



# HP Verified Reference Architecture for Hortonworks HDP 2.2 on HP ProLiant DL380 Gen9 with RHEL

HP Converged Infrastructure with Hortonworks Data Platform 2.2 for Apache Hadoop

## Table of contents

Executive summary .....	2
Introduction .....	3
Hortonworks Data Platform solution overview.....	3
Key highlights of HDP 2.2.....	4
Hortonworks Data Platform .....	5
Solution components .....	6
High-availability considerations.....	6
Pre-deployment considerations/system selection .....	7
Reference architectures .....	9
Multi-rack reference architecture.....	11
Capacity and sizing.....	12
Analysis and recommendations .....	13
Guidance deployment steps .....	16
Server selection.....	16
Worker nodes.....	21
Networking .....	23
Switch selection.....	23
Management .....	24
Bill of materials .....	26
Summary.....	29
Implementing a proof-of-concept .....	30
Appendix A: Cluster design – heat map for server platforms.....	30
Appendix B: Hadoop cluster tuning/optimization .....	32
Appendix C: Alternate parts .....	34
Appendix D: HP value-added services and support.....	35
For more information .....	37

## Executive summary

HP and Apache Hadoop allow you to derive new business insights from all of your data by providing a platform to store, manage and process data at scale. As organizations start to capture and collect these big datasets, increased storage becomes a necessity. With a centralized repository, increased compute power and storage are requirements to respond to the sheer scale of Big Data and its projected growth. Enterprises need Big Data platforms that are purpose-built to support their big data strategies. This white paper provides several performance optimized configurations for deploying Hortonworks Data Platform (HDP) clusters of varying sizes on HP infrastructure that provide a significant reduction in complexity and increase in value and performance.

The configurations are based on Hortonworks Data Platform (HDP), 100% open source distribution of Apache Hadoop, specifically HDP 2.2 and the HP ProLiant DL380 Gen9 server platform. The configurations reflected in this document have been jointly designed and developed by HP and Hortonworks to provide optimum computational performance for Hadoop and are also compatible with other HDP 2.x releases.

HP Big Data solutions provide best-in-class performance and availability, with integrated software, services, infrastructure, and management – all delivered as one proven configuration as described at [hp.com/go/hadoop](http://hp.com/go/hadoop). In addition to the benefits described above, the reference architecture in this white paper also includes the following features that are unique to HP:

**Servers** – HP ProLiant DL360 Gen9 and DL380 Gen9 include:

- The HP Smart Array P440ar controller provides increased<sup>1</sup> I/O throughput performance resulting in a significant performance increase for I/O bound Hadoop workloads (a common use case) and the flexibility for the customer to choose the desired amount of resilience in the Hadoop Cluster with either JBOD or various RAID configurations. The HP H240ar Smart Host Bus Adapter is a 2X4 internal ports flexible Smart HBA for DL360, DL380 and ML350 Gen9 servers which provide customers with the flexibility and speed they have come to expect from HP. HP Smart Host Bus Adapters are well-suited for workloads that require reliable high performance and scalability for direct attached storage with basic RAID functionality and shared storage support at a lower cost. Typical uses include: OS boot, Hadoop.
- Two sockets with 10 core processors, using Intel® Xeon® E5-2600 v3 product family, provide the high performance required for CPU bound Hadoop workloads. Alternative processors are provided in Appendix C. For Management and Head nodes, the same CPU family with 8 core processors are used.
- The HP iLO Management Engine on the servers contains HP Integrated Lights-Out 4 (iLO 4) and features a complete set of embedded management features for HP Power/Cooling, Agentless Management, Active Health System, and Intelligent Provisioning which reduces node and cluster level administration costs for Hadoop.

**Cluster Management** – HP Insight Cluster Management Utility (CMU) provides push-button scale out and provisioning with industry leading provisioning performance, reducing deployments from days to hours. HP CMU provides real-time and historical infrastructure and Hadoop monitoring with 3D visualizations allowing customers to easily characterize Hadoop workloads and cluster performance. This allows customers to further reduce complexity and improve system optimization leading to improved performance and reduced cost. In addition, HP Insight Management and HP Service Pack for ProLiant, allow for easy management of firmware and servers.

**Networking** – The HP 5900AF-48XGT-4QSFP+ 10GbE Top of Rack switch has 48 RJ-45 1/10GbE ports and 4 QSFP+ 40GbE ports. It provides IRF bonding and sFlow for simplified management, monitoring and resiliency of Hadoop network. The 512MB flash, 2GB SDRAM and packet buffer size of 9MB provide excellent performance of 952 million pps throughput and switching capacity of 1280Gb/s with very low 10Gb/s latency of less than 1.5  $\mu$ s (64-byte packets).

HP FlexFabric 5930-32QSFP+ 40GbE Aggregation switch provides IRF Bonding and sFlow which simplifies the management, monitoring and resiliency of the customer's Hadoop network. The 1GB flash, 4GB SDRAM memory and packet buffer size of 12.2MB provide excellent performance of 1429 million pps throughput and routing/switching capacity of 2560Gb/s with very low 10Gb/s latency of less than 1  $\mu$ s (64-byte packets). The switch seamlessly handles burst scenarios such as shuffle, sort and block replication which are common in Hadoop clusters.

**Analytics database** – The HP Vertica connectors for Hadoop allow seamless integration of both structured and unstructured data providing end-to-end analytics thereby simplifying bi-directional data movement for Hadoop and reducing customer integration costs. Vertica is a leading real-time, scalable, analytical platform for structured data.

All of these features reflect HP balanced building blocks of servers, storage and networking, along with integrated management software.

<sup>1</sup> Compared to the previous generation of Smart Array controllers

**Target audience:** This document is intended for decision makers, system and solution architects, system administrators and experienced users who are interested in reducing the time to design or purchase an HP and Hortonworks solution. An intermediate knowledge of Apache Hadoop and scale out infrastructure is recommended. Those already possessing expert knowledge about these topics may proceed directly to [Solution components](#).

**Document purpose:** The purpose of this document is to describe a reference architecture, highlighting recognizable benefits to technical audiences and providing guidance for end users on selecting the right configuration for building their Hadoop cluster needs.

This white paper describes testing performed in May-June 2015.

## Introduction

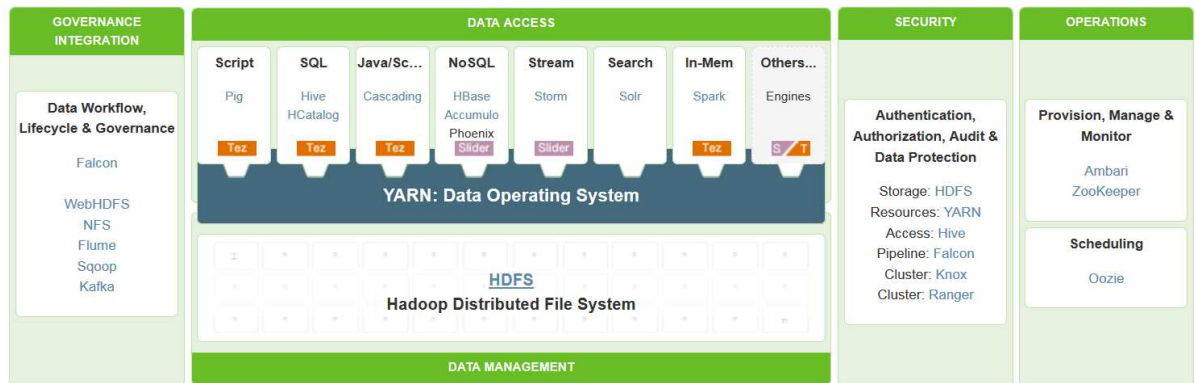
This white paper has been created to assist in the rapid design and deployment of Hortonworks Data Platform software on HP infrastructure for clusters of various sizes. It is also intended to identify the software and hardware components required in a solution to simplify the procurement process. The recommended HP Software, 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 configurations included provide the best value for optimum MapReduce, YARN, Hive, HBase and Solr computational performance, resulting in a significant performance increase at an optimum cost.

## Hortonworks Data Platform solution overview

Hortonworks is a major contributor to Apache Hadoop, the world’s most popular big data platform. Hortonworks focuses on further accelerating the development and adoption of Apache Hadoop by making the software more robust and easier to consume for enterprises and more open and extensible for solution providers. The Hortonworks Data Platform (HDP), powered by Apache Hadoop, is a massively scalable and 100% open source platform for storing, processing and analyzing large volumes of data. It is designed to deal with data from many sources and formats in a very quick, easy and cost-effective manner.

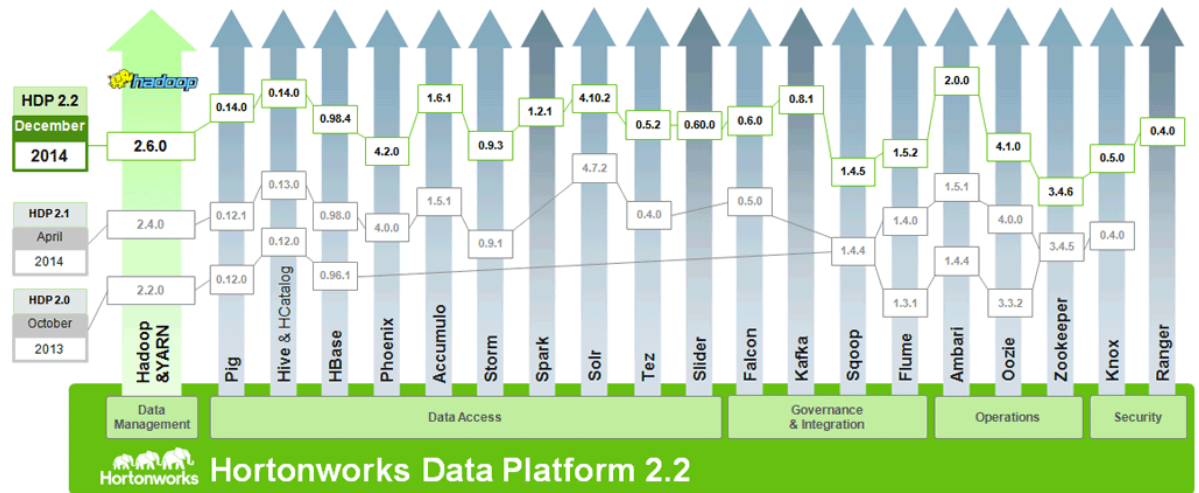
HDP is a platform for multi-workload data processing across an array of processing methods – from batch through interactive and real-time – all supported with solutions for governance, integration, security and operations. As the only completely open Hadoop data platform available, HDP integrates with and augments your existing best-of-breed applications and systems so you can gain value from your enterprise Big Data, with minimal changes to your data architectures. Finally, HDP allows you to deploy Hadoop wherever you want it – from cloud or on-premises as an appliance, and across both Linux® and Microsoft® Windows®. Figure 1 shows the Hortonworks Data Platform.

**Figure 1.** Hortonworks Data Platform: A full Enterprise Hadoop Data Platform



Hortonworks Data Platform Version 2.2 represents yet another major step forward for Hadoop as the foundation of a Modern Data Architecture. This release incorporates the most recent innovations that have happened in Hadoop and its supporting ecosystem of projects. HDP 2.2 packages more than a hundred new features across all Apache Hadoop open source existing projects. Every component is updated and Hortonworks has added some key technologies and capabilities to HDP 2.2. Figure 2 shows Hortonworks Data Platform 2.2.

**Figure 2.** Hortonworks Data Platform 2.2



Hortonworks Data Platform enables Enterprise Hadoop: the full suite of essential Hadoop capabilities that are required by the enterprise and that serve as the functional definition of any data platform technology. This comprehensive set of capabilities is aligned to the following functional areas: Data Management, Data Access, Data Governance, and Integration, Security, and Operations.

### Key highlights of HDP 2.2

**Enterprise SQL at scale in Hadoop** – While YARN has allowed new engines to emerge for Hadoop, the most popular integration point with Hadoop continues to be SQL and Apache Hive is still the defacto standard.

**Updated SQL semantics for hive transactions for update and delete** – ACID transactions provide atomicity, consistency, isolation, and durability. This helps with streaming and baseline update scenarios for Hive such as modifying dimension tables or other fact tables.

**Improved performance of hive with a cost based optimizer** – The cost based optimizer for Hive uses statistics to generate several execution plans and then chooses the most efficient path as it relates system resources required to complete the operation. This presents a major performance increase for Hive.

For detailed information on Hortonworks Data Platform, please see [hortonworks.com/hdp](http://hortonworks.com/hdp)

## Hortonworks Data Platform

The platform functions within Hortonworks Data Platform are provided by two key groups of services, namely the Management and Worker services. Management services manage the cluster and coordinate the jobs whereas Worker services are responsible for the actual execution of work on the individual scale out nodes. Tables 1 and 2 below specify which services are management services and which services are worker services. Each table contains two columns. The first column is the description of the service and the second column specifies the maximum number of nodes a worker service can be distributed to on the platform. The Reference Architectures (RAs) we provide in this document will map the Management and Worker services onto HP infrastructure for clusters of varying sizes. The RAs factor in the scalability requirements for each service. Tables 1 and 2 list the maximum distribution across nodes for the HDP base management services and HDP base enterprise worker services.

