

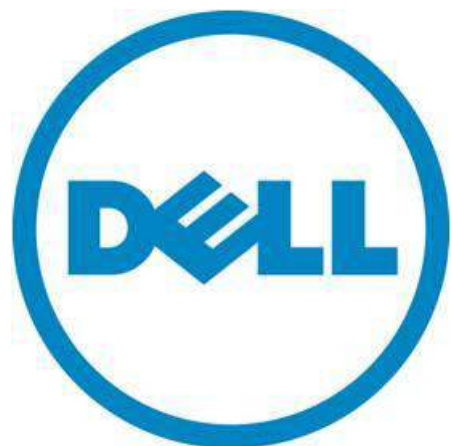
# Dell Reference Configuration for Hortonworks Data Platform

A Quick Reference Configuration Guide

Armando Acosta  
Hadoop Product Manager  
Dell Revolutionary Cloud and Big Data Group

Kris Applegate  
Solution Architect  
Dell Solution Centers

Rob Wilbert  
Solution Architect  
Dell Solution Centers



## Executive Summary

This document details the configuration set-up for Hortonworks Data Platform (HDP) software on the PowerEdge R720XD. The intended audiences for this document are customers and system architects looking for information on configuring Apache Hadoop clusters within their information technology environment for big data analytics.

The reference configuration introduces the server set-up that can run the Hortonworks stack. The document will only focus on configuration; it will not go into detail about Hadoop solution components or resiliency, performance, or software considerations. This document does not focus on best practices or complete architecture for a Hortonworks Data Platform Solution.

Dell developed this document to help streamline configuration for the Hortonworks Data Platform software.

THIS WHITE PAPER IS FOR INFORMATIONAL PURPOSES ONLY, AND MAY CONTAIN TYPOGRAPHICAL ERRORS AND TECHNICAL INACCURACIES. THE CONTENT IS PROVIDED AS IS, WITHOUT EXPRESS OR IMPLIED WARRANTIES OF ANY KIND.

© 2013 Dell Inc. All rights reserved. Reproduction of this material in any manner whatsoever without the express written permission of Dell Inc. is strictly forbidden. For more information, contact Dell. Dell, the DELL logo, and the DELL badge are trademarks of Dell Inc. Intel and Xeon are registered trademarks of Intel Corp. Red Hat is a registered trademark of Red Hat Inc. Linux is a registered trademark of Linus Torvalds. Other trademarks and trade names may be used in this document to refer to either the entities claiming the marks and names or their products. Dell Inc. disclaims any proprietary interest in trademarks and trade names other than its own.



# Reference Configuration

Hortonworks Data Platform is available on both Linux and, in partnership with Microsoft, on Windows. This initial configuration will target deployment on bare-metal servers running RedHat Linux 6.x.

## Server Roles

**Name Node(s)**<sup>1</sup> – Name nodes serve as control nodes for the HDFS, MapReduce, and HBase processes. For HDFS, name nodes own the block map and directory tree for all the data on the cluster. With MapReduce, the name node owns the job tracking daemon (JobTracker) that handles job execution and monitoring. Lastly, with HBase, name nodes are responsible for running the monitoring processes as well as owning any metadata operations. In addition to a primary name node, a secondary name node is strongly recommended for any deployment beyond a proof-of-concept.

**Data Node(s)** – Data nodes are the nodes that hold the data as well as execute MapReduce jobs. Data nodes are generally filled with large amounts of local disk, enabling the parallel processing and distributed storage features of Hadoop. The number of data nodes is dictated by use case. Adding additional data nodes increases both performance and capacity simultaneously. Maintaining a 1:1 ratio of CPU cores to disk spindles can be important in many high I/O workloads.

**Edge Node(s)** – Edge nodes lie on the perimeter of the dedicated Hadoop network and bridge the Hadoop environment with the production IT environment. Edge nodes enable external users and business processes to interact with the cluster. Additional edge nodes may be added to the Hadoop cluster as external access requirements increase.

