**Hewlett Packard Enterprise**

# HPE Reference Architecture for Hortonworks HDP 2.4 on HPE ProLiant DL380 Gen9 servers

HPE Converged Infrastructure with Hortonworks
HDP 2.4 for Apache Hadoop

**HORTONWORKS®**

**Technical white paper**

# Contents

## Executive summary

Hewlett Packard Enterprise 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 HPE infrastructure that provide a significant reduction in complexity and a recognized increase in value and performance.

The configurations are based on Hortonworks Data Platform (HDP), 100% open source distribution of Apache Hadoop, specifically HDP 2.4 and the HPE ProLiant DL380 Gen9 server platform. The configurations reflected in this document have been jointly designed and developed by HPE and Hortonworks to provide optimum computational performance for Hadoop and are also compatible with other HDP 2.x releases.

HPE Big Data solutions provide excellent performance and availability, with integrated software, services, infrastructure, and management – all delivered as one proven configuration, described in more detail at hpe.com/info/hadoop

**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 HPE and Hortonworks Data Platform (HDP) 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 the Solution components section.

**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 April 2016.

## Introduction

This white paper has been created to assist in the rapid design and deployment of Hortonworks software on HPE 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 HPE software, HPE ProLiant servers, and HPE 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 optimum MapReduce, YARN, Spark, Hive and HBase computational performance, resulting in a significant performance at an optimum cost.

## Solution overview

HPE and HDP allow you to derive new business insights from Big Data by providing a platform to store, manage and process data at scale. This RA provides several performance optimized configurations for deploying HDP clusters on HPE infrastructure.

The configurations are based on Hortonworks' Distribution of Apache Hadoop (HDP), specifically HDP 2.4 and Ambari 2.2, and the HPE ProLiant DL380 Gen9 server platform.

The reference architecture in this white paper also includes the following features that are unique to HPE:

- Servers

  The HPE Smart Array P440ar and P840 controllers provide increased[1] 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 Just a Bunch of Disks (JBOD) or various RAID configurations.

  Two sockets with 12 core processors, using Intel® Xeon® E5-2600 v3 product family, provide the high performance required for CPU bound Hadoop workloads. For Management and Head nodes, the same CPU family with 8 core processors are used. Alternate processors are provided in Appendix C.

---

[1] Compared to the previous generation of Smart Array controllers

The HPE iLO management engine on the servers contains HPE Integrated Lights-Out 4 (iLO 4) and features a complete set of embedded management features for HPE Power/Cooling, Agentless Management, Active Health System, and Intelligent Provisioning which reduces node and cluster level administration costs for Hadoop.

- Cluster management

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

- Networking

The HPE FlexFabric 5900AF-48XGT-4QSFP+ 10GbE Top of Rack (ToR) switch has 48 RJ-45 1/10GbE ports and 4 QSFP+ 40GbE ports. It provides Intelligent Resilient Fabric (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 up to 952 million pps throughput and switching capacity of 1280Gb/s with very low 10Gb/s latency of less than 1.5µs (64-byte packets).

HPE FlexFabric 5940 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 up to 1904 Mpps throughput and routing/switching capacity of 2560Gb/s with very low 10Gb/s latency of less than 1µs (64-byte packets). The switch seamlessly handles burst scenarios such as shuffle, sort and block replication which are common in Hadoop clusters.

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

In order to simplify the build for customers, HPE provides the exact bill of materials in this document to allow a customer to purchase this complete solution. HPE recommends that customers purchase the option in which HPE Technical Services Consulting will install the prebuilt operating system images, verify all firmware and versions are correctly installed, and run a suite of tests that verify that the configuration is performing optimally. Once this has been done, the customer can perform a standard HDP installation using the recommended guidelines in this document.

## Hortonworks Data Platform (HDP) 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.
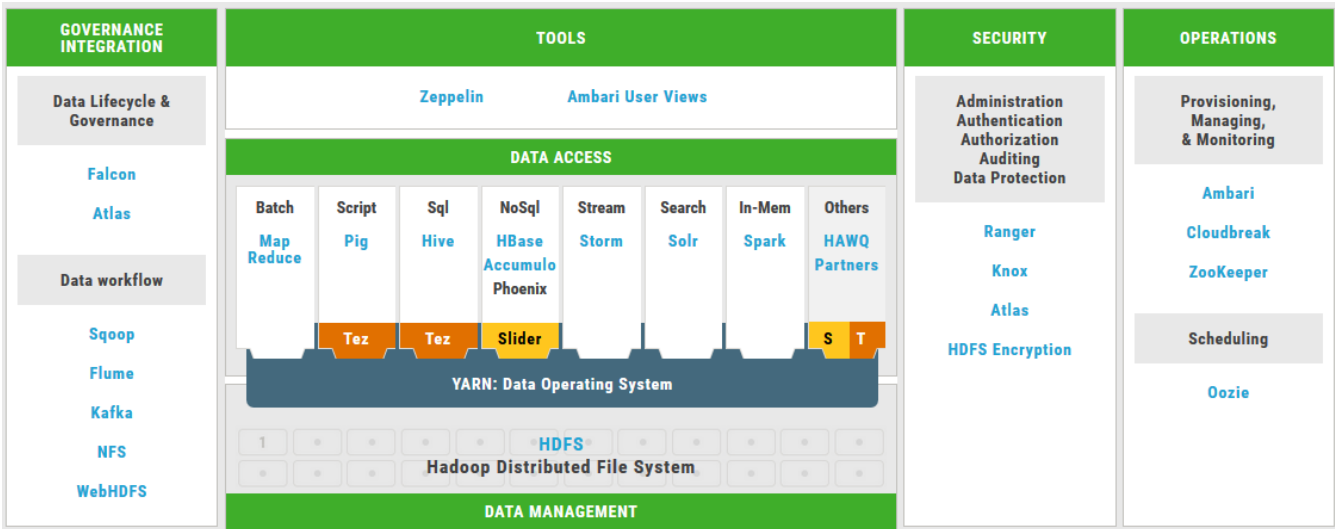
| GOVERNANCE INTEGRATION | TOOLS | SECURITY | OPERATIONS |
|---|---|---|---|

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

Hortonworks Data Platform Version 2.4 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 2.4 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 2.4. Figure 2 shows Hortonworks Data Platform 2.4.
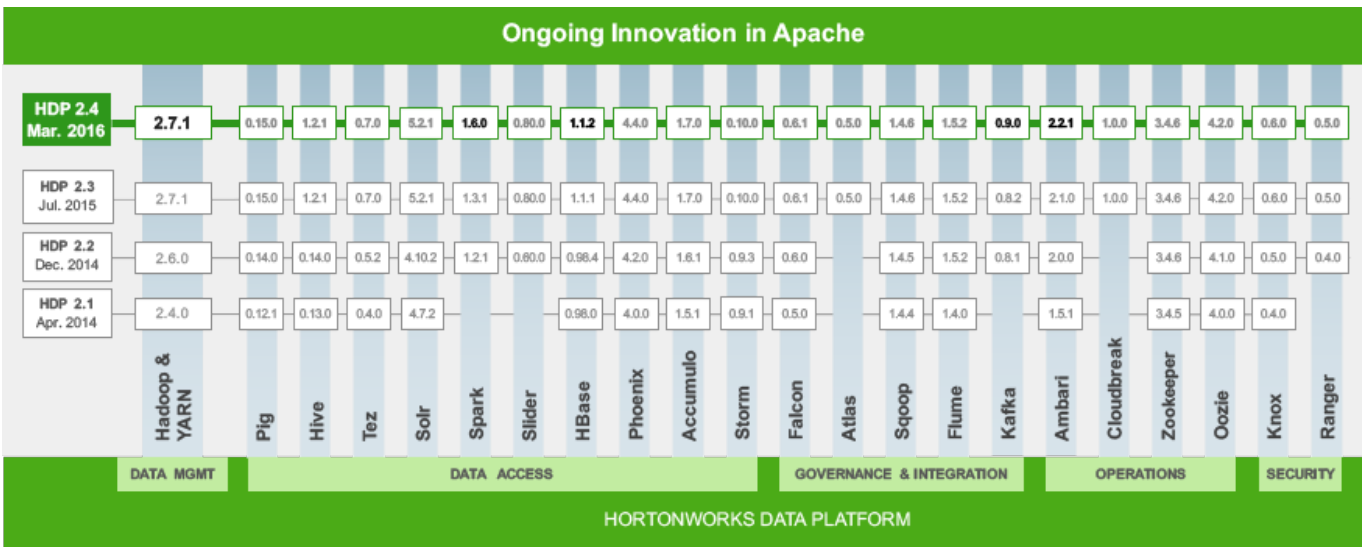
**Ongoing Innovation in Apache**