### Management services

**Table 1.** HDP Base Management services

Service	Maximum distribution across nodes
Ambari	1
HueServer	1
ResourceManager	2
JobHistoryServer	1
HBaseMaster	Varies
NameNode	2
Oozie	1
ZooKeeper	Varies

### Worker services

**Table 2.** HDP Base Enterprise Worker services

Service	Maximum distribution across nodes
DataNode	Typically, all except management nodes
NodeManager	Typically, all except management nodes
ApplicationMaster	One for each job
HBaseRegionServer	Varies

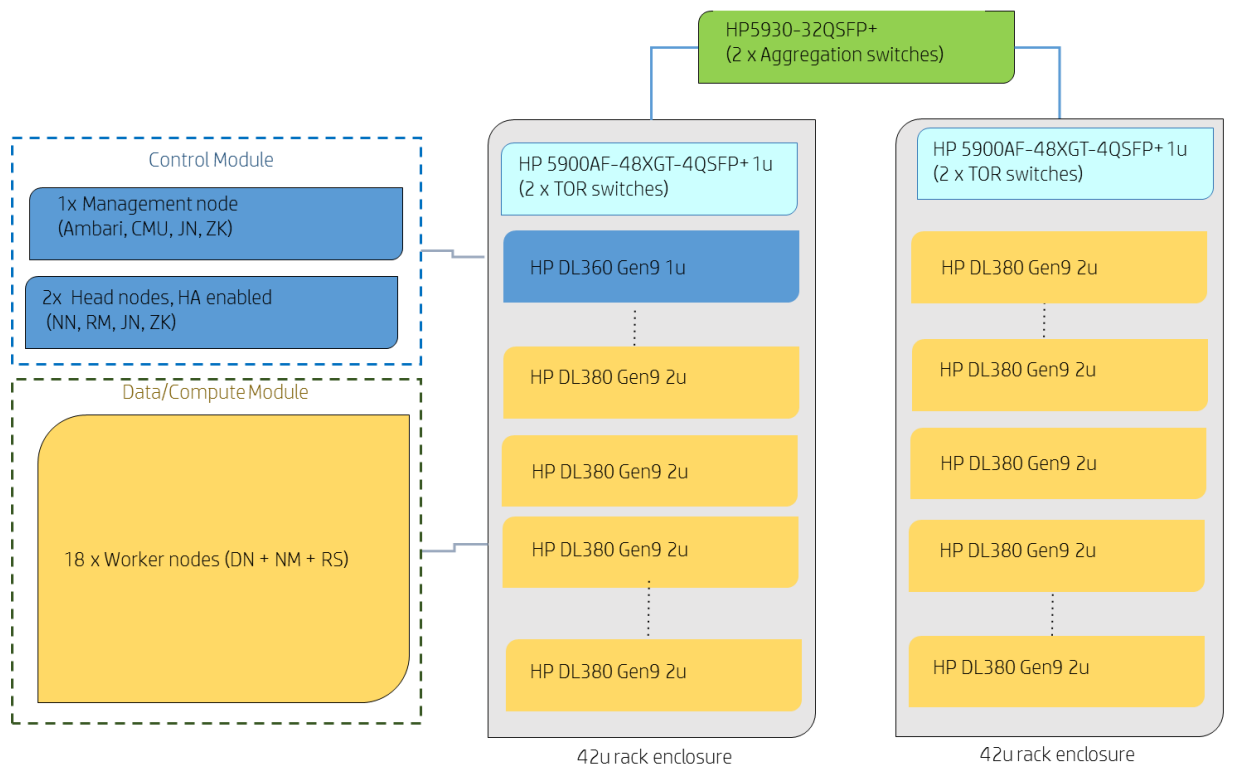
#### Note

A full Hadoop cluster will have component-specific services on each worker node and management node.

## Solution components

Figure 3 shows a conceptual diagram of an HP DL380 Gen9 Reference Architecture.

**Figure 3.** HP DL380 Gen9 Reference Architecture conceptual diagram



Legend: JN – Journal Node, ZK – ZooKeeper, NN – NameNode, RM – ResourceManager, DN – DataNode, NM – NodeManager, AM – ApplicationManager, JH – JobHistoryServer, RS - RegionServer

Each of the components is discussed at length below. For full BOM listing on products selected, please refer to the Bill of Materials section of this white paper.

### High-availability considerations

The following are some of the high-availability features considered in this reference architecture configuration:

**Hadoop NameNode HA** – The configurations in this white paper utilize quorum-based journaling high-availability features in HDP 2.2. For this feature, servers should have similar I/O subsystems and server profiles so that each NameNode server could potentially take the role of another. Another reason to have similar configurations is to ensure that ZooKeeper’s quorum algorithm is not affected by a machine in the quorum that cannot make a decision as fast as its quorum peers.

**ResourceManager HA** – To make a YARN cluster highly-available (similar to JobTracker HA in MR1), the underlying architecture of an Active/Standby pair is configured – hence the completed tasks of in-flight MapReduce jobs are not re-run on recovery after the ResourceManager is restarted or failed over. One ResourceManager is Active and one or more ResourceManagers are in standby mode waiting to take over should anything happen to the Active ResourceManager.

**OS availability and reliability** – For the reliability of the server, the OS disk is configured in a RAID 1+0 configuration thus preventing failure of the system from OS hard disk failures.

**Network reliability** – The reference architecture configuration uses two HP 5900AF-48XGT switches for redundancy, resiliency and scalability through using Intelligent Resilient Framework (IRF) bonding. We recommend using redundant power supplies.

To ensure the servers and racks have adequate power redundancy we recommend that each server have a backup power supply, and each rack have at least two Power Distribution Units (PDUs).

## Pre-deployment considerations/system selection

There are a number of key factors you should consider prior to designing and deploying a Hadoop Cluster. The following subsections articulate the design decisions in creating the baseline configurations for the reference architectures. The rationale provided includes the necessary information for you to take the configurations and modify them to suit a particular custom scenario. Table 3 lists the items to be taken into consideration before deployment.

**Table 3.** Pre-deployment considerations

Functional Component	Value
Operating system	Improves availability and reliability
Computation	Ability to balance price with performance
Memory	Ability to balance price with capacity and performance
Storage	Ability to balance price with capacity and performance
Network	Ability to balance price with performance

### Operating system

Hortonworks HDP 2.2 supports the following 64-bit operating systems:

- For Red Hat® Enterprise Linux® (RHEL) 6.6 systems, Hortonworks provides 64-bit packages for RHEL 6.6, the current release as of this reference architecture publication, or later is required.
- For RHEL previous to 6.6, Hortonworks provides 64-bit packages for RHEL 5 and RHEL 6.
- For Ubuntu systems, Hortonworks provides 64-bit packages for Precise (12.04) LTS. Full details on supported OS, databases and JDK versions are available at the Hortonworks site.

### Best practices

HP recommends using a 64-bit operating system to avoid constraining the amount of memory that can be used on worker nodes. The 64-bit version of RHEL 6.6 is recommended due to its superior filesystem, performance and scalability characteristics, plus the comprehensive, certified support of the Hortonworks and HP supplied software used in Big Data clusters. The Reference Architectures listed in this document were tested with 64-bit RHEL 6.6. For Ambari and its supporting services database, HP best practice is to use a single database such as PostgreSQL.

### Computation

Unlike MR1, where the processing or computational capacity of a Hadoop cluster is determined by the aggregate number of Resource Containers available across all the worker nodes, under YARN/MR2, the notion of slots has been discarded and resources are now configured in terms of amounts of memory and CPU (virtual cores). There is no distinction between resources available for map, and resources available for reduces – all MR2 resources are available for both in the form of **containers**. Employing Hyper-Threading increases your effective core count, potentially allowing ResourceManager to assign more cores as needed.

Resource request tasks that will use multiple threads can request more than one core with the `mapreduce.map.cpu.vcores` and `mapreduce.reduce.cpu.vcores` properties. HDP2.x only supports Capacity Scheduler with YARN. The default scheduler in HDP2.2 is the Capacity Scheduler.

### Key point

When computation performance is of primary concern, HP recommends higher CPU powered DL380 servers for worker nodes with 256GB RAM. HDP 2.2 components, such as HBase, HDFS Caching, Storm and Solr, benefit from large amounts of memory.

### Memory

Use of Error Correcting Memory (ECC) is a practical requirement for Apache Hadoop and is standard on all HP ProLiant servers. Memory requirements differ between the management nodes and the worker nodes. The management nodes typically run one or more memory intensive management processes and therefore have higher memory requirements. Worker nodes need sufficient memory to manage the NodeManager and Container processes. If you have a memory bound

YARN job we recommend that you increase the amount of memory on all the worker nodes. In addition, a high memory cluster can also be used for Spark, HBase or interactive Hive, which could be memory intensive.

### Storage

Fundamentally, Hadoop is designed to achieve performance and scalability by moving the compute activity to the data. It does this by distributing the Hadoop job to worker nodes close to their data, ideally running the tasks against data on local disks.

---

### Best practice

HP recommends choosing LFF drives over SFF drives (same angular speed) as the LFF drives have better performance than the SFF drives due to faster tangential speed that leads to higher disk I/O. Given the architecture of Hadoop, the data storage requirements for the worker nodes are best met by direct attached storage (DAS) in a Just a Bunch of Disks (JBOD) configuration.

---

There are several factors to consider and balance when determining the number of disks a Hadoop worker node requires.

**Storage capacity** – The number of disks and their corresponding storage capacity determines the total amount of the HDFS storage capacity for your cluster.

**Redundancy** – Hadoop ensures that a certain number of block copies are consistently available. This number is configurable in the block replication factor setting, which is typically set to three. If a Hadoop worker node goes down, Hadoop will replicate the blocks that had been on that server onto other servers in the cluster to maintain the consistency of the number of block copies. For example, if the NIC (Network Interface Card) on a server with 16TB of block data fails, 16TB of block data will be replicated between other servers in the cluster to ensure the appropriate number of replicas exist. Furthermore, the failure of a non-redundant ToR (Top of Rack) switch will generate even more replication traffic. Hadoop provides data throttling capability in the event of a node/disk failure so as to not overload the network.

**I/O performance** – The more disks you have, the less likely it is that you will have multiple tasks accessing a given disk at the same time. This avoids queued I/O requests and incurring the resulting I/O performance degradation.

**Disk configuration** – The management nodes are configured differently from the worker nodes because the management processes are generally not redundant and as scalable as the worker processes. For management nodes, storage reliability is therefore important and SAS drives are recommended. For worker nodes, one has the choice of SAS or SATA and as with any component there is a cost/performance tradeoff. If performance and reliability are important, we recommend SAS MDL disks; otherwise, we recommend SATA MDL disks. Specific details around disk and RAID configurations will be provided in the [Server selection](#) section of this white paper.

### Network

Configuring a single Top of Rack (ToR) switch per rack introduces a single point of failure for each rack. In a multi-rack system such a failure will result in a very long replication recovery time as Hadoop rebalances storage, and in a single-rack system such a failure could bring down the whole cluster. Consequently, configuring two ToR switches per rack is recommended for all production configurations as it provides an additional measure of redundancy. This can be further improved by configuring link aggregation between the switches. The most desirable way to configure link aggregation is by bonding the two physical NICs on each server. Port1 wired to the first ToR switch and Port2 wired to the second ToR switch, with the two switches IRF bonded. When done properly, this allows the bandwidth of both links to be used. If either of the switches fail, the servers will still have full network functionality, but with the performance of only a single link. Not all switches have the ability to do link aggregation from individual servers to multiple switches; however, the HP 5900AF-48XGT switch supports this through HP Intelligent Resilient Framework (IRF) technology. In addition, switch failures can be further mitigated by incorporating dual power supplies for the switches.

Hadoop is rack-aware and tries to limit the amount of network traffic between racks. The bandwidth and latency provided by two bonded 10 Gigabit Ethernet (GbE) connections from the worker nodes to the ToR switch is more than adequate for most Hadoop configurations. Multi-Rack Hadoop clusters, that are not using IRF bonding for inter-rack traffic, will benefit from having ToR switches connected by 40GbE uplinks to core aggregation switches. Large Hadoop clusters introduce multiple issues that are not typically present in small to medium sized clusters. To understand the reasons for this, it is helpful to review the network activity associated with running Hadoop jobs and with exception events such as server failure.

A more detailed white paper for Hadoop Networking best practices is available. For more information, please refer to [hp.com/go/hadoop](http://hp.com/go/hadoop).



---

**Best practice**

HP recommends ToR switches with packet buffering and connected by 40GbE uplinks to core aggregation switches for large clusters. For MapReduce jobs, during the shuffle phase, the intermediate data has to be pulled by the reduce tasks from mapper output files across the cluster. While network load can be reduced if partitioners and combiners are used, it is possible that the shuffle phase will place the core and ToR switches under a large traffic load.

Each reduce task can concurrently request data from a default of five mapper output files. Thus, there is the possibility that servers will deliver more data than their network connections can handle which will result in dropped packets and can lead to a collapse in traffic throughput. ToR switches with packet buffering protect against this event.

---

**Reference architectures**

The following sections illustrate a reference progression of Hadoop clusters from a single rack to a multi-rack configuration. Best practices for each of the components within the configurations specified have been articulated earlier in this document.

**Single rack reference architecture**

The Single Rack Hortonworks Enterprise Reference Architecture (RA) is designed to perform well as a single rack cluster design but also form the basis for a much larger multi-rack design. When moving from the single rack to multi-rack design, one can simply add racks to the cluster without having to change any components within the single rack. The Reference Architecture reflects the following.

*Single rack network*

As previously described in the [Network](#) section, two IRF Bonded HP 5900AF-48XGT ToR switches are specified for performance and redundancy. The HP 5900AF-48XGT includes four 40GbE uplinks which can be used to connect the switches in the rack into the desired network or to the 40GbE HP 5930-32QSFP+ aggregation switch. Keep in mind that if IRF bonding is used, it requires 2x 40GbE ports per switch, which would leave 2x 40GbE ports on each switch for uplinks.

*Cluster isolation and access configuration*

It is important to isolate the Hadoop Cluster on the network so that external network traffic does not affect the performance of the cluster. In addition, this also allows the Hadoop cluster to be managed independently from that of its users, which ensures that the cluster administrator is the only one capable of making changes to the cluster configurations. To achieve this, we recommend isolating the ResourceManager, NameNode and Worker nodes on their own private Hadoop Cluster subnet.

---

**Key point**

Once a Hadoop cluster is isolated, the users of the cluster will still need a way to access the cluster and submit jobs to it. To achieve this we recommend multi-homing the Management node so that it participates in both the Hadoop Cluster subnet and a subnet belonging to the users of the cluster. Ambari is a web application that runs on the Management node and allows users to manage and configure the Hadoop cluster (including seeing the status of jobs) without being on the same subnet, provided the Management node is multi-homed. Furthermore, this allows users to shell into the Management node and run the Apache Pig or Apache Hive command line interfaces and submit jobs to the cluster that way.

