Technical white paper

# HP Reference Architecture for Cloudera Enterprise

HP Converged Infrastructure with Cloudera Enterprise for Apache Hadoop

## Table of contents

# Executive summary

HP and Apache Hadoop allow you to derive new business insights from Big Data by providing a platform to store, manage and process data at scale. However, Apache Hadoop is complex to deploy, configure, manage and monitor. This white paper provides several performance optimized configurations for deploying Cloudera Enterprise 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 Cloudera's Distribution including Apache Hadoop (CDH), specifically CDH3u3, Cloudera Manager 3.7 and the HP ProLiant DL Gen8 server platform. The configurations reflected in this document have been jointly designed and developed by HP and Cloudera to provide optimum computational performance for Hadoop and are also compatible with CDH4 upon its release.

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

- **For Vertica**, the 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.

- **For networking**, the HP 5830AF-48G 1GbE Top of Rack switch and the HP 5920AF-24XG 10GbE Aggregation switch provide IRF Bonding and sFlow which simplifies the management, monitoring and resiliency of the customer's Hadoop network. In addition, the 1GB and 3.6GB respective packet buffers increase Hadoop network performance by seamlessly handling burst scenarios such as Shuffle, Sort and Block Replication which are common in a Hadoop network.

- **For servers**, the HP ProLiant Gen8 DL360p and DL380p include:
  - The HP Smart Array P420i controller which provides 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 JBOD or various RAID configurations.
  - Two sockets with the fastest 6 core processors and the Intel® C600 Series Chipset, providing the performance required for fastest time to completion for CPU bound Hadoop workloads.
  - The HP FlexibleLOM network technologies provide customers the flexibility to easily move between 1GbE and 10GbE network interconnects on the server.
  - The HP iLO Management Engine on the servers contains HP Integrated Lights-Out 4 (iLO 4) which 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.

- **For management**, HP Insight Cluster Management Utility (CMU) provides push-button scale out and provisioning with industry leading provisioning performance (deployment of 800 nodes in 30 minutes), reducing deployments from days to hours. In addition, CMU provides real-time and historical infrastructure and Hadoop monitoring with 3D visualizations allowing customers to easily characterize Hadoop workloads and cluster performance reducing complexity and improving system optimization leading to improved performance and reduced cost. HP Insight Management and HP Service Pack for ProLiant, allow for easy management of firmware and the server.

All of these features reflect HP's 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 Cloudera Enterprise software on HP infrastructure for clusters of various sizes. It is also intended to concretely identify the software and hardware components required in a solution in order 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 and HBase computational performance, resulting in a 3x performance advantage over the closest competition[2].

**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 Cloudera solution. An

---

[1] Compared to the previous generation of Smart Array controllers
[2] http://www.hp.com/hpinfo/newsroom/press_kits/2012/HPDiscover2012/Hadoop_Appliance_Fact_Sheet.pdf

intermediate knowledge of Apache Hadoop and scale out infrastructure is recommended. Those already possessing expert knowledge about these topics may proceed directly to

# Cloudera Enterprise overview

Apache Hadoop is an open-source project administered by the Apache Software Foundation. Hadoop's contributors work for some of the world's biggest technology companies. That diverse, motivated community has produced a genuinely innovative platform for consolidating, combining and understanding large-scale data in order to better comprehend the data deluge. Enterprises today collect and generate more data than ever before. Relational and data warehouse products excel at OLAP and OLTP workloads over structured data. Hadoop, however, was designed to solve a different problem: the fast, reliable analysis of both structured data and complex data. As a result, many enterprises deploy Hadoop alongside their legacy IT systems, which allows them to combine old data and new data sets in powerful new ways.

Technically, Hadoop consists of two key services: reliable data storage using the Hadoop Distributed File System (HDFS) and high-performance parallel data processing using a technique called MapReduce. Hadoop runs on a collection of commodity, shared-nothing servers. You can add or remove servers in a Hadoop cluster at will; the system detects and compensates for hardware or system problems on any server. Hadoop, in other words, is self-healing. It can deliver data – and can run large-scale, high-performance processing jobs – in spite of system changes or failures. Originally developed and employed by dominant web companies like Yahoo and Facebook, Hadoop is now widely used in finance, technology, telecom, media and entertainment, government, research institutions and other markets with significant data. With Hadoop, enterprises can easily explore complex data using custom analyses tailored to their information and questions.
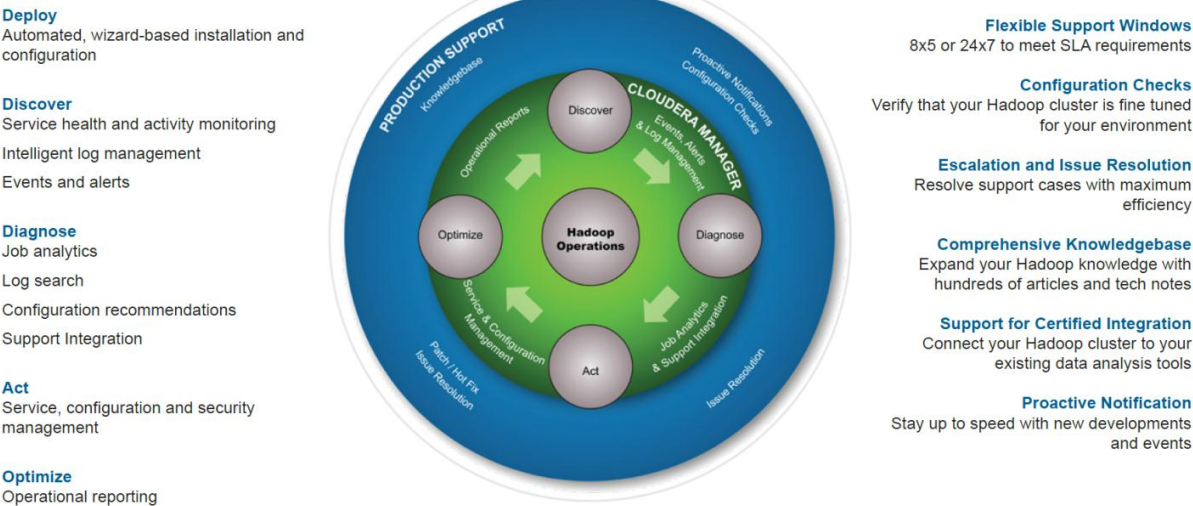
Cloudera is an active contributor to the Hadoop project and provides an enterprise-ready, 100% open source distribution that includes Hadoop and related projects. Cloudera's distribution bundles the innovative work of a global open-source community; this includes critical bug fixes and important new features from the public development repository and applies all this to a stable version of the source code. In short, Cloudera integrates the most popular projects related to Hadoop into a single package, which is run through a suite of rigorous tests to ensure reliability during production. In addition, Cloudera Enterprise is a subscription offering which enables data-driven enterprises to run Apache Hadoop environments in production cost effectively with repeatable success. Comprised of Cloudera Support and Cloudera Manager, a software layer that delivers deep visibility into and across Hadoop clusters, Cloudera Enterprise gives Hadoop operators an efficient way to precisely provision and manage cluster resources. It also allows IT shops to apply familiar business metrics – such as measurable SLAs and chargebacks – to Hadoop environments so they can run at optimal utilization. Built-in predictive capabilities anticipate shifts in the Hadoop infrastructure, ensuring reliable operation.

Cloudera Enterprise makes it easy to run open source Hadoop in production:

- Simplify and accelerate Hadoop deployment
- Reduce the costs and risks of adopting Hadoop in production
- Reliably operate Hadoop in production with repeatable success
- Apply SLAs to Hadoop
- Increase control over Hadoop cluster provisioning and management

For detailed information on Cloudera Enterprise, please see cloudera.com/products-services/enterprise/

Figure 1. Cloudera Enterprise

**Deploy**
Automated, wizard-based installation and configuration

**Discover**
Service health and activity monitoring

Intelligent log management

Events and alerts

**Diagnose**
Job analytics

Log search

Configuration recommendations

Support Integration

**Act**
Service, configuration and security management

**Optimize**
Operational reporting

**Flexible Support Windows**
8x5 or 24x7 to meet SLA requirements

**Configuration Checks**
Verify that your Hadoop cluster is fine tuned for your environment

**Escalation and Issue Resolution**
Resolve support cases with maximum efficiency

**Comprehensive Knowledgebase**
Expand your Hadoop knowledge with hundreds of articles and tech notes

**Support for Certified Integration**
Connect your Hadoop cluster to your existing data analysis tools

**Proactive Notification**
Stay up to speed with new developments and events

Typically, Hadoop clusters are either used for batch MapReduce analysis of data or they are used to run HBase, which is an online distributed store for reading and writing structured data. It is up to the user to choose which services to install and configure. We recommend that you run either MapReduce or HBase on your worker nodes as running both the master (JobTracker and HBaseMaster) and worker services (TaskTracker and HBaseRegionServer) will result in both services competing for the same resources, thereby resulting in degraded performance.

The platform functions within Cloudera Enterprise 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. The two tables below specify which services are management services and which services are workers services. Each table contains two columns. The first column is the description of the service and the second column specifies the number of nodes the service can be distributed to. 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 so this is not something you will need to manage.

# Management services

Table 1. Cloudera Management Services

| Service | Maximum Distribution across Nodes |
| --- | --- |
| Cloudera Manager | 1 |
| Hue Server | 1 |
| JobTracker | 1 |
| HBase Master | Varies |
| NameNode | 1 |
| Secondary NameNode | 1 |

## Worker services

Table 2. Cloudera Enterprise Worker Services

| Service | Maximum Distribution across Nodes |
|---|---|
| DataNode | Most or all nodes |
| TaskTracker | Most or all nodes |
| HBase RegionServer | Varies |

# Pre-deployment considerations

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 suite a particular custom scenario.

| 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

Cloudera Manager 3.7 supports only the following 64-bit operating systems:

- For Red Hat systems, Cloudera provides 64-bit packages for Red Hat Enterprise Linux 5 and Red Hat Enterprise Linux 6. Cloudera recommends using update 5 or later for Red Hat Enterprise Linux 5.
- For SUSE systems, Cloudera provides 64-bit packages for SUSE Linux Enterprise Server 11 (SLES 11). Service pack 1 or later is required.

CDH3u3 supports the following 32-bit and 64-bit operating systems:

- For Ubuntu systems, Cloudera provides 32-bit and 64-bit packages for Lucid (10.04) and Maverick (10.10).
- For Debian systems, Cloudera provides 32-bit and 64-bit packages for Squeeze (6.0.2) and Lenny (5.0.8).
- For Red Hat systems, Cloudera provides 32-bit and 64-bit packages for Red Hat Enterprise Linux 5 and CentOS 5, and 64-bit packages for Red Hat Enterprise Linux 6 and CentOS 6. Cloudera recommends using update 5 or later of Red Hat Enterprise Linux 5.
- For SUSE systems, Cloudera provides 64-bit packages for SUSE Linux Enterprise Server 11 (SLES 11). Service pack 1 or later is required.

HP recommends using a 64-bit operating system to avoid constraining the amount of memory that can be used on worker nodes. 64-bit Red Hat Enterprise Linux 5.5 update 5 or greater is recommended due to better ecosystem support, more comprehensive functionality for components such as RAID controllers and compatibility with HP Insight CMU. The Reference Architectures listed in this document were tested with 64-bit Red Hat Enterprise Linux 6.2.

## Computation

MapReduce slots are configured on a per server basis and are decided upon via an examination of the resources available on the server and how they can cater to the requirements of the tasks involved in a Hadoop Job. The processing or computational capacity of a Hadoop cluster is determined by the aggregate number of MapReduce slots available across all the worker nodes. Employing Hyper-Threading increases your effective core count, potentially allowing you to configure more MapReduce slots. Refer to the Storage section below to see how I/O performance issues arise from sub-optimal disk to core ratios (too many slots and too few disks).

To remove the bottleneck for CPU bound workloads, for the best cost/performance tradeoff, we recommend buying 6 core processors with faster clock speeds as opposed to buying 8 core processors.

## 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 TaskTracker and DataNode processes in addition to the sum of all the memory assigned to each of the MapReduce slots. If you have a memory bound MapReduce job we recommend that you increase the amount of memory on all the worker nodes. In addition, the cluster can also be used for HBase which is very memory intensive.

It is important to saturate all the memory channels available to ensure optimal use of the memory bandwidth. For example, on a two socket processor with eight memory channels available per server one would typically fully populate the channels with either 4GB DIMMs or 8GB DIMMs resulting in a configuration of 32GB or 64GB of memory per server, respectively.

## 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.

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 and not as DAS with RAID or Network Attached Storage (NAS).

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 16 TB of block data fails, 16 TB of block data will be replicated between other servers in the cluster to ensure the appropriate amount of replicas exist. Furthermore, the failure of a non-redundant TOR (Top of Rack) switch will generate even more replication traffic. One needs to ensure that the performance of the network is sufficient to adequately handle MapReduce shuffle and sort phases occurring at the same time as block replication.