| | Hadoop & YARN | Pig | Hive | Tez | Solr | Spark | Slider | HBase | Phoenix | Accumulo | Storm | Falcon | Atlas | Sqoop | Flume | Kafka | Ambari | Cloudbreak | Zookeeper | Oozie | Knox | Ranger |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **HDP 2.4** Mar. 2016 | 2.7.1 | 0.15.0 | 1.2.1 | 0.7.0 | 5.2.1 | 1.6.0 | 0.80.0 | 1.1.2 | 4.4.0 | 1.7.0 | 0.10.0 | 0.6.1 | 0.5.0 | 1.4.6 | 1.5.2 | 0.9.0 | 2.2.1 | 1.0.0 | 3.4.6 | 4.2.0 | 0.6.0 | 0.5.0 |
| **HDP 2.3** Jul. 2015 | 2.7.1 | 0.15.0 | 1.2.1 | 0.7.0 | 5.2.1 | 1.3.1 | 0.80.0 | 1.1.1 | 4.4.0 | 1.7.0 | 0.10.0 | 0.6.1 | 0.5.0 | 1.4.6 | 1.5.2 | 0.8.2 | 2.1.0 | 1.0.0 | 3.4.6 | 4.2.0 | 0.6.0 | 0.5.0 |
| **HDP 2.2** Dec. 2014 | 2.6.0 | 0.14.0 | 0.14.0 | 0.5.2 | 4.10.2 | 1.2.1 | 0.60.0 | 0.98.4 | 4.2.0 | 1.6.1 | 0.9.3 | 0.6.0 | | 1.4.5 | 1.5.2 | 0.8.1 | 2.0.0 | | 3.4.6 | 4.1.0 | 0.5.0 | 0.4.0 |
| **HDP 2.1** Apr. 2014 | 2.4.0 | 0.12.1 | 0.13.0 | 0.4.0 | 4.7.2 | | | 0.98.0 | 4.0.0 | | 1.5.1 | 0.9.1 | 0.5.0 | | 1.4.4 | 1.4.0 | | 1.5.1 | | 3.4.5 | 4.0.0 | | 0.4.0 |

| DATA MGMT | DATA ACCESS | GOVERNANCE & INTEGRATION | OPERATIONS | SECURITY |
|---|---|---|---|---|

**HORTONWORKS DATA PLATFORM**
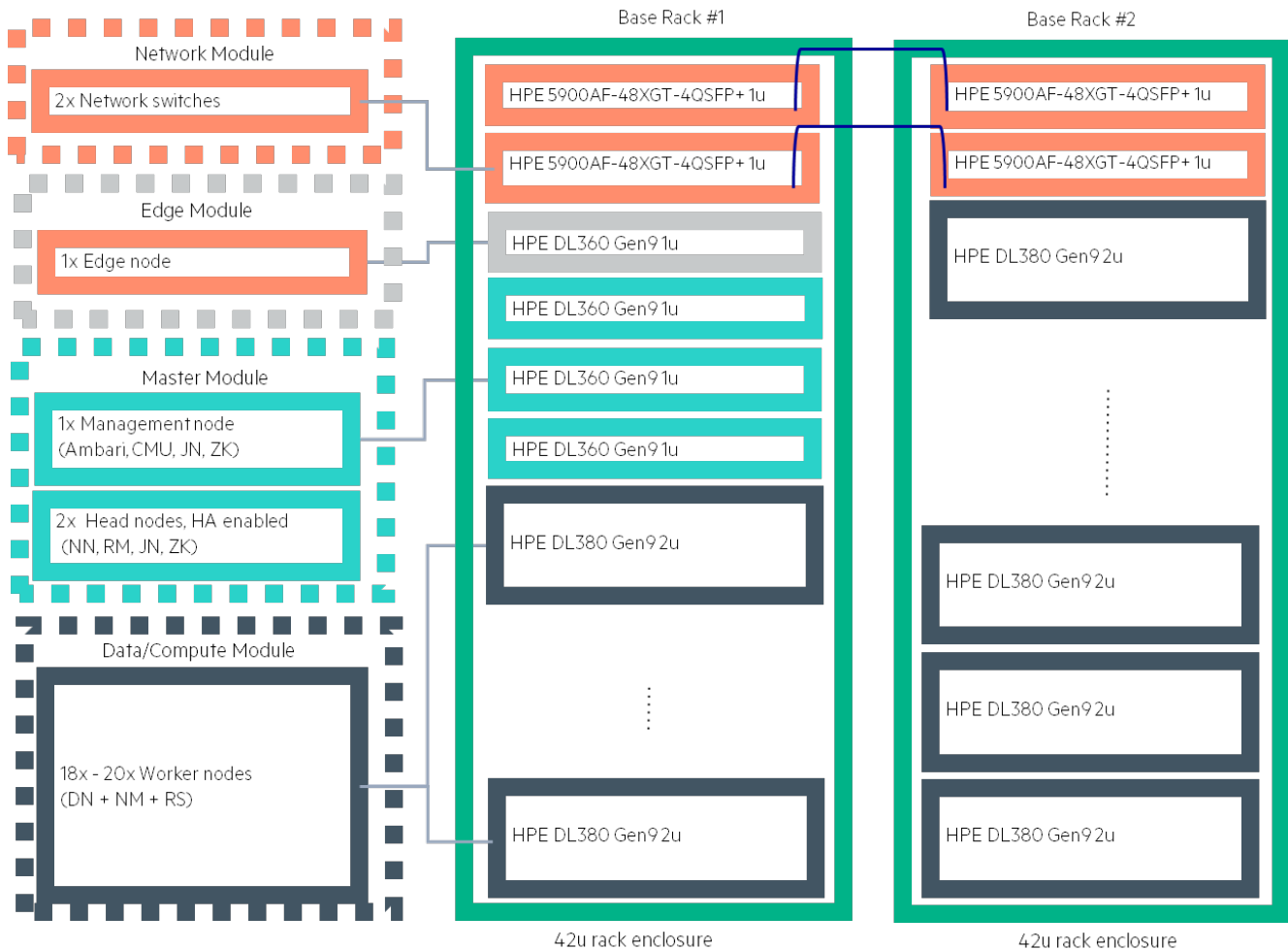
**Figure 2.** Hortonworks Data Platform 2.4

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.

For detailed information on Hortonworks Data Platform, see hortonworks.com/hdp

## Solution components

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 to 5 below specify different services that run on Management, ResourceManager, NameNode and Worker nodes. The Reference Architectures (RAs) we provide in this document will map the Management and Worker services onto HPE infrastructure for clusters of varying sizes. The RAs factor in the scalability requirements for each service.

Each of the HDP solution components shown in Figure 3 is discussed at length in the section below. For a full BOM listing on the products selected, refer to the Bill of materials section of this white paper.



**Figure 3.** Basic conceptual diagram of an HPE ProLiant DL380 Gen9 HDP reference architecture

### High Availability considerations

The following are some of the High Availability (HA) features considered in this reference architecture configuration:

- **Hadoop NameNode HA** – The configurations in this white paper utilize quorum-based journaling high-availability feature. 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. Ambari provides a simple wizard to enable HA for YARN ResourceManager.

- **OS availability and reliability** – For the reliability of the server, the OS disk is configured in a RAID1+0 configuration thus preventing failure of the system from OS hard disk failures.

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

- **Power supply** – 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).

---

**Note**

For High Availability considerations, it is recommended to have a minimum configuration of three Master nodes (Management, NameNode and ResourceManager) and five data nodes.

---

## Pre-deployment considerations

The operating system and the network are key factors you need to 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.

### Operating system

Hortonworks HDP 2.4 supports the following 64-bit operating systems, visit http://docs.hortonworks.com/HDPDocuments/HDP2/HDP-2.4.2/bk_installing_manually_book/content/meet-min-system-requirements.html for the minimum requirements:

- For Red Hat® systems, HDP provides 64-bit packages for Red Hat Enterprise Linux® (RHEL) 6 update 6 or later and Red Hat Enterprise Linux 7.1 or later.

- For SUSE systems, HDP provides 64-bit packages for SUSE Linux Enterprise Server 11 (SLES 11). Service pack 2 or later is required.

### Computation

Unlike MR1, under YARN/MR2, there is no distinction between resources available for maps, and resources available for reduces – all MR2 resources are available for both in the form of containers. Employing Hyper-Threading increases 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. HDP supports the use of the Fair and FIFO schedulers in MR2.

---

**Key point**

When computation performance is of primary concern, HPE recommends high CPU powered HPE ProLiant DL380 servers for Worker nodes with 256GB RAM. HDP 2.4 components such as Spark, HDFS Caching and Solr benefit from large amounts of memory.

---

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** – For Management nodes, storage reliability is 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. Specific details around disk and RAID configurations will be provided in the next section.

### Network
Configuring a single 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 HPE FlexFabric 5900AF-48XGT switch supports this through HPE Intelligent Resilient Fabric (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.

A more detailed white paper for Hadoop Networking best practices is available. For more information refer to http://h20195.www2.hpe.com/V2/GetDocument.aspx?docname=4AA5-3279ENW

## Best practices and configuration guide for the solution
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. This section specifies which server to use and the rationale behind it.

### Management/Head/Edge 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**
HPE 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 HDP 2.4. 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 five subsections:

- Server platform

- Management node

- ResourceManager server

- NameNode server

- Edge node

**Server platform: HPE ProLiant DL360 Gen9**

The HPE ProLiant DL360 Gen9 (1U) server platform shown in Figure 4 below, is an excellent choice as the server platform for the Management nodes and Head nodes.



**Figure 4.** HPE ProLiant DL360 Gen9 server

**Processor configuration**

The configuration features two sockets with 8-core processors of the Intel E5-2600 v3 processor 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 Intel Xeon E5-2640 v3 processors for the Management servers with the ResourceManager, NameNode and Ambari services. Alternate CPU configurations are available in <u>Appendix C: Alternate parts</u> section.

**Drive configuration**

The HPE Smart Array P440ar controller is specified to drive eight 1TB 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 HPE Smart Array P440ar controller to apply the following RAID schemes:

- **Management node:** eight disks, with OS on (2) RAID1 set, one JBOD or RAID0 for ZooKeeper, one JBOD or RAID0 for QJN (Quorum Journal Node), 4 disks in RAID1+0 for database.

- **ResourceManager and NameNode servers:** four disks with RAID1+0 for OS and Hadoop software, one JBOD or RAID0 for ZooKeeper, one JBOD or RAID0 for QJN (Quorum Journal Node), two disks are available for other filesystems.

- **Edge node:** A multi-homed Edge node can be added with eight disks to provide a sufficient amount of disk capacity for a staging area for ingesting data into the Hadoop cluster from another subnet.