---

**Staging data**

In addition, once the Hadoop Cluster is on its own private network one needs to think about how to be able to reach the HDFS in order to ingest data. The HDFS client needs the ability to reach every Hadoop DataNode in the cluster in order to stream blocks of data onto the HDFS. The Reference Architecture provides two options to do this.

The first option is to use the already multi-homed Management node, which can be configured with 4 additional disks to provide twice the amount of disk capacity (an additional 3.6TB) compared to the other management servers in order to provide a staging area for ingesting data into the Hadoop Cluster from another subnet.

The other option is to make use of the open ports that have been left available in the switch. This Reference Architecture has been designed such that if both NICs are used on each worker node and 2 NICs are used on each management node it leaves 26 ports still available across both the switches in the rack. These 26 10GbE ports or the remaining 40GbE ports on the switches can be used by other multi-homed systems outside of the Hadoop cluster to move data into the Hadoop Cluster.

---

**Note**

The benefit of using dual-homed edge node(s) to isolate the in-cluster Hadoop traffic from the ETL traffic flowing to the cluster is often debated. One benefit of doing so is better security. However, the downside of a dual-homed network architecture is ETL performance/connectivity issues, since a relatively few number of nodes in the cluster are capable of ingesting data. For example, the customer may want to kick off Sqoop tasks on the worker nodes to ingest data from external RDBMs, which will maximize the ingest rate. However this requires that the worker nodes be exposed to the external network to parallelize data ingestion, which is less secure. The customer has to weigh their options before committing to an optimal network design for their environments.

---

One can leverage WebHDFS which provides an HTTP proxy to securely read and write data to and from the Hadoop Distributed File System. For more information on WebHDFS, please see [http://docs.hortonworks.com/HDPDocuments/HDP2/HDP-2.2.0/HDFS\\_Admin\\_Tools\\_v22/web\\_hdfs\\_admin/index.html](http://docs.hortonworks.com/HDPDocuments/HDP2/HDP-2.2.0/HDFS_Admin_Tools_v22/web_hdfs_admin/index.html)

**Rack enclosure**

The rack contains eighteen HP ProLiant DL380 servers, three HP ProLiant DL360 servers and two HP 5900AF-48XGT switches within a 42U rack. This leaves 1U open for a 1U KVM switch.

**Network**

As later described in the [Switch selection](#) section, two HP 5900AF-48XGT switches are specified for performance and redundancy. The HP 5900AF-48XGT includes up to four 40GbE uplinks which can be used to connect the switches in the rack.

**Management nodes**

Three ProLiant DL360 Gen9 management nodes are specified:

- The Management node
- The ResourceManager/NameNode HA
- The NameNode/ResourceManager HA

**Worker nodes**

As specified in this design, eighteen ProLiant DL380 Gen9 worker nodes will fully populate a rack.

---

**Best practice**

Although it is possible to deploy with as few nodes as a single worker node, HP recommends starting with a minimum of three worker nodes in order to provide the redundancy that comes with the default replication factor of 3 for availability. Performance improves with additional worker nodes as ResourceManager can leverage idle nodes to land jobs on servers that have the appropriate blocks, leveraging data locally rather than pulling data across the network. These servers are homogenous and run the DataNode and the NodeManager (or HBaseRegionServer) processes.

---

**Power and cooling**

In planning for large clusters, it is important to properly manage power redundancy and distribution. To ensure the servers and racks have adequate power redundancy we recommend that each server have a backup power supply, and each rack have at least two Power Distribution Units (PDUs). There is an additional cost associated with procuring redundant power supplies. This is less important for larger clusters as the inherent redundancy within the Hortonworks Distribution of Hadoop will ensure there is less impact.

---

**Best practice**

For each server, HP recommends that each power supply is connected to a different PDU than the other power supply on the same server. Furthermore, the PDUs in the rack can each be connected to a separate data center power line to protect the infrastructure from a data center power line failure.

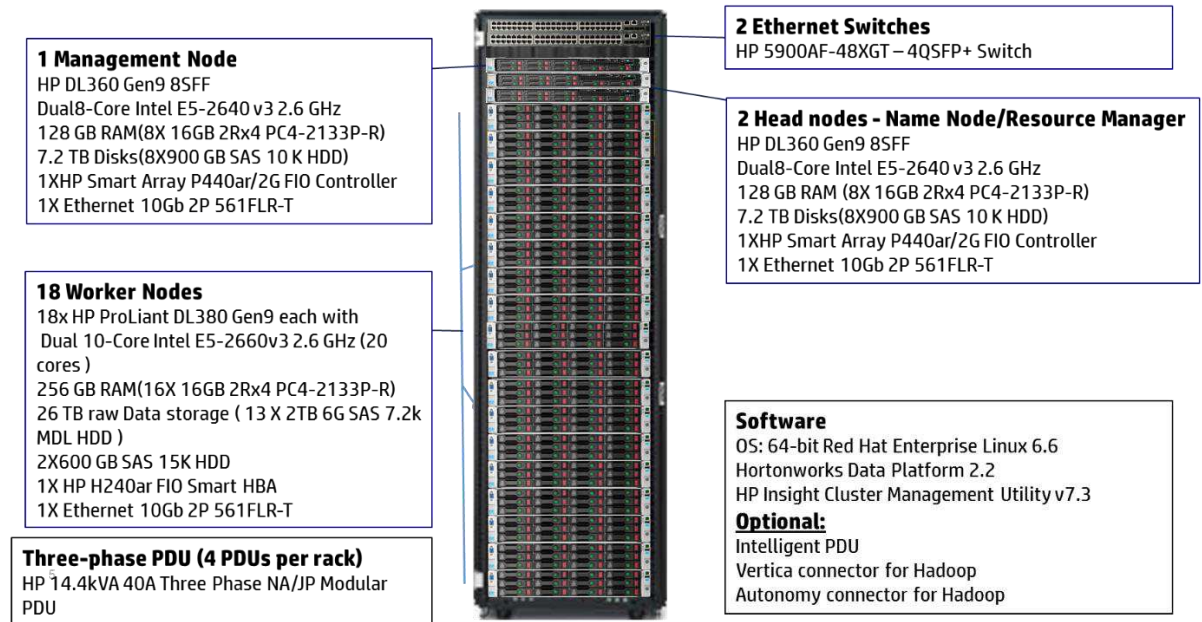
Additionally, distributing the server power supply connections evenly to the in-rack PDUs, as well as distributing the PDU connections evenly to the data center power lines ensures an even power distribution in the data center and avoids overloading any single data center power line. When designing a 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.

---

### Open rack space

The design leaves 1U open in the rack allowing for a KVM switch when using a standard 42U rack. Figure 4 shows the rack level view for single rack reference architecture.

**Figure 4.** Single rack reference Architecture – Rack level view



### Multi-rack reference architecture

The Multi-Rack design assumes the Single Rack RA Cluster design is already in place and extends its scalability. The Single Rack configuration ensures the required amount of management services are in place for large scale out. For Multi-Rack clusters, one simply adds more racks of the configuration provided below to the Single Rack configuration. This section reflects the design of those racks, and Figures 5 and 6 show the rack level view of the multi-rack architecture.

#### Rack enclosure

The rack contains eighteen HP ProLiant DL380 Gen9 servers and two HP 5900AF-48XGT switches within a 42U rack. 4U remains open and can accommodate an additional 2x DL380 servers (2U each) or a KVM switch (1U); or, install two HP 5930-32QSFP+ (1U) aggregation switches in the first expansion rack.

#### Multi-Rack Network

As later described in the [Switch selection](#) section, two HP 5900AF-48XGT ToR switches are specified per each expansion rack for performance and redundancy. The HP 5900AF-48XGT includes up to four 40GbE uplinks which can be used to connect the switches in the rack into the desired network, via a pair of HP 5930-32QSFP+ aggregation switches.

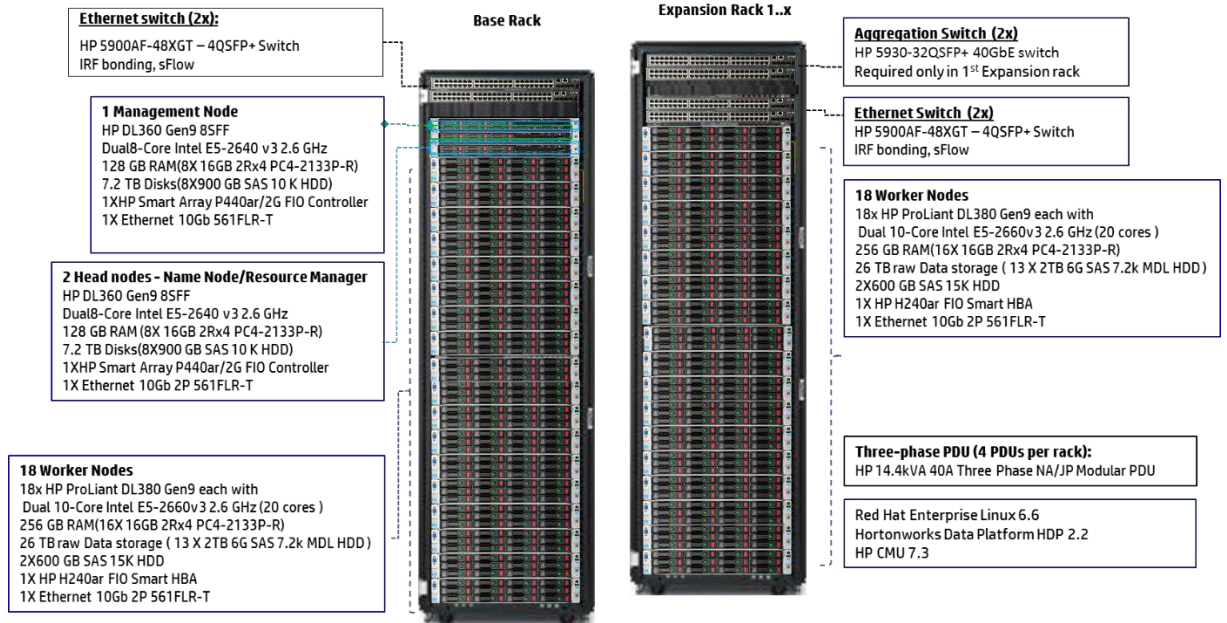
#### Software

The ProLiant DL380 servers in the rack all are configured as Worker nodes in the cluster, as all required management processes are already configured in the Single Rack RA. Aside from the OS, each worker node typically runs DataNode, NodeManager (and HBaseRegionServer if you are using HBase).

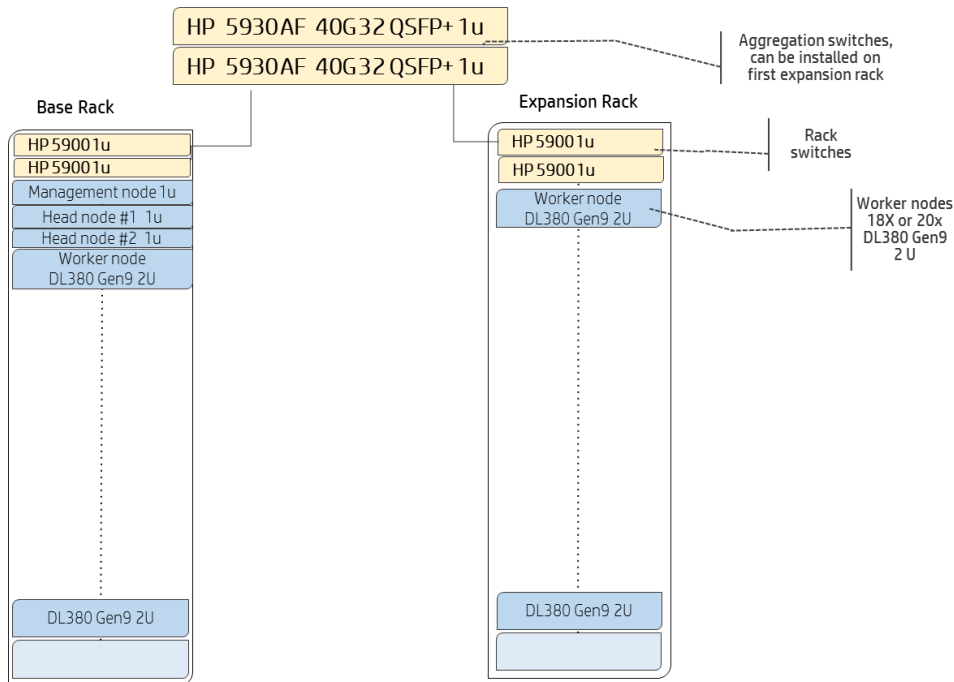
#### Note

While much of the architecture for Multi-Rack Hadoop cluster was borrowed from the Single Rack design, the architecture suggested here for multi-rack is based on previous iterations of testing on the DL380e platform. It is provided here as a general guideline for designing multi-rack Hadoop clusters.

**Figure 5.** Multi-rack reference architecture – Rack level view



**Figure 6.** Multi-rack reference architecture (extension of the single rack reference architecture)



## Capacity and sizing

Hadoop cluster storage sizing requires careful planning and identifying the current and future storage and compute needs.

Here is a general guideline on data inventory:

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

- Default compression turned on
- Space for intermediate files
- How to calculate storage needs – guidelines

To calculate the storage needs find the number of TB of data per day, week, month and year. Then add the ingestion rate of all data sources. It would make sense to identify storage requirements for short term, medium term and long term.

Another important consideration is data retention both size and duration. What data is required to be kept and for how long. Consider the maximum fill-rate and file system format space requirements on hard drives while estimating size of storage.

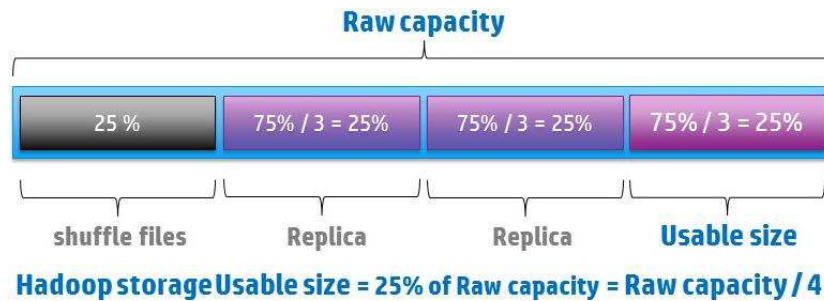
## Analysis and recommendations

### Usable HDFS capacity sizing

Some customers may have a specific requirement for usable HDFS capacity. Generally speaking there is a 10% deduction (decimal to binary conversion) when converting from raw to usable.

For example, a full rack cluster (18 x DL380 Gen9) that has 13x 2TB drives per node, the raw capacity would be  $2TB * 13 * 18 = 468TB$ , subtract 10% from the raw space to calculate usable HDFS space of 421TB. 25% of usable space should be set aside for MapReduce leaving 315TB. If the replication factor is **3**, then you would get  $315TB/3$  or approximately 105TB of usable space. Following this rule of thumb would help decide how many data nodes to order. Compression would provide additional usable space, depending on the type of compression, such as snappy or gzip. Figure 7 shows the usable storage for a replication factor of 3.

**Figure 7.** Usable storage for replication factor of 3.



### Key point

HP recommends using compression as it reduces file size on disks and speeds up data transfer to disks and network. The parameter to set to configure job output files to be compressed:

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

### System configuration guidance

#### Workload matters

Hadoop distributes data across a cluster of balanced machines and uses replication to ensure data reliability and fault tolerance. Because data is distributed on machines with compute power, processing can be sent directly to the machines storing the data. Since each machine in a Hadoop cluster stores as well as processes data, those machines need to be configured to satisfy both data storage and processing requirements. Table 4 shows examples of CPU and I/O bound workloads.

**Table 4.** Examples of CPU and I/O bound workloads

I/O bound jobs	CPU bound job
Sorting	Classification
Grouping	Clustering
Data import and export	Complex text mining
Data Movement and transformation	Natural language processing
	Feature extraction

Based on feedback from the field, most users looking to build a Hadoop cluster are not aware of the eventual profile of their workload. Often the first jobs that an organization runs with Hadoop differ very much from the jobs that Hadoop is ultimately used for as proficiency increases.

Building a cluster appropriate for the workload is key to optimizing the Hadoop cluster.

#### *Processor options*

For workloads that are CPU intensive it is recommended to choose higher capacity processors with more cores. Typically workloads such as interactive Hive, Spark and Solr Search will benefit from higher capacity CPUs. Table 5 shows alternate CPUs for the selected ProLiant DL380.

**Table 5.** CPU recommendations

<b>CPU</b>	<b>Description</b>
2 x E5-2660 v3	Base configuration (10 cores/2.6GHz)
2 x E5-2670 v3	Enhanced (12 cores/2.3GHz)
2 x E5-2680 v3	High Performance (12 cores/ 2.5GHz)

See [Appendix C: Alternate parts](#), Table C-1, for BOM details on various CPU choices.

#### *Memory options*

When calculating memory requirements, remember that Java uses up to 10 percent of memory to manage the virtual machine. HP recommends to configure Hadoop to use strict heap size restrictions to avoid memory swapping to disk.

It is important to optimize RAM for the memory channel width. For example, when using dual-channel memory, each machine should be configured with pairs of DIMMs. With triple-channel memory each machine should have triplets of DIMMs. Similarly, quad-channel DIMMs should be in groups of four. Table 6 shows the memory configurations that we recommend.

#### **Key point**

Interactive Hive, Spark, Storm workloads are generally more memory intensive workloads; hence, it is suggested that 256GB RAM be used on each DL380 server.

**Table 6.** Memory recommendations

<b>Memory</b>	<b>Description</b>
128GB – 8 x HP 16GB 2Rx4 PC4-2133P-R Kit	Base configuration
256GB – 16 x HP 16GB 2Rx4 PC4-2133P-R Kit	High capacity configuration

See [Appendix C: Alternate parts](#), Table C-2 for BOM details on alternate memory configurations.

### Storage options

For workloads such as ETL and similar long running queries where the amount of storage is likely to grow, it is recommended to pick higher capacity and faster drives. For a performance oriented solution, SAS drives are recommended as they offer a significant read and write throughput performance enhancement over SATA disks. Table 7 shows alternate storage options for the selected ProLiant DL380.

**Table 7.** HDD recommendations

HDD	Description
2/3/4TB LFF SATA	Base configuration
2/3/4TB LFF SAS	Performance configuration

See [Appendix C: Alternate parts](#), Table C-3 for BOM details of alternate hard disks.

### Network planning

Hadoop is sensitive to network speed in addition to CPU, memory and disk I/O, hence 10GbE is ideal, considering current and future needs for additional workloads. Switches with deep buffer caching are very useful. Network redundancy is a must for Hadoop. Hadoop shuffle phase does generate a lot of network traffic.

Generally accepted oversubscription ratios are around 2-4:1. Lower/higher oversubscription ratios can be considered if higher performance is required. A separate white paper on networking best practices for Hadoop is available at <http://h20195.www2.hp.com/V3/GetDocument.aspx?docname=4AA5-3279ENW>

See [Appendix C: Alternate parts](#), Table C-4 for BOM details of alternate network cards.

### Key point

HP recommends all ProLiant systems be upgraded to the latest BIOS and firmware versions before installing the OS. HP Service Pack for ProLiant (SPP) is a comprehensive systems software and firmware update solution, which is delivered as a single ISO image. The minimum SPP version recommended is 2015.04.0 (B). This updated version includes a fix for the OpenSSL Heartbleed Vulnerability. The latest version of SPP can be obtained from: [http://h18004.www1.hp.com/products/servers/service\\_packs/en/index.html](http://h18004.www1.hp.com/products/servers/service_packs/en/index.html)

### Vertica and Hadoop

Relational database management systems such as HP Vertica excel at analytic processing for big volumes of structured data including call detail records, financial tick streams and parsed weblog data. HP Vertica is designed for high speed load and query when the database schema and relationships are well defined. Hortonworks Distribution for Hadoop, built on the popular open source Apache Software Foundation project, addresses the need for large-scale batch processing of unstructured or semi-structured data. When the schema or relationships are not well defined, Hadoop can be used to employ massive MapReduce style processing to derive structure out of data. The Hortonworks Distribution simplifies installation, configuration, deployment and management of the powerful Hadoop framework for enterprise users.

Each can be used standalone – HP Vertica for high-speed loads and ad-hoc queries over relational data, Hortonworks Distribution for general-purpose batch processing, for example from log files. Combining Hadoop and Vertica creates a nearly infinitely scalable platform for tackling the challenges of big data.

### Note

Vertica was the first analytic database company to deliver a bi-directional Hadoop Connector enabling seamless integration and job scheduling between the two distributed environments. With Vertica's Hadoop and Pig Connectors, users have unprecedented flexibility and speed in loading data from Hadoop to Vertica and querying data from Vertica in Hadoop as part of MapReduce jobs for example. The Vertica Hadoop and Pig Connectors are supported by Vertica, and available for download.

For more information, please see [vertica.com/the-analytics-platform/native-bi-etl-and-hadoop-mapreduce-integration/](http://vertica.com/the-analytics-platform/native-bi-etl-and-hadoop-mapreduce-integration/)

### HP IDOL and Hadoop

HP IDOL for Hadoop from HP Autonomy is designed to make Hadoop data come alive by helping manage, control, and enable greater insight so businesses can profit from information. Based on the technology of HP IDOL 10 (Intelligent Data Operating Layer), a market-leading information analytics platform currently being used to solve the toughest information challenges across a wide range of industries, HP IDOL for Hadoop tightly integrates with Hadoop to enrich its unstructured content capabilities. By using HP IDOL for Hadoop, components in an existing Hadoop environment and a programming framework (MapReduce), you can reduce coding and implementation costs while performing complex analysis that improve business outcomes. Data analysis techniques include inquire (search), investigate (explore), interact (engage), and improve (filter) data at scale to extract timely and actionable insights from big data in Hadoop.

HP IDOL for Hadoop leverages the market-leading file filtering technology of HP IDOL KeyView to extract text, metadata, and security information from all forms of data. The product can then apply its entity extraction function to identify sensitive fields such as social security and credit card numbers. For instance, if you are a healthcare provider with hundreds of thousands of patient forms stored in your Hadoop farms, you can ensure that no sensitive information is being stored. HP IDOL for Hadoop helps you remain compliant with your company policies, as well as HIPAA, PCI DSS, Sarbanes-Oxley, and more.

For more information, please see [autonomy.com/idolhadoop](http://autonomy.com/idolhadoop)

### Use cases

Figure 8 shows the use cases for IDOL and Hadoop.

**Figure 8.** Use cases for IDOL and Hadoop

Finance	Government	Telecom	Manufacturing	Energy	Healthcare
Fraud detection	Fraud detection	Broadcast monitoring	Supply chain optimization	Weather forecasting	Drug development
Anti-money laundering	Anti-money laundering	Churn prevention	Defect tracking	Natural resource exploration	Scientific research
Risk management	Risk management	Advertising optimization	RFID Correlation		Evidence based medicine
			Warranty management		Healthcare outcomes analysis

Horizontal Use Cases

Sentiment analysis	Marketing campaign optimization	Logistics optimization
Social CRM / network analysis	Brand management	Clickstream analysis
Churn mitigation	Social media analytics	Influencer analysis
Brand monitoring	Pricing optimization	IT infrastructure analysis
Cross and Up sell	Internal risk assessment	Legal discovery
Loyalty & promotion analysis	Customer behavior analysis	Equipment monitoring
Web application optimization	Revenue assurance	Enterprise search

## Guidance deployment steps

### Server selection

This section will provide topologies for the deployment of management and worker nodes for single and multi-rack clusters. Depending on the size of the cluster, a Hadoop deployment consists of one or more nodes running management services and a quantity of worker nodes. We have designed this reference architecture so that regardless of the size of the cluster, the server used for the management nodes and the server used for the worker nodes remains consistent. This section specifies which server to use and the rationale behind it.

### Management nodes

Management services are not distributed redundantly across as many nodes as the services that run on the worker nodes and therefore benefit from a server that contains redundant fans and power supplies. In addition, Management nodes require storage only for local management services and the OS, unlike worker nodes that store data in HDFS, and so do not require large amounts of internal storage. However, as these local management services are not replicated across multiple



servers, an array controller supporting a variety of RAID schemes and SAS direct attached storage is required. In addition, the management services are memory and CPU intensive; therefore, a server capable of supporting a large amount of memory is also required.

---

### Best practice

HP recommends that all Management nodes in a cluster be deployed with either identical or highly similar configurations. The configurations reflected in this white paper are also cognizant of the high availability feature in Hortonworks HDP 2.2. For this feature, servers should have similar I/O subsystems and server profiles so that each management server could potentially take the role of another. Similar configurations will also ensure that ZooKeeper's quorum algorithm is not affected by a machine in the quorum that cannot make a decision as fast as its quorum peers.

---

This section contains 4 subsections:

- Server platform
- Management node
- ResourceManager server
- NameNode server

### Server platform: HP ProLiant DL360 Gen9

The HP ProLiant DL360 Gen9 (1U), shown in Figure 9 below, is an excellent choice as the server platform for the management nodes and head nodes.

**Figure 9.** HP ProLiant DL360 Gen9 Server



#### Processor configuration

The configuration features two sockets with 8 core processors of the Intel E5-2600 v3 product family, which provide 16 physical cores and 32 Hyper-Threaded cores per server. We recommend that Hyper-Threading be turned on.

The reference architecture was tested using the Intel Xeon E5-2640 v3 processors for the management servers with the ResourceManager, NameNode and Ambari services. The configurations for these servers are designed to be able to handle an increasing load as your Hadoop cluster grows. Choosing a powerful processor such as this to begin with sets the stage for managing growth seamlessly. An alternative CPU option is a 10-core E5-2660 v3 for head nodes to support large numbers of services/processes.

#### Drive configuration

The Smart Array P440ar Controller is specified to drive eight 900GB 2.5" SAS disks on the Management node, ResourceManager and NameNode servers. Hot pluggable drives are specified so that drives can be replaced without restarting the server. Due to this design, one should configure the P440ar controller to apply the following RAID schemes:

- Management node: 8 disks with RAID 1+0 for OS and PostgreSQL database, and management stack software.
- ResourceManager and NameNode Servers: 8 disks with RAID 1+0 for OS and Hadoop software.

---

### Best practice

For a performance oriented solution HP recommends SAS drives as they offer a significant read and write performance enhancement over SATA disks. The Smart Array P440ar controller provides two port connectors per controller with each containing 4 SAS links. The drive cage for the DL360 Gen9 contains 8 disks slots and thus each disk slot has a dedicated SAS link which ensures the server provides the maximum throughput that each drive can give you.

---

*Memory configuration*

Servers running management services such as the HBaseMaster, ResourceManager, NameNode and Ambari should have sufficient memory as they can be memory intensive. When configuring memory, one should always attempt to populate all the memory channels available to ensure optimum performance. The dual Intel Xeon E5-2600 v3 series processors in the HP ProLiant DL380 Gen9 have 4 memory channels per processor which equates to 8 channels per server. The configurations for the management and head node servers were tested with 128GB of RAM, which equated to eight 16GB DIMMs.