- **I/O performance** – Each worker node has a certain number of MapReduce slots available for processing Hadoop tasks. Each slot operates on one block of data at a time. 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 disks. Specific details around disk and RAID configurations will be provided in the Server selection section.

## 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 flood of network traffic as Hadoop rebalances storage, and in a single-rack system such a failure brings down the whole cluster. Consequently, configuring two TOR switches per rack is recommended for all production configurations.

Hadoop is rack-aware and tries to limit the amount of network traffic between racks. The bandwidth and latency provided by a 1 Gigabit Ethernet (GbE) connection from worker nodes to the TOR switch is 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 10 GbE 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.

During the map phase of Hadoop jobs that utilize the HDFS, the majority of tasks reference data on the server that executes the task (node-local). For those tasks that must access data remotely, the data is usually on other servers in the same rack (rack-local). Only a small percentage of tasks need to access data from remote racks. Although the amount of remote-rack accesses increases for larger clusters, it is expected to put a relatively small load on the TOR and core switches.

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. Consequently, large clusters will benefit from having TOR switches with packet buffering and connected by 10 GbE uplinks to core aggregation switches in order to accommodate this 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. This is why we recommend TOR switches with deep packet buffering.

# Switches

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

## Aggregation switches

The HP 5920AF-24XG 10GbE switch is an ideal aggregation switch as it is well suited to handle large volumes of inter-rack traffic and scenarios such as block replication occurring at the same time as a MapReduce shuffle and sort phase. The switch has a 3.6 GB packet buffer depth for very deep buffering, aggregation switch redundancy and better high availability (HA) support with IRF bonding and 24 10GbE ports. For more information on the HP 5920AF-24XG, please see http://h17007.www1.hp.com/us/en/products/switches/HP_5920_Switch_Series/index.aspx

The configuration for the HP 5920AF-24XG switch is provided below.

Figure 2. HP 5920AF-24XG Aggregation switch



Table 3. HP 5920AF-24XG Single Aggregation Switch options

| Qty | Description |
| --- | --- |
| 1 | HP 5920AF-24XG Switch |
| 2 | HP 58x0AF 650W AC Power Supply |
| 2 | HP 5920AF-24XG Front (port-side) to Back (power-side) Airflow Fan Tray |
| 24 | HP X130 10G SFP+ LC SR Transceiver |
| 24 | HP 15m Premier Flex LC/LC Optical Cable |

## Top of Rack (TOR) switches

The HP 5830AF-48G is an ideal TOR switch and has a 1 GB packet buffer depth for very deep buffering, resiliency and high availability, scalability support, forty eight 1GbE ports, two 10GbE uplinks, and the option for adding two more 10GbE ports. A dedicated management switch for iLO traffic is not required as the ProLiant DL360p Gen8 and DL380p Gen8 are able to share iLO traffic over NIC1. The volume of iLO traffic is minimal and does not degrade performance over that port.

For more information on the HP 5830AF-48G switch, please see http://h17007.www1.hp.com/us/en/products/switches/HP_5830_Switch_Series/index.aspx

The configuration for the HP 5830AF-48G switch is provided below.

Figure 3. HP 5830AF-48G Top of Rack (TOR) switch



Table 4. HP 5830AF-48G Single Switch options

| Qty | Description |
| --- | --- |
| 1 | HP 5830AF-48G Switch with 1 Interface Slot |
| 2 | HP 58x0AF 650W AC Power Supply |
| 1 | HP 5500/5120 2-port 10GbE SFP+ Module |
| 1 | HP 5830AF-48G Back(power)-Front(prt) Fan Tray |
| 1 | HP X240 10G SFP+ SFP+ 0.65m DAC Cable |

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

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.

Insight CMU is highly flexible and customizable, offers both GUI and CLI interfaces, and is being used to deploy a range of software environments, from simple compute farms to highly customized, application-specific configurations. Insight CMU is available for HP ProLiant and HP BladeSystem servers, with Red Hat Enterprise Linux and Novell SUSE Linux operating systems, including Red Hat Enterprise Linux, SUSE Linux Enterprise, CentOS, and Ubuntu. Insight CMU also includes options for monitoring GPUs and for installing GPU drivers and software.

For more information, please see hp.com/go/cmu.

Table 5. HP Insight CMU options

| Qty | Description |
| --- | --- |
| 1 | HP Insight CMU 1yr 24x7 Flex Lic |
| 1 | HP Insight CMU 1yr 24x7 Flex E-LTU |
| 1 | HP Insight CMU 3yr 24x7 Flex Lic |
| 1 | HP Insight CMU 3yr 24x7 Flex E-LTU |
| 1 | HP Insight CMU Media |