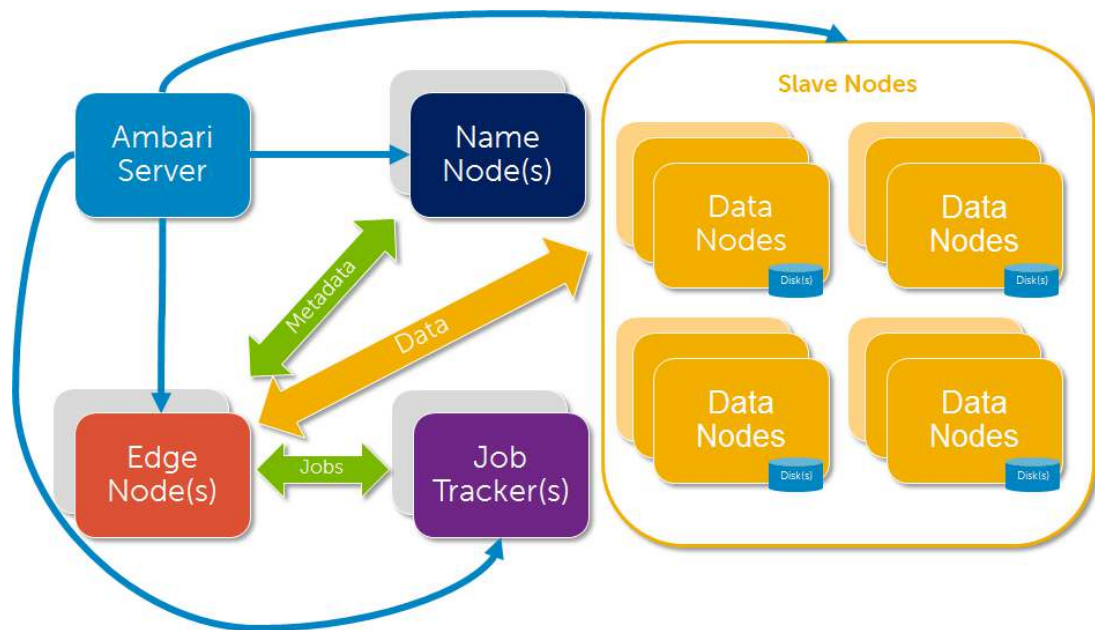
**Ambari Manager Node** – The Ambari management node is where the Ambari server resides. The Ambari management node runs the configuration management processes, web server software, monitoring software (open-source project Nagios) and performance monitoring (open-source project Ganglia) software. In a production environment, the Ambari server should run on a dedicated node; however, for the purposes of this document, Ambari server was installed on the edge node.

---

<sup>1</sup> In Hortonworks terminology the Name Node can be referred to as the Master Node



Figure 1. Dell Big Data Cluster Logical Diagram



## Node Count Recommendations

Dell recognizes that use-cases for Hadoop range from small development clusters all the way through large multi petabyte production installations. Dell has a Professional Services team that sizes Hadoop clusters for a customer's particular use. As a starting point, three cluster configurations can be defined for typical use:

**Minimum Development Cluster** – The minimum development cluster is targeted at functional testing and may even be built from existing equipment; however, the performance of these types of clusters can be significantly less as development clusters typically do not benefit from the highly distributed nature of HDFS.

**Recommended Small Cluster** – The recommended small cluster is a good starting point for customers taking the initial steps for running HDP in production. A small cluster provides some layers of basic resiliency that is expected in today's production IT world.

**Recommended Production Cluster** – The recommended production cluster configuration provides dense storage and compute capacity, coupled with high degree of resiliency. The production cluster allows for an adequate number of data nodes to demonstrate the performance benefits of distributed storage and parallel computing.

Table 1. Recommended Cluster Sizes

	Minimum Development Cluster <sup>3</sup>	Recommended SmallCluster <sup>3</sup>	Recommended Production Cluster
Name Node(s)	1 <sup>1</sup>	1 <sup>2</sup>	2
Job Tracker(s)	0 <sup>1</sup>	0 <sup>2</sup>	1
Edge Node(s)	0 <sup>1</sup>	1 <sup>2</sup>	1
Data Node(s)	3	6	14
Ambari Management Node	0 <sup>1</sup>	0 <sup>2</sup>	1
1 GbE Switches	1	1	2
10 GbE Switches	0	2	2
Rack Units	9U	19U	42U

<sup>1</sup> In this case a single node serves as the name node, job tracker, edge node and Ambari management node.  
<sup>2</sup> In this case the Ambari management node, job tracker, and edge node roles are combined.  
<sup>3</sup> Configurations include high availability and resiliency which is recommended for production clusters, proof of concepts and small cluster can exclude high availability and resiliency