**Best practice**

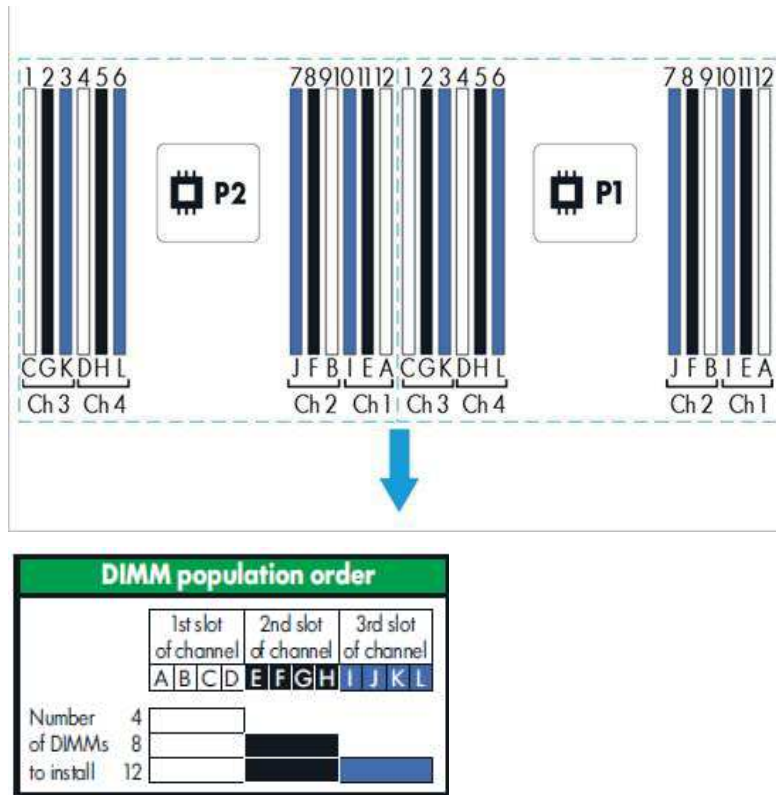
Configure all the memory channels according to recommended guidelines to ensure optimal use of memory bandwidth.

For example, on a two socket processor with eight memory channels available per server one would typically populate channels with 16GB DIMMs resulting in a configuration of 128GB in sequential alphabetical order balanced between the two processors: P1-A, P2-A, P1-B, P2-B, P1-C, P2-C, P1-D, P2-D, as shown in Figure 10.

Each Intel Xeon E5-2600 v3 family processor socket contains four memory channels per installed processor with three DIMMs per channel for a total of twelve (12) DIMMs or a grand total of twenty-four (24) DIMMs for the server. Figure 9 shows the memory configuration.

- There are four (4) Memory channels per processor; eight (8) channels per 2 processor server.
- There are three DIMM slots for each memory channel; twenty-four total slots for 2 processor server.
- Memory channels 1 and 3 consist of the three DIMMs that are furthest from the processor.
- Memory channel 2 and 4 consist of the three DIMMs that are closest to the processor.

**Figure 10.** Memory configuration



### Network configuration

The HP ProLiant DL380 Gen9 is designed for network connectivity to be provided via a FlexibleLOM. The FlexibleLOM can be ordered as a 2 x 10GbE NIC configuration. This Reference Architecture was tested using the 2 x 10GbE NIC configuration (as specified in the server configuration below).

### Best practice

For each management server HP recommends bonding and cabling two 10GbE NICs to create a single bonded pair which will provide 20GbE of throughput as well as a measure of NIC redundancy. In the reference architecture configurations later in the document you will notice that we use two IRF Bonded switches. In order to ensure the best level of redundancy we recommend cabling NIC 1 to Switch 1 and NIC 2 to Switch 2.

### Management node components

The Management node hosts the applications that submit jobs to the Hadoop Cluster. We recommend that you install with the software components shown in Table 8.

**Table 8.** Management node basic software components

Software	Description
Red Hat Enterprise Linux 6.6	Recommended Operating System
HP Insight CMU 7.3	Infrastructure Deployment, Management, and Monitoring
Oracle JDK 1.7.0_45	Java Development Kit
PostgreSQL 8.4	Database Server for Ambari
Ambari 1.6.x	Hortonworks Hadoop Cluster Management Software
Hue Server	Web Interface for Applications
NameNode HA	NameNode HA (Journal Node)
Apache Pig and Apache Hive	Analytical interfaces to the Hadoop Cluster
HiveServer2	Hue application to run queries on Hive with authentication
ZooKeeper	Cluster coordination service

Please see the following link for the Ambari Installation guide [http://docs.hortonworks.com/HDPDocuments/Ambari-1.7.0.0/Ambari\\_Doc\\_Suite/ADS\\_v170.html#Installing\\_HDP\\_Using\\_Ambari](http://docs.hortonworks.com/HDPDocuments/Ambari-1.7.0.0/Ambari_Doc_Suite/ADS_v170.html#Installing_HDP_Using_Ambari)

The management node and head nodes, as tested in the reference architecture, contain the following base configuration:

- 2 x Eight-Core Intel Xeon E5-2640 v3 Processors
- Smart Array P440ar Controller with 2GB FBWC
- 7.2TB – 8 x 900GB SFF SAS 10K RPM disks
- 128GB DDR4 Memory – 8 x HP 16GB 2Rx4 PC4-2133P-R Kit
- 10GbE 2P NIC 561FLR-T card adapter

A BOM for the Management node is available in the BOM section of this white paper.

### Head node 1 – ResourceManager server

The ResourceManager server contains the following software components. Table 9 lists the server base hardware components. Please see the following link for more information on installing and configuring the ResourceManager and

NameNode HA: [http://docs.hortonworks.com/HDPDocuments/HDP2/HDP-2.1.5/bk\\_system-admin-guide/content/ch\\_hadoop-ha-rm.html](http://docs.hortonworks.com/HDPDocuments/HDP2/HDP-2.1.5/bk_system-admin-guide/content/ch_hadoop-ha-rm.html).

**Table 9.** ResourceManager Server Base Software components

Software	Description
Red Hat Enterprise Linux 6.6	Recommended Operating System
Oracle JDK 1.7.0_45	Java Development Kit
ResourceManager	YARN ResourceManager
NameNode HA	NameNode HA (Failover Controller, Journal Node, NameNode Standby)
Oozie	Oozie Workflow scheduler service
HBaseMaster	The HBase Master for the Hadoop Cluster (Only if running HBase)
ZooKeeper	Cluster coordination service
Flume	Flume

### Head node 2 – NameNode server

The NameNode server contains the following software components. Table 10 lists the NameNode server base software components. Please see the following link for more information on installing and configuring the NameNode and ResourceManager HA. [http://docs.hortonworks.com/HDPDocuments/HDP2/HDP-2.1.5/bk\\_system-admin-guide/content/ch\\_hadoop-ha-1.html](http://docs.hortonworks.com/HDPDocuments/HDP2/HDP-2.1.5/bk_system-admin-guide/content/ch_hadoop-ha-1.html)

**Table 10.** NameNode server base software components

Software	Description
Red Hat Enterprise Linux 6.6	Recommended Operating System
Oracle JDK 1.7.0_45	Java Development Kit
NameNode	The NameNode for the Hadoop Cluster (NameNode Active, Journal node, Failover Controller)
JobHistoryServer	Job History for ResourceManager
ResourceManager HA	Failover, Passive mode
Flume	Flume agent (if required)
HBMaster	HBase Master (Only if running HBase)
ZooKeeper	Cluster coordination service

## Worker nodes

The worker nodes run the NodeManager and YARN container processes and thus storage capacity and performance are important factors.

### Server platform: HP ProLiant DL380 Gen9

The HP ProLiant DL380 Gen9 (2U), shown below in Figure 11, is an excellent choice as the server platform for the worker nodes. For ease of management we recommend a homogenous server infrastructure for your worker nodes.

**Figure 11.** HP ProLiant DL380 Gen9 Server



#### Processor selection

The configuration features two processors from the Intel Xeon E5-2600 v3 family. The Base configuration provides 20 physical or 40 Hyper-Threaded cores per server. Hadoop manages the amount of work each server is able to undertake via the ResourceManager configured for that server. The more cores available to the server, the better for ResourceManager Utilization (see the [Computation](#) section for more detail). We recommend that Hyper-Threading be turned on. For this RA, we chose 2 x E5-2660 v3 (10 cores/2.6GHz) CPUs.

#### Memory selection

Servers running the worker node processes should have sufficient memory for either HBase or for the amount of MapReduce Slots configured on the server. The dual Intel Xeon E5-2600 v3 series processors in the HP ProLiant DL380 Gen9 have 4 memory channels per processor which equates to 8 channels per server. When configuring memory, one should always attempt to populate all the memory channels available to ensure optimum performance.

With the advent of YARN in HDP 2.2, the memory requirement has gone up significantly to support a new generation of Hadoop applications. A base configuration of 128GB is recommended, and for certain high memory capacity applications 256GB is recommended. For this RA, we chose 256GB memory (16 x HP 16GB 2Rx4 PC4-2133P-R Kit).

---

### Best practice

To ensure optimal memory performance and bandwidth, HP recommends using 16GB DIMMs to populate each of the 4 memory channels on both processors which will provide an aggregate of 128GB of RAM. For 256GB capacity, we recommend adding an additional 8 x 16GB DIMMs, one in each memory channel. For any applications requiring more than 256GB capacity, we recommend going with 32GB DIMMs, not to populate the third slot on the memory channels to maintain full memory channel speed.

---

#### Drive configuration

Redundancy is built into the Apache Hadoop architecture and thus there is no need for RAID schemes to improve redundancy on the worker nodes as it is all coordinated and managed by Hadoop. Drives should use a Just a Bunch of Disks (JBOD) configuration, which can be achieved with the H240ar Smart Host Bus Adapter or HP Smart Array P440ar controller by configuring each individual disk as a separate RAID 0 volume. Additionally array acceleration features on the P440ar should be turned off for the RAID 0 data volumes. The first two positions on the rear drive cage with 600GB disks allow the OS to be placed in RAID1.

The H240ar Smart Host Bus Adapter or HP Smart Array P440ar controller provides two port connectors per controller with each containing 4 SAS links. For a performance oriented solution, we recommend SAS drives as they offer a significant read and write throughput performance enhancement over SATA disks. An HP 12G SAS Expander Card is needed to access more than 8 internal drives. For this RA, we chose 2TB LFF SATA drives.

**Best practice**

For OS drives, HP recommends using two 600GB SATA MDL LFF disks with the H240ar Smart Host Bus Adapter or HP Smart Array P440ar controller configured as one RAID1 mirrored 600GB logical drive for the OS. The other 13 disks are each configured for HDFS as 2TB in RAID0, for a total of 13 logical drives for Hadoop data. Protecting the OS provides an additional measure of redundancy on the worker nodes.

*DataNode settings*

By default, the failure of a single `dfs.data.dir` or `dfs.datanode.data.dir` will cause the HDFS DataNode process to shut down, which results in the NameNode scheduling additional replicas for each block that is present on the DataNode. This causes needless replications of blocks that reside on disks that have not failed. To prevent this, you can configure DataNodes to tolerate the failure of `dfs.data.dir` or `dfs.datanode.data.dir` directories; use the `dfs.datanode.failed.volumes.tolerated` parameter in `hdfs-site.xml`. For example, if the value for this parameter is 3, the DataNode will only shut down after four or more data directories have failed. This value is respected on DataNode startup; in this example the DataNode will start up as long as no more than three directories have failed.

**Note**

For configuring YARN, update the default values of the following attributes with ones that reflect the cores and memory available on a worker node.

```
yarn.nodemanager.resource.cpu-vcores
yarn.nodemanager.resource.memory-mb
```

Similarly, specify the appropriate size for map and reduce task heap sizes using the following attributes:

```
mapreduce.map.java.opts
mapreduce.reduce.java.opts
```

*Network configuration*

For 10GbE networks, we recommend that the two 10GbE NICs be bonded to improve throughput performance to 20GbE/s and thereby improve performance. In addition, in the reference architecture configurations later on in the document you will notice that we use two IRF Bonded switches. In order to ensure the best level of redundancy we recommend cabling NIC 1 to Switch 1 and NIC 2 to Switch 2.

**Worker node components**

Table 11 lists the Worker Node software components. Please see the following link for more information on installing and configuring the NodeManager (or HBaseRegionServer) and DataNode manually: For adding worker nodes: [http://docs.hortonworks.com/HDPDocuments/HDP2/HDP-2.1.5/bk\\_system-admin-guide/content/admin\\_add-nodes-2.html](http://docs.hortonworks.com/HDPDocuments/HDP2/HDP-2.1.5/bk_system-admin-guide/content/admin_add-nodes-2.html).

For adding HBase RegionServer: [http://docs.hortonworks.com/HDPDocuments/HDP2/HDP-2.1.5/bk\\_system-admin-guide/content/admin\\_add-nodes-3.html](http://docs.hortonworks.com/HDPDocuments/HDP2/HDP-2.1.5/bk_system-admin-guide/content/admin_add-nodes-3.html)

**Table 11.** Worker Node base software components

Software	Description
Red Hat Enterprise Linux 6.6	Recommended Operating System
Oracle JDK 1.7.0_45	Java Development Kit
NodeManager	The NodeManager process for MR2/YARN
DataNode	The DataNode process for HDFS
<i>HBaseRegionServer</i>	<i>The HBaseRegionServer for HBase (Only if running HBase)</i>

The ProLiant DL380 Gen9 (2U) as configured for the reference architecture as a worker node has the following configuration:

- Dual 10-Core Intel Xeon E5-2660 v3 Processors with Hyper-Threading
- Thirteen 2TB 3.5" 7.2K LFF SATA SC MDL (26TB for Data)
- Two HP 600GB 12G SAS 15K 3.5in ENT SCC HDD (for OS)
- HP DL380 Gen9 3LFF Rear SAS/SATA Kit (2x 600GB for OS and 1x 2TB for Data)
- 256GB DDR4 Memory (16 x HP 16GB), 4 channels per socket
- 1 x 10GbE 2 Port NIC FlexibleLOM (Bonded)
- 1 x H240ar Smart Host Bus Adapter or Smart Array P440ar Controller with 2GB FBWC

---

### Note

Customers also have the option of purchasing a second power supply for additional power redundancy. This is especially appropriate for single rack clusters where the loss of a node represents a noticeable percentage of the cluster.

---

The BOM for the Worker node is provided in BOM section of this white paper.