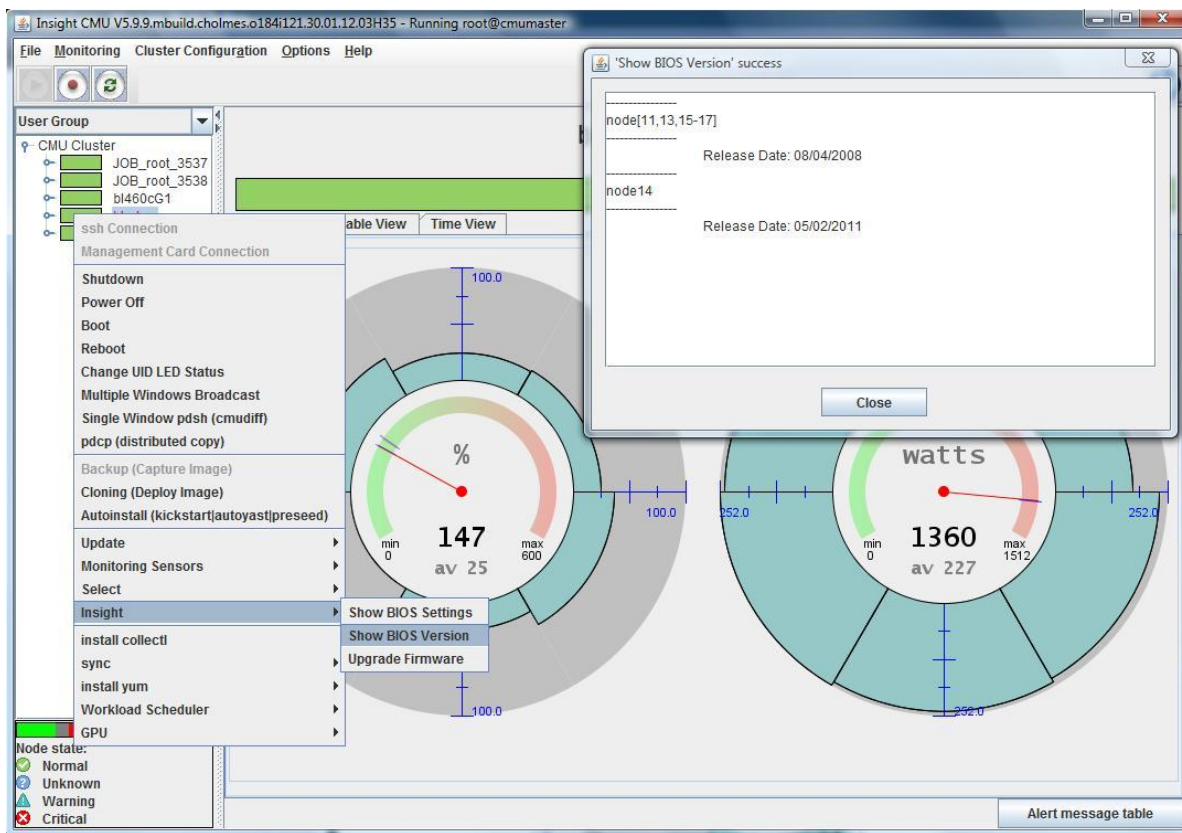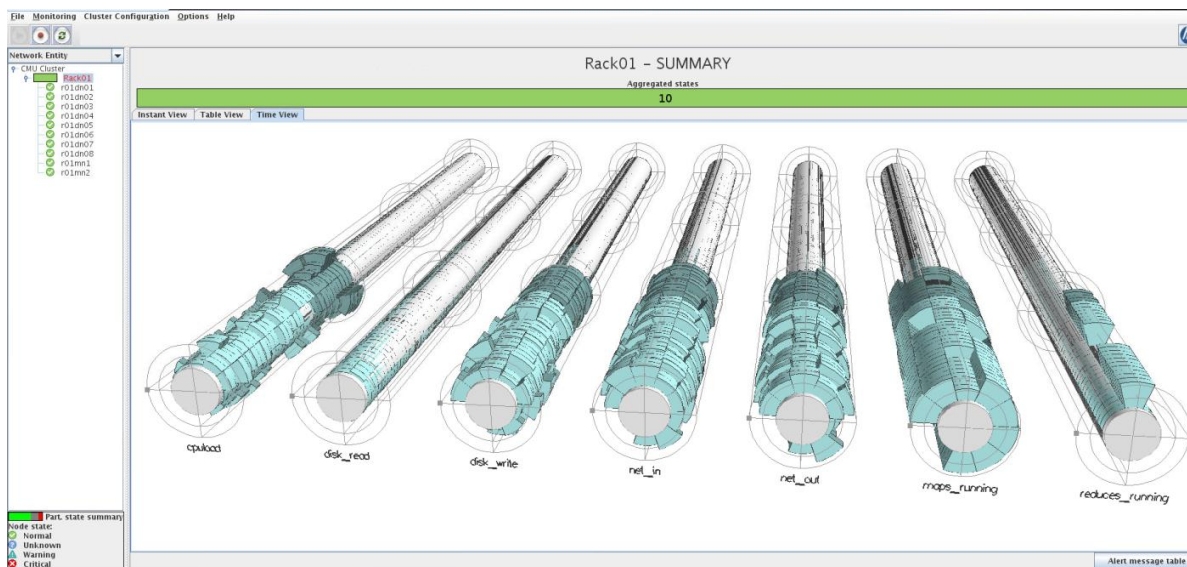Figure 4. HP Insight CMU Interface – real-time view

Figure 5. HP Insight CMU Interface – Time View



# Server selection

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. The Reference Architecture section will provide topologies for the deployment of management and worker nodes for single and multi-rack clusters.

## Management nodes

Management services are not distributed 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, as well as an array controller supporting a variety of RAID schemes and SAS direct attached storage. In addition, the management services are memory and CPU intensive; therefore, a server capable of supporting a large amount of memory is also required. Management nodes do not participate in storing data for the HDFS and have much lower storage capacity requirements than worker nodes and thus a 2U server with a large amount of disks is not required.

The configurations reflected in this white paper are also cognizant of the upcoming features in Cloudera CDH4, most noticeably, high availability. 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. 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.

This section contains 4 subsections:

- Server platform
- Management node
- JobTracker server
- NameNode server

## Server platform: HP ProLiant DL360p Gen8

The HP ProLiant DL360p Gen8 (1U) is an excellent choice as the server platform for the management nodes.

Figure 6. HP ProLiant DL360p Gen8 Server



### Processor configuration

The configuration features two sockets with the fastest 6 core processors and the Intel C600 Series, which provide 12 physical cores and 24 Hyper-Threaded cores per server at the fastest clock speeds available. We recommend that Hyper-Threading be turned on.

We recommend processors with 6 cores for the management servers because the JobTracker, NameNode and Cloudera Manager are CPU intensive and multi-threaded and will fully utilize all the cores available. Furthermore, the configurations for these servers are designed to be able to handle an increasing load as your Hadoop cluster grows so one needs to ensure the right processing capacity is available to begin with.

### Drive configuration

Apache Hadoop does not provide software redundancy for the management servers of a Hadoop cluster the same way it does for the workers and thus RAID is appropriate. The P420i Smart Array Controller is specified to drive eight 900GB 2.5" SAS disks on the Management node and four 900GB 2.5" SAS disks on the JobTracker and NameNode servers. The Management node has more disks than the JobTracker and NameNode servers due to the fact that the Management node needs to have extra storage capacity for RAID Mirroring, the Cloudera Manager Databases, and logs, as well as to act as a multi-homed staging server for data import and export out of the HDFS. Hot pluggable drives are specified so that drives can be replaced without restarting the server. Due to this design, one should configure the Gen8 P420i controller to apply the following RAID schemes:

- Management node: 4 Disks with RAID 1+0 for OS and MySQL database, 4 Disks with RAID 5 for data staging

- JobTracker and NameNode Servers: 4 Disks with RAID 1+0 for OS

The Gen8 P420i controller provides two port connectors per controller with each containing 4 SAS links. Each drive cage for the DL360p contains 8 disks and thus each disk has a dedicated SAS link which ensures the server provides the maximum throughput that each drive can give you. For a performance oriented solution, we recommend SAS drives as they offer a significant read and write throughput performance enhancement over SATA disks.