Figure 2. Reference Configuration Diagram

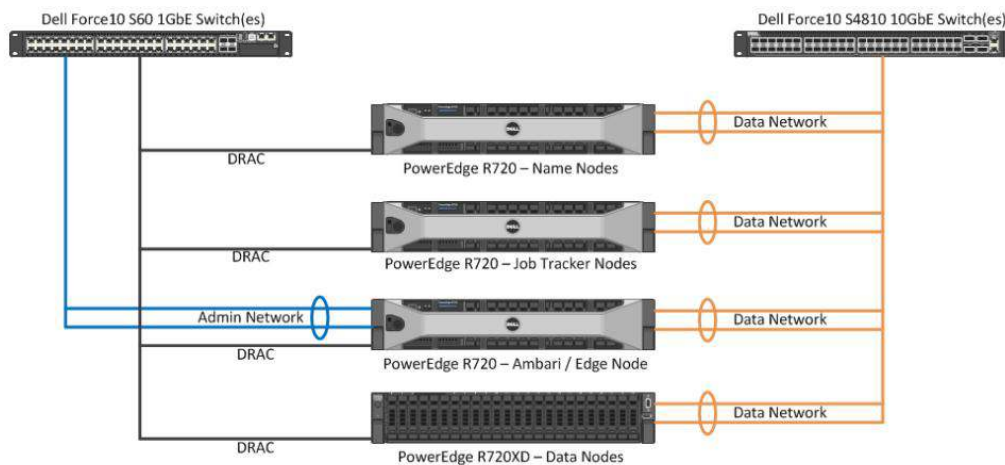
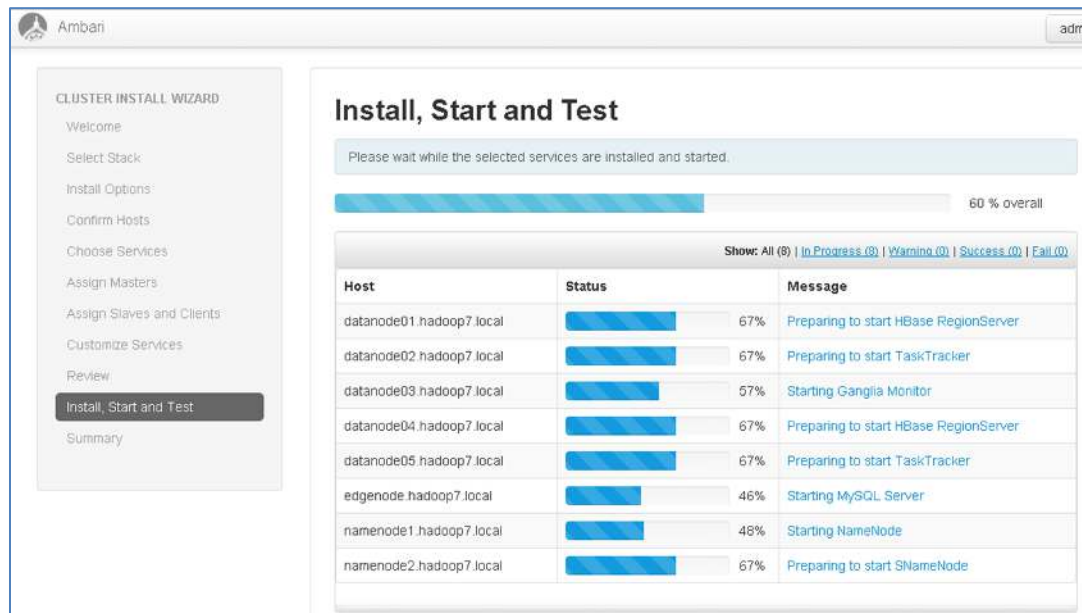


Figure 3. Ambari Manager - Node Installation



## Tested Configuration

For the purposes of this document, a small Hadoop cluster was deployed as recommended in **Table 1**. The specific software revisions used in the test are shown in **Table 2**. The PowerEdge R720 and R70XD hardware configurations we tested are shown in **Table 3** and **Table 4**. The hardware listed should be used as initial guidance only. Additional configurations are possible and will likely be required as each customer's environment and use-case is unique. Common parameters that could differ include:

1. **Processors** – Higher frequencies and core counts may improve performance while lower voltage/TDP processors, such as the Intel Xeon E5-2630L processor, can improve power efficiency
2. **Local Storage** – Disk capacity, drive technology, and spindle speed can be matched to budget and performance requirements as necessary
3. **Memory** – Depending on the usage of various services (Hbase versus Map Reduce) more or less memory may be necessary on both the infrastructure and data nodes

## Teragen / Terasort

These two HDFS / MapReduce benchmarks are used in conjunction with each other to stress Hadoop systems and provide valuable metrics with regards to network, disk and CPU utilization. By starting with these benchmarks as a baseline, Hadoop administrators can tune Hadoop's wide variety of parameters to achieve the desired performance. Teragen starts by generating flat text files that contain pseudo-random data that Terasort then sorts. This type of sort / shuffle exercise simulates customer workloads as they manipulate data through MapReduce jobs.

