Figure 7.** Multi-rack HP BDRA, extending the capabilities of a single rack



**Note**
While much of the architecture for the multi-rack solution was borrowed from the single-rack design, the architecture suggested here for multi-rack solutions is based on previous iterations of Hadoop testing on the HP ProLiant DL380e platform rather than cartridge servers. It is provided here as a general guideline for designing multi-rack Hadoop clusters.

**Extending networking in a multi-rack configuration**
For performance and redundancy, two HP 5930-32QSFP+ ToR switches are specified per expansion rack. The HP 5930-32QSFP+ switch includes up to six 40GbE uplinks that can be used to connect these switches to the desired network via a pair of HP 5930-32QSFP+ aggregation switches.

# Capacity and sizing

Hadoop cluster storage sizing requires careful planning, including the identification of current and future storage and compute needs.

## Guidelines for calculating storage needs

The following are general considerations for data inventory:

• Sources of data
• Frequency of data
• Raw storage
• Processed HDFS storage
• Replication factor

- Default compression turned on
- Space for intermediate files

To calculate your storage needs, determine the number of TB of data per day, week, month, and year; and then add the ingestion rates of all data sources.

It makes sense to identify storage requirements for the short-, medium-, and long-term.

Another important consideration is data retention – both size and duration. Which data must you keep? For how long?

In addition, consider maximum fill-rate and file system format space requirements on hard drives when estimating storage size.
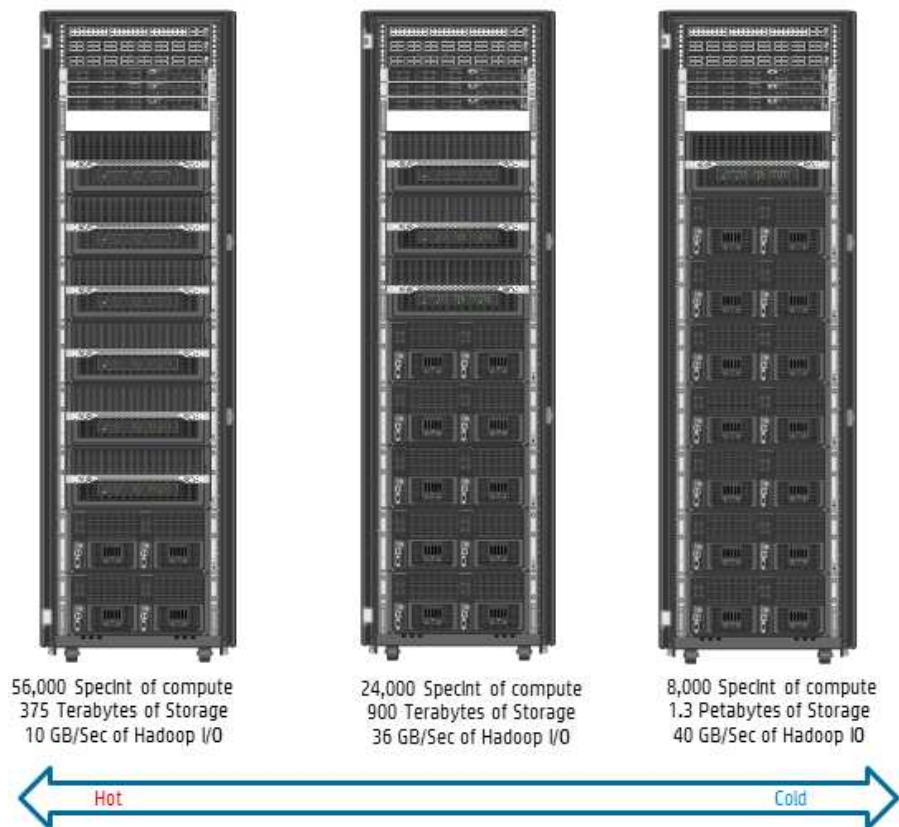
## HP BDRA cluster sizing

HP BDRA consists of two tiers of servers that deliver storage and compute resources for a modern Hadoop cluster; tiers are interconnected by high-performance HP networking components. De-coupling storage and compute resources allows each tier to be scaled independently, providing maximum flexibility in the design, procurement, and growth of the cluster, as shown in Figure 8.

Three different HP BDRA single-rack configurations are shown side-by-side. The racks can be characterized as follows:

- Left rack – This is a hot configuration, where compute resources are relatively high and storage capability is lower.
- Middle rack – This configuration has a balanced compute/storage mix.
- Right rack – This is a cold configuration, where storage resources are relatively high and compute capability is lower.