### Memory configuration

Servers running management services such as the HBaseMaster, JobTracker, NameNode and Cloudera Manager 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® E5-2667 2.9 GHz processors in the HP ProLiant DL360p Gen8 have 4 memory channels per processor which equates to 8 channels per server. The configurations for the management servers were tested with 64GB of RAM, which equated to eight 8GB DIMMS.

### Network configuration

The HP ProLiant DL360p Gen8 is designed for network connectivity to be provided via a FlexibleLOM. The FlexibleLOM can be ordered in a 4 x 1GbE NIC configuration or a 2 x 10GbE NIC configuration. This Reference Architecture was tested using the 4 x 1GbE NIC configuration (as specified in the server configuration below).

For each management server we recommend bonding and cabling only two of the 1GbE NICs to create a single bonded pair which will provide 2GbE of throughput as well as a measure of NIC redundancy. In the reference 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.

**Bill of materials**

Table 6. The HP ProLiant DL360p Gen8 Server Configuration

| Qty | Description |
|-----|-------------|
| 1 | HP DL360p Gen8 8-SFF CTO Chassis |
| 1 | HP DL360p Gen8 E5-2667 FIO Kit |
| 1 | HP DL360p Gen8 E5-2667 Kit |
| 8 | HP 8GB 1Rx4 PC3-12800R-11 Kit |
| 4 | HP 900GB 6G SAS 10K 2.5in SC ENT HDD (Note: Management node needs 8) |
| 1 | HP Ethernet 1GbE 4P 331FLR FIO Adapter |
| 1 | HP 512MB FBWC for P-Series Smart Array |
| 2 | HP 460W CS Gold Hot Plug Power Supply Kit |
| 1 | HP 1U SFF BB Gen8 Rail Kit |
| 1 | ProLiant DL36x(p) HW Support |

### Management node

The Management node hosts the applications that submit jobs to the Hadoop Cluster. We recommend that you install with the following software components:

Table 7. Management node Software

| Software | Description |
|----------|-------------|
| RHEL 6.2 | Recommended Operating System |
| HP Insight CMU 7.0 | Infrastructure Deployment, Management, and Monitoring |
| Oracle JDK 1.6.0_26 | Java Development Kit |
| MySQL | Database Server for Cloudera Manager |
| Cloudera Manager 3.7 | Cloudera Hadoop Cluster Management Software |
| Cloudera Hue Server | Web Interface for Cloudera Applications |

| Software | Description |
| --- | --- |
| NFS Server | Provides an NFS Mount for the NameNode Edit Log |
| Apache Pig and/or Apache Hive from CDH3u3 | Analytical interfaces to the Hadoop Cluster |
| *ZooKeeper* | *Synchronization service (Only if running HBase)* |

Please see the following link for the Cloudera Manager and MySQL Installation guide, https://ccp.cloudera.com/display/ENT/Cloudera+Manager+Installation+Guide.

Please see the following link for the Cloudera Hue, Apache Pig and Apache Hive installation guides, https://ccp.cloudera.com/display/CDHDOC/CDH3+Documentation

The Management node contains the following base configuration:

- Dual Six-Core Intel E5-2667 2.9 GHz Processors
- P420i Smart Array Controller
- Eight 900GB SFF SAS 10K RPM disks
- 64 GB DDR3 Memory
- 4 x 1GbE FlexibleLOM NICs

### JobTracker server

The JobTracker server contains the following software. Please see the following link for more information on installing and configuring the JobTracker and secondary NameNode, https://ccp.cloudera.com/display/CDHDOC/CDH3+Installation+Guide

Table 8. JobTracker Server Software

| Software | Description |
| --- | --- |
| RHEL 6.2 | Recommended Operating System |
| Oracle JDK 1.6.0_26 | Java Development Kit |
| JobTracker | The JobTracker for the Hadoop Cluster |
| *HBaseMaster* | *The HBase Master for the Hadoop Cluster (Only if running HBase)* |
| *ZooKeeper* | *Synchronization service (Only if running HBase)* |
| Secondary NameNode | Process to handle Check pointing from the NameNode |

The JobTracker server contains the following base configuration:

- Dual Six-Core Intel E5-2667 2.9 GHz Processors
- Four 900GB SFF SAS 10K RPM disks
- 64 GB DDR3 Memory
- 4 x 1GbE FlexibleLOM NICs
- 1 x P420i Smart Array Controller

### NameNode server

The NameNode server contains the following software. Please see the following link for more information on installing and configuring the NameNode, https://ccp.cloudera.com/display/CDHDOC/CDH3+Installation+Guide

Table 9. NameNode Server Software

| Software | Description |
| --- | --- |
| RHEL 6.2 | Recommended Operating System |
| Oracle JDK 1.6.0_26 | Java Development Kit |
| NameNode | The NameNode for the Hadoop Cluster |
| NFS Client | Allows the NameNode to write logs to the Management node |
| *ZooKeeper* | *Synchronization service (Only if running HBase)* |

The NameNode server contains the following base configuration:

- Dual Six-Core Intel E5-2667 2.9 GHz Processors
- Four 900GB SFF SAS 10K RPM disks
- 64 GB DDR3 Memory
- 4 x 1GbE FlexibleLOM NICs
- 1 x P420i Smart Array Controller

## Worker nodes

The worker nodes run the TaskTracker (or HBaseRegionServer) and DataNode processes and thus storage capacity and performance are important factors.

### Server platform: HP ProLiant DL380p Gen8

The HP ProLiant DL380p Gen8 (2U) 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 7. HP ProLiant DL380p Gen8 Server



### Processor configuration

The configuration features two sockets with the fastest 6 core processors and the Intel C600 Series which provide 12 physical or 24 Hyper-Threaded cores per server at the fastest clock speeds available. Hadoop manages the amount of work each server is able to undertake via the amount of Map/Reduce slots configured for that server. The more cores available to the server, the more MapReduce slots can be configured for the server (see the Computation section for more detail). We recommend that Hyper-Threading be turned on.

**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 HP Smart Array P420i controller by configuring each individual disk as a separate RAID 0 volume. Additionally array acceleration features on the P420i should be turned off for the RAID 0 data volumes. The worker node design includes a second P420 controller connected to a second drive cage for a total of 16 drives per server.

Customers also have the option of keeping all the drives in the JBOD configured described above, but removing two of the 1TB SAS MDL disks in the server and replacing them with two 500GB SAS MDL disks and using the HP Smart Array P420i Controller to configure those with a RAID 1 mirrored OS and Hadoop runtime. This provides additional measures of redundancy on the worker nodes. We do not recommend sharing drives that contain the OS and Hadoop runtimes with drives that contain the temporary MapReduce data and the HDFS block data as it results in degraded I/O performance.

*Performance*
The HP Smart Array P420i controller provides two port connectors per controller with each containing 4 SAS links. Each drive cage for the DL380p contains 8 disks and thus each disk has a dedicated SAS link which ensures the server provides the maximum throughput that each drive can give you. For a performance oriented solution, we recommend SAS drives as they offer a significant read and write throughput performance enhancement over SATA disks.

*Core to disk ratio*
The more drives a server contains, the more efficiently it can service I/O requests because it reduces the likelihood of multiple threads contending for the same drive which can result in interleaved I/O and degraded performance.

*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.

**Memory configuration**
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 Intel Xeon E5-2667 has 4 memory channels per processor. When configuring memory, one should always attempt to populate all the memory channels available to ensure optimum performance.

To ensure optimal memory performance and bandwidth, we recommend using 4 or 8GB DIMMs to populate each of the 4 memory channels on the processor which will provide an aggregate of 32 or 64GB of RAM respectively. We used 8GB DIMMs (which gave us an aggregate of 64GB of RAM per server) in our testing.

**Network configuration**
For 1GbE networks we recommend that the four 1GbE NICs be bonded to improve throughput performance to 4 Gb/s and thereby improve performance. In addition, in the reference 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 and 3 to Switch 1 and NIC 2 and 4 to Switch 2.

The worker node contains the following software. Please see the following link for more information on installing and configuring the TaskTracker (or HBaseRegionServer) and DataNode,
https://ccp.cloudera.com/display/CDHDOC/CDH3+Installation+Guide

Table 10. Worker Node Software

| Software | Description |
|---|---|
| RHEL 6.2 | Recommended Operating System |
| Oracle JDK 1.6.0_26 | Java Development Kit |
| TaskTracker | The TaskTracker process for Map/Reduce (Only if running Map/Reduce) |
| DataNode | The DataNode process for HDFS |
| *HBaseRegionServer* | *The HBaseRegionServer for HBase (Only if running HBase)* |

A worker node server contains the following base configuration:

- Dual Six-Core Xeon E5-2667 2.9 GHz Processors with Hyper-Threading
- 16 x 1TB 2.5" SAS MDL 7.2K RPM disks
- 64 GB DDR3 Memory
- 4 x 1GbE NICs FlexibleLOM
- 2 x P420 Smart Array Controllers

> The DL380p ships with an onboard controller (P420i) so only one additional controller (P420) is required. Customers also have the option of purchasing a second power supply for additional power redundancy.

Table 11. The HP ProLiant DL380p Gen8 Server Configuration

| Qty | Description |
|---|---|
| 1 | HP DL380p Gen8 8-SFF CTO Chassis |
| 1 | HP DL380p Gen8 E5-2667 FIO Kit |
| 1 | HP DL380p Gen8 E5-2667 Kit |
| 8 | HP 8GB 1Rx4 PC3-12800R-11 Kit |
| 16 | HP 1TB 6G SAS 7.2K 2.5in SC MDL HDD |
| 1 | HP 380/385 Gen8 8-SFF Cage/Backplane Kit |
| 1 | HP Smart Array P420/1GB FBWC Controller |
| 1 | HP Ethernet 1GbE 4P 331FLR FIO Adapter |

| Qty | Description |
| --- | --- |
| 1 | HP 512MB FBWC for P-Series Smart Array |
| 1 | HP 750W CS Gold Hot Plug Power Supply Kit |
| 1 | ProLiant DL38x(p) HW Support |

# Reference Architectures

The following illustrates a reference progression of Hadoop clusters from a single rack to a multi-rack configuration.

## Single Rack Reference Architecture

The Single Rack Cloudera 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 5830AF-48G TOR switches are specified for performance and redundancy. The HP 5830AF-48G includes up to four 10GbE uplinks which can be used to connect the switches in the rack into the desired network or the 10GbE HP 5920AF Aggregation switch. Keep in mind that if IRF bonding is used, it requires up to 2 10GbE ports per switch, which would leave between 2 to 3 10GbE 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 for 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 JobTracker, NameNode and Worker nodes on their own private Hadoop Cluster subnet.

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. Cloudera Manager is a web application that runs on the Management node and allows users to be able 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 be able 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 move data onto it. The HDFS client needs to potentially be able to reach every Hadoop DataNode in the cluster in order to stream blocks onto it to move data onto the HDFS. The Reference Architecture provides two ways to do this.

The first option is to use the already multi-homed Management node. This server has been configured with twice the amount of disk capacity (an additional 3.6 TB) compared to the other management servers in order to provide an immediate solution to move 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 all 4 NICs are used on each worker node and 2 NICs are used on each management node it

leaves 16 ports still available across both the switches in the rack. These 16 1GbE ports or the remaining 10GbE ports on the switches can be used by other multi-homed systems outside of the Hadoop cluster to move data into the Hadoop Cluster.

## Configuring backup for the Hadoop NameNode

In Cloudera's CDH3u3, the Apache Hadoop NameNode is not highly available. In order to facilitate swift recovery of a NameNode, Hadoop has the NameNode checkpoint the FSImage which is configured on the Secondary NameNode i.e. JobTracker server.

> To ensure ease of NameNode recovery we recommend that you persist the NameNode metadata and edit logs for Hadoop to at least two directories, one of which is located on an NFS mount point on the Management node.

To do this, one needs to create the NFS mount and then reference it in the dfs.name.dir property in the hdfs-site.xml file, as referenced below:

**hdfs-site.xml:**

```
<property>
  <name>dfs.name.dir</name>
  <value>/data/1/dfs/nn,/nfsmount/dfs/nn</value>
</property>
```

In the following instructions, local path examples are used to represent Hadoop parameters. Change the path examples to match your configuration.