Figure 4. Ambari Manager Monitoring

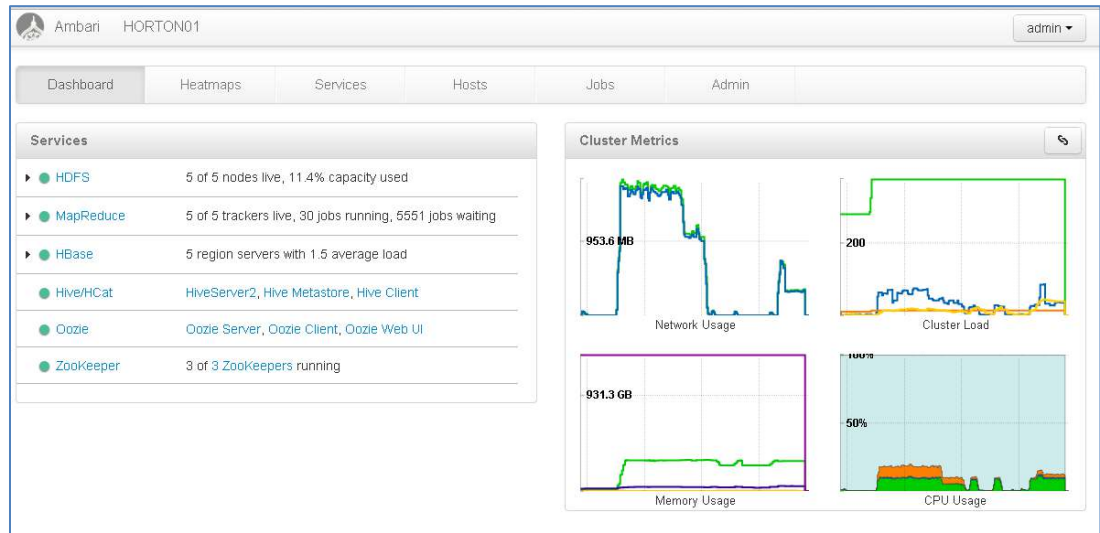


Table 2. Software Revisions (As Tested)

Component	Revision
Redhat Enterprise Linux	6.4
Hortonworks Data Platform Hadoop	1.3
Hadoop	1.2.0

Table 3. PowerEdge R720 Infrastructure Node Configuration (As Tested)

Component	Specification
Height	2 Rack Units (3.5")
Processor	2x Intel Xeon E5-2650 2 GHz 8-core processors
Memory	128 GB
Disk	6x 600 GB 15K SAS Drives
Network	4x 1GbE Intel LOMs, 2x 10GbE Intel NICs
RAID Controller	PowerEdge RAID Controller H710 (PERC)
Management Card	Integrated Dell Remote Access Controller (iDRAC)

Table 4. PowerEdge R720XD Data Node Configuration (As Tested)

Component	Specification
Height	2 Rack Units (3.5")
Processor	2x Intel Xeon E5-2667 2.9 GHz 6-core processors
Memory	64 GB
Disk	24x 500GB or 1TB 7200 RPM Nearline SAS drives
Network	4x 1GbE Intel LOMs, 2x 10GbE Intel NICs
RAID Controller	PowerEdge RAID Controller H710 (PERC)
Management Card	Integrated Dell Remote Access Controller (iDRAC)



## Dell Solution Centers

The Dell Solution Centers are a global network of connected labs that allow Dell to help customers architect, validate and build solutions. With multiple footprints in every region, they help customers understand anything from simple hardware platforms, to more complex solutions. These engagements range from an informal 30-60 minute briefing, through a longer half-day workshop, and on to a proof-of-concept that allow customers to kick the tires of their solution prior to signing on the dotted line. Customers may engage with their account team and have them submit a request to take advantage of these free services.

## Links

Hortonworks – <http://hortonworks.com>

Hortonworks Data Platform - <http://hortonworks.com/products/hdp/>

Hortonworks Sandbox - <http://hortonworks.com/products/hortonworks-sandbox/>