Thus, the flexibility of HP BDRA allows you to create a solution that meets your particular needs.

**Figure 8.** Three typical HP BDRA solutions, with estimated performance guidelines



56,000 SpecInt of compute
375 Terabytes of Storage
10 GB/Sec of Hadoop I/O

24,000 SpecInt of compute
900 Terabytes of Storage
36 GB/Sec of Hadoop I/O

8,000 SpecInt of compute
1.3 Petabytes of Storage
40 GB/Sec of Hadoop IO

Hot ⟷ Cold

Performance metrics listed for these HP BDRA configurations are based on a Hadoop cluster utilizing HP ProLiant DL380p Gen8 servers in the same footprint.

# Configuring an HP BDRA solution

The instructions provided here assume you have already created your Hadoop cluster on an HP BDRA solution. They are intended to help you:

1. Optimize settings for the Mellanox Unstructured Data Accelerator for Hadoop (UDA).
2. Set up HDP.

## Optimizing UDA settings

UDA is a software plug-in designed to accelerate Hadoop networks and improve the scaling of Hadoop clusters that execute intensive applications such as data analytics. UDA facilitates efficient data shuffles and merges over Mellanox InfiniBand and 10GbE/40GbE/56GbE RDMA over Converged Ethernet (RoCE) adapter cards. UDA is integrated into Apache Hadoop 2.0.x and 3.0.x; patches are available for all binary-compatible distributions. UDA installation is a simple RPM addition to the execution library.

UDA software is installed in /usr/lib64/uda and should be installed on every node, including the management node.

UDA documentation is available from the following sources:

• Mellanox Unstructured Data Acceleration (UDA) Release Notes Rev 3.4

https://www.dropbox.com/s/dp7onn0lhjefkxt/Mellanox_Unstructured_Data_Acceleration_%28UDA%29_Release_Notes_v3_4_0.pdf?dl=1

• Mellanox Unstructured Data Acceleration (UDA) Quick Start Guide Rev 3.4

https://www.dropbox.com/s/02lg4x16a60gx8g/Mellanox_Unstructured_Data_Acceleration_%28UDA%29_Quick_Start_Guide_v3_4_0.pdf?dl=1

This section provides general assistance for configuring Mellanox Ethernet cards and related software to work with Hadoop.

• **CompletedMaps** – This parameter gives all MapReduce resources (RAM/CPU) to mappers at the map phase and all MapReduce resources to reducers at the reduce phase. UDA's levitated merge causes reducers to wait until all mappers complete; thus, there is no need to give reducers resources while mappers are running and vice versa.

```
mapreduce.job.reduce.slowstart.completedmaps=1.00
```

• **Invoking the UDA shuffle** – The UDA plug-in for shuffles is not configured as the default shuffle plug-in. To invoke the UDA shuffle, you need to specify the plug-in when starting the job, as in the following example:

```
hadoop jar /usr/lib/hadoop-mapreduce/hadoop-mapreduce-examples.jar
terasort  -Dmapreduce.job.reduce.shuffle.consumer.plugin.class=
com.mellanox.hadoop.mapred.UdaShuffleConsumerPlugin  -
Dmapred.reduce.child.log.level=WARN -Dyarn.app.mapreduce.am.job.cbd-
mode.enable=false -Dyarn.app.mapreduce.am.job.map.pushdown=false -
Dmapreduce.job.maps=538 -Dmapreduce.job.reduces=178 /moonshot/terasort/input
/moonshot/terasort/output
```

• **Adding a symbolic link** – Unless you add a symbolic link to the UDA jar file, NodeManagers do not start. The following example must be updated whenever the Hadoop version changes unless you create a link that is common across releases pointing to the latest version of jar files:

```
ln -s /usr/lib64/uda/uda-hadoop-2.x.jar  /usr/lib/hadoop-mapreduce/uda-
hadoop-2.x.jar
```

## Setting up HDP

### Compute node components

Compute nodes contain the software shown in Table 1.

**Table 1.** Compute node base software components

| Software | Description |
|---|---|
| Red Hat® Enterprise Linux® (RHEL) 6.5 | Recommended operating system |
| Oracle Java Development Kit (JDK) 1.7.0_45 | JDK |
| NodeManager | NodeManager process for MR2/YARN |

Visit docs.hortonworks.com/HDPDocuments/HDP2/HDP-2.1.5/bk_system-admin-guide/content/admin_add-nodes-2.html for more information on the following topics:

- Manually installing and configuring NodeManager (or HBaseRegionServer) and DataNode
- Adding nodes

### Storage node components

Storage nodes contain the software shown in Table 2.

**Table 2.** Storage node base software components

| Software | Description |
|---|---|
| RHEL 6.5 | Recommended operating system |
| Oracle JDK 1.7.0_45 | JDK |
| DataNode | DataNode process for HDFS |

Visit docs.hortonworks.com/HDPDocuments/HDP2/HDP-2.1.5/bk_system-admin-guide/content/admin_add-nodes-2.html for more information on the following topics:

- Manually installing and configuring DataNode
- Adding nodes

### Configuration changes

Whenever possible, make the changes described here using Ambari.

*Hosts*

```
Memory Overcommit Validation Threshold .75     (default .80)
```

*HDFS*

Make the following changes to the HDFS configuration:

- Increase the dfs.blocksize value to allow more data to be processed by each map task, thus reducing the total number of mappers and NameNode memory consumption.
```
dfs.ha.fencing.method shell (true)
dfs.block.size, dfs.blocksize 256   (default 128)
```

- Increase the dfs.namenode.handler.count value to better manage multiple HDFS operations from multiple clients.
```
dfs.namenode.handler.count 120    (default 30)
```

- Increase the dfs.namenode.service.handler.count value to better manage multiple HDFS service requests from multiple clients.
```
dfs.namenode.service.handler.count 120 (default 30)
```

- Increase the Java heap size of the NameNode to provide more memory for NameNode metadata and ensure Maximum Process File Descriptors is set to 65,536:

```
# Maximum Process File descripts can be set in many services.
Maximum Process File Descriptors 65536

Java Heap Size of NameNode in Bytes 2138 Mib   (default 1 Gib)
Java Heap Size of Secondary NameNode in Bytes 2138 Mib  (default 1 Gib)
```

- Increase certain timeout values and numbers of retries for network-intensive environments. Make the following changes in the HDFS Service Advanced Configuration Snippet (Safety Valve) for hdfs-site.xml:

```
<property>
        <name>dfs.socket.timeout</name>
        <value>6000000</value>
</property>

<property>
        <name>dfs.datanode.socket.write.timeout</name>
        <value>6000000</value>
</property>

<property>
        <name>dfs.client.block.write.locateFollowingBlock.retries</name>
        <value>30</value>
</property>
```

*YARN*
Both yarn-site.xml and mapreduce-site.xml appear in several groups. The following values, which are used in all locations, are used to configure the UDA shuffle plug-in.

Make the following changes in the YARN Service Advanced Configuration Snippet (Safety Valve) for yarn-site.xml:

---

**Note**
The uda_shuffle can only be configured after UDA installation and configuration is complete. NodeManagers may fail to start if not properly set up.

---

- Replace the default MapReduce shuffle handler with the Mellanox UDA shuffle handler:

```
<property>
        <name>yarn.nodemanager.aux-services</name>
        <value>mapreduce_shuffle,uda_shuffle</value>
</property>
```

- Specify the Java class file for the Mellanox UDA shuffle handler:

```
<property>
        <name>yarn.nodemanager.aux-services.uda_shuffle.class</name>
        <value>com.mellanox.hadoop.mapred.UdaShuffleHandler</value>
</property>
```

- Set the block size for Snappy Compression, which is used in HP BDRA, helping to reduce the amount of I/O:

```
<property>
        <name>io.compression.codec.snappy.buffersize</name>
        <value>32768</value>
</property>
```

*Mellanox UDA shuffle handler*
Make the following changes in YARN Service MapReduce Advanced Configuration Snippet (Safety Valve). Refer to Optimizing UDA settings for more information.

- Specify Java class files for the Mellanox shuffle provider:

```
<property>
        <name>mapreduce.shuffle.provider.plugin.classes</name>
        <value>com.mellanox.hadoop.mapred.UdaShuffleProviderPlugin,org.apache
        .hadoop.mapred.TaskTracker$DefaultShuffleProvider</value>
</property>
```

- Replace the default MapReduce shuffle handler with the Mellanox shuffle handler:

```
<property>
        <name>mapreduce.job.shuffle.provider.services</name>
         <value>uda_shuffle</value>
</property>
```

- Specify the total amount of memory per reducer allocated for UDA's RDMA shuffle phase and levitated merge algorithm. More memory is preferable.

```
<property>
        <name>mapred.rdma.mem.use.contig.pages</name>
        <value>1</value>
</property>

<property>
        <name>mapred.rdma.shuffle.total.size</name>
        <value>4g</value>
</property>
```

- Specify the replication factor for the MapReduce job jar file. Whenever a client submits a job, the jar file is copied to many nodes (the default is 10) so that it is available to TaskTrackers and NodeManagers. However, since there are only six storage nodes, HP sets the replication factor to three to avoid under-replication errors:

```
<property>
        <name>mapred.submit.replication</name>
        <value>3</value>
</property>
```

*Optimizing compute nodes*
The changes described here apply to compute nodes, which have less memory than storage nodes.

- Specify the amount of memory for the Reduce task. Note that UDA needs a lot of memory because it works with RAM, not disk; as a result, Mellanox recommends configuring more memory per reducer but fewer reducers. Thus, HP increased the amount of memory from the default 1GB to 6GB:

```
yarn.app.mapreduce.am.resource.mb 2048 Mib  (default 1 Gib)
mapreduce.reduce.memory.mb 6144 Mib  (default 1 Gib)
```

- Specify the amount of memory for all YARN containers on a node. Since compute nodes each have 32GB of memory, HP recommends allocated 24GB for YARN:

```
yarn.nodemanager.resource.memory-mb 24 Gib   (default 8 Gib)
```

- Specify the number of virtual CPU cores for all YARN containers on a node; HP recommends using discretion:

```
yarn.nodemanager.resource.cpu-vcores 14   (default 8)
```

- Increase the minimum memory allocation by ResourceManager for each container from the default 1GB to 1,536MB to match the map task memory setting:

```
yarn.scheduler.minimum-allocation-mb 1536 Mib
```

- Reduce the maximum memory allocation by ResourceManager for each container from the default 64GB to 24GB:

```
yarn.scheduler.maximum-allocation-mb 24 Gib  (default 64 Gib)

#below three settings not valid for capacity scheduler
yarn.scheduler.fair.assignmultiple true  (default false)
yarn.scheduler.fair.locality-delay-node-ms 0
yarn.scheduler.fair.locality-delay-rack-ms 0
```

- Increase the maximum number of job counters to accommodate workloads such as Hive metrics:

```
mapreduce.job.counters.max 1000   (default 120)
```

- Give all MapReduce resources (RAM/CPU) to mappers during the map phase and to reducers during the reduce phase. Since UDA's levitated merge causes reducers to wait until all mappers complete, there is no need to give reducers resources while mappers are running and vice versa.

```
resource manager java heap size 4  GiB
```

*Compute node configuration*
The compute node configuration is as follows:

- System memory: 32GB

- Logical cores: 8

- Disks: 1

- Reserved memory: 12GB

- Available memory: 20GB

- Minimum container size: 2GB

- Number of containers: 10

- RAM per container: 2GB

Based on this configuration, the following settings are recommended in mapred-site.xml:

```
yarn.app.mapreduce.am.command-opts "-Xmx3277"
mapreduce.map.memory.mb 1536
mapreduce.reduce.memory.mb 6144
mapreduce.map.java.opts -Xmx1228
mapreduce.reduce.java.opts -Xmx4915
```

## Configuring compression

HP recommends using compression since it reduces file size on disk and speeds up disk and network I/Os. Set the following MapReduce parameter to compress job output files:

```
mapreduce.output.fileoutputformat.compress=true
```

# Summary

HP BDRA is a modern, flexible architecture for big data solutions that improves overall access to information and time-to-solution, while providing the inherent flexibility to support the rapidly changing requirements of big data applications. HP BDRA leverages HP's innovative big data building blocks of servers, storage, and networking, along with integrated management software and bundled support.

HP BDRA and HDP allow you to derive new business insights from big data by providing a platform to store, manage, and process data at scale. However, the design, procurement and deployment of a Hadoop cluster can be both complex and time consuming. Thus, this white paper outlines reference architectures for heterogeneous clusters of varying sizes with HDP 2.1 on HP infrastructure and management software. Guidelines for optimizing HDP software settings are provided.

## Advantages of HP BDRA

### Extremely elastic
Because HP BDRA is built around a two-tier system that de-couples compute and storage resources, the earlier commitment to fixed compute/storage ratios has been discarded; and tiers can now scale independently. You only need to grow your cluster where needed, thus maximizing solution density while minimizing Total Cost of Ownership (TCO). Moreover, nodes can be allocated by time-of-day or even for a single job without redistributing data.

### Consolidation
HP BDRA provides a distributed file system that has enough performance and can scale to large enough capacities to become the single source for big data within your organization. You no longer need to maintain multiple pools of data for various big data projects; now, you can consolidate that data into a single solution that can be shared by a broad range of data management tools. This consolidation can reduce your operational and management costs, while ensuring that the data does not become fragmented.

### Optimization and unification of big data projects
HP BDRA provides an ideal platform to deliver support for different types of workloads. Compute cores can be allocated as needed to provide better capacity management; and, since storage nodes are a smaller subset of the cluster, they are less costly to overprovision.

# Appendix A – Bill of Materials

Bills of Materials (BOMs) provided here are based on the tested configuration for a single-rack HP BDRA solution featuring the following key components:

- One management node
- Two head nodes
- 90 compute nodes
- Six storage nodes (refer to Appendix B – Alternate parts for storage nodes for other processor, memory, and hard drive options)
- Two ToR switches
- One HP iLO switch
- Hortonworks Data Platform (HDP 2.1 was tested)

**Note**
Part numbers are correct at the time of publication and subject to change.
BOMs do not include complete support options or support for non-North American rack and power requirements.
If you have questions on ordering, consult your HP Reseller or HP Sales Representative for more information.
Alternatively, to contact an HP sales expert, visit hp.com/large/contact/enterprise/index.html.

**Management node and head nodes**

**Important**
Table A-1 provides a BOM for one HP ProLiant DL360p Gen8 server. The tested solution featured one management server and two head nodes, requiring a total of three HP ProLiant DL360p Gen8 servers.

**Table A-1.** BOM for a single HP ProLiant DL360p Gen8 server

| Quantity | Part number | Description |
|----------|-------------|-------------|
| 1 | 654081-B21 | HP DL360p Gen8 8-SFF CTO Server |
| 1 | 712726-L21 | HP DL360p Gen8 E5-2650v2SDHS FIO Kit |
| 1 | 712726-B21 | HP DL360p Gen8 E5-2650v2SDHS Kit |
| 8 | 708641-B21 | HP 16GB 2Rx4 PC3-14900R-13 Kit |
| 8 | 652589-B21 | HP 900GB 6G SAS 10K 2.5in SC ENT HDD |
| 1 | 661069-B21 | HP 512MB FBWC for P-Series Smart Array |
| 1 | 649281-B21 | HP IB FDR/EN 10/40Gb 2P 544QSFP Adptr |
| 1 | 684208-B21 | HP Ethernet 1GbE 4P 331FLR FIO Adptr |
| 1 | 663201-B21 | HP 1U SFF BB Gen8 Rail Kit |
| 2 | 655874-B21 | HP QSFP/SFP+ Adaptor Kit |
| 2 | 656362-B21 | HP 460W CS Plat PL Ht Plg Pwr Supply Kit |
| 1 | C6N36A | HP Insight Control ML/DL/BL FIO Bndl Lic |
| 1 | G3J28AAE | RHEL Svr 2 Sckt/2 Gst 1yr 24x7 E-LTU |
| 2 | SG508A | HP C13 - C14 WW 250V 10Amp IPD 1.37m 1pc Jumper Cord |
| 3 | 498385-B23 | HP 3M 4X DDR/QDR QSFP IB Cu Cable |

**Compute nodes**

---

**Important**

Table A-2 provides a BOM for one HP Moonshot 1500 chassis with 45 m710 server cartridges. The tested solution featured two Moonshot chassis with a total of 90 server cartridges.