1.  Create the dfs.name.dir local directories:

    ```
    $ sudo mkdir -p /data/1/dfs/nn /nfsmount/dfs/nn
    ```

2.  Configure the owner of the dfs.name.dir and dfs.data.dir directories to be the hdfs user:

    ```
    $ sudo chown -R hdfs:hadoop /data/1/dfs/nn /nfsmount/dfs/nn
    ```

Use the chmod command to reset permissions for these dfs.name.dir directories to drwx------ (700); for example:

```
$sudo chmod 700 /data/1/dfs/nn /nfsmount/dfs/nn
```

In addition, to keep NameNode processes from hanging when the NFS server is unavailable, configure the NFS mount as a soft mount (so that I/O requests that time out fail rather than hang), and set other options as follows:
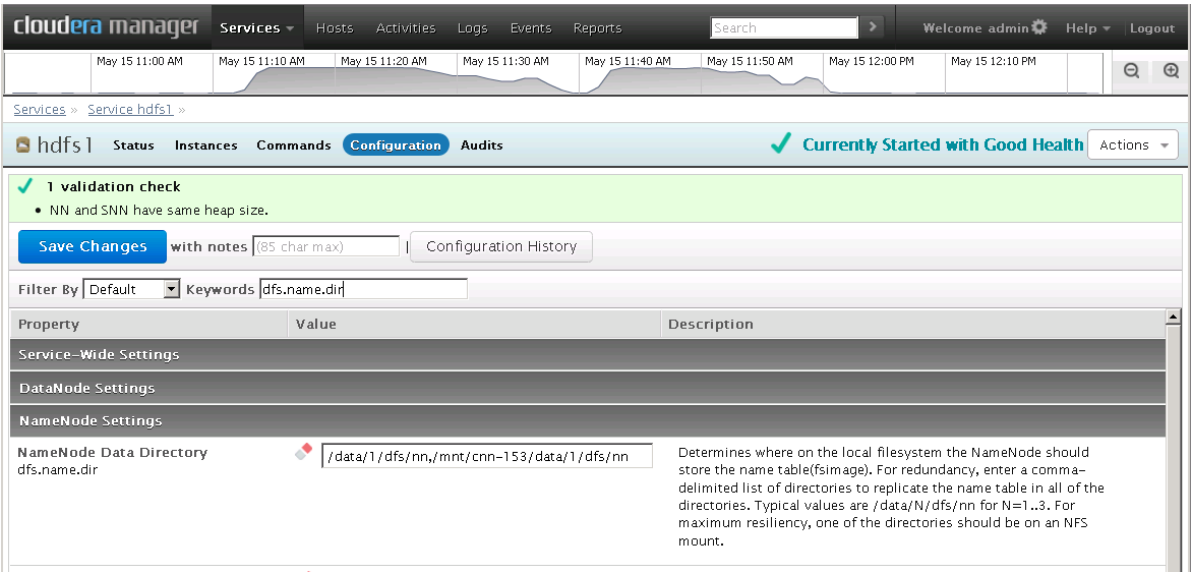
```
tcp,soft,intr,timeo=10,retrans=10
```

These options configure a soft mount over TCP; transactions will be retried ten times (retrans=10) at 1-second intervals (timeo=10) before being deemed to have failed. For example:

```
mount -t nfs -o tcp,soft,intr,timeo=10,retrans=10, <server>:<export>
<mount_point>
```

where `<server>` is the remote host, `<export>` is the exported file system, and `<mount_point>` is the local mount point. For more information, see the man pages for mount and nfs.

The following screenshot is from a cluster configuration where /mnt/cnn-153/data/1 is the local mount point of the exported filesystem on the JobTracker server.

Figure 8. Cloudera Manager Services Configuration



## Management nodes

Three ProLiant DL360p Gen8 management nodes are specified:

- The Management Node
- The JobTracker Node
- The NameNode

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

## Worker nodes

As specified in this design, eighteen ProLiant DL380p Gen8 worker nodes will fully populate a rack.

> One can have as few nodes as a single worker node, however starting with at least three worker nodes is recommended to provide the redundancy that comes with the default replication factor of 3. Performance improves with additional worker nodes as the JobTracker can leverage idle nodes to land jobs on servers that have the appropriate blocks, leveraging data locality rather than pulling data across the network. These servers are homogenous and run the DataNode and the TaskTracker (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).
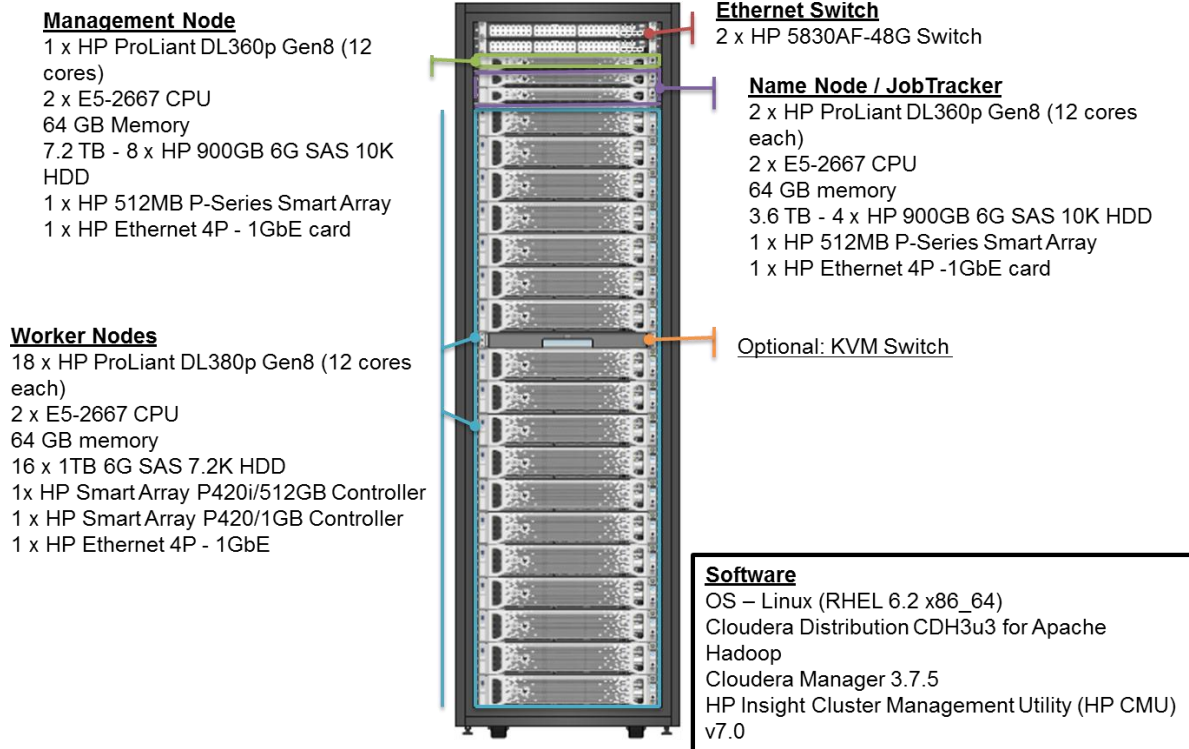
For each server, we recommend that each power supply is connected to a different PDU than the other power supply on the same server. Furthermore, the PDUs in the rack can each be connected to a separate data center power line to protect the infrastructure from a data center power line failure.

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 9. Single Rack Reference Architecture – Rack Level View



**Management Node**
1 x HP ProLiant DL360p Gen8 (12 cores)
2 x E5-2667 CPU
64 GB Memory
7.2 TB - 8 x HP 900GB 6G SAS 10K HDD
1 x HP 512MB P-Series Smart Array
1 x HP Ethernet 4P - 1GbE card

**Worker Nodes**
18 x HP ProLiant DL380p Gen8 (12 cores each)
2 x E5-2667 CPU
64 GB memory
16 x 1TB 6G SAS 7.2K HDD
1x HP Smart Array P420i/512GB Controller
1 x HP Smart Array P420/1GB Controller
1 x HP Ethernet 4P - 1GbE

**Ethernet Switch**
2 x HP 5830AF-48G Switch

**Name Node / JobTracker**
2 x HP ProLiant DL360p Gen8 (12 cores each)
2 x E5-2667 CPU
64 GB memory
3.6 TB - 4 x HP 900GB 6G SAS 10K HDD
1 x HP 512MB P-Series Smart Array
1 x HP Ethernet 4P -1GbE card

Optional: KVM Switch

**Software**
OS – Linux (RHEL 6.2 x86_64)
Cloudera Distribution CDH3u3 for Apache Hadoop
Cloudera Manager 3.7.5
HP Insight Cluster Management Utility (HP CMU) v7.0

Figure 10. Single Rack Reference Architecture – Software Distribution

| Rack | | Software Distribution |
|---|---|---|
| HP 5830 1u | | 2 x TOR Switches |
| HP 5830 1u | | 1 x Management Node |
| DL360p Gen8 1u | | 1 x JobTracker Node |
| DL360p Gen8 1u | | 1 x NameNode |
| DL360p Gen8 1u | | 18 x Worker Nodes |
| DL380p Gen8 2u | | Open for KVM Switch |
| DL380p Gen8 2u | | |
| DL380p Gen8 2u | | |
| DL380p Gen8 2u | | |
| DL380p Gen8 2u | | |
| Open 1u | | |

## Multi-Rack Reference Architecture

The Multi-Rack design assumes the Single Rack RA Cluster design is already in place and extends it allowing for a multi-rack design.

### Multi-Rack Network

As previously described in the Network section, two HP 5830AF-48G TOR switches are specified per rack for performance and redundancy. The HP 5830AF-48G includes up to four 10GbE uplinks, which can be used to connect the TOR switches into a 10GbE aggregation switch such as the HP 5920AF which can then be connected to your network of choice. Keep in mind that if IRF bonding is used, it requires up to 2 10GbE ports per switch, which would leave between 2 to 3 10GbE ports on each switch for uplinks.

### Worker nodes

Nineteen HP ProLiant DL380p worker nodes are specified. These servers are homogenous and run the DataNode and TaskTracker (or HBaseRegionServer) processes.

### Open rack space

The design leaves 2U open in the rack allowing for an additional Worker node or KVM switch.

Figure 11. Multi Rack Reference Architecture – Rack Level View



One or More Racks

**Management Node**
1 x HP ProLiant DL360p Gen8 (12 cores)
2 x E5-2667 CPU
64 GB Memory
7.2 TB - 8 x HP 900GB 6G SAS 10K HDD
1 x HP 512MB P-Series Smart Array
1 x HP Ethernet 4P - 1GbE card

**Worker Nodes**
19 x HP ProLiant DL380p Gen8 (12 cores each)

**Worker Nodes**
18 x HP ProLiant DL380p Gen8 (12 cores each)
2 x E5-2667 CPU
64 GB memory
16 x 1TB 6G SAS 7.2K HDD
1 x HP Smart Array P420/1GB Controller
1x HP Smart Array P420i/512GB Controller
1 x HP Ethernet 4P - 1GbE

**Software**
OS – Linux (RHEL 6.2 x86_64)
Cloudera Distribution CDH3u3 for Apache Hadoop
Cloudera Manager 3.7.5
HP Insight Cluster Management Utility v7.0

**Ethernet Switch**
2 x HP 5830AF-48G Switch

**Name Node / JobTracker**
2 x HP ProLiant DL360p Gen8 (12 cores each)
2 x E5-2667 CPU
64 GB memory
3.6 TB - 4 x HP 900GB 6G SAS 10K HDD
1 x HP 512MB P-Series Smart Array
1 x HP Ethernet 4P -1GbE card
Optional: KVM Switch

Figure 12. Multi-Rack Reference Architecture (extension of the single rack reference architecture)



HP 5920AF

Aggregation Switch

Single Rack RA

1 or More 42u Racks

HP 5830 1u
HP 5830 1u
DL380p Gen8 2u
DL380p Gen8 2u
DL380p Gen8 2u
DL380p Gen8 2u
DL380p Gen8 2u
DL380p Gen8 2u
DL380p Gen8 2u
DL380p Gen8 2u
DL380p Gen8 2u
⋮
DL380p Gen8 2u
Open 2u

2 x TOR Switches
19 x DL380p Worker Nodes
Open for KVM Switch

# Vertica and Hadoop

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

Each can be used standalone – Vertica for high-speed loads and ad-hoc queries over relational data, Cloudera's 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.

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/

# Summary

HP and Cloudera 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 configurations for deploying clusters of varying sizes with Cloudera Enterprise on HP infrastructure and management software. These configurations leverage HP's 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 Cloudera Enterprise software on HP infrastructure for clusters of various sizes.

**For more information**

Cloudera, cloudera.com

Hadoop on HP, hp.com/go/hadoop

Hadoop and Vertica, vertica.com/the-analytics-platform/native-bi-etl-and-hadoop-mapreduce-integration

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

HP 5830 Switch Series, hp.com/hpinfo/newsroom/press_kits/2011/HPatVMworld2011/Datasheet_A5830.pdf

HP ProLiant servers, hp.com/go/proliant

HP Enterprise Software, hp.com/go/software

HP Networking, hp.com/go/networking

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

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

HP Services, hp.com/go/services

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

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


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

**cloudera**

---

**Get connected**

hp.com/go/getconnected

Current HP driver, support, and security alerts
delivered directly to your desktop