HP iLO Management Engine is a complete set of embedded features, standard on all ProLiant Gen9 servers. The engine is easy to use and includes HP iLO, HP Agentless Management, HP Active Health System, HP Intelligent Provisioning, and HP Embedded Remote Support.

HP Insight Online with HP Insight Remote Support provides 24x7 remote monitoring and anywhere, anytime personalized access to your IT and support status.

HP Smart Update, including HP Smart Update Manager (HP SUM), HP Service Pack for ProLiant (SPP) and other products, reduces deployment time and update complexity by systematically and securely updating server infrastructure in the data center; in most cases the downtime is limited to a single reboot.

## Networking

Hadoop clusters contain two types of switches, namely Top of Rack (ToR) switches and Aggregation switches. Top of Rack switches route the traffic between the nodes in each rack and Aggregation switches route the traffic between the racks.

### Switch selection

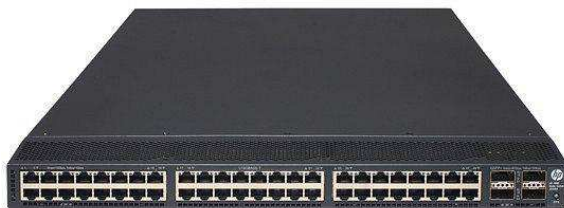
#### Top of Rack (ToR) switches

The HP 5900AF-48XGT-4QSFP+, 10GbE, is an ideal ToR switch with forty eight 10GbE ports and four 40GbE uplinks providing resiliency, high availability and scalability support. In addition this model comes with support for CAT6 cables (copper wires) and Software Defined Networking (SDN). A dedicated management switch for iLO traffic is not required as the ProLiant DL360 Gen9 and DL380 Gen9 are able to share iLO traffic over NIC 1. The volume of iLO traffic is minimal and does not degrade performance over that port. For more information on the 5900AF-48XGT-4QSFP+, 10GbE switch, please see [hp.com/networking/5900](http://hp.com/networking/5900)

The BOM for the HP 5900AF-48XGT switch is provided in the BOM section of this white paper.

Customers who would like to separate iLO and PXE traffic from the data/hadoop network traffic can add 1GbE HP 5900AF-48G-4XG-2QSFP+ Switch (JG510A, shown in Figure 12) network switch. Modify the BOM to include HP Ethernet 10Gb 2-port 561T Adapter NIC instead of HP Ethernet 10Gb 2P 561FLR-T FIO Adptr FlexibleLOM NIC to separate iLO and PXE network traffic on all the systems. The BOMs for 1GbE switch and network cards are provided in the BOM section of this white paper.

**Figure 12.** 5900AF-48G-4XG-2QSFP+ switch



### Aggregation switches

The HP FlexFabric 5930-32QSFP+, 40GbE, switch is an ideal aggregation switch as it is well suited to handle very large volumes of inter-rack traffic such as can occur during shuffle and sort operations, or large scale block replication to recreate a failed node. The switch has better connectivity with 32 40GbE ports, supporting up to 104 ports of 10GbE via breakout cables and six 40GbE uplink ports, aggregation switch redundancy and high availability (HA) support with IRF bonding ports. SDN ready with OpenFlow 1.3 and overlay networks with VXLAN and NVGRE support. Figure 13 shows the HP 5930 Aggregation switch. For more information on the HP 5930-32QSFP+ please see [hp.com/networking/5930](http://hp.com/networking/5930)

The BOM for the HP 5930-32QSFP+ switch is provided in the BOM section of this white paper.

**Figure 13.** HP 5930-32QSFP+ 40GbE Aggregation switch



## Management

### HP Insight Cluster Management Utility

HP Insight Cluster Management Utility (CMU) is an efficient and robust hyper-scale cluster lifecycle management framework and suite of tools for large Linux clusters such as those found in High Performance Computing (HPC) and Big Data environments. A simple graphical interface enables an “at-a-glance” real-time or 3D historical view of the entire cluster for both infrastructure and application (including Hadoop) metrics, provides frictionless scalable remote management and analysis, and allows rapid provisioning of software to all nodes of the system. HP Insight CMU makes the management of a cluster more user friendly, efficient, and error free than if it were being managed by scripts, or on a node-by-node basis. HP Insight CMU offers full support for iLO 2, iLO 3, iLO 4 and LO100i adapters on all ProLiant servers in the cluster.

---

### Best practice

HP recommends using HP Insight CMU for all Hadoop clusters. HP Insight CMU allows one to easily correlate Hadoop metrics with cluster infrastructure metrics, such as CPU Utilization, Network Transmit/Receive, Memory Utilization and I/O Read/Write. This allows characterization of Hadoop workloads and optimization of the system thereby improving the performance of the Hadoop Cluster. CMU Time View Metric Visualizations will help you understand, based on your workloads, whether your cluster needs more memory, a faster network or processors with faster clock speeds. In addition, Insight CMU also greatly simplifies the deployment of Hadoop, with its ability to create a Golden Image from a node and then deploy that image to up to 4000 nodes. Insight CMU is able to deploy 800 nodes in 30 minutes.

---

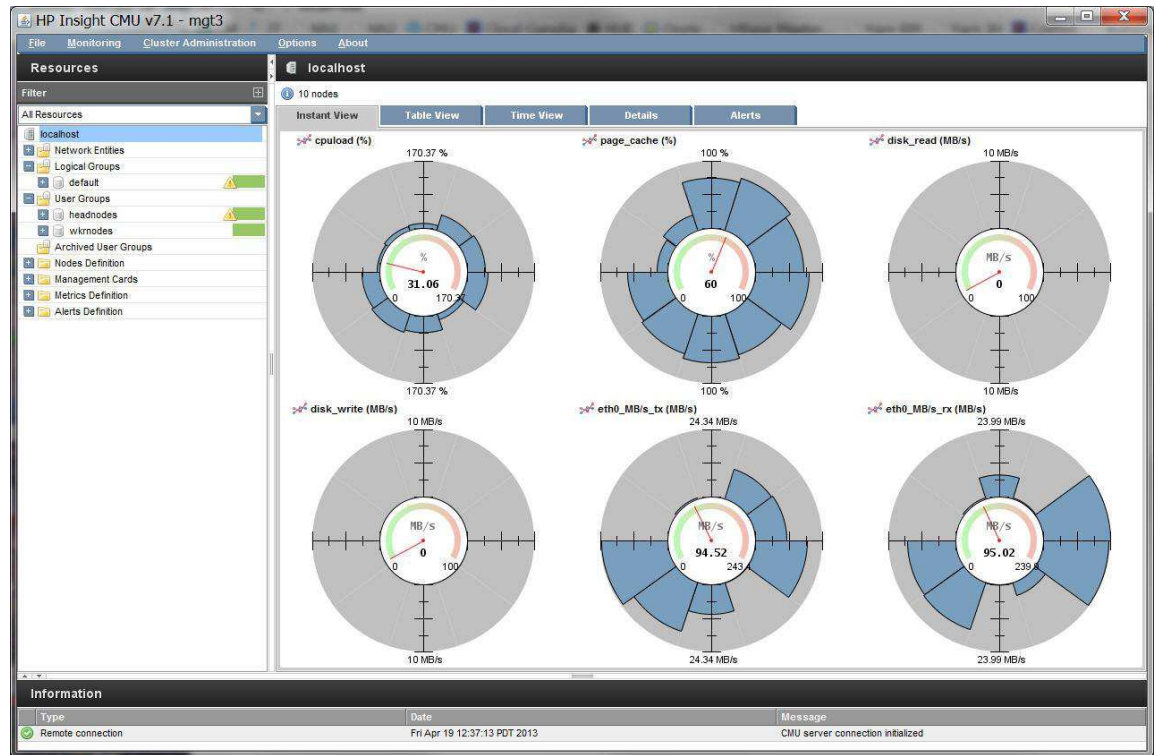
HP Insight CMU is highly flexible and customizable, offers both GUI and CLI interfaces, and can be used to deploy a range of software environments, from simple compute farms to highly customized, application-specific configurations. HP Insight CMU is available for HP ProLiant and HP BladeSystem servers, and is supported on a variety of Linux operating systems, including Red Hat Enterprise Linux, SUSE Linux Enterprise Server, CentOS, and Ubuntu. HP Insight CMU also includes options for monitoring graphical processing units (GPUs) and for installing GPU drivers and software. Figures 14, 15 and 16 show views of the HP Insight CMU.

For more information, please see [hp.com/go/cmu](http://hp.com/go/cmu).

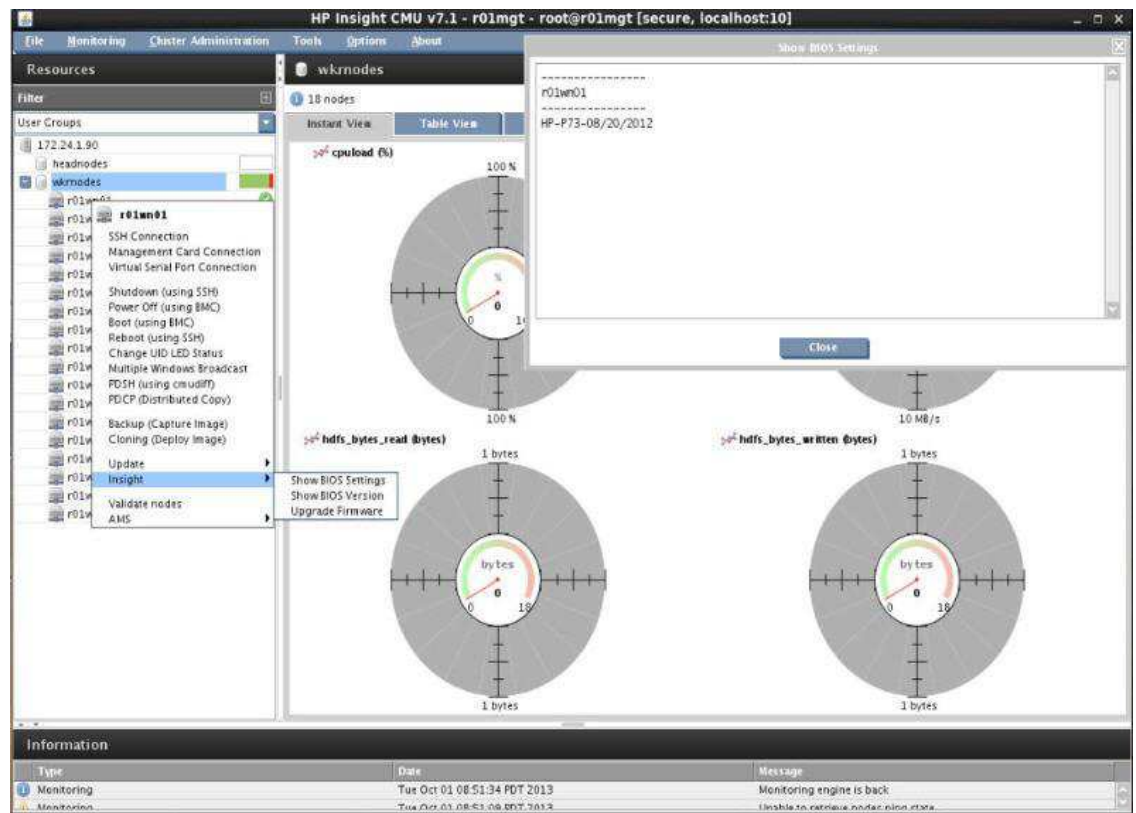
For the CMU BOM please see the BOM section of this white paper.



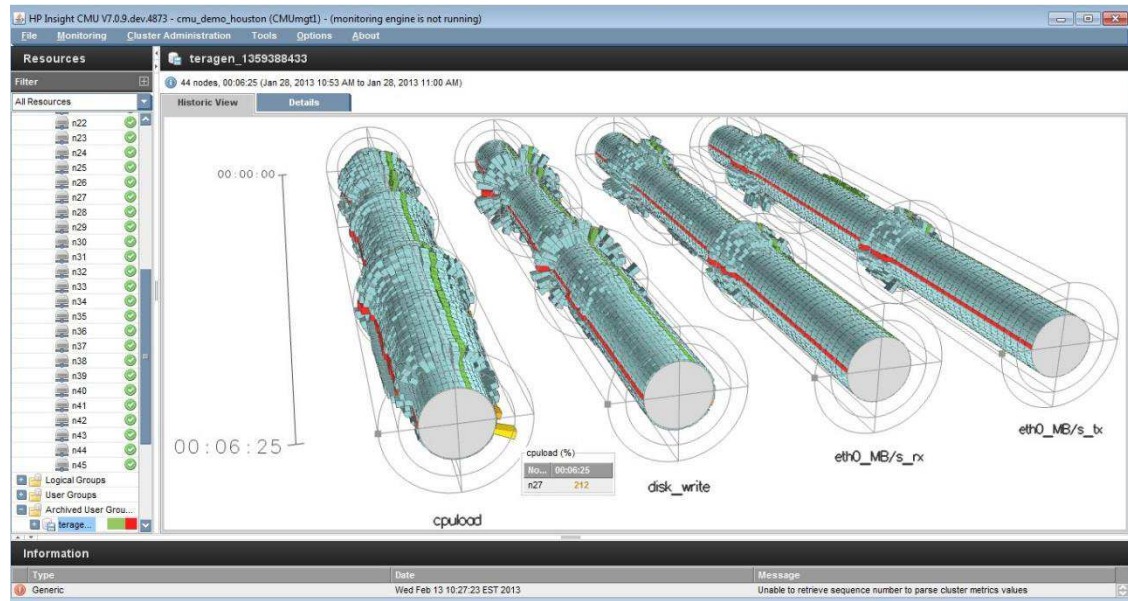
**Figure 14.** HP Insight CMU Interface – real-time view



**Figure 15.** HP Insight CMU Interface – real-time view – BIOS settings



**Figure 16.** HP Insight CMU Interface – Time View



**Key point**

CMU can be configured to support High Availability with an Active-Passive cluster

**Bill of materials**

The BOMs outlined in Tables 12-21 below are based on the tested configuration for a Single-Rack Reference Architecture with 1 management node, 2 Head nodes, 18 worker nodes and 2 ToR switches. Quantities specified on each table are based on per server configuration.

The following BOMs contain electronic license to use (E-LTU) parts. Electronic software license delivery is now available in most countries. HP recommends purchasing electronic products over physical products (when available) for faster delivery and for the convenience of not tracking and managing confidential paper licenses. For more information, please contact your reseller or an HP representative.

**Note**

Part numbers used are based on the part numbers that were used for testing and subject to change. The bill of materials does not include complete support options or other rack and power requirements. If you have questions regarding ordering, please consult with your HP Reseller or HP Sales Representative for more details.  
<http://www8.hp.com/us/en/business-services/it-services/it-services.html>

**Management Node and Head Node BOM**

**Table 12.** The HP ProLiant DL360 Gen9 Server Configuration

Qty	Part Number	Description
1	755258-B21	HP DL360 Gen9 8-SFF CTO Chassis
1	755386-L21	HP DL360 Gen9 E5-2640v3 FIO Kit
1	755386-B21	HP DL360 Gen9 E5-2640v3 Kit
8	726719-B21	HP 16GB 2Rx4 PC4-2133P-R Kit
8	652589-B21	HP 900GB 6G SAS 10K 2.5in SC ENT HDD

Qty	Part Number	Description
1	700699-B21	HP Ethernet 10Gb 2P 561FLR-T Adptr
1	749974-B21	HP Smart Array P440ar/2G FIO Controller
2	720478-B21	HP 500W FS Plat Ht Plg Pwr Supply Kit
1	SG506A	HP C13 - C14 WW 250V 10Amp IPD 0.76m 1pc Jumper Cord
1	SG508A	HP C13 - C14 WW 250V 10Amp IPD 1.37m 1pc Jumper Cord
1	663201-B21	HP 1U SFF Ball Bearing Rail Kit
1	C6N36ABE	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)
1	G3J28AAE	RHEL Svr 2 Sckt/2 Gst 1yr 24x7 E-LTU

### Worker Node BOM

**Table 13.** The HP ProLiant DL380 Gen9 Server Configuration

Qty	Part Number	Description
1	719061-B21	HP DL380 Gen9 12-LFF CTO Server
1	762764-L21	HP DL380 Gen9 E5-2660v3 FIO Kit
1	762764-B21	HP DL380 Gen9 E5-2660v3 Kit
16	726719-B21	HP 16GB 2Rx4 PC4-2133P-R Kit
13	652757-B21	HP 2TB 6G SATA 7.2k 3.5in SC MDL HDD
2	765424-B21	HP 600GB 12G SAS 15K 3.5in ENT SCC HDD
1	700699-B21	HP Ethernet 10Gb 2P 561FLR-T FIO Adptr
1	749976-B21	HP H240ar FIO Smart HBA
1	727250-B21	HP 12Gb DL380 Gen9 SAS Expander Card
1	768856-B21	HP DL380 Gen9 3LFF Rear SAS/SATA Kit
1	720864-B21	HP 2U LFF BB Gen9 Rail Kit
2	720479-B21	HP 800W FS Plat Ht Plg Pwr Supply Kit
1	SG506A	HP C13 - C14 WW 250V 10Amp IPD 0.76m 1pc Jumper Cord
1	SG508A	HP C13 - C14 WW 250V 10Amp IPD 1.37m 1pc Jumper Cord
1	C6N36ABE	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)
1	G3J28AAE	RHEL Svr 2 Sckt/2 Gst 1yr 24x7 E-LTU

**Network BOMs****Table 14.** Network – Top of Rack Switch

Qty	Part Number	Description
2	JG336A	HP 5900AF-48XGT-4QSFP+ Switch
2	JG326A	HP X240 40G QSFP+ QSFP+ 1m DAC Cable
4	JC680A	HP A58x0AF 650W AC Power Supply
4	JG553A	HP X712 Bck(pwr)-Frt(prt) HV Fan Tray

**Table 15.** Network – Aggregation/Spine Switch (Only required for first Expansion Rack. Not required for Single Rack Architecture.)

Qty	Part Number	Description
2	JG726A	HP FF 5930-32QSFP+
4	JC680A	HP 58x0AF 650W AC Power Supply
4	JG553A	HP X712 Bck(pwr)-Frt(prt) HV Fan Tray
2	JG326A	HP X240 40G QSFP+ QSFP+ 1m DAC Cabl
8	JG328A	HP X240 40G QSFP+ QSFP+ 5m DAC Cable

**10GbE PCI NIC to replace 10GbE FlexibleLOM NIC:****Table 16.** Modified BOM for HP ProLiant DL360/DL380 Gen9 Server Configuration

Qty	Part Number	Description
1	716591-B21	HP Ethernet 10Gb 2-port 561T Adapter

**Table 17.** Network – Top of Rack Switch for separate iLO and PXE network

Qty	Part Number	Description
1	JG510A	HP 5900AF-48G-4XG-2QSFP+ Switch
1	JC680A	HP A58x0AF 650W AC Power Supply
2	JC682A	HP A58x0AF Back (power side) to Front (port side) Airflow Fan Tray

**Other hardware and software BOMs****Table 18.** Hardware – Rack and PDU

**Note:** The quantity specified below is for a full rack with 2 switches, 3x DL360 and 18x DL380.

Qty	Part Number	Description
4	AF520A	HP Intelligent Mod PDU 24a Na/Jpn Core
8	AF547A	HP 5xC13 Intelligent PDU Extension Bar G2 Kit
1	BW946A	HP 42U Location Discovery Kit
1	BW904A	HP 642 1075mm Shock Intelligent Series Rack
1	BW932A	HP 600mm Rack Stabilizer Kit

Qty	Part Number	Description
1	BW930A	HP Air Flow Optimization Kit
1	BW906A	HP 42U 1075mm Side Panel Kit
1	BW891A	HP Rack Grounding Kit

**Table 19.** Software – HP Insight Cluster Management Utility (CMU) options**Note**

The quantity specified below is for a single node.

Qty	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

**Table 20.** Software – Red Hat Enterprise Linux

Qty	Part Number	Description
21	G3J28AAE	RHEL Svr 2 Sckt/2 Gst 1yr 24x7 E-LTU

**Note**

While HP is a certified reseller of Hortonworks software subscription, all application support (L1-L3) for Hortonworks software is provided by Hortonworks. The HP ProLiant DL380 platform is Hortonworks Certified.

**Table 21.** Software – Hortonworks Subscription options

Qty	Part Number	Description
5	F5Z52A	Hortonworks Data Platform Enterprise 4 Nodes or 50TB Raw Storage 1 year 24x7 Support LTU.

## Summary

HP and Hortonworks allow one to derive new business insights from Big Data by providing a platform to store, manage and process data at scale. However, designing and ordering Hadoop Clusters can be both complex and time consuming. This white paper provided several reference architecture configurations for deploying clusters of varying sizes with Hortonworks Data Platform 2.2 on HP infrastructure and management software. These configurations leverage HP balanced building blocks of servers, storage and networking, along with integrated management software and bundled support. In addition, this white paper has been created to assist in the rapid design and deployment of Hortonworks Data Platform software on HP infrastructure for clusters of various sizes.

## Implementing a proof-of-concept

As a matter of best practice for all deployments, HP recommends implementing a proof-of-concept using a test environment that matches as closely as possible the planned production environment. In this way, appropriate performance and scalability characterizations can be obtained. For help with a proof-of-concept, contact an HP Services representative (<http://www8.hp.com/us/en/business-services/it-services/it-services.html>) or your HP partner.

## Appendix A: Cluster design – heat map for server platforms

### Tiered compute/storage deployment

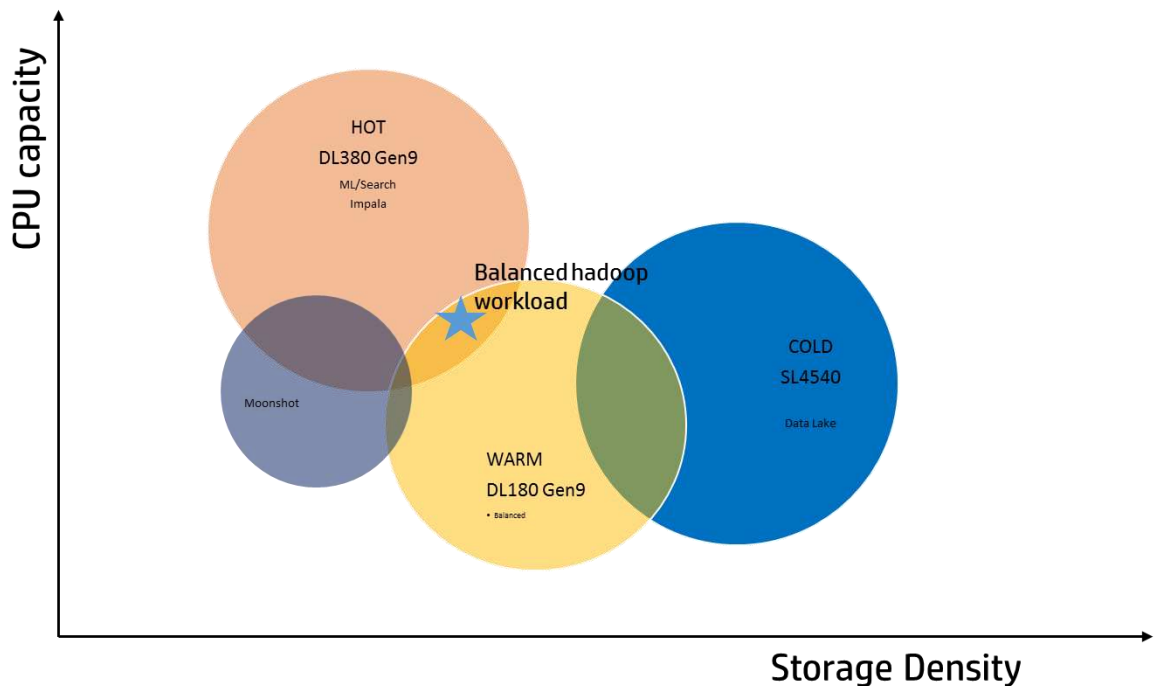
Generally speaking, a data center has tiered compute/storage deployments. For simplicity reasons, dissect the data center into three major tiers according to their computation characteristics: high performance (hot), high capacity (cold) and balanced (warm) zones. Please refer to Figure A-1 below.

**HP ProLiant SL4540 Gen8 Servers:** A highly efficient converged design that delivers the right combination of capacity and performance, in the least amount of space and at low cost, with the reliability and manageability you expect from HP ProLiant Gen8 servers. The ProLiant SL4540 Server has a converged and balanced architecture ideal for Hadoop and scale-out parallel processing applications. For workloads that require a balanced 1:1 spindle to core ratio, the HP ProLiant SL4540 Server 3x15 supports three Intel Xeon E5-2400 v2 processor-based nodes with up to 45 hot-plug hard drives in just 4.3U with the enterprise-class reliability and manageability of HP ProLiant. Visit <http://www8.hp.com/us/en/products/proliant-servers/product-detail.html?oid=6464822> for additional details.

**HP ProLiant DL180 Gen9 Server:** A cost-effective, enterprise-class server that delivers essential performance for data centers to meet their compute and storage needs today with scalability to grow as business requirements change. The HP ProLiant DL180 Gen9 Server provides ample storage capacity on demand with drive configurations ranging from (4) to (12) Large Form Factor (LFF) or (8) to (16) Small Form Factor (SFF) HP SmartDrives, with flexibility of LFF non-hot plug (NHP) options at the right economics. It supports up to two Intel Xeon processor E5-2600 v3 processors with up to 12 cores and 16 DIMM slots of HP DDR4 Smart Memory offering improved performance and efficiency. The platform is ideal for the most demanding Big Data applications like Apache Hadoop, which require the right mix of compute and storage. Visit <http://www8.hp.com/us/en/products/proliant-servers/product-detail.html?oid=7252820> for additional details.

**HP ProLiant m710 Server Cartridge (HP Moonshot):** HP Moonshot delivers the right compute power with breakthrough economics for your data center. And with server solutions designed for specific workloads, HP Moonshot helps you get more out of your infrastructure with less space, power consumption, and complexity. HP ProLiant m710 servers (HP Moonshot) are used as computational modules in HP Big Data Reference Architecture (BDRA). Visit [hp.com/go/moonshot](http://hp.com/go/moonshot) for additional details.

**Figure A-1.** Compute/storage tiers in the data center



**HP Big Data Reference Architecture:** HP BDRA is a modern, flexible 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. This big data reference architecture 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. Hosting big data on a single node can be an issue if a server goes down and the time to redistribute all data to the other nodes could be significant. 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. HP ProLiant SL4540 and m710 Server Cartridge (HP Moonshot) servers are used in BDRA architecture. For additional information <http://h20195.www2.hp.com/V2/GetDocument.aspx?docname=4AA5-6136ENW>

#### **High performance – interactive/SQL (hot)**

Compute power optimized systems such as the HP ProLiant DL380 are ideal candidates when price/performance with high I/O bandwidth workloads are of primary importance. The data crunching rate in this tier is extremely high, so the demand for CPU power in this tier is high, as well as the memory size. Typical applications in this tier could be ETL apps, real time streaming ([Spark Streaming](#)), NoSQL applications such as [MongoDB](#), [Cassandra](#), [Hbase](#), and in-memory processing such as [Spark](#). Additionally, machine learning apps like K-means clustering and search indexing apps such as [Solr](#) could make best use in this tier. Both require significant computing power. The DL380 Gen9 server is a perfect selection for this tier as it has high CPU clock rate and large memory configuration. The data footprint in this tier ranges from hundreds of gigabytes to multi terabytes. If the data volume grows over time, it may not be economical to store data in this tier, we recommend migrating infrequently accessed data (cold or warm data) to high capacity or balanced tiers.

#### **High capacity – storage optimized (cold)**

Storage optimized systems, such as HP ProLiant SL4540 servers, are an ideal consideration when low cost per terabyte is the prime factor. This is a perfect place for data warehousing analytics as the bigger the data volume is, the higher accuracy of the analytics outcome. Typical workloads are web log analytics, sentiment and click stream analysis. Applications such as Hive, Pig, and MapReduce often find their home in this tier. The SL4540 (3x15) server platform is an ideal candidate for this tier as it is slanted toward higher storage density to tackle the large data volume. The other variants of the SL4540, 2x25 and 1x60, are also very good for mirroring cold data for backup/DR purpose.

#### **Balanced performance and capacity – batch focused (warm)**

A mixture of workloads typically found in high performance and high capacity tiers are also seen in this tier. The compute and storage density requirement for the server is equally important. The data footprint in this tier is from multi terabytes to a couple of petabytes. The HP ProLiant DL180 Gen9 server is a good selection in this tier as its compute and storage architecture is more balanced. Typical workloads are similar to cold tier, such as web log analytics, sentiment and click stream analysis, but compute and storage needs are more balanced.

#### **Overlapped tiers**

As one can see in Figure A-1 there are no distinct boundaries between these tiers. When there is a variety or mixture of workloads a hybrid approach may be useful, a combination of server platforms can help fulfillment of the performance and data storage requirements. It is recommended to have early discussions with IT infrastructure and other related teams to fully discuss the implications of workload requirements vs. platform selection, in order to arrive at the right Hadoop system.

## Appendix B: Hadoop cluster tuning/optimization

### Server tuning

Below are some general guidelines for tuning the server OS and the storage controller for a typical Hadoop proof of concept (POC). Please note that these parameters are recommended for MapReduce workloads which are most prevalent in Hadoop environments. Please note that there is no silver bullet performance tuning. Modifications will be needed for other types of workloads.

#### a) OS tuning

- As a general recommendation, update to the latest patch level available to improve stability and optimize performance
- The recommended Linux file system is ext4, 64 bit OS:
  - Enable defaults, nodiratime, noatime (/etc/fstab)
  - Do not use logical volume management (LVM)
- Tune OS block readahead to 8K (/etc/rc/local):  
**blockdev --setra 8192 <storage device>**
- Turn off disk swappiness or to min=5:  
**Set sysctl vm.swappiness=0 in /etc/sysctl.conf**
- Tune ulimits for number of open files to a high number:  
Example: in /etc/security/limits.conf:  
**soft nofile 65536**  
**hard nofile 65536**
- Set **nproc = 65536**  
Add it to end of (/etc/security/limits.conf)
- Set IO scheduler policy to deadline on all the data drives  
**echo deadline > /sys/block/<device>/queue/scheduler**  
  
For persistency across boot, append the following to kernel boot line in /etc/grub.conf:  
**elevator=deadline**
- Configure network bonding on a minimum two 10GbE server ports, for up to a max 20GbE throughput.
- Ensure forward and reverse DNS is working properly.
- Install and configure ntp to ensure clocks on each node are in sync to management node.
- For good performance improvements, disable transparent huge page compaction:  
**echo never > /sys/kernel/mm/transparent\_hugepage/enabled**

#### b) Storage controller tuning

- Tune array controller stripe size to 1024MB:  
**hpssacli ctrl slot=<slot number> ld <ld number> modify ss=1024**
- Disable array accelerator(caching) (aa=disable):  
**hpssacli ctrl slot=<slot number> ld <ld number> modify aa=disable**

#### c) Power settings

Please note for a performance driven POC, we recommend using settings that help boost performance but could have negative impact on power consumption measurement:

- HP Power Profile → Maximum Performance
- HP Power regulator → Static High Performance mode
- Intel\_QPI\_Link\_Mgt → Disabled
- Min\_Proc\_Idle\_power\_Core\_State → No C-states
- Mem\_Power\_Saving → Max Perf
- Thermal Configuration → increased cooling



- Min\_Proc\_Idle Power Package state → No Package state
- Energy/Performance Bios → Disabled
- Collaborative Power Control → Disabled
- Dynamic Power Capping Functionality → Disabled
- DIMM Voltage Preference → Optimized for Performance

**d) CPU tuning**

The default BIOS settings for CPU should be adequate for most Hadoop workloads. Make sure that Hyper-Threading is turned on as it will help with additional performance gain.

**e) HP ProLiant BIOS**

- SPP version >= 2014.09.0 (B)
- Update System BIOS version to be >= P89
- Update Integrated Lights-Out (iLO) version to be >= 2.03
- Intel Virtualization Technology → Disabled
- Intel VT-d → Disabled

**f) HP Smart Array P440ar / Smart HBA H240ar**

- Update controller firmware to be >= v1.34
- Configure each Hadoop data drive as a separate RAID 0 array with stripe size of 1024KB
- Turn Off "Array Acceleration" / "Caching" for all data drives

Example: (with two controllers)

- ctrl.slot=0 ld all modify caching=disable ← disable caching on all logical drives on 1<sup>st</sup> ctrlr
- ctrl.slot=0 ld1 modify caching=enable ← enable caching on the OS logical drive on 1<sup>st</sup> ctrlr

**g) Network cards**

- Ethernet driver ixgbe, version >= 3.23.0.79 and firmware version >= 0x800005ac, 1.949.0

**h) Oracle Java**

- java.net.preferIPv4Stack set to true

**i) Patch common security vulnerabilities**

- Bash 9740 'Shellshock' fix for bash shell vulnerability as per instructions on <https://access.redhat.com/articles/1200223>
- For Heartbleed vulnerability, all versions of OpenSSL 1.0.1 prior to 1.0.1g need to be updated to 1.0.1g. <https://access.redhat.com/solutions/781793>

## Appendix C: Alternate parts

**Table C-1.** Alternate Processors – DL380

Qty/Node	Part Number	Description
1	719048-L21	HP DL380 Gen9 E5 2650v3 FIO Kit 10 cores at 2.3GHz
1	719048-B21	HP DL380 Gen9 E5 2650v3 Kit
1	762764-L21	HP DL380 Gen9 E5 2660v3 FIO Kit 10 cores at 2.6GHz
1	762764-B21	HP DL380 Gen9 E5 2660v3 Kit
1	762766-L21	HP DL380 Gen9 E5 2680v3 FIO Kit 12 cores at 2.5GHz
1	762766-B21	HP DL380 Gen9 E5 2680v3 Kit

**Table C-2.** Alternate Memory – DL380

Qty/Node	Part Number	Description
8	728629-B21	HP 32GB 2Rx4 PC4-2133P-R Kit for 256GB of Memory
16	728629-B21	HP 32GB 2Rx4 PC4-2133P-R Kit for 512GB of Memory

**Table C-3.** Alternate Disk Drives – DL380

Qty/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
12	658079-B21	HP 2TB 6G SATA 7.2k 3.5in SC MDL HDD
12	628061-B21	HP 3TB 6G SATA 7.2k 3.5in SC MDL HDD
12	693687-B21	HP 4TB 6G SATA 7.2k 3.5in SC MDL HDD

**Table C-4.** Alternate Network Cards – DL380

Qty/Node	Part Number	Description
1	665243-B21	HP Ethernet 10GbE 560FLR SFP+ FIO Adapter for 10Gb networking only
1	652503-B21	HP Ethernet 10Gb 2P 530SFP+ Adapter

### Note

The SFP+ network cards are used with DAC cabling and will not work with CAT6 cabling. If SFP+ network cards are used the 5900 SFP+ equivalent ToR network switches are required (HP 5900AF-48XG-4QSFP+ Part number JC772A).

**Table C-5.** Alternate Controller Cards – DL380 and DL360

Qty/Node	Part Number	Description
1	726821-B21	HP Smart Array P440/4GB FBWC 12Gb 1-port Int SAS Controller
1	726897-B21	HP Smart Array P840/4GB FBWC 12Gb 2-port Int SAS Controller
1/Drive	D8S84A	HP Secure Encryption No Media Flexible License per Drive
1/Drive	D8S85AAE	HP Secure Encryption No Media E-LTU per Drive

## Appendix D: HP value-added services and support

In order to help customers jump-start their Hadoop solution development, HP offers several Big Data services, including Factory Express and Technical Services (TS) Consulting. 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. TS Consulting offers specialized Hadoop design, implementation, and installation and setup services. HP offers a variety of support levels to meet your needs.

### 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. Please engage TS Consulting for details and quoting assistance.

### TS Consulting – Reference Architecture Implementation Service for Hadoop

With HP Reference Architecture Implementation Service for Hadoop, experienced HP Big Data consultants install, configure, deploy, and test your Hadoop environment based on the HP Reference Architecture. We'll 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, we'll work with you to configure the environment according to your goals and needs. We'll also conduct an acceptance test to validate and prove that the system is operating to your satisfaction.

### TS Consulting – Big Data services

HP Big Data Services can help you reshape your IT infrastructure to corral these increasing volumes of bytes – from e-mails, social media, and website downloads – and convert them into beneficial information. Our Big Data solutions encompass strategy, design, implementation, protection and compliance. We deliver these solutions in three steps.

1. **Big Data Architecture Strategy:** We'll define 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. **Big Data System Infrastructure:** HP experts will design and implement a high-performance, integrated platform to support a strategic architecture for Big Data. Choose form 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. **Big Data Protection:** Ensure availability, security and compliance of Big Data systems. Our consultants can help you safeguard your data, achieve regulatory compliance and lifecycle protection across your Big Data landscape, as well as improve your backup and continuity measures.

For additional information, please visit:

[hp.com/services/bigdata](http://hp.com/services/bigdata)

## HP Support options

HP offers a variety of support levels to meet your needs. HP Proactive Care helps prevent problems, resolve problems faster, and improve productivity. It helps customers identify and address IT problems before they cause performance issues or outages through analysis, reports, and update recommendations. Customers experiencing any performance issues are rapidly connected to experts for faster resolution. HP recommends adding HP Personalized Support option and HP Proactive Select with Proactive Care. Personalized support option is delivered by a local Account Support Manager who helps customers plan support options specific to their environment, deliver services, and review results. Proactive Select credits can be purchased upfront and used for optimization and improvement services related to health-checks, availability, and firmware updates throughout the year.

HP Datacenter Care provides a more personalized, customized approach for large, complex environments, with one solution for reactive, proactive, and multi-vendor support needs.

### *HP Support Plus 24*

For a higher return on your server and storage technology, our combined reactive support service delivers integrated onsite hardware/software support services available 24x7x365, including access to HP technical resources, 4-hour response onsite hardware support and software updates.

### *HP Proactive Care*

**HP Proactive Care** – HP Proactive Care begins with providing all of the benefits of proactive monitoring and reporting along with rapid reactive care. You also receive enhanced reactive support, through access to HP's expert reactive support specialists. You can customize your reactive support level by selecting either 6 hour call-to-repair or 24x7 with 4 hour onsite response. You may also choose DMR (Defective Media Retention) option.

**HP Proactive Care with the HP Personalized Support Option** – Adding the Personalized Support Option for HP Proactive Care is highly recommended. The Personalized Support option builds on the benefits of HP Proactive Care Service, providing you an assigned Account Support Manager who knows your environment and delivers support planning, regular reviews, and technical and operational advice specific to your environment. These proactive services will be coordinated with Microsoft's proactive services that come with Microsoft Premier Mission Critical, if applicable.

### *HP Proactive Select*

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

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

For additional information, please visit:

HP Education Services: <http://h10076.www1.hp.com/education/bigdata.htm>

HP Technology Consulting Services: [hp.com/services/bigdata](http://hp.com/services/bigdata)

HP Services: [hp.com/services](http://hp.com/services)

## For more information

Hortonworks, [hortonworks.com](http://hortonworks.com)

HP Solutions for Apache Hadoop, [hp.com/go/hadoop](http://hp.com/go/hadoop)

Hadoop and Vertica, [vertica.com/the-analytics-platform/native-bi-etl-and-hadoop-mapreduce-integration](http://vertica.com/the-analytics-platform/native-bi-etl-and-hadoop-mapreduce-integration)

HP Insight Cluster Management Utility (CMU), [hp.com/go/cmu](http://hp.com/go/cmu)

HP 5900 Switch Series, [hp.com/networking/5900](http://hp.com/networking/5900)

HP FlexFabric 5930 Switch Series, [hp.com/networking/5930](http://hp.com/networking/5930)

HP ProLiant servers, [hp.com/go/proliant](http://hp.com/go/proliant)

HP Enterprise Software, [hp.com/go/software](http://hp.com/go/software)

HP Networking, [hp.com/go/networking](http://hp.com/go/networking)

HP Integrated Lights-Out (iLO), [hp.com/servers/ilo](http://hp.com/servers/ilo)

HP Product Bulletin (QuickSpecs), [hp.com/go/quickspecs](http://hp.com/go/quickspecs)

HP Services, [hp.com/go/services](http://hp.com/go/services)

HP Support and Drivers, [hp.com/go/support](http://hp.com/go/support)

HP Systems Insight Manager (HP SIM), [hp.com/go/hpsim](http://hp.com/go/hpsim)

Red Hat, [redhat.com](http://redhat.com)

To help us improve our documents, please provide feedback at [hp.com/solutions/feedback](http://hp.com/solutions/feedback)



### Sign up for updates

[hp.com/go/getupdated](http://hp.com/go/getupdated)

---

© Copyright 2015 Hewlett-Packard Development Company, L.P. The information contained herein is subject to change without notice. The only warranties for HP products and services are set forth in the express warranty statements accompanying such products and services. Nothing herein should be construed as constituting an additional warranty. HP shall not be liable for technical or editorial errors or omissions contained herein.

Microsoft and Windows are trademarks of the Microsoft group of companies. Intel and Xeon are trademarks of Intel Corporation in the U.S. and other countries. Oracle and Java are registered trademarks of Oracle and/or its affiliates. Red Hat is a registered trademark of Red Hat, Inc. in the United States and other countries. Linux is the registered trademark of Linus Torvalds in the U.S. and other countries.