---

**Table A-2.** BOM for a single HP Moonshot chassis with 45 server cartridges

| Quantity | Part number | Description |
|---|---|---|
| 1 | 755371-B21 | HP Moonshot 1500 Chassis |
| 4 | 684532-B21 | HP 1500W Ht Plg Pwr Supply Kit |
| 2 | 704654-B21 | HP Moonshot-45XGc Switch Kit |
| 2 | 704652-B21 | HP Moonshot 4QSFP Uplink Kit |
| 45 | 755860-B21 | HP ProLiant m710 Server Cartridge |
| 45 | 765483-B21 | HP Moonshot 480GB M.2 2280 FIO Kit |
| 1 | 681254-B21 | HP 4.3U Rail Kit |
| 45 | C6N36AAE | HP Insight Control ML/DL/BL Bundle E-LTU |
| | C6N36A | HP Insight Control ML/DL/BL FIO Bndl Lic (optional if E-LTU is not available) |
| 45 | G3J28AAE | RHEL Svr 2 Sckt/2 Gst 1yr 24x7 E-LTU |

**Storage nodes**

---

**Important**

Table A-3 provides a BOM for one HP SL4500 chassis with two SL4540 servers. The tested solution featured three chassis and a total of six servers.

---

**Table A-3.** BOM for a single HP SL4500 chassis with two nodes

| Quantity | Part number | Description |
|---|---|---|
| 1 | 663600-B22 | HP 2xSL4500 Chassis |
| 4 | 512327-B21 | HP 750W CS Gold Ht Plg Pwr Supply Kit |
| 1 | 681254-B21 | HP 4.3U Rail Kit |
| 1 | 681260-B21 | HP 0.66U Spacer Blank Kit |
| 2 | 664644-B22 | HP 2xSL4540 Gen8 Tray Node Svr |
| 2 | 740695-L21 | HP SL4540 Gen8 E5 2450v2 FIO Kit |
| 2 | 740695-B21 | HP SL4540 Gen8 E5-2450v2 Kit |
| 12 | 713985-B21 | HP 16GB 2Rx4 PC3L-12800R-11 Kit |
| 2 | 631681-B21 | HP 2GB FBWC for P-Series Smart Array |
| 4 | 655708-B21 | HP 500GB 6G SATA 7.2k 2.5in SC MDL HDD |
| 50 | 652766-B21 | HP 3TB 6G SAS 7.2K 3.5in SC MDL HDD |
| 2 | 692276-B21 | HP Smart Array P420i Mezz Ctrllr FIO Kit |
| 2 | 682632-B21 | HP SL4500 Storage Mezz to PCIe Opt Kit |
| 2 | 668943-B21 | HP 12in Super Cap for Smart Array |
| 2 | G3J28AAE | RHEL Svr 2 Sckt/2 Gst 1yr 24x7 E-LTU |
| 2 | C6N27A | HP Insight Control Lic |
| 2 | 649281-B21 | HP IB FDR/EN 10/40Gb 2P 544QSFP Adptr |
| 4 | 498385-B23 | HP 3M 4X DDR/QDR QSFP IB Cu Cable |

**Networking**
Table A-4 provides a BOM for two ToR switches and one HP iLO switch, as featured in the tested configuration.

**Table A-4.** BOM for two HP 5930 switches (ToR) switches and one HP 5900 switch (HP iLO)

| Quantity | Part number | Description |
|---|---|---|
| 2 | JG726A | HP FF 5930-32QSFP+ Switch |
| 4 | JG553A | HP X712 Bck(pwr)-Frt(prt) HV Fan Tray |
| 4 | JC680A | HP A58x0AF 650W AC Power Supply |
| 4 | JC680A#B2B | Power PDU Cable |
| 1 | JG510A | HP 5900AF-48G-4XG-2QSFP+ Switch |
| 2 | JC680A | HP A58x0AF 650W AC Power Supply |
| 2 | JC682A | HP 58x0AF Bck(pwr)-Frt(ports) Fan Tray |
| 4 | JG330A | QSFP+ to 4SFP+ 3m DAC cable |
| 20 | JG327A | HP X240 40G QSFP+ QSFP+ 3m DAC Cable |
| 10 | JG326A | HP X240 40G QSFP+ QSFP+ 1m DAC Cable |
| 12 | AF595A | HP 3.0M,Blue,CAT6 STP,Cable Data |

**Other hardware**

**Important**
Quantities listed in Table A-5 are based on a full rack with three switches, 90 compute nodes, and six storage nodes.

**Table A-5.** BOM for a single rack with four PDUs

| Quantity | Part number | Description |
|---|---|---|
| 1 | BW908A | HP 642 1200mm Shock Intelligent Rack |
| 1 | BW908A   001 | HP Factory Express Base Racking Service |
| 1 | BW946A | HP 42U Location Discovery Kit |
| 1 | BW930A | HP Air Flow Optimization Kit |
| 1 | BW930A   B01 | Include with complete system |
| 1 | BW909A | HP 42U 1200mm Side Panel Kit |
| 1 | BW891A | HP Rack Grounding Kit |
| 4 | AF520A | HP Intelligent Mod PDU 24a Na/Jpn Core |
| 6 | AF547A | HP 5xC13 Intlgnt PDU Ext Bars G2 Kit |

**HP Insight Cluster Management Utility options**

**Important**
Options listed in Table A-6 are based on a single node.

**Table A-6.** BOM for HP Insight Cluster Management Utility (Insight CMU) options, per-node

| Quantity | Part number | Description |
|----------|-------------|-------------|
| 1 | QL803B | HP Insight CMU 1yr 24x7 Flex Lic |
| 1 | QL803BAE | HP Insight CMU 1yr 24x7 Flex E-LTU |
| 1 | BD476A | HP Insight CMU 3yr 24x7 Flex Lic |
| 1 | BD476AAE | HP Insight CMU 3yr 24x7 Flex E-LTU |
| 1 | BD477A | HP Insight CMU Media |

**Hortonworks software**

**Important**
Options listed in Table A-7 are based on a single node.
While HP is a certified reseller of Hortonworks software subscriptions, all application support (level-one through level-three) is provided by Hortonworks.

**Table A-7.** BOM for Hortonworks software

| Quantity | Part number | Description |
|----------|-------------|-------------|
| 5 | F5Z52A | Hortonworks Data Platform Enterprise 4 Nodes or 50TB Raw Storage 1 year 24x7 Support LTU |

# Appendix B – Alternate parts for storage nodes

This appendix provides BOMs for alternate processors, memory, and disk drives for the SL4540 servers used as storage nodes.

**Table B-1.** Alternate processors

| Quantity per node | Part number | Description |
| --- | --- | --- |
| 1 | 740674-L21 | HP SL4540 Gen8 E5 2470v2 FIO Kit 10 cores at 2.4GHz |
| 1 | 740674-B21 | HP SL4540 Gen8 E5-2470v2 Kit |
| 1 | 740695-L21 | HP SL4540 Gen8 E5 2450v2 FIO Kit 8 cores at 2.5GHz |
| 1 | 740695-B21 | HP SL4540 Gen8 E5-2450v2 Kit |
| 1 | 740677-L21 | HP SL4540 Gen8 E5 2440v2 FIO Kit 8 cores at 1.9GHz |
| 1 | 740677-B21 | HP SL4540 Gen8 E5-2440v2 Kit |

**Table B-2.** Alternate memory – SL4540

| Quantity per node | Part number | Description |
| --- | --- | --- |
| 6 | 713985-B21 | HP 16GB 2Rx4 PC3L-12800R-11 Kit for 96GB of Memory |
| 12 | 713985-B21 | HP 16GB 2Rx4 PC3L-12800R-11 Kit for 192GB of Memory |

**Table B-3.** Alternate disk drives

| Quantity per node | Part number | Description |
| --- | --- | --- |
| 12 | 652757-B21 | HP 2TB 6G SAS 7.2K 3.5in SC MDL HDD |
| 12 | 652766-B21 | HP 3TB 6G SAS 7.2K 3.5in SC MDL HDD |
| 12 | 695510-B21 | HP 4TB 6G SAS 7.2K 3.5in SC MDL HDD |

# Appendix C – HP value added services and support

In order to help you jump-start your Hadoop solution development, HP offers a range of big data services, which are outlined in this appendix.

**Factory Express Services**
Factory-integration services are available for customers seeking a streamlined deployment experience. With the purchase of Factory Express services, your Hadoop cluster will arrive racked and cabled, with software installed and configured per an agreed-upon custom statement of work, for the easiest deployment possible. You should contact HP Technical Services for more information and for assistance with a quote.

**Technical Services Consulting – Reference Architecture Implementation Service for Hadoop (Hortonworks)**
With HP Reference Architecture Implementation Service for Hadoop, HP can install, configure, deploy, and test a Hadoop cluster that is based on HP BDRA. Experienced consultants implement all the details of the original Hadoop design: naming, hardware, networking, software, administration, backup, disaster recovery, and operating procedures. Where options exist, or the best choice is not clear, HP works with you to configure the environment to meet your goals and needs. HP also conducts an acceptance test to validate that the system is operating to your satisfaction.

**Technical Services Consulting – Big Data Services**
HP Big Data Services can help you reshape your IT infrastructure to corral increasing volumes of data – from e-mails, social media, and website downloads – and convert them into beneficial information. These services encompass strategy, design, implementation, protection, and compliance. Delivery is in the following three steps:

1. **Architecture strategy:** HP defines the functionalities and capabilities needed to align your IT with your big data initiatives. Through transformation workshops and roadmap services, you'll learn to capture, consolidate, manage and protect business-aligned information, including structured, semi-structured, and unstructured data.
2. **System infrastructure:** HP designs and implements a high-performance, integrated platform to support a strategic architecture for big data. Choose from design and implementation services, reference architecture implementations, and integration services. Your flexible, scalable infrastructure will support big data variety, consolidation, analysis, share, and search on HP platforms.
3. **Data protection:** Ensure the availability, security, and compliance of your big data systems. HP can help you safeguard your data and achieve regulatory compliance and lifecycle protection across your big data landscape, while also enhancing your approach to backup and business continuity.

For additional information, visit hp.com/services/bigdata.

## HP Support options

HP offers a variety of support levels to meet your needs. More information is provided below.

**HP Support Plus 24**
HP can provide integrated onsite hardware/software support services, available 24x7x365, including access to HP technical resources, four-hour response onsite hardware support and software updates.

**HP Proactive Care**
HP Proactive Care provides all of the benefits of proactive monitoring and reporting, along with rapid reactive support through HP's expert reactive support specialists. You can customize your reactive support level by selecting either six-hour call-to-repair or 24x7 with four-hour onsite response.

HP Proactive Care helps prevent problems, resolve problems faster, and improve productivity. Through analysis, reports, and update recommendations, you are able to identify and address IT problems before they can cause performance issues or outages.

**HP Proactive Care with the HP Personalized Support Option**
Adding the Personalized Support Option for HP Proactive Care is highly recommended. This option builds on the benefits of HP Proactive Care Service, providing you an assigned Account Support Manager who knows your environment and can deliver support planning, regular reviews, and technical and operational advice specific to your environment.

**HP Proactive Select**
To address your ongoing/changing needs, HP recommends adding Proactive Select credits to provide tailored support options from a wide menu of services that can help you optimize the capacity, performance, and management of your environment. These credits may also be used for assistance in implementing solution updates. As your needs change over time, you have the flexibility to choose the services best suited to address your current challenges.

**HP Datacenter Care**

HP Datacenter Care provides a more personalized, customized approach for large, complex environments, providing a single solution for reactive, proactive, and multi-vendor support needs. You may also choose the Defective Media Retention (DMR) option.

**Other offerings**

HP highly recommends HP Education Services (customer training and education) and additional Technical Services, as well as in-depth installation or implementation services when needed.

**More information**

For additional information, visit:

- HP Education Services: http://h10076.www1.hp.com/education/bigdata.htm
- HP Technology Consulting Services: hp.com/services/bigdata
- HP Deployment Services: hp.com/services/deployment

# For more information

Hortonworks: hortonworks.com

Hortonworks partner site: hortonworks.com/partner/hp/

HP Solutions for Apache Hadoop: hp.com/go/hadoop

HP Insight Cluster Management Utility (Insight CMU): hp.com/go/cmu

HP 5930 Switch Series: hp.com/networking/5930

HP ProLiant servers: hp.com/go/proliant

HP Moonshot system: hp.com/go/moonshot

HP Enterprise Software: hp.com/go/software

HP Networking: hp.com/go/networking

HP Integrated Lights-Out (HP iLO): hp.com/servers/ilo

HP Product Bulletin (QuickSpecs): hp.com/go/quickspecs

HP Services: hp.com/go/services

HP Support and Drivers: hp.com/go/support

HP Systems Insight Manager (HP SIM): hp.com/go/hpsim

HP Trafodion: http://wiki.trafodion.org


To help us improve our documents, please provide feedback at hp.com/solutions/feedback.


**About Hortonworks**

Hortonworks develops, distributes and supports the only 100% open source Apache Hadoop data platform. Our team comprises the largest contingent of builders and architects within the Hadoop ecosystem who represent and lead the broader enterprise requirements within these communities. The Hortonworks Data Platform provides an open platform that deeply integrates with existing IT investments and upon which enterprises can build and deploy Hadoop-based applications. Hortonworks has deep relationships with the key strategic data center partners that enable our customers to unlock the broadest opportunities from Hadoop. For more information, visit hortonworks.com.

**Sign up for updates**
**hp.com/go/getupdated**