---

**Best practice**

For a performance oriented solution, HPE recommends 15K SAS or SSD drives as they offer a significant read and write performance enhancement over SATA disks. The HPE Smart Array P440ar controller provides two port connectors per controller with each containing four SAS links. The drive cage for the HPE ProLiant DL360 Gen9 contains eight disk slots and thus each disk slot has a dedicated SAS link which ensures the server provides the maximum throughput that each drive can give. These nodes require highly reliable storage for databases, namespace storage, and edit-log journaling (ZooKeeper) services running on the servers.

---

**Cluster isolation and access configuration**

It is important to isolate the Hadoop cluster so that external network traffic does not affect the performance of the cluster. In addition, isolation allows the Hadoop cluster to be managed independently from its users, ensuring that the cluster administrator is the only person able to make changes to the cluster configuration. Thus, HPE recommends deploying ResourceManager, NameNode, and Worker nodes on their own private Hadoop cluster subnet.

## Key point

Once a Hadoop cluster is isolated, the users still need a way to access the cluster and submit jobs to it. To achieve this, HDP recommends multi-homing the Edge nodes, so that they can participate in both the Hadoop cluster subnet and a subnet belonging to users.

## Performance critical Hadoop clusters

For performance critical clusters, HPE and Hortonworks recommend 2 disks in RAID1 for OS, 2 disks in RAID1 for HDFS metadata, 1 JBOD or RAID0 for ZooKeeper, 1 JBOD or RAID0 for QJN (Quorum Journal Node), and 4 disks in RAID10 for database. The HPE ProLiant DL360 Gen9 supports 10 Small Form Factor disks with the addition of the optional 2 SFF SAS/SATA kit (p/n 764630-B21). In the case of an 8-disk configuration, such as the one tested for this reference architecture, the recommendation is to co-locate the HDFS metadata with the OS on the RAID1 set.

### 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 HPE ProLiant DL360 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.

### Management node

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 1.

**Table 1.** Management node basic software components

| SOFTWARE | DESCRIPTION |
| --- | --- |
| Red Hat Enterprise Linux 7.1 | Recommended Operating System |
| HPE Insight CMU 8.0 | Infrastructure Deployment, Management, and Monitoring |
| Oracle JDK 1.7.0_75 | Java Development Kit |
| PostgreSQL 8.4 | Database Server for Ambari |
| Ambari 2.2.2 | Ambari Management Software |
| Hue Server | Web Interface for Hadoop Applications |
| ZooKeeper | Cluster coordination service |
| NameNode HA | NameNode HA (Journal Node) |

See the following link for the Ambari and PostgreSQL Installation guide,
http://docs.hortonworks.com/HDPDocuments/Ambari/Ambari-2.2.2.0/index.html

The Management node and Head nodes, as tested in the reference architecture, contain the following configuration:

- 2 x eight-core Intel Xeon E5-2640 v3 processors

- HPE Smart Array P440ar controller with 2GB FBWC

- 8TB – 8 x 1TB SFF SAS 7.2K RPM disks

- 128GB DDR4 memory – 8 x HPE 16GB 2Rx4 PC4-2133P-R kit

- HPE Ethernet 10GbE 2P 561FLR-T network adapter

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

## Head node 1 – ResourceManager server

The ResourceManager server contains the following software components. See the following link for more information on installing and configuring the ResourceManager and NameNode HA, http://docs.hortonworks.com/index.html

Table 2 shows the ResourceManager server base software components.

**Table 2.** ResourceManager server base software components

| SOFTWARE | DESCRIPTION |
| --- | --- |
| Red Hat Enterprise Linux 7.1 | Recommended Operating System |
| Oracle JDK 1.7.0_75 | 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. See the following link for more information on installing and configuring the NameNode and ResourceManager HA at http://docs.hortonworks.com/index.html

Table 3 shows the NameNode server base software components.

**Table 3.** NameNode server base software components

| SOFTWARE | DESCRIPTION |
| --- | --- |
| RHEL 7.1 | Recommended Operating System |
| Oracle JDK 1.7.0_67 | 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) |
| HBaseMaster | HBase Master (Only if running HBase) |
| ZooKeeper | Cluster coordination service |
| HiveServer2 | Hue application to run queries on Hive with authentication |
| Apache Pig and Apache Hive | Analytical interfaces to the Hadoop cluster |

### Edge node

The Edge node hosts the client configurations that submit jobs to the Hadoop cluster. We recommend that you install with the software components shown in Table 4.

**Table 4.** Edge node basic software components

| SOFTWARE | DESCRIPTION |
| --- | --- |
| Red Hat Enterprise Linux 7.1 | Recommended Operating System |
| Oracle JDK 1.7.0_75 | Java Development Kit |
| Gateway Services | Hadoop Gateway Services (HDFS, YARN, MapReduce, HBase, and others) |

The Edge nodes, as tested in the reference architecture, contain the following configuration:

- 2 x eight-core Intel Xeon E5-2640 v3 processors

- HPE Smart Array P440ar controller with 2GB FBWC

- 8TB – 8 x 1TB SFF SAS 7.2K RPM disks

- 64GB DDR4 memory – 8 x HPE 8GB 2Rx4 PC4-2133P-R kit

- HPE Ethernet 10GbE 2P 561FLR-T network adapter

### Worker nodes

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

#### Server platform: HPE ProLiant DL380 Gen9

The HPE ProLiant DL380 Gen9 (2U) shown below in Figure 5, 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 5.** HPE ProLiant DL380 Gen9 server

#### Processor selection

The configuration features two processors from the Intel Xeon E5-2600 v3 family. The high performance architecture configuration provides 24 physical or 48 Hyper-Threaded cores per server with 2 x E5-2680 v3 (12 cores/2.5 GHz) 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 HPE 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 Spark and YARN, 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 reference architecture, we chose 256GB memory (16 x HPE 16GB 2Rx4 PC4-2133P-R Kit).

## Best practice

To ensure optimal memory performance and bandwidth, HPE 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 HPE Smart Array P840 controller by configuring each individual disk as a separate RAID0 volume. Additionally array acceleration features on the HPE Smart Array P840 should be turned off for the RAID0 data volumes. The 340GB M.2 SSD drives are configured for OS.

The HPE Smart Array P840 controller provides two port connectors per controller with each containing 8 SAS links.

## Best practice

For a performance oriented solution HPE recommends 15K SAS or SSD drives as they offer a significant read and write performance enhancement over SATA disks, which improves the performance on I/O intensive workloads.

## 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. Enable CPU scheduler in Ambari.

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

While configuring YARN for MapReduce jobs make sure that the following attributes have been specified with sufficient vcores and memory. They represent resource allocation attributes for map and reduce containers. Note that the optimum values for these attributes depend on the nature of workload/use case.

```
mapreduce.map.memory.mb
mapreduce.reduce.memory.mb
mapreduce.map.cpu.vcores
mapreduce.reduce.cpu.vcores
```

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

```
mapreduce.map.java.opts.max.heap
mapreduce.reduce.java.opts.max.heap
```

**Worker node software components**

Table 5 lists the Worker node software components. See the following link for more information on installing and configuring the NodeManager (or HBaseRegionServer) and DataNode: http://docs.hortonworks.com/index.html

Table 5 shows the Worker node base software components.

**Table 5.** Worker node base software components

| SOFTWARE | DESCRIPTION |
|---|---|
| Red Hat Enterprise Linux 7.1 | Recommended Operating System |
| Oracle JDK 1.7.0_75 | Java Development Kit |
| NodeManager | The NodeManager process for MR2/YARN |
| DataNode | The DataNode process for HDFS |
| HBaseRegionServer | The HBaseRegionServer for HBase (Only if running HBase) |

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

- Dual 12-core Intel Xeon E5-2680 v3 processors with Hyper-Threading enabled

- 15 x 3TB 3.5" 7.2K LFF SATA MDL (42TB for data)

- HPE ProLiant DL380 Gen9 3LFF rear SAS/SATA kit

- 256GB DDR4 memory (16 x HPE 16GB), 4 channels per socket

- 1 x dual 340GB RI-2 Solid State M.2 enablement kit

- 1 x HPE Ethernet 10GbE 2-port FlexibleLOM (bonded)

- 1 x HPE Smart Array P840/4G FIO controller

**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 the Bill of materials section of this white paper.

## Networking/switch selection
Hadoop clusters contain two types of switches, namely ToR switches and aggregation switches. ToR switches route the traffic between the nodes in each rack and aggregation switches route the traffic between the racks.

### Top of Rack (ToR) switches
The HPE FlexFabric 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). For more information on the HPE FlexFabric 5900AF-48XGT-4QSFP+, 10GbE, switch, visit hpe.com/networking/5900. The BOM for the HPE FlexFabric 5900AF-48XGT switch is provided in the Bill of materials section of this white paper.

A dedicated management switch for iLO traffic is not required as the HPE ProLiant DL360 Gen9 and HPE ProLiant 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. Customers who would like to separate iLO and PXE traffic from the data/Hadoop network traffic can add a 1GbE HPE FlexFabric 5900AF-48G-4XG-2QSFP+ network switch (JG510A, shown in Figure 6). The BOMs for 1GbE switch and network cards are provided in the Bill of materials section of this white paper.

**Figure 6.** HPE FlexFabric 5900AF-48G-4XG-2QSFP+ switch

**Aggregation switches**

The HPE FlexFabric 5940 switch series is a family of high performance and low-latency 10GbE, 40GbE top-of-rack (ToR) data center switches. The switch series also include 100G uplink technology and is part of the HPE FlexFabric data center solution, which is a cornerstone of the FlexNetwork architecture. The switch has better connectivity with 32 40GbE ports, which may be split into four 10GbE ports each for a total of 96 10GbE ports with 8 40GbE uplinks per switch, 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 7 shows the HPE FlexFabric 5940 aggregation switch. For more information on the HPE FF 5940 32QSFP+, see
http://www8.hp.com/us/en/products/networking-switches/product-detail.html?oid=1009148840

The BOM for the HPE FF 5940 32QSFP+ switch is provided in the Bill of materials section of this white paper.



**Figure 7.** HPE FlexFabric 5940 32QSFP+ 40GbE aggregation switch

## 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 HDP 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 RA reflects the following.

#### Single-rack network

As described in the Networking/switch selection section, two IRF bonded HPE FlexFabric 5900AF-48XGT ToR switches are specified for performance and redundancy. The HPE FlexFabric 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 HPE FlexFabric 5940 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 HPE FlexFabric 5900AF-48XGT switch for uplinks.

**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.

One option is to use the already multi-homed Edge node, which can be configured to provide a staging area for ingesting data into the Hadoop cluster from another subnet.

---

**Note**

The benefit of using dual-homed Edge nodes 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.

---

Another option is to 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, see
https://docs.hortonworks.com/HDPDocuments/HDP2/HDP-2.4.0/bk_hdfs_admin_tools/content/ch11.html

**Rack**

The rack contains eighteen HPE ProLiant DL380 servers, four HPE ProLiant DL360 servers and two HPE FlexFabric 5900AF-48XGT switches within a 42U rack.

**Management/Head/Edge nodes**

Four HPE ProLiant DL360 Gen9 Management nodes are specified:

- The Management node

- The ResourceManager/NameNode HA

- The NameNode/ResourceManager HA

- The Edge node

Detailed information on the hardware and software configurations is available in the configuration guide section of this document.

**Worker nodes**

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

**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 Data Platform will ensure there is less impact.

---

**Best practice**

For each server, HPE 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.

---

**Single-rack reference architecture**

Refer to Figure 8 for a rack-level view of the single-rack HDP reference architecture for this solution.

1 Edge Node
HPE DL360 Gen9 8SFF
Dual 8-Core Intel E5-2640 v3 2.6 GHz
64GB RAM(8x 8GB 2Rx4 PC4-2133P-R)
8TB Disks (8x 1TB SAS 7.2K HDD)
1xHPE Smart Array P440ar/2G FIO Controller
1x Ethernet 10Gb 2P 561FLR-T

2 Ethernet Switches
HPE 5900AF-48XGT – 4QSFP+ Switch

1 Management Node
HPE DL360 Gen9 8SFF
Dual 8-Core Intel E5-2640 v3 2.6 GHz
128 GB RAM (8x 16GB 2Rx4 PC4-2133P-R)
8TB Disks (8x 1TB SAS 7.2K HDD)
1x HPE Smart Array P440ar/2G FIO Controller
1x Ethernet 10Gb 2P 561FLR-T

2 Head nodes - NameNode/ResourceManager
HPE DL360 Gen9 8SFF
Dual 8-Core Intel E5-2640 v3 2.6 GHz
128 GB RAM (8x 16GB 2Rx4 PC4-2133P-R)
8TB Disks (8x 1TB SAS 7.2K HDD)
1x HPE Smart Array P440ar/2G FIO Controller
1x Ethernet 10Gb 2P 561FLR-T

18 Worker Nodes
18x HPE ProLiant DL380 Gen9
Dual 12-Core Intel E5-2680 v3 2.5 GHz
256GB RAM (16x 16GB 2Rx4 PC4-2133P-R)
Dual 340GB RI-2 Solid State M.2 Enablement Kit
45TB raw storage (15x 3TB 6G SAS 7.2k HDD )
1x HPE Smart Array P840/4G FIO Controller
1x Ethernet 10Gb 2P 561FLR-T

Software
Operating System: 64-bit Red Hat Enterprise Linux 7.1
Hortonworks Data Platform 2.4
HPE Insight Cluster Management Utility v8.0
Optional:
Intelligent PDU
Vertica connector for Hadoop
Autonomy connector for Hadoop

Three-phase PDU (4 PDUs per rack):
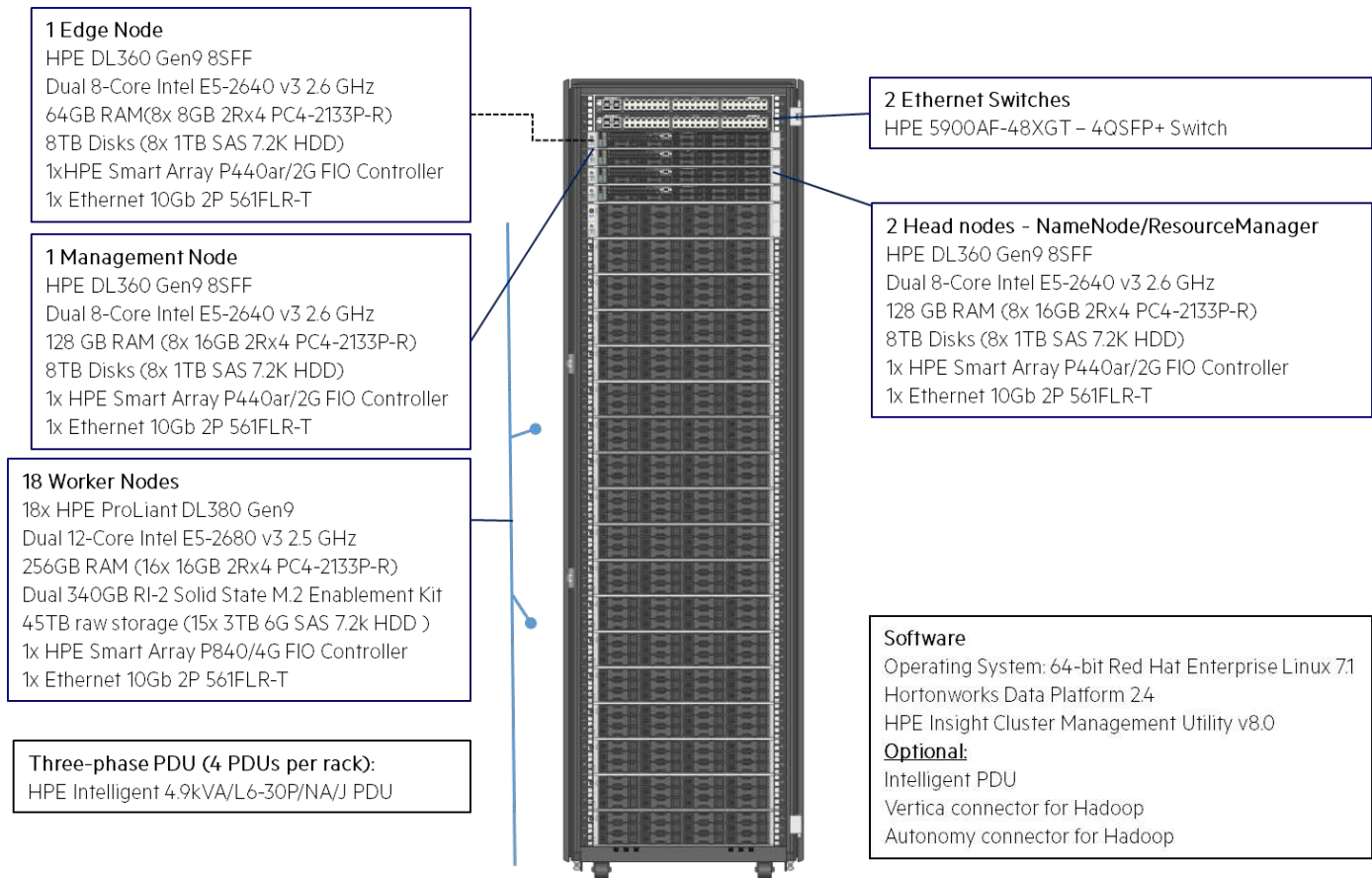HPE Intelligent 4.9kVA/L6-30P/NA/J PDU

**Figure 8**. Single-rack HDP 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 extension racks of a similar configuration to the single-rack configuration as shown in Figure 8. This section reflects the design of those racks, and Figures 9 and 10 show the rack-level view of the multi-rack architecture.

**Extension rack**

The extension rack contains eighteen HPE ProLiant DL380 Gen9 servers and two HPE FlexFabric 5900AF-48XGT switches within a 42U rack; an additional 2x HPE ProLiant DL380 servers (2U each) or, two HPE FlexFabric 5940 32QSFP+ (1U) aggregation switches can be installed in the first expansion rack.

**Multi-Rack network**

As described in the Networking/switch selection section of this white paper, two HPE FlexFabric 5900AF-48XGT ToR switches are specified per each expansion rack for performance and redundancy. The HPE FlexFabric 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 HPE FlexFabric 5940 32QSFP+ aggregation switches.

**Multi-rack architecture**

The HPE 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 base rack. Aside from the OS, each Worker node typically runs DataNode, NodeManager (and HBaseRegionServer if you are using HBase).

**Note**

In a multi-rack configuration we recommend moving the Master nodes onto different racks for better resiliency. Also move one of the aggregation switches to a different rack for better resiliency still using IRF.
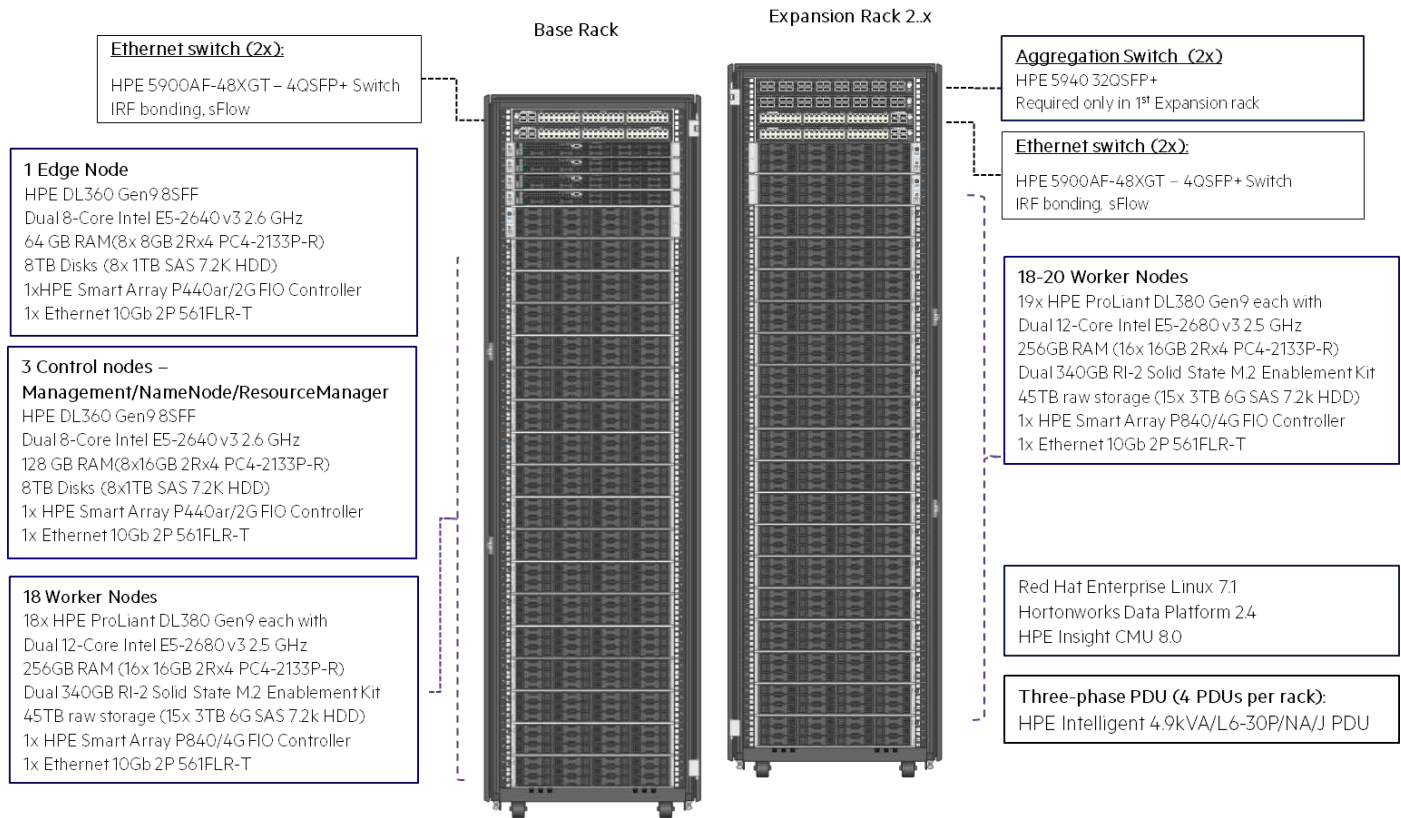
Ethernet switch (2x):
HPE 5900AF-48XGT – 4QSFP+ Switch
IRF bonding, sFlow

1 Edge Node
HPE DL360 Gen9 8SFF
Dual 8-Core Intel E5-2640 v3 2.6 GHz
64 GB RAM (8x 8GB 2Rx4 PC4-2133P-R)
8TB Disks (8x 1TB SAS 7.2K HDD)
1x HPE Smart Array P440ar/2G FIO Controller
1x Ethernet 10Gb 2P 561FLR-T

3 Control nodes –
Management/NameNode/ResourceManager
HPE DL360 Gen9 8SFF
Dual 8-Core Intel E5-2640 v3 2.6 GHz
128 GB RAM (8x16GB 2Rx4 PC4-2133P-R)
8TB Disks (8x1TB SAS 7.2K HDD)
1x HPE Smart Array P440ar/2G FIO Controller
1x Ethernet 10Gb 2P 561FLR-T

18 Worker Nodes
18x HPE ProLiant DL380 Gen9 each with
Dual 12-Core Intel E5-2680 v3 2.5 GHz
256GB RAM (16x 16GB 2Rx4 PC4-2133P-R)
Dual 340GB RI-2 Solid State M.2 Enablement Kit
45TB raw storage (15x 3TB 6G SAS 7.2k HDD)
1x HPE Smart Array P840/4G FIO Controller
1x Ethernet 10Gb 2P 561FLR-T

Base Rack

Expansion Rack 2..x

Aggregation Switch (2x)
HPE 5940 32QSFP+
Required only in 1st Expansion rack

Ethernet switch (2x):
HPE 5900AF-48XGT – 4QSFP+ Switch
IRF bonding, sFlow

18-20 Worker Nodes
19x HPE ProLiant DL380 Gen9 each with
Dual 12-Core Intel E5-2680 v3 2.5 GHz
256GB RAM (16x 16GB 2Rx4 PC4-2133P-R)
Dual 340GB RI-2 Solid State M.2 Enablement Kit
45TB raw storage (15x 3TB 6G SAS 7.2k HDD)
1x HPE Smart Array P840/4G FIO Controller
1x Ethernet 10Gb 2P 561FLR-T

Red Hat Enterprise Linux 7.1
Hortonworks Data Platform 2.4
HPE Insight CMU 8.0

Three-phase PDU (4 PDUs per rack):
HPE Intelligent 4.9kVA/L6-30P/NA/J PDU

**Figure 9.** Multi-rack reference architecture - Rack-level view

Figure 10 shows how a single-rack reference architecture can be extended to a multi-rack reference architecture.

**Note**

There is no need for an aggregation switch when a second rack is added to the base rack. The existing HPE FlexFabric 5900 switches in the racks can be used to network between both (base racks) as shown in Figure 3. Figure 10 shows the network connection, using an additional aggregation switch, when a third rack (expansion rack) is added to the base racks. For simplicity, Figure 9 depicts only one of the base racks and the third expansion rack. The other base rack connects to the aggregation switch in the same way as the first base rack as shown in Figure 10.
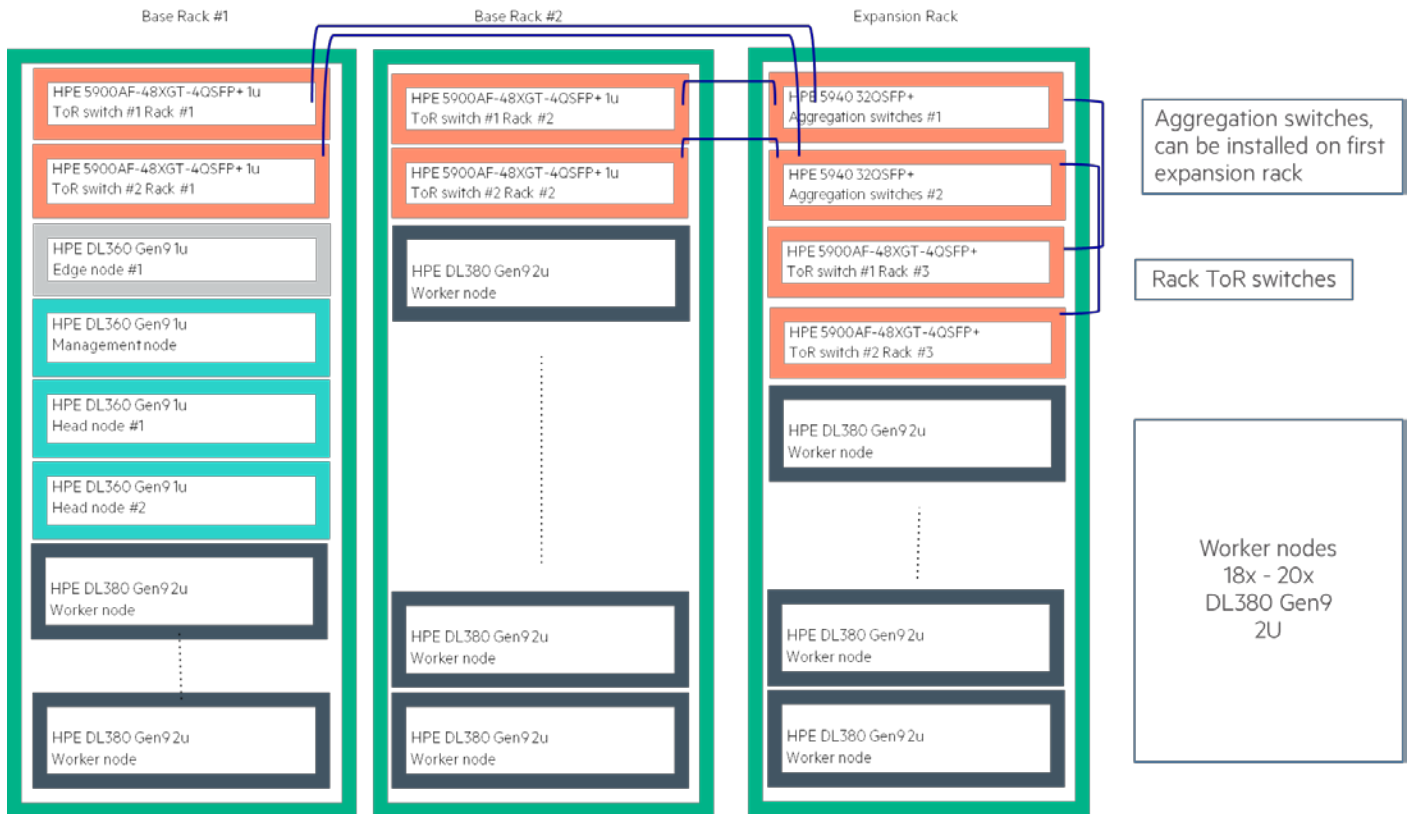
**Figure 10.** Multi-rack reference architecture (extension of the single-rack reference architecture)

## Capacity and sizing

Hadoop cluster storage sizing requires careful planning and identifying the current and future storage and compute needs. Use the following as general guidelines for data inventory:

- Sources of data

- Frequency of data

- Raw storage

- Processed HDFS storage

- Replication factor

- Default compression turned on

- Space for intermediate files

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

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

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

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

### Processor options

For workloads that are CPU intensive it is recommended to choose higher capacity processors with more cores. Typically workloads such as Spark and Solr Search will benefit from higher capacity CPUs. Table 6 shows alternate CPUs for the selected HPE ProLiant DL380 servers.

**Table 6.** CPU recommendations

| CPU | DESCRIPTION |
| --- | --- |
| 2 x E5-2660 v3 | Min configuration (10 cores/2.6GHz) |
| 2 x E5-2680 v3 | Base configuration (12 cores/2.5GHz) |
| 2 x E5-2697 v3 | Enhanced (14 cores/2.6GHz) |

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. HPE 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 7 shows the recommended memory configurations.

**Table 7.** Memory recommendations

| MEMORY | DESCRIPTION |
| --- | --- |
| 128GB 8 x HPE 16GB 2Rx4 PC4-2133P-R Kit | Base configuration |
| 256GB 16 x HPE 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. Table 8 shows alternate storage options for the selected HPE ProLiant DL380.

**Table 8.** HDD recommendations

| HDD | DESCRIPTION |
| --- | --- |
| 2/3/4/6/8TB LFF SATA | Base configuration |
| 2/3/4/6/8TB LFF SAS | Performance configuration |

See Appendix C: Alternate parts Table C-3 for BOM details of alternate hard disks.

### Key point

HPE recommends all HPE ProLiant systems be upgraded to the latest BIOS and firmware versions before installing the OS. HPE 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 2015.10.0 (B). The latest version of SPP can be obtained from:
http://h18004.www1.hp.com/products/servers/service_packs/en/index.html

# HPE Vertica and Hadoop

Relational database management systems such as HPE Vertica excel at analytic processing for big volumes of structured data including call detail records, financial tick streams and parsed weblog data. HPE Vertica is designed for high speed load and query when the database schema and relationships are well defined. Hortonworks Data Platform, 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 Data Platform simplifies installation, configuration, deployment and management of the powerful Hadoop framework for enterprise users.

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

**Note**

HPE 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 HPE Vertica's Hadoop and Pig Connectors, users have unprecedented flexibility and speed in loading data from Hadoop to HPE Vertica and querying data from HPE Vertica in Hadoop as part of MapReduce jobs for example. The HPE Vertica Hadoop and Pig Connectors are supported by HPE Vertica, and available for download.

For more information, see hpe.com/info/vertica

# HPE Insight Cluster Management Utility

HPE 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. HPE 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. HPE Insight CMU offers full support for HPE iLO 2, iLO 3, iLO 4 and LO100i adapters on all HPE ProLiant servers in the cluster.

**Best practice**

HPE recommends using HPE Insight CMU for all Hadoop clusters. HPE 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. HPE Insight 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, HPE 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. HPE Insight CMU is able to deploy 800 nodes in 30 minutes.

HPE 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. HPE Insight CMU is available for HPE ProLiant and HPE BladeSystem servers, and is supported on a variety of Linux operating systems, including Red Hat Enterprise Linux, SUSE Linux Enterprise Server, CentOS, and Ubuntu. HPE Insight CMU also includes options for monitoring graphical processing units (GPUs) and for installing GPU drivers and software. Figures 11 and 12 show views of the HPE Insight CMU.

HPE Insight CMU can be configured to support High Availability with an active-passive cluster. For more information, see hpe.com/info/cmu

For the HPE Insight Cluster Management Utility BOM, see the Bill of materials section of this paper.

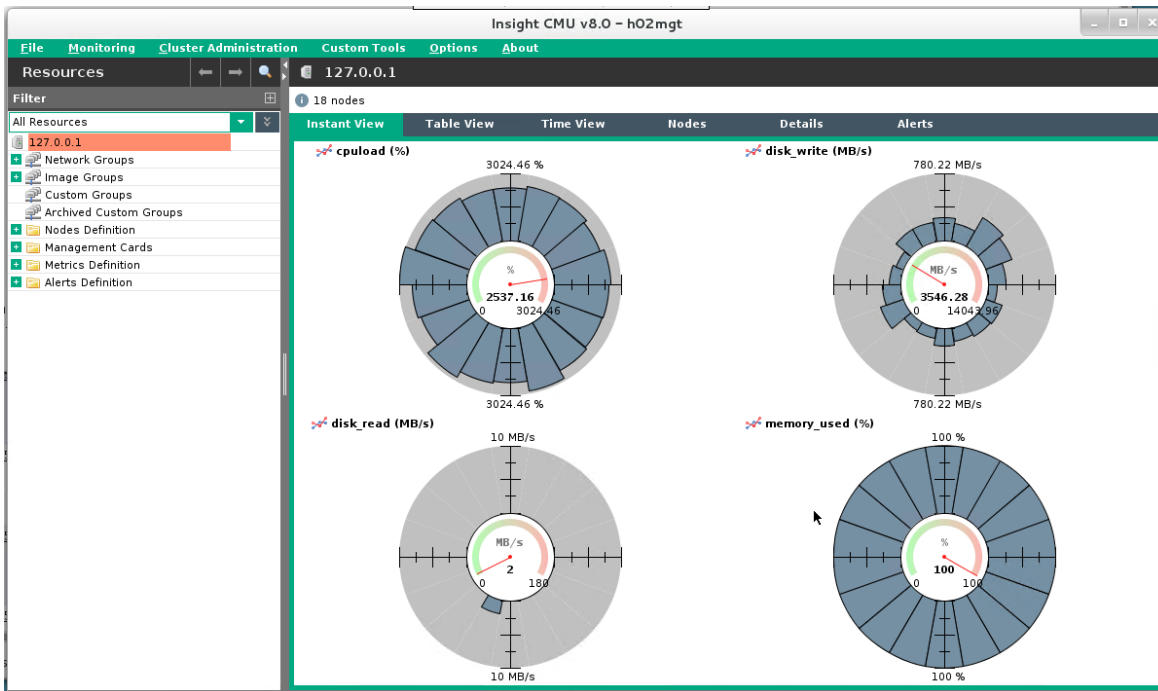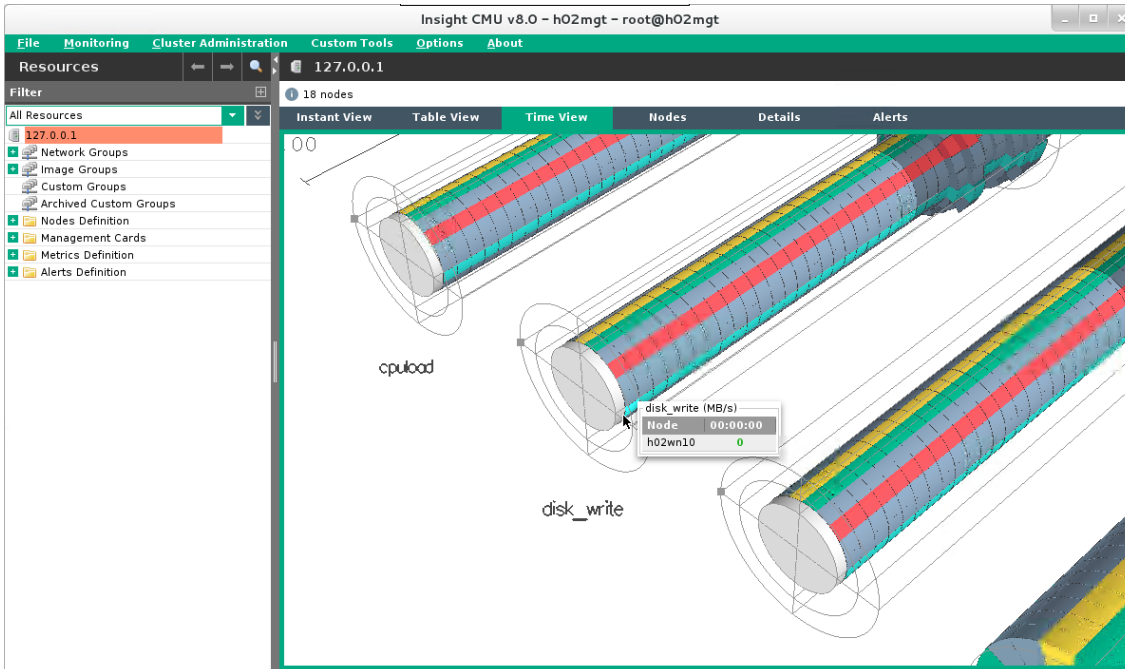**Figure 11.** HPE Insight CMU Interface - real-time view



**Figure 12.** HPE Insight CMU Interface - Time view

## Summary

HPE 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 (HDP) 2.4 on HPE infrastructure and management software. These configurations leverage HPE 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 HPE infrastructure for clusters of various sizes.

## Implementing a proof-of-concept

As a matter of best practice for all deployments, HPE 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 HPE Services representative (hpe.com/us/en/services/consulting.html) or your HPE partner.

## Appendix A: Bill of materials

The BOMs outlined in this section are based on a tested configuration for a single-rack reference architecture with 1 Management node, 2 Head nodes, 18 Worker nodes and 2 ToR switches. The following tables show the Bill of materials for nodes and switches.

The following BOMs contain electronic license to use (E-LTU) parts. Electronic software license delivery is now available in most countries. HPE 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, contact your reseller or an HPE representative.

---

**Note**

Part numbers are at time of publication 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 HPE Reseller or HPE Sales Representative for more details. hpe.com/us/en/services/consulting.html

---

**Management node and Head node BOM**

**Table 9.** BOM for the HPE ProLiant DL360 Gen9 server configuration

| QTY | PART NUMBER | DESCRIPTION |
|---|---|---|
| | | **Management and Head nodes** |
| 1 | 755258-B21 | HPE DL360 Gen9 8-SFF CTO Chassis |
| 1 | 755836-L21 | HPE DL360 Gen9 E5-2640v3 FIO Kit |
| 1 | 755836-B21 | HPE DL360 Gen9 E5-2640v3 Kit |
| 8 | 726719-B21 | HPE 16GB 2Rx4 PC4-2133P-R Kit |
| 8 | 652749-B21 | HPE 1TB 6G SAS 7.2K 2.5in SC MDL HDD |
| 1 | 700699-B21 | HPE Ethernet 10Gb 2P 561FLR-T Adptr |
| 1 | 749974-B21 | HPE Smart Array P440ar/2G FIO Controller |
| 2 | 720478-B21 | HPE 500W FS Plat Ht Plg Pwr Supply Kit |
| 1 | 663201-B21 | HPE 1U SFF Ball Bearing Rail Kit |
| 1 | C6N36ABE | HPE Insight Control ML/DL/BL Bundle E-LTU |
| 1 | G3J28AAE | RHEL Svr 2 Sckt/2 Gst 1yr 24x7 E-LTU |

**Edge node BOM**

**Table 10.** BOM for the HPE ProLiant DL360 Gen9 server configuration

| QTY | PART NUMBER | DESCRIPTION |
| --- | --- | --- |
| 1 | | **Edge nodes** |
| 1 | 755258-B21 | HPE DL360 Gen9 8-SFF CTO Chassis |
| 1 | 755836-L21 | HPE DL360 Gen9 E5-2640v3 FIO Kit |
| 1 | 755836-B21 | HPE DL360 Gen9 E5-2640v3 Kit |
| 8 | 759934-B21 | HPE 8GB 2Rx4 PC4-2133P-R Kit |
| 8 | 652749-B21 | HPE 1TB 6G SAS 7.2K 2.5in SC MDL HDD |
| 1 | 700699-B21 | HPE Ethernet 10Gb 2P 561FLR-T Adptr |
| 1 | 749974-B21 | HPE Smart Array P440ar/2G FIO Controller |
| 2 | 720478-B21 | HPE 500W FS Plat Ht Plg Pwr Supply Kit |
| 1 | 663201-B21 | HPE 1U SFF Ball Bearing Rail Kit |
| 1 | C6N36ABE | HPE Insight Control ML/DL/BL Bundle E-LTU |
| | C6N36A | HPE 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 11.** BOM for the HPE ProLiant DL380 Gen9 server configuration

| QTY | PART NUMBER | DESCRIPTION |
| --- | --- | --- |
| 1 | | **Worker node** |
| 1 | 719061-B21 | HPE DL380 Gen9 12-LFF CTO Server |
| 1 | 762766-L21 | HPE DL380 Gen9 E5-2680v3 FIO Kit |
| 1 | 762766-B21 | HPE DL380 Gen9 E5-2680v3 Kit |
| 8 | 726719-B21 | HPE 16GB 2Rx4 PC4-2133P-R Kit |
| 1 | 835565-B21 | HPE Dual 340GB Read Intensive-2 Solid State M.2 Enablement Kit |
| 15 | 628061-B21 | HPE 3TB 6G SATA 7.2k 3.5in SC MDL HDD |
| 1 | 700699-B21 | HPE Ethernet 10Gb 2P 561FLR-T FIO Adptr |
| 1 | 761874-B21 | HPE Smart Array P840/4G FIO Controller |
| 1 | 768856-B21 | HPE DL380 Gen9 3LFF Rear SAS/SATA Kit |
| 1 | 727250-B21 | HPE 12Gb DL380 Gen9 SAS Expander Card |
| 1 | 720864-B21 | HPE 2U LFF BB Gen9 Rail Kit |
| 2 | 720479-B21 | HPE 800W FS Plat Ht Plg Pwr Supply Kit |
| 1 | C6N36ABE | HPE Insight Control ML/DL/BL Bundle E-LTU |
| | C6N36A | HPE 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 12.** Network - Top of Rack switch

| QTY | PART NUMBER | DESCRIPTION |
| --- | --- | --- |
| | | **Top-of-Rack network switches** |
| 2 | JG336A | HPE 5900AF-48XGT-4QSFP+ Switch |
| 2 | JG326A | HPE X240 40G QSFP+ QSFP+ 1m DAC Cable |
| 4 | JC680A | HPE A58x0AF 650W AC Power Supply |
| 4 | JG553A | HPE X712 Bck(pwr)-Frt(prt) HV Fan Tray |

**Table 13.** Network - Aggregation/Spine switch (Only required for first expansion rack. Not required for single-rack architecture)

| QTY | PART NUMBER | DESCRIPTION |
| --- | --- | --- |
| | | **Aggregation switches** |
| 2 | JH396A | HPE FF 5940 32QSFP+ |
| 4 | JC680A | HPE 58x0AF 650W AC Power Supply |
| 4 | JG553A | HPE X712 Bck(pwr)-Frt(prt) HV Fan Tray |
| 2 | JG326A | HPE X240 40G QSFP+ QSFP+ 1m DAC Cabl |
| 8 | JG328A | HPE X240 40G QSFP+ QSFP+ 5m DAC Cable |

**Table 14.** 10GbE PCI NIC instead of FlexibleLOM NIC

| QTY | PART NUMBER | DESCRIPTION |
| --- | --- | --- |
| 1 | 716591-B21 | HPE Ethernet 10Gb 2-port 561T Adapter |

**Table 15.** Network - Top of Rack switch for separate iLO and PXE network

| QTY | PART NUMBER | DESCRIPTION |
| --- | --- | --- |
| 1 | JG510A | HPE 5900AF-48G-4XG-2QSFP+ Switch |
| 1 | JC680A | HPE A58x0AF 650W AC Power Supply |
| 2 | JC682A | HPE A58x0AF Back (power side) to Front (port side) Airflow Fan Tray |

## Other hardware and software BOM

**Table 16.** Hardware - Rack and PDU

| QTY | PART NUMBER | DESCRIPTION |
|---|---|---|
| | | **Rack infrastructure** |
| 4 | AF520A | HPE Intelligent 4.9kVA/L6-30P/NA/J PDU |
| 8 | AF547A | HPE 5xC13 Intelligent PDU Extension Bar G2 Kit |
| 1 | BW946A | HPE 42U Location Discovery Kit |
| 1 | BW904A | HPE 642 1075mm Shock Intelligent Series Rack |
| 1 | BW932A | HPE 600mm Rack Stabilizer Kit |
| 1 | BW930A | HPE Air Flow Optimization Kit |
| 1 | BW906A | HPE 42U 1075mm Side Panel Kit |
| 1 | BW891A | HPE Rack Grounding Kit |

**Note**

The quantity specified below is for a single node.

**Table 17.** Software - HPE Insight Cluster Management Utility (CMU) options

| QTY | PART NUMBER | DESCRIPTION |
|---|---|---|
| | | **CMU-options** |
| 1 | QL803B | HPE Insight CMU 1yr 24x7 Flex Lic |
| 1 | QL803BAE | HPE Insight CMU 1yr 24x7 Flex E-LTU |
| 1 | BD476A | HPE Insight CMU 3yr 24x7 Flex Lic |
| 1 | BD476AAE | HPE Insight CMU 3yr 24x7 Flex E-LTU |
| 1 | BD477A | HPE Insight CMU Media |

**Table 18.** Software - Red Hat Enterprise Linux

| QTY | PART NUMBER | DESCRIPTION |
|---|---|---|
| 21 | G3J28AAE | RHEL Svr 2 Sckt/2 Gst 1yr 24x7 E-LTU |

**Note**

While HPE is a certified reseller of Hortonworks HDP software subscription, all application support for HDP software is provided by Hortonworks.

**Table 19.** Software - Hortonworks subscription

| QTY | PART NUMBER | DESCRIPTION |
|---|---|---|
| 5 | F5Z52A | Hortonworks Data Platform Enterprise 4 Nodes or 50TB Raw Storage 1 year 24x7 Support LTU. |

# 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. For additional HDP performance tuning, visit https://docs.hortonworks.com/HDPDocuments/HDP2/HDP-2.4.0/bk_installing_manually_book/content/determine-hdp-memory-config.html

- 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>`

  – Decrease disk swappiness to minimum 1:

    `Set sysctl vm.swappiness=1 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 two 10GbE server ports, for 20GbE throughput.

  – Ensure forward and reverse DNS is working properly.

  – Install and configure ntp to ensure clocks on each node are in sync to the Management node.

  – For good performance improvements, disable transparent huge page compaction:

    `echo never > /sys/kernel/mm/transparent_hugepage/enabled`

- 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`

- 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:

  – HPE Power Profile ➔ Maximum Performance

  – HPE 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

- 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.

- HPE ProLiant BIOS

  – SPP version >= 2015.04.0 (B)

  – Update System BIOS version to be >= P89

  – Update HPE Integrated Lights-Out (iLO) version to be >= 2.03

  – Intel Virtualization Technology ➔ Disabled

  – Intel VT-d ➔ Disabled

- Embedded HPE Dynamic Smart Array B140i controller

  – The embedded HPE Dynamic Smart Array B140i controller used for m.2 OS drives will operate in UEFI only mode. The HPE Dynamic Smart Array B140i defaults to AHCI off the chipset. As shown below, the HPE Smart Array needs to be enabled in BIOS on the SATA-only models, if required.

    Embedded SATA Configuration ➔ SATA_RAID_ENABLED

  – By default, the HPE Dynamic Smart Array B140i will not operate in Legacy mode. For legacy support, additional controllers will be needed; and, for CTO orders, select the Legacy mode settings part, 758959-B22.

  – Enable caching on the OS logical drive

    ```
    ctrl slot=<slot number> create type=ld drives=<drives> caching=enable
    ```

- HPE Smart Array P440ar / P840

  - Update controller firmware to be >= v1.34

  - Configure each Hadoop data drive as a separate RAID0 array with stripe size of 1024KB

  - Turn Off "Array Acceleration" / "Caching" for all data drives

    Example:

    ```
    ctrl slot=<slot number> ld all modify caching=disable ← disable caching on all logical drives on
    1st ctrlr
    ```

    ```
    ctrl slot=<slot number> ld 1 modify caching=enable ← enable caching on the OS logical drive on
    1st ctrlr
    ```

  - Tune array controller stripe size to 1024MB:

    ```
    hpssacli ctrl slot=<slot number> ld <ld number> modify ss=1024
    ```

- Network cards

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

- Oracle Java

  ```
  java.net.preferIPv4Stack set to true
  ```

- Patch common security vulnerabilities

  - Check Red Hat Enterprise Linux and SUSE security bulletins for more information.

## Appendix C: Alternate parts

**Table C-1.** Alternate processors - HPE ProLiant DL380

| QTY | PART NUMBER | DESCRIPTION |
| --- | --- | --- |
| 1 | 762764-L21 | HPE DL380 Gen9 E5 2660v3 FIO Kit (10 cores/2.6GHz) |
| 1 | 762764-B21 | HPE DL380 Gen9 E5 2660v3 Kit |
| 1 | 719054-L21 | HPE DL380 Gen9 E5-2697v3 FIO Kit (14 cores/2.3GHz) |
| 1 | 719054-B21 | HPE DL380 Gen9 E5 2697v3 Kit |
| 1 | 817945-L21 | HPE DL380 Gen9 E5-2660v4 FIO Kit (14 cores/2.0GHz) |
| 1 | 817945-B21 | HPE DL380 Gen9 E5 2660v4 Kit |
| 1 | 817951-L21 | HPE DL380 Gen9 E5-2680v4 FIO Kit (14 cores/2.4GHz) |
| 1 | 817951-B21 | HPE DL380 Gen9 E5 2680v4 Kit |
| 1 | 817959-L21 | HPE DL380 Gen9 E5-2690v4 FIO Kit (14 cores/2.6GHz) |
| 1 | 817959-B21 | HPE DL380 Gen9 E5 2690v4 Kit |

**Table C-2.** Alternate memory - HPE ProLiant DL380

| QTY | PART NUMBER | DESCRIPTION |
|-----|-------------|-------------|
| 8 | 726719-B21 | HPE 16GB 2Rx4 PC4-2133P-R Kit for 128GB of Memory |
| 8 | 728629-B21 | HPE 32GB 2Rx4 PC4-2133P-R Kit for 256GB of Memory |
| 16 | 728629-B21 | HPE 32GB 2Rx4 PC4-2133P-R Kit for 512GB of Memory |
| 8 | 836220-B21 | HPE 16GB 2Rx4 PC4-2400P-R Kit for 128GB of Memory |
| 8 | 805351-B21 | HPE 32GB 2Rx4 PC4-2400P-R Kit for 256GB of Memory |
| 16 | 805353-B21 | HPE 32GB 2Rx4 PC4-2400P-R Kit for 512GB of Memory |

**Table C-3.** Alternate disk drives - HPE ProLiant DL380

| QTY | PART NUMBER | DESCRIPTION |
|-----|-------------|-------------|
| 15 | 652757-B21 | HPE 2TB 6G SAS 7.2K 3.5in SC MDL HDD |
| 15 | 652766-B21 | HPE 3TB 6G SAS 7.2K 3.5in SC MDL HDD |
| 15 | 695510-B21 | HPE 4TB 6G SAS 7.2K 3.5in SC MDL HDD |
| 15 | 761477-B21 | HPE 6TB 6G SAS 7.2K 3.5in SC MDL HDD |
| 15 | 793703-B21 | HPE 8TB 12G SAS 7.2K 3.5in 512e SC HDD |
| 15 | 658079-B21 | HPE 2TB 6G SATA 7.2k 3.5in SC MDL HDD |
| 15 | 628061-B21 | HPE 3TB 6G SATA 7.2k 3.5in SC MDL HDD |
| 15 | 693687-B21 | HPE 4TB 6G SATA 7.2k 3.5in SC MDL HDD |
| 15 | 765255-B21 | HPE 6TB 6G SATA 7.2k 3.5in SC MDL HDD |
| 15 | 793695-B21 | HPE 8TB 6G SATA 7.2K 3.5in 512e SC HDD |

**Table C-4.** Alternate network cards - HPE ProLiant DL380

| QTY | PART NUMBER | DESCRIPTION |
|-----|-------------|-------------|
| 1 | 665243-B21 | HPE Ethernet 10GbE 560FLR SFP+ FIO Adapter for 10Gb networking only |
| 1 | 665249-B21 | HPE Ethernet 10Gb 2P 560SFP+ 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 HPE FlexFabric 5900 SFP+ equivalent ToR network switches are required (HPE FlexFabric 5900AF-48XG-4QSFP+ Part number JC772A).

**Table C-5.** Alternate controller cards - HPE ProLiant DL380 and HPE ProLiant DL360

| QTY | PART NUMBER | DESCRIPTION |
|-----|-------------|-------------|
| 1 | 726821-B21 | HPE Smart Array P440/4GB FBWC 12Gb 1-port Int SAS Controller |
| 1 | 726897-B21 | HPE Smart Array P840/4GB FBWC 12Gb 2-port Int SAS Controller |
| 1 | 749976-B21 | HPE H240ar FIO Smart HBA |
| 1 | 764630-B21 | HPE DL360 Gen9 2SFF HDD Kit (Required when upgrading an 8SFF DL360 model to a 10SFF model) |
| 1/Server | C9A82AAE | HPE Secure Encryption per Svr Entitlement |

# Appendix D: HPE value-added services and support

In order to help customers jump-start their Hadoop solution development, HPE 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. HPE 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 (Hortonworks)
With HPE Reference Architecture Implementation Service for Hadoop, experienced HPE Big Data consultants install, configure, deploy, and test your Hadoop environment based on the HPE 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.

HPE Factory Express Level 4 Service (HA454A1) is the recommended Factory Integration service for Big Data covering hardware and software integration, as well as end-to-end delivery project management.

For more information and assistance on Factory Integration services, you can contact:

- AMS: easy.solutions.americas@hpe.com
- APJ: ap.fe-engagement@hpe.com
- EMEA: sol_eng_support@hpe.com

### TS Consulting – Big Data services
HPE Big Data Services can help you reshape your IT infrastructure to corral 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:** HPE experts will design and implement a high-performance, integrated platform to support a strategic architecture for Big Data. Choose from design and implementation services, reference architecture implementations and integration services. Your flexible, scalable infrastructure will support Big Data variety, consolidation, analysis, share and search on HPE 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, visit: hpe.com/us/en/services/consulting/big-data.html

### HPE Support options
HPE offers a variety of support levels to meet your needs.

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

### HPE 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 HPE technical resources, 4-hour response onsite hardware support and software updates.

### HPE Proactive Care

HPE 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 HPE'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.

### HPE Proactive Care with the HPE Personalized Support Option

Adding the Personalized Support Option for HPE Proactive Care is highly recommended. The Personalized Support option builds on the benefits of HPE 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.

### HPE Proactive Select

And to address your ongoing/changing needs, HPE 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.

### Other offerings

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

### More information

For additional information, visit:

HPE Education Services: http://www8.hp.com/us/en/trainingww/bigdataoverview.html

HPE Big Data Solutions: hpe.com/info/bigdata

HPE Services: hpe.com/services

## Resources and additional links

Hortonworks, hortonworks.com

HPE Solutions for Apache Hadoop, hpe.com/info/hadoop

Hadoop and Vertica, hpe.com/info/vertica

HPE Insight Cluster Management Utility (CMU), hpe.com/info/cmu

HPE FlexFabric 5900 switch series, hpe.com/networking/5900

HPE FlexFabric 5940 switch series, http://www8.hp.com/us/en/products/networking-switches/product-detail.html?oid=1009161439

HPE ProLiant servers, hpe.com/info/proliant

HPE Networking, hpe.com/networking

HPE Services, hpe.com/services

Red Hat, redhat.com


To help us improve our documents, please provide feedback at hpe.com/contact/feedback.


**About Hortonworks**

Hortonworks is a leading innovator in the industry, creating, distributing and supporting enterprise-ready open data platforms and modern data applications. Our mission is to manage the world's data. We have a single-minded focus on driving innovation in open source communities such as Apache Hadoop, NiFi, and Spark. We along with our 1600+ partners provide the expertise, training and services that allow our customers to unlock transformational value for their organizations across any line of business. Our connected data platforms powers modern data applications that deliver actionable intelligence from all data: data-in-motion and data-at-rest.

**Sign up for updates**

**Hewlett Packard Enterprise**