

Dell | Cloudera Apache Hadoop Solution Crowbar Administration User Guide

A Dell User Guide for Apache Hadoop Deployment

Revision 2.0



Table of Contents

Tables and Figures 4

Notes, Cautions, and Warnings 5

Abbreviations 5

Overview 6

Introduction 6

Document Scope 7

 Opscode Chef Server..... 7

 Dell Specific Options..... 8

Dell | Cloudera Apache Hadoop Solution..... 8

 Network Setup 8

 Managing Growth..... 8

 Rack..... 8

 Pod..... 8

 Cluster..... 8

 Default Networks..... 9

 Layout..... 9

 IP Addressing 10

 Rack Awareness..... 10

 Hadoop Basics..... 11

 Apache Hadoop Component Deployment..... 12

Crowbar User Interface 13

Cloudera Manager Overview..... 13

 Functionality Outline..... 13

Barclamps 14

 Cloudera Manager Barclamp..... 15

 Cloudera Manager Installation Overview 15

 Cloudera Manager Node Inventory Page..... 15

 Cloudera Manager Administration Console..... 16

 Login Screen..... 18

 License Key Upload Screen..... 19

 License Key Restart Screen 20

 License Key Confirmation Screen..... 21

 Cloudera Registration Screen..... 22

 Node Search Screen..... 23

 Node Search Results Screen..... 24

 Repository Configuration Screen..... 25

 SSH Credentials Screen 26

 Package Install Screen 27

 Host Inspector Screen 28

 Service Selection Screen 29

 Inspect Role Assignments Screen # 1..... 30

 Inspect Role Assignments Screen # 2..... 31

 Monitoring Database Setup Screen 32

Dell | Cloudera Apache Hadoop Solution

- Review Configuration Changes Screen33
- Cluster Services Initialization Screen..... 34
- Configuration Completion Screen35
- Service Display Screen 36
- Pig Barclamp.....37
- Hive Barclamp38
- Sqoop Barclamp 40
- Support.....41**
 - Cloudera Support 41
- Appendix A: Dell | Hadoop Solution Components41**
- Appendix B: External References41**
- To Learn More.....41**

Tables and Figures

TABLE 1-1: DEFAULT NETWORKS.....	9
TABLE 1-2: MASTER/SECONDARY (ADMIN) NAME NODES NETWORK CONNECTIONS.....	9
TABLE 1-3: EDGE NODES NETWORK CONNECTIONS.....	9
TABLE 1-4: SLAVE NODES NETWORK CONNECTIONS.....	9
TABLE 1-5: IP ADDRESSING SCHEMA.....	10
TABLE 1-6: POD 1 IP EXAMPLE ADDRESSING LAYOUT.....	10
TABLE 1-7: POD 2 IP EXAMPLE ADDRESSING LAYOUT.....	10
TABLE 2 SUPPORTED APACHE HADOOP COMPONENTS.....	12
TABLE 3 USER INTERFACE SERVICE URLS.....	13
TABLE 4 DIFFERENCES BETWEEN CLOUDERA MANAGER FREE EDITION AND CLOUDERA MANAGER.....	13
TABLE 5 BARCLAMP DESCRIPTIONS.....	14
TABLE 6 OPEN FILE HANDLES CONFIGURATION PARAMETERS.....	15
TABLE 7 HADOOP HIGH AVAILABILITY CONFIGURATION PARAMETERS.....	15
FIGURE 4-1: CLOUDERA NODE INVENTORY PAGE.....	16
FIGURE 1 LOGIN PAGE.....	18
FIGURE 2 LICENSE KEY UPLOAD SCREEN.....	19
FIGURE 3 LICENSE KEY RESTART SCREEN.....	20
FIGURE 4 LICENSE KEY CONFIRMATION SCREEN.....	21
FIGURE 5 REGISTRATION SCREEN.....	22
FIGURE 6 CLOUDERA CLUSTER NODE SEARCH SCREEN.....	23
FIGURE 7 NODE SEARCH RESULTS SCREEN.....	24
FIGURE 8 SSH CREDENTIALS SCREEN.....	26
FIGURE 9 PACKAGE INSTALL SCREEN.....	27
FIGURE 10 SERVICE SELECTION SCREEN.....	29
FIGURE 11 INSPECT ROLE ASSIGNMENTS SCREEN # 1.....	30
FIGURE 12 INSPECT ROLE ASSIGNMENTS SCREEN #2.....	31
FIGURE 13 MONITORING DATABASE SETUP SCREEN.....	32
FIGURE 14 REVIEW CONFIGURATION CHANGES SCREEN.....	33
FIGURE 15 CLUSTER SERVICES INITIALIZATION SCREEN.....	34
FIGURE 16 CONFIGURATION COMPLETION SCREEN.....	35
FIGURE 17 SERVICE DISPLAY SCREEN.....	36
TABLE 4-34: PIG BARCLAMP PARAMETERS.....	37
TABLE 4-41: HIVE BARCLAMP PARAMETERS.....	38
TABLE 4-48: SQOOP BARCLAMP PARAMETERS.....	40

Dell | Cloudera Apache Hadoop Solution

THIS 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.

Information in this document is subject to change without notice.

© 2012 Dell Inc. All rights reserved.

Reproduction of these materials is allowed under the Apache 2 license.

Information in this document is subject to change without notice.

© 2012 Dell Inc. All rights reserved.

Dell, the DELL logo, and the DELL badge, PowerConnect, and PowerEdge are trademarks of Dell Inc. Cloudera, CDH, Cloudera Enterprise are trademarks of Cloudera and its affiliates in the US and other countries. Intel and Xeon are registered trademarks of Intel Corporation in the U.S. and other countries. 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.

Other Trademarks used in this text: Dell™, the DELL logo, Cloudera™, Nagios™, Ganglia™, Opscode Chef™, OpenStack™, Canonical Ubuntu™, VmWare™, Dell Precision™, OptiPlex™, Latitude™, PowerEdge™, PowerVault™, PowerConnect™, OpenManage™, EqualLogic™, KACE™, FlexAddress™ and Vostro™ are trademarks of Dell Inc. Intel®, Pentium®, Xeon®, Core™ and Celeron® are registered trademarks of Intel Corporation in the U.S. and other countries. AMD® is a registered trademark and AMD Opteron™, AMD Phenom™, and AMD Sempron™ are trademarks of Advanced Micro Devices, Inc. Microsoft®, Windows®, Windows Server®, MS-DOS® and Windows Vista® are either trademarks or registered trademarks of Microsoft Corporation in the United States and/or other countries. Red Hat Enterprise Linux® and Enterprise Linux® are registered trademarks of Red Hat, Inc. in the United States and/or other countries. Novell® is a registered trademark and SUSE™ is a trademark of Novell Inc. in the United States and other countries. Oracle® is a registered trademark of Oracle Corporation and/or its affiliates. Citrix®, Xen®, XenServer® and XenMotion® are either registered trademarks or trademarks of Citrix Systems, Inc. in the United States and/or other countries. VMware®, Virtual SMP®, vMotion®, vCenter®, and vSphere® are registered trademarks or trademarks of VMWare, Inc. in the United States or other countries.

Other trademarks and trade names may be used in this publication 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.

April 2012

Notes, Cautions, and Warnings



Note: A NOTE indicates important information that helps you make better use of your computer.



CAUTION: A CAUTION indicates potential damage to hardware or loss of data if instructions are not followed.



WARNING: A WARNING indicates a potential for property damage, personal injury, or death.

Abbreviations

Abbreviation	Definition
BMC	Baseboard management controller.
DBMS	Database management system.
EDW	Enterprise data warehouse.
EoR	End-of-row switch/router.
HDFS	Hadoop Distributed File System.
IPMI	Intelligent Platform Management Interface.
LAG	Link aggregation group.
LOM	Local Area Network on Motherboard.
NIC	Network interface card.
ToR	Top-of-rack switch/router.

Overview

Hadoop is an Apache project being built and used by a global community of contributors, written in the Java programming language. Yahoo! has been the largest contributor to the project, and uses Hadoop extensively across its businesses. Other contributors and users include Facebook, LinkedIn, eHarmony, and eBay. Cloudera has created a quality controlled distribution of Hadoop and offers commercial management software, support, and consulting services.

Dell developed a solution for Hadoop that includes optimized hardware, software, and services to streamline deployment and improve the customer experience.

Introduction

This document provides instructions you to use when deploying Cloudera Manager and Apache Hadoop Eco-System components with Crowbar. This guide is for use with the Crowbar Users Guide, and is *not* a stand-alone document. It specifically covers Cloudera Manager, Apache Hadoop and the deployment steps from a Crowbar prospective. Please refer to the Crowbar User Guide for assistance with installing common Crowbar components and configuring the target systems.

 **Note:** Concepts beyond the scope of this guide are introduced as needed in notes and references to other documentation.

The Dell | Cloudera Apache Hadoop Solution is based on the Cloudera CDH Enterprise distribution of Hadoop. Dell's solution includes:

- Dell Reference architecture (RA) and best practices documentation.
- Optimized hardware and network infrastructure.
- Cloudera CDH software (CDH Community-provided for customer-deployed solutions).
- Cloudera Manager free edition with the ability to upgrade to enterprise level via Cloudera issued license key.
- Cloudera Manager provided Hadoop infrastructure management tools.
- Dell Crowbar software framework.

This solution provides Dell a foundation to offer additional solutions as the Hadoop environment evolves and expands.

Document Scope

The focus of this guide is the use of Crowbar, *not* Cloudera Manager or Apache Hadoop. While Crowbar includes substantial components to assist in the deployment of Cloudera Manager and Apache Hadoop, its operational aspects are completely independent. For more detailed information, please refer to the links below;

Cloudera Manager Free Edition Documentation


- [Cloudera Manager Free Edition 4.0.x Release Notes](#)
- [Cloudera Manager Free Edition Frequently Asked Questions](#)
- [Cloudera Manager Free Edition Installation Guide](#)
- [Cloudera Manager Free Edition User Guide](#)
- [Configuring Ports for Cloudera Manager Free Edition](#)
- [Configuring TLS Security for Cloudera Manager Free Edition](#)

Cloudera Manager Enterprise Edition Documentation

- [Cloudera Manager 4.0.x Release Notes](#)
- [Cloudera Manager Frequently Asked Questions](#)
- [Cloudera Manager Installation Guide](#)
- [Cloudera Manager User Guide](#)
- [Configuring Hadoop Security with Cloudera Manager](#)
- [Configuring Ports for Cloudera Manager](#)
- [Configuring TLS Security for Cloudera Manager](#)

Other related documentation

- [CDH4 Documentation](#)
- [High Availability for the Hadoop Distributed File System - HDFS](#)
- [CDH4 High Availability Overview](#)
- [CDH4 High Availability Guide](#)
- [Hadoop at Apache.org](#)
- [Hue 2.0 User Guide](#)

 This guide provides this additional information about Cloudera as notes flagged with the Cloudera logo. For detailed operational support for Hadoop, we suggest visiting the Cloudera documentation web site at <http://www.cloudera.com>.

Opscode Chef Server

Crowbar makes extensive use of Opscode Chef Server, <http://opscode.com>. To explain Crowbar actions, you should understand the underlying Chef implementation. This guide provides this additional Chef information as notes flagged with the Opscode logo.



To use Crowbar, it is not necessary to log into the Chef Server; consequently, use of the Chef UI is not covered in this guide. Supplemental information about Chef is included.

Dell Specific Options

The Dell EULA version of Crowbar provides additional functionality and color pallets than the open source version. When divergences are relevant, they are identified.



To perform some configuration options and provide some integration, we use libraries that cannot be distributed using open source.

Crowbar is not limited to managing Dell servers and components. Due to driver requirements, some barclamps, for example: BIOS and RAID must be targeted to specific hardware; however, those barclamps are not required for system configuration.

Dell | Cloudera Apache Hadoop Solution

Network Setup

The network configuration assumes a flat L2 wiring. All network connections should be accessible at that layer. Where isolation between different logical networks is required, VLANs are used.

Managing Growth

The system architecture is organized into three components, for sizing as the Hadoop environment grows. From smallest to largest, they are:

- Rack
- Pod
- Cluster

Each has specific characteristics and sizing considerations. You can scale the environment by adding additional capacity as needed, without the need to replace any existing components.

Rack

A rack is the smallest component in a Hadoop environment, and consists of all of the power, network cabling, and two Ethernet switches required to support up to 20 data nodes. These nodes should utilize their own power connectivity and data center space – separate from other racks – and be treated as a fault zone.

Pod

A pod is a single set of stacked Ethernet switches. For the Dell | Cloudera Reference Architecture, both the maximum and minimum are six. A pod consists of the administration and operation infrastructure to support three racks.

Cluster

A cluster is a set of greater than one pod, up to a maximum of 12 pods. A cluster is a set of Hadoop nodes that share the same Network Node and management tools for operating the Hadoop environment.




Note: Please see the Dell | Cloudera Solution Reference Architecture Guide for more detailed information.

Default Networks

The default networks are presented in the following table.

Table 1-1: Default Networks

Usage	Description	Default reserved vLAN tag	Tagged
Admin/Internal vLAN	Used for administrative functions such as Crowbar node installation, TFTP booting, DHCP assignments, KVM, system logs, backups, and other monitoring. There is only one vLAN set up for this function and it is spanned across the entire network.	100	Not tagged
BMC vLAN	Used for connecting to the BMC of each node.	100	Not tagged
Storage vLAN	Used by the Swift storage system for replication of data between machines, monitoring of data integrity, and other storage specific functions (802.1q Tagged).	200	Tagged
Edge vLANs	Used for connections to devices external to the Hadoop cluster; these include externally visible services such as load balancers and web servers. Use one or many of these networks, dependent on the need to segregate traffic among groups of servers (802.1q Tagged).	300	Tagged

 **Note:** The admin and BMC networks are expected to be in the same L2 network.

Layout

Due to the nature of Crowbar's network layout, addresses are assigned to a whole network based upon interface, Network Type (Production, Management, and External) and teaming type.

Table 1-2: Master/Secondary (Admin) Name Nodes Network Connections

Interface	Network Type	Teaming Type
BMC	Management LAN	Single
LOM1	Production LAN	Teamed
LOM2	Production LAN	Teamed
Eth1	Production LAN	Teamed
Eth2	Management LAN	Single

Table 1-3: Edge Nodes Network Connections

Interface	Network Type	Teaming Type
BMC	Management LAN	Single
LOM1	Production LAN	Teamed 1
LOM2	Production LAN	Teamed 1
Eth1	External LAN	Teamed 2
Eth2	External LAN	Teamed 2

Table 1-4: Slave Nodes Network Connections

Interface	Network Type	Teaming Type
BMC	Management LAN	Single
LOM1	Production LAN	Teamed 1
LOM2	Production LAN	Teamed 1

IP Addressing

The IP address can be assigned in this fashion, using large subnets to support many machines on the production network. The management network is a Class C network with 254 IP addresses. The Production network is what is known as a /23 with 512 IP addresses. In each network, the first 10 IP addresses are reserved for switches, routers, and firewalls.


 **Note:** Each network's ".1" address is reserved for the network gateway.

Table 1-5: IP Addressing Schema

LAN	Network	Subnet	Gateway	Reserved
Management LAN	172.16.0.0	255.255.255.0	172.16.0.1	0.1 – 0.10
Production LAN	172.16.2.0	255.255.254.0	172.16.2.1	2.1-2.20
Name Nodes	DHCP Allocated			
Slave Nodes	DHCP Allocated			
External LAN	TBD by Customer			

Rack Awareness

With the network set up using Top of Rack (ToR) switches, Rack Awareness can be programmed using the Chef information about which switch the LOM1 is plugged into. A simple script has been added to the Hadoop configuration to pull the information out of Chef, and then use it for Rack Awareness.

Table 1-6: Pod 1 IP Example Addressing Layout

Network: 172.16.0.0	Netmask: 255.255.252.0
Multicast: 172.16.0.0	Broadcast 172.16.3.255

Pod	Rack Number	Network	Server Type	IP Range	Subnet Mask	Gateway
1	1	Production	Slave	172.16.0.1-42	255.255.252.0	172.16.0.1
1	2	Production	Slave	172.16.1.1-42	255.255.252.0	172.16.0.1
1	3	Production	Slave	172.16.2.1-42	255.255.252.0	172.16.0.1
1		Production	Master Name	172.16.3.1-19	255.255.252.0	172.16.0.1
1		Production	Secondary Name	172.16.3.20-30	255.255.252.0	172.16.0.1
1		Production	Edge	172.16.3.41-50	255.255.252.0	172.16.0.1
1	1	BMC	Slave	172.16.0.200-242	255.255.252.0	172.16.0.1
1	2	BMC	Slave	172.16.1.200-242	255.255.252.0	172.16.0.1
1	3	BMC	Slave	172.16.2.200-242	255.255.252.0	172.16.0.1
1		BMC	Master Name	172.16.3.201-219	255.255.252.0	172.16.0.1
1		BMC	Secondary Name	172.16.3.220-230	255.255.252.0	172.16.0.1
1		BMC	Edge	172.16.3.231-250	255.255.252.0	172.16.0.1

Table 1-7: Pod 2 IP Example Addressing Layout

Network: 172.16.0.0	Netmask: 255.255.252.0
Multicast: 172.16.0.0	Broadcast: 172.16.3.255

Pod	Rack Number	Network	Server Type	IP Range	Subnet Mask	Gateway
2	1	Production	Slave	172.16.4.1-42	255.255.252.0	172.16.4.1
2	2	Production	Slave	172.16.5.1-42	255.255.252.0	172.16.4.1

Dell | Cloudera Apache Hadoop Solution

2	3	Production	Slave	172.16.6. 1-42	255.255.252.0	172.16.4.1
2		Production	Master Name	172.16.7.1-19	255.255.252.0	172.16.4.1
2		Production	Secondary Name	172.16.7.20-30	255.255.252.0	172.16.4.1
2		Production	Edge	172.16.7.41-50	255.255.252.0	172.16.4.1
2	1	BMC	Slave	172.16.4.200-242	255.255.252.0	172.16.4.1
2	2	BMC	Slave	172.16.5.200-242	255.255.252.0	172.16.4.1
2	3	BMC	Slave	172.16.6.200.242	255.255.252.0	172.16.4.1
2		BMC	Master Name	172.16.7.201-219	255.255.252.0	172.16.4.1
2		BMC	Secondary Name	172.16.7.220-230	255.255.252.0	172.16.4.1
2		BMC	Edge	172.16.7.231-250	255.255.252.0	172.16.4.1
2		External	Edge	TBD by Customer	TBD	TBD

Hadoop Basics

The Hadoop software library is a framework that allows for the distributed processing of large data sets across clusters of computers using a simple programmatic driven processing model. Hadoop is designed to scale up from a minimum of three servers to thousands of machines, each offering local computation and storage.

Rather than rely on hardware to deliver high-availability, the Hadoop library itself is designed to detect and handle failures at the application layer, so delivering a highly-available service on a cluster of computers, each of which may be prone to failures.

Hadoop is ideal for organizations with a growing need to store and process massive application datasets. It enables applications to work with thousands of nodes and petabytes of data. This Crowbar barclamp provides the ability to deploy and maintain Hadoop cluster Admin, Master, Slave and Edge nodes. It also provides the capability to configure and deliver Hadoop HDFS and MapReduce components.

- **Hadoop Core:** The common libraries and utilities that provide the basic Hadoop runtime environment. A set of components and interfaces which implement a distributed filesystem and provide general I/O access for the Hadoop framework (serialization, Java RPC and persistent data storage).
- **Hadoop Distributed File System (HDFS):** A distributed file system that provides redundant, high-throughput access to application data.
- **MapReduce:** A software framework for distributed processing of large data sets on compute clusters.

Apache Hadoop Component Deployment

For Hadoop (Cloudera Manager) and Eco-System components (Hive, Sqoop and Pig), employ Crowbar tools to construct a starting proposal and then edit any parameters to fit the specific needs of your environment. Once the proposal is ready, apply the proposal to deploy each system components.


 **Note:** The Base Hadoop system (HDFS and Map Reduce), Zookeeper, Hbase, Oozie and Hue are deployed using the Cloudera Manager administration console. Crowbar also provided some supplemental Hadoop Eco-System Barclamps (Hive, Sqoop and Pig) and you must install the base Hadoop system (HDFS and Map Reduce) using Cloudera Manager before deploying any of these add-ons.

Table 2 Supported Apache Hadoop Components

Component	Deployment Method	Description
Hadoop Core (HDFS/Map Reduce)	Cloudera Manager	Common libraries and utilities that provides the basic Hadoop runtime environment (HDFS/map reduce), a set of components and interfaces which implements a distributed filesystem and provides general i/o access for the Hadoop framework (serialization, java rpc and persistent data storage).
HUE	Cloudera Manager	HUE (Hadoop User Experience) is a user interface framework and SDK platform for visual Hadoop applications. It delivers a suite of web base UI applications which can be used to access and modify the Hadoop Distributed File System (HDFS) and Map Reduce job queue. HUE provides UI application portals for HDFS file browsing, Map/Reduce job control, user account administration and web based on-line help.
HBase	Cloudera Manager	HBase is an open source, non-relational, distributed database modeled after Google's BigTable and is written in Java. It is developed as part of Apache Software Foundation's Apache Hadoop project and runs on top of HDFS (Hadoop Distributed Filesystem), providing BigTable -like capabilities for Hadoop. That is, it provides a fault-tolerant way of storing large quantities of sparse data. HBase features compression, in-memory operation, and Bloom filters on a per-column basis as outlined in the original BigTable paper. Tables in HBase can serve as the input and output for MapReduce jobs run in Hadoop, and may be accessed through the Java API but also through REST, Avro or Thrift gateway APIs. HBase is not a direct replacement for a classic SQL Database, although recently its performance has improved, and it is now serving several data-driven websites, including Facebook's Messaging Platform
ZooKeeper	Cloudera Manager	High-performance coordination service for distributed applications. ZooKeeper provides primitives such as distributed locks which can be used for building large scale distributed processing applications.
Oozie	Cloudera Manager	Oozie is an open-source workflow/coordination service to manage data processing jobs for Apache Hadoop™. It is an extensible, scalable and data-aware service to orchestrate dependencies between jobs running on Hadoop (including HDFS, Pig and MapReduce).
Hive	Crowbar Barclamp	Data warehouse infrastructure that provides SQL based data summarization and ad hoc querying.
Pig	Crowbar Barclamp	Platform for analyzing large data sets that consists of a high-level language for expressing data algorithms.
Sqoop	Crowbar Barclamp	SQL based command-line tool to assist with HDFS data import/export (SQL-to-Hadoop).

For more information about Hadoop, please visit <http://hadoop.apache.org/>.

Crowbar User Interface

Crowbar is delivered as a Web application available on the admin node using HTTP on port 3000. By default, you can access it using <http://192.168.124.10:3000>. Additionally, the default installation contains an implementation of Hadoop specific components (see table below).



 **Note:** Crowbar has been tested on the following browsers: FireFox 3.5+, FireFox 4.0, Internet Explorer 7, and Safari 5. A minimum screen resolution of 1024x768 or higher is recommended.

Table 3 User Interface Service URLs

User Interface Service	Default Location	Port	Example URL
Crowbar	Crowbar Admin Node	3000	<a href="http://<crowbar_admin_node>:3000">http:// <crowbar_admin_node >:3000
Cloudera Manager	Crowbar Admin Node	7180	<a href="http://<Cloudera_web_app_node>:7180">http:// <Cloudera_web_app_node >:7180
Hadoop Name Node	Hadoop Name Node	50070	<a href="http://<master_name_node>:50070">http:// <master_name_node >:50070
Hadoop Secondary Name Node	Hadoop Secondary Name Node	50090	<a href="http://<secondary_name_node>:50090">http:// <secondary_name_node >:50090
Hadoop Data Node	Hadoop Data Node	50075	<a href="http://<data_node>:50075">http:// <data _node>: 50075
Hadoop Job Tracker Web	Hadoop Job Tracker Node	50030	<a href="http://<job_tracker_node>:50030">http:// <job_tracker_node>: 50030
Hadoop Task Tracker Web	Task Tracker Node	50060	<a href="http://<task_tracker_node>:50060">http://<task_tracker_node>:50060

 **Note:** The crowbar_admin_node IP address (192.168.124.10) is the default address. Replace it with the address assigned to the Crowbar Admin node. Nagios, Ganglia and Chef can be accessed directly from a web browser or via selecting one of the links on the Crowbar Dashboard.



Cloudera Manager Overview



































Cloudera Manager deploys and centrally operates a complete Hadoop stack. The application automates the installation process, reducing deployment time from weeks to minutes, gives you a cluster-wide, real time view of the services running and the status of their hosts, provides a single, central place to enact configuration changes across your cluster; and incorporates a full range of reporting and diagnostic tools to help you optimize cluster performance and utilization. Cloudera Manager provides full lifecycle management for Hadoop deployments.

Functionality Outline

- Installs the complete Hadoop stack in minutes via a wizard-based interface.
- Gives you complete, end-to-end visibility and control over your Hadoop cluster from a single interface.
- Let's you set server roles and configure services across the cluster.
- Let's you gracefully start, stop and restart of services as needed.
- Shows information pertaining to hosts in your cluster including status, resident memory, virtual memory and roles.

Table 4 Differences between Cloudera Manager Free Edition and Cloudera Manager

Feature	Cloudera Manager Free Edition	Cloudera Manager
Max number of hosts supported	50	Unlimited
Automated installer		

Host Monitoring		
Secure communication between server and agents (TLS)		
Service and Configuration Management		
Manage HDFS, MapReduce , HBase, Hue, Oozie, and Zookeeper		
Configuration audit trails		
Workflows to start, stop, restart, add, delete, and decommission role instances		
Configuration versioning and history		
Support for Kerberos		
Service Monitoring		
Status and health summary		
Proactive health tests		
Log Search and Management		
Events Management		
Alerts		
Activity Monitoring		
Operational Reports		
Global Time Control for historical diagnosis		

Barclamps

Table 5 Barclamp Descriptions

Barclamp	Description
Cloudera Manager	Provides end-to-end management for apache Hadoop with the ability to deploy and centrally operate a complete Hadoop stack gives you a cluster wide, real time view of nodes and services running and provides a single central place to enact configuration changes across your cluster. Cloudera Manager incorporates a full range of reporting and diagnostic tools to help you optimize cluster performance and utilization.
Hive	Data warehouse that infrastructure provides SQL based data summarization and ad hoc querying.
Pig	Platform for analyzing large data sets that consists of a high-level language for expressing data algorithms.
Sqoop	SQL based command-line tool to assist with HDFS data import/export (SQL-to-Hadoop).
ZooKeeper	High-performance coordination service for distributed applications. ZooKeeper provides primitives such as distributed locks which can be used for building large scale distributed processing applications.

Cloudera Manager Barclamp

The Cloudera Manager Barclamp performs all the low level operating system configuration setup for the Hadoop cluster and installs the Cloudera Manager server setup in order to prepare for Hadoop cluster deployment.

Table 6 Open File Handles Configuration Parameters

Name	Description	Required	Default
Map/Reduce	Maximum number of Map/Reduce open file handles.	True	32768
HDFS	Maximum number of HDFS open file handles.	True	32768
HBASE	Maximum number of HBASE open file handles.	True	32768

Table 7 Hadoop High Availability Configuration Parameters

Name	Description	Required	Default
Shared Edits Directory	Specifies the HA shared edits directory.	True	/dfs/ha
Shared Edits Export Options	Specifies the HA shared edits export options.	True	rw,async,no_root_squash,no_subtree_check
Shared Edits Mount Options	Specifies the HA shared edits mount options.	True	rsync=65536,wsync=65536,intr,soft,bg

Cloudera Manager Installation Overview

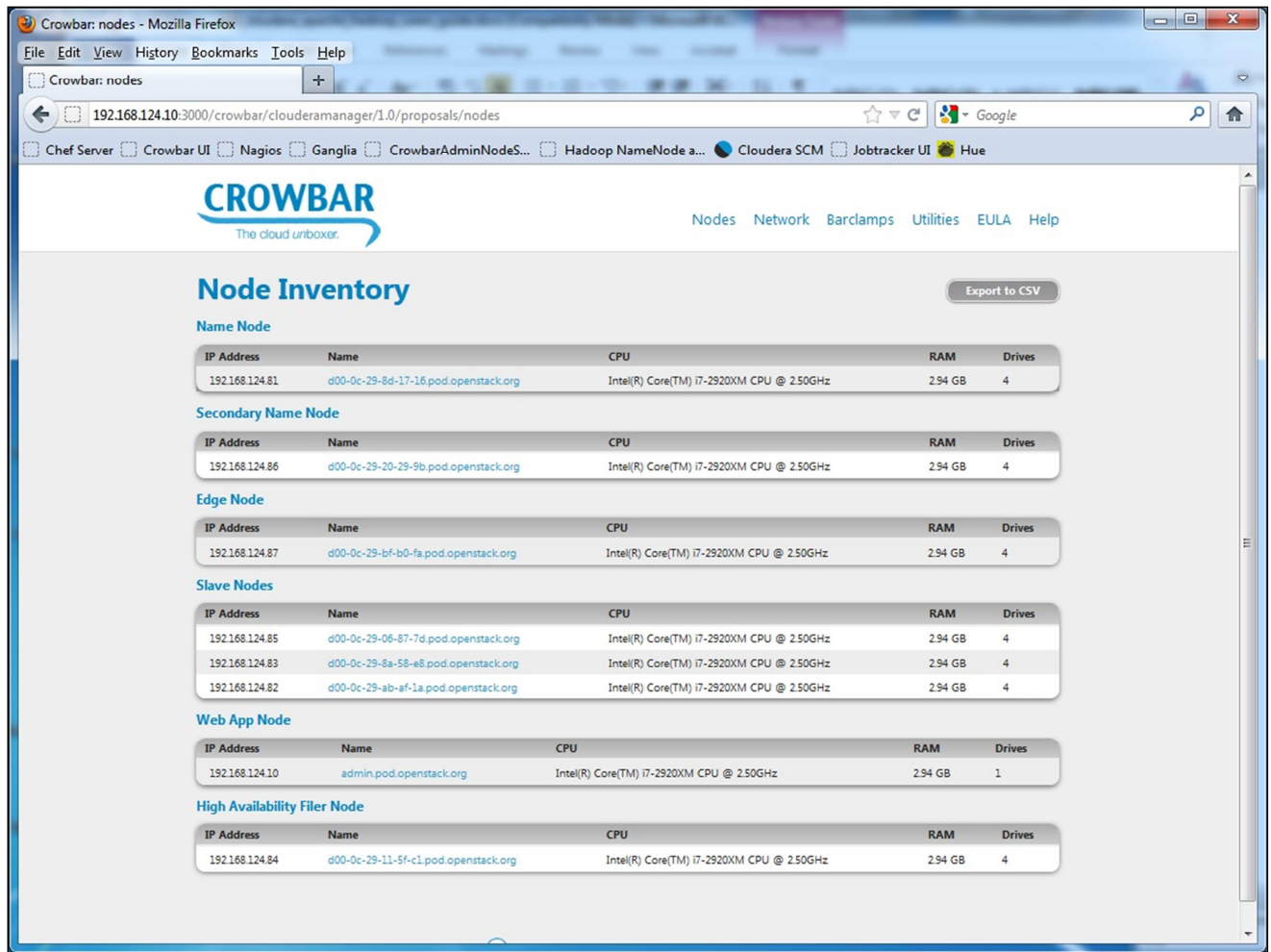
After the Cloudera Manager Barclamp has been deployed from Crowbar, you must run the Cloudera Manager configuration wizard in order to fully deploy the Hadoop cluster. This operation will perform the following tasks:

- Using SSH, discover the cluster hosts you specify via IP address ranges or hostnames.
- Install the Cloudera Manager Agent and CDH4 (including Hue) on the cluster data nodes.
- Install the Oracle JDK if it's not already installed on the cluster hosts.
- Configures the package repositories for Cloudera Manager, CDH4 and the Oracle JDK.
- Allow you to select and configure optional eco-system components.
- Determine mapping of services to host.
- Suggest a Hadoop configuration and automatically starts the Hadoop services.

You can choose to abort the Cloudera Manager Agent and CDH installation process and Cloudera Manager wizard will automatically revert and completely rollback the installation process for any uninstalled components. Installed components are not uninstalled during an abort.

Cloudera Manager Node Inventory Page

Once the Cloudera barclamp has been deployed, from the Edit Proposal page, there is a link below the Proposal Attributes section called "Cloudera Manager Nodes." Clicking on this link will display a page titled "Cloudera Node Inventory." This screen is pictured in Figure 4-1 below. You may print this page as it will be very useful during the Cloudera Manager installation to ensure the correct nodes are selected for their intended Cloudera Manager roles.



Note: You may also export this data to a comma separated value file by selecting the "Export to CSV" button at the top of the page.

Cloudera Manager Administration Console

Note: The Cloudera Manager Administration console supports the following web browsers:


- Internet Explorer 8 and 9.
- Google Chrome.
- Safari 5.
- Firefox 3.6 and later.

To start the Cloudera Manager Administration Console;

- In a web browser, type the following URL: http(s):// IP_ADDRESS: PORT_NUMBER.
 - IP_ADDRESS is the name or IP address of the host machine where the Cloudera Manager Web Server is installed.
 - PORT_NUMBER is the default port number (7180).
 - Crowbar Installation defaults are Crowbar Admin Node on port 7180 (http://192.168.124.10:7180).
- Log into the Cloudera Manager Admin Console. The default login credentials are;
 - Username: admin
 - Password: admin

Dell | Cloudera Apache Hadoop Solution

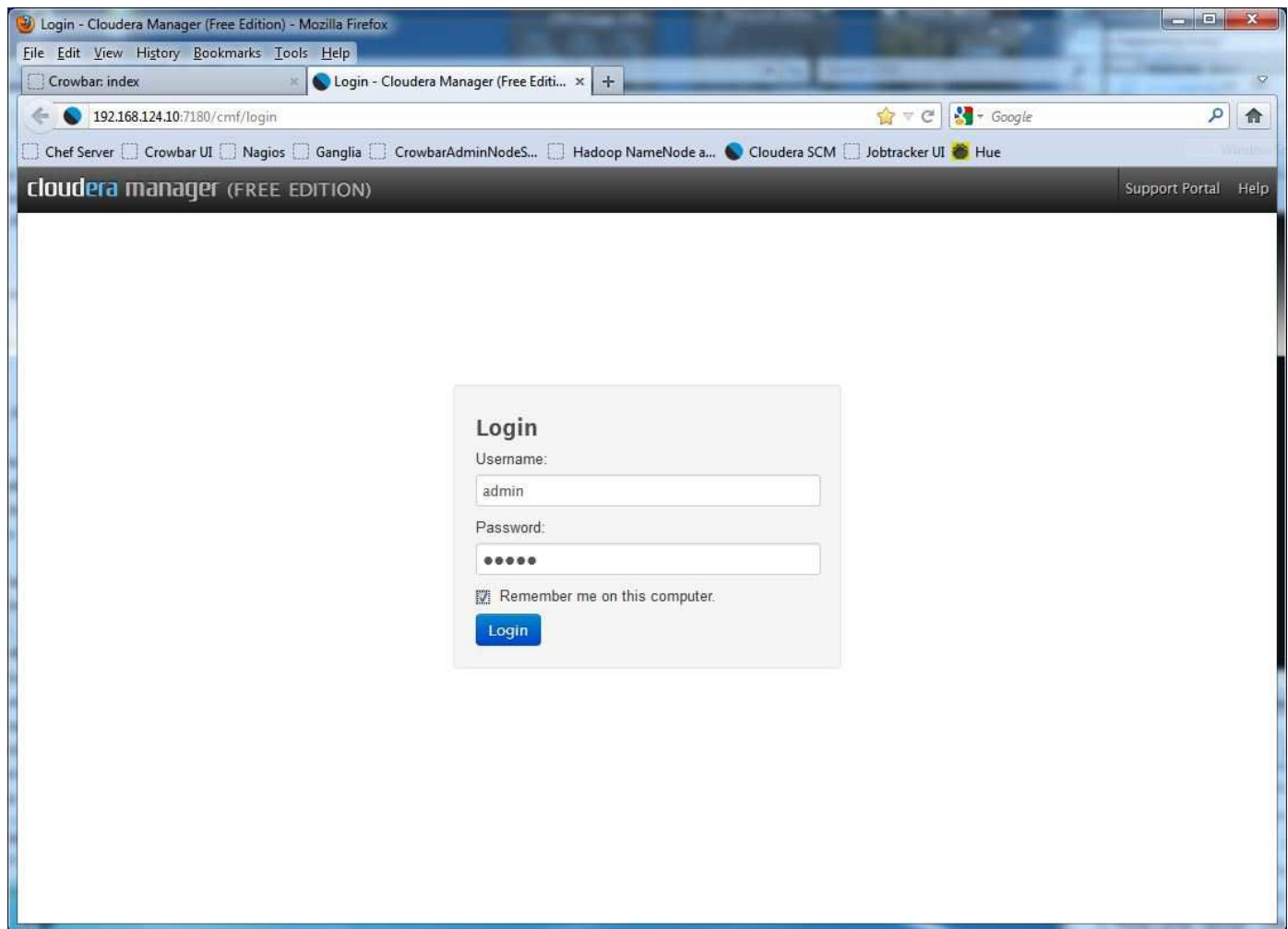
- You can also access the Cloudera Manager Administration Console from the Crowbar User Interface using the link located on the crowbar admin node view page (Cloudera Manager).

 **Note:** For security, you should change the password for the default admin user account as soon as possible. This option from the Cloudera Manager application under the Administration->Password tab.

Login Screen

- Enter the user login name and password (default=admin, admin).
- Check **Remember me on this computer** to store the browser session data if desired.
- Click the **Login** button to proceed.

Figure 1 Login Page

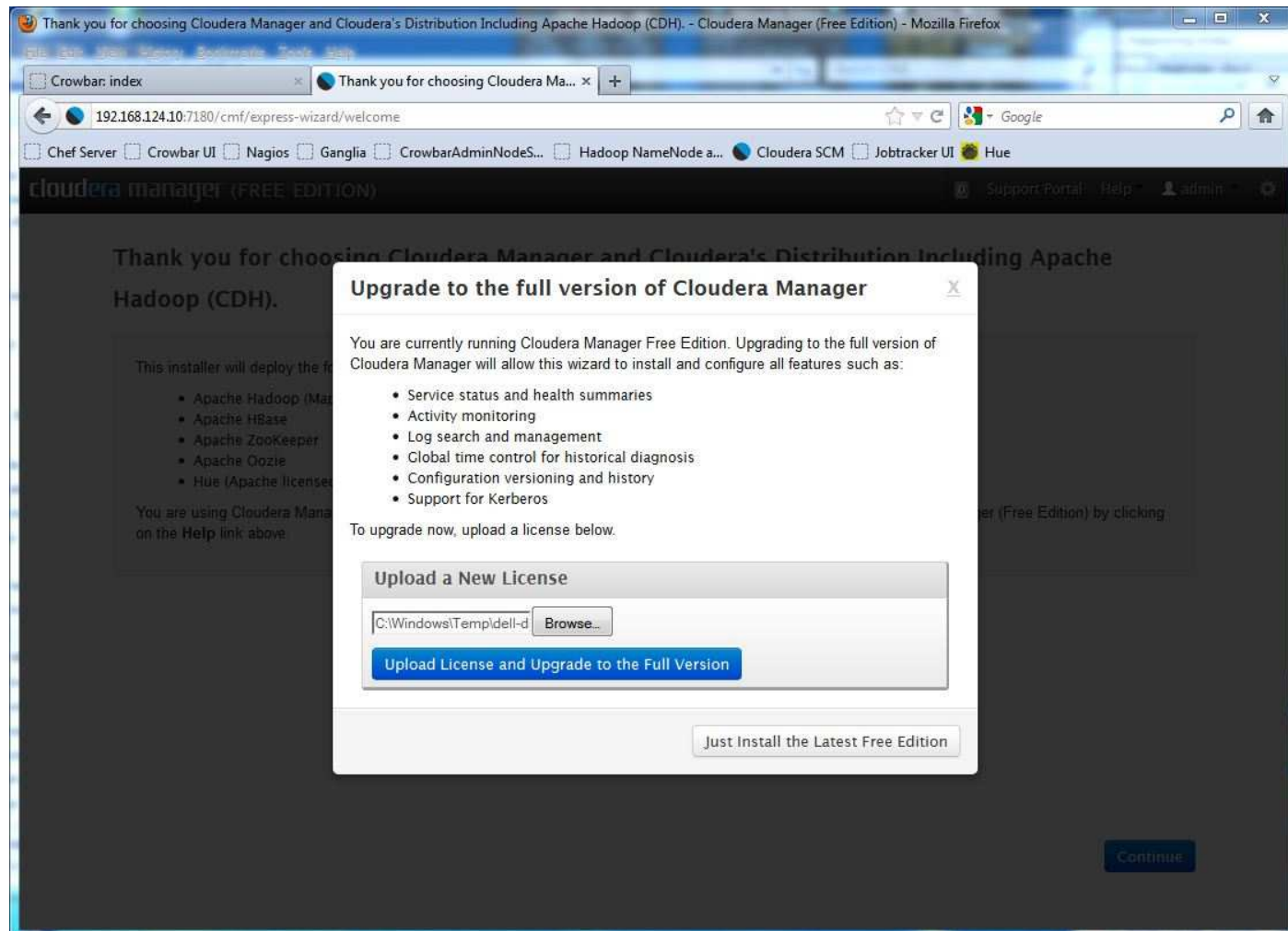


License Key Upload Screen

Note: Applying the license key is an optional step and you can always enter the license key later on in the process by clicking on the **welcome->admin** link in the Cloudera Manager administration console. This menu option is located at the top right side of the Cloudera Manager administration console screen.

- If you have obtained a Cloudera Manager License key and you wish to upgrade to the Cloudera Manager Enterprise edition, you can enter the license key at this point in time.
- Hit the **Browse** button and select the license key file.
- Hit the **Upload License and Upgrade to the Full Version** button to apply the license key.
- Hit the **Continue** Button to proceed after the license key has been applied.

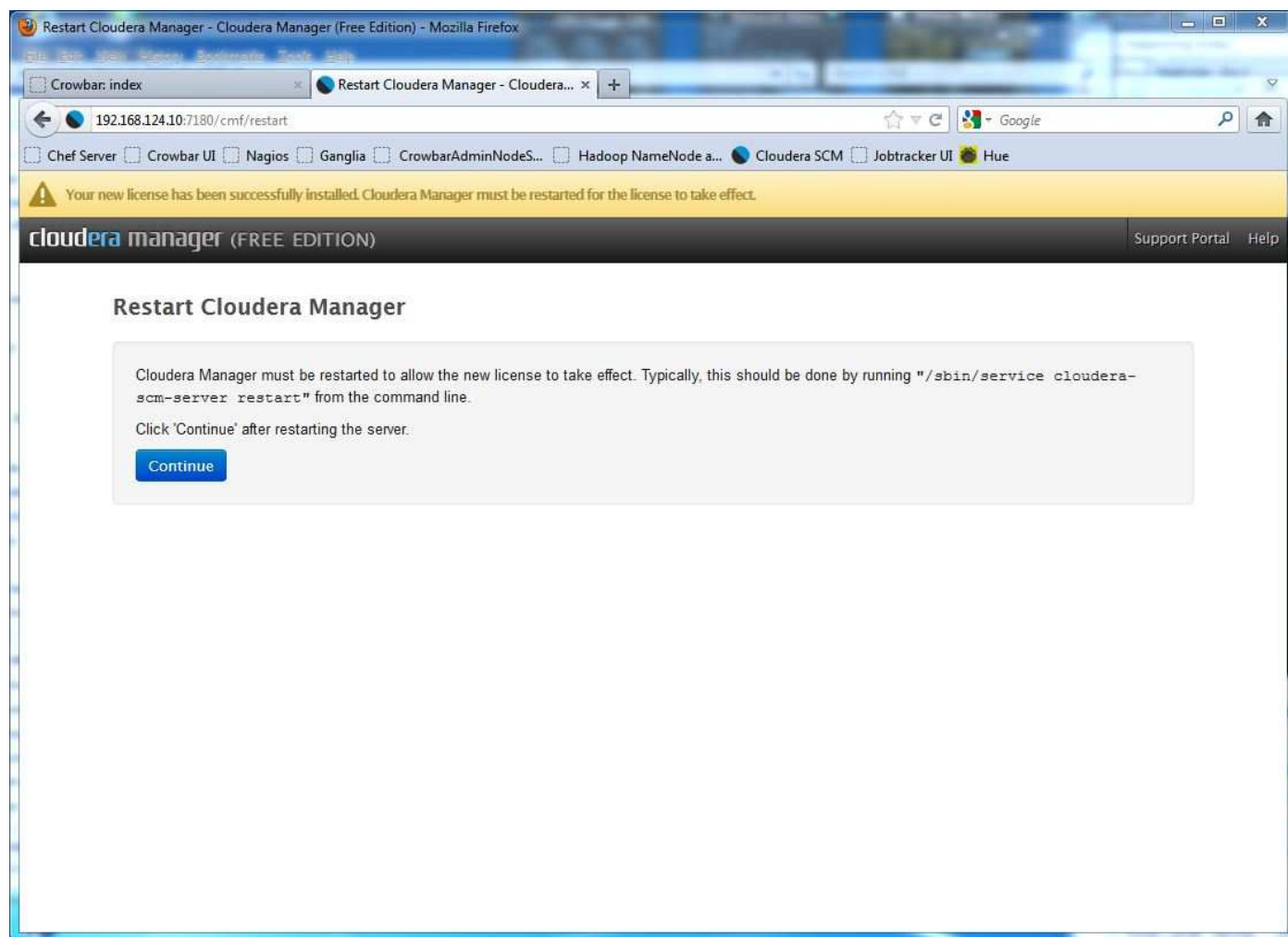
Figure 2 License Key Upload Screen



License Key Restart Screen

- Once the license key has been uploaded, the Cloudera Manager application will ask you to restart the Cloudera Manager server in order for it to take effect. You need to open an SSH console on the node which has the Cloudera Web Application role applied to it (login=root/crowbar) and execute the following commands;
 - `cd /etc/init.d`
 - `service cloudera-scm-server restart`
- Once the Cloudera manager server has been restarted, you need to restart the browser and log back into the Cloudera Manager Administration console.
- Once the Cloudera manager server has been restarted, you need to restart the browser and log back into the Cloudera Manager Administration console.

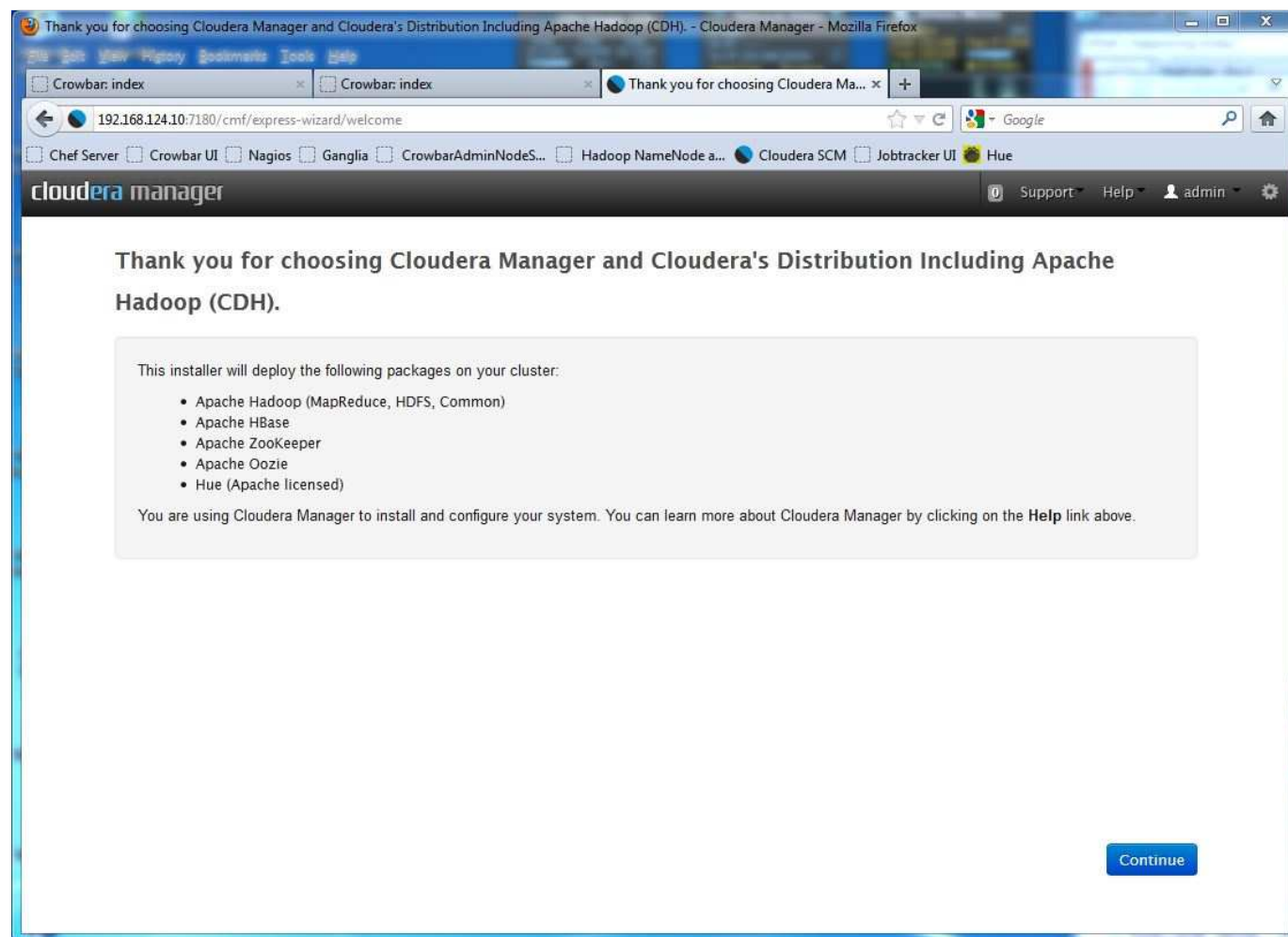
Figure 3 License Key Restart Screen



License Key Confirmation Screen

- If you have entered the Cloudera Manager License key, you will see this additional screen.
- Hit the **Continue** Button to proceed.

Figure 4 License Key Confirmation Screen



Cloudera Registration Screen

- If you wish to register with Cloudera, fill out the registration form and click the **Submit Registration** button.
- If you do not want to register, click the **Skip Registration** button to proceed.

Figure 5 Registration Screen

The screenshot shows a web browser window titled "Register your CDH installation. - Cloudera Manager - Mozilla Firefox". The address bar shows the URL "192.168.124.10:7180/cmf/express-wizard/registration?submit=Continue". The browser tabs include "Crowbar: index" and "Register your CDH installation. - Clou...". The Cloudera Manager interface has a top navigation bar with links for Chef Server, Crowbar UI, Nagios, Ganglia, CrowbarAdminNodeS..., Hadoop NameNode a..., Cloudera SCM, Jobtracker UI, and Hue. The main content area is titled "Register your CDH installation." and contains a registration form. The form includes a message: "Please help us improve our product by registering. Cloudera Manager will send your registration information to Cloudera. Fields marked * are required." The form fields are: First Name *, Last Name *, Email *, Company *, Position, Phone, City *, and State * (a dropdown menu with "Select a state..." as the current selection). Below the form is a checkbox labeled "Send me information about Cloudera products and services." which is checked. At the bottom right of the form are three buttons: "Back", "Skip Registration", and "Submit Registration".

Node Search Screen

- Enter the IP address range for all nodes on Hadoop cluster. This should be an alphanumeric search pattern and the crowbar admin node should be included in the search. For example;
 - **192.168.124.10, 192.168.124.8[1-3]** will attempt to discover 192.168.124.10, 192.168.124.81, 192.168.124.82 and 192.168.124.83,
- Click the **Search** button to search for available hosts within the specified IP address range.

Figure 6 Cloudera Cluster Node Search Screen

Specify hosts for your CDH cluster installation. - Cloudera Manager - Mozilla Firefox

192.168.124.10:7180/cmf/express-wizard/hosts

cloudera manager

Support Help admin

Specify hosts for your CDH cluster installation.

Hint: Search for hostnames and/or IP addresses using [patterns](#).

192.168.124.[10-90]

[View Currently Managed Hosts](#)

Search

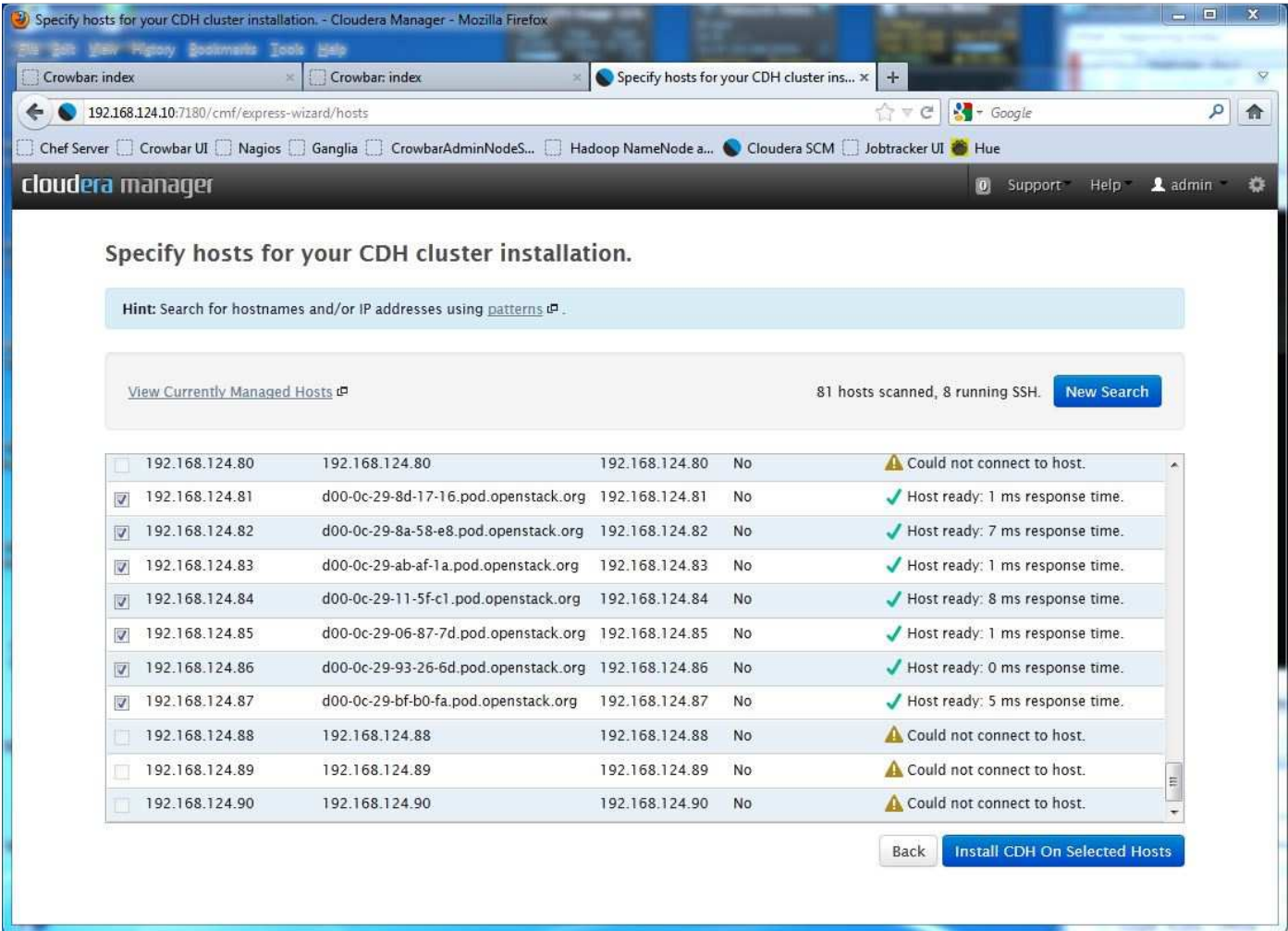
Back Install CDH On Selected Hosts

Waiting for www.google-analytics.com...

Node Search Results Screen


- Verify that all your hosts are discovered.
- Make any cluster configuration adjustments by selecting or deselecting nodes.
- Click the **Install CDH On Selected Hosts** button to proceed.

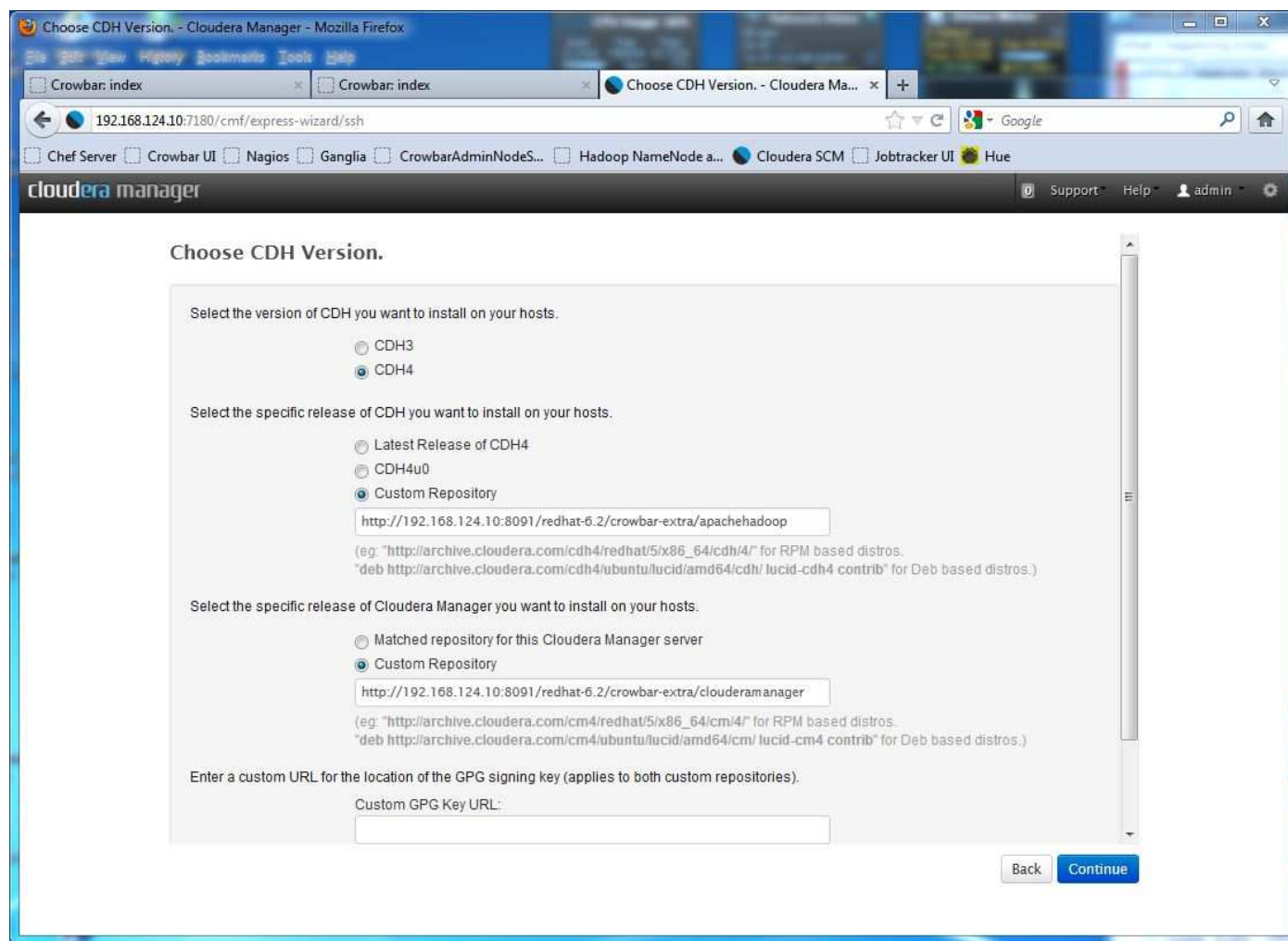
Figure 7 Node Search Results Screen



Repository Configuration Screen

- Select **CDH4** for installation.
- Select **Custom Repository** for CDH.
- Enter this URL - <http://192.168.124.10:8091/redhat-6.2/crowbar-extra/apachehadoop>
- Select **Custom Repository** for Cloudera Manager.
- Enter this URL - <http://192.168.124.10:8091/redhat-6.2/crowbar-extra/clouderamanager>
- Leave the GPG Key URL field blank.
- Click the **Continue** button to proceed.

 RPM based packages are served from the crowbar admin node. By default, the IP address is 192.168.124.10 on port 8091 (http:// 192.168.124.10:8091). If you configure the crowbar admin node to be on another IP address, you will have to make the appropriate adjustments to the URL.



SSH Credentials Screen

- Select **Login to all hosts as root**.
- Select **All hosts accept same password**.
- Enter the SSH login credentials for the cluster (default=root, crowbar).
- Click the **Start Installation** button to proceed.

Figure 8 SSH Credentials Screen

Provide SSH login credentials. - Cloudera Manager - Mozilla Firefox

192.168.124.10:7180/cmf/express-wizard/ssh#ssh

cloudera manager

Support Help admin

Provide SSH login credentials.

Root access to your hosts is required to install the Cloudera packages. This installer will connect to your hosts via SSH and log in either directly as root or as another user with password-less sudo privileges to become root.

Login to all hosts as: ☒ root
☐ Another User:

You may connect via password or public-key authentication for the user selected above.

Authentication Method: ☒ All hosts accept same password
☐ All hosts accept same public key

Enter Password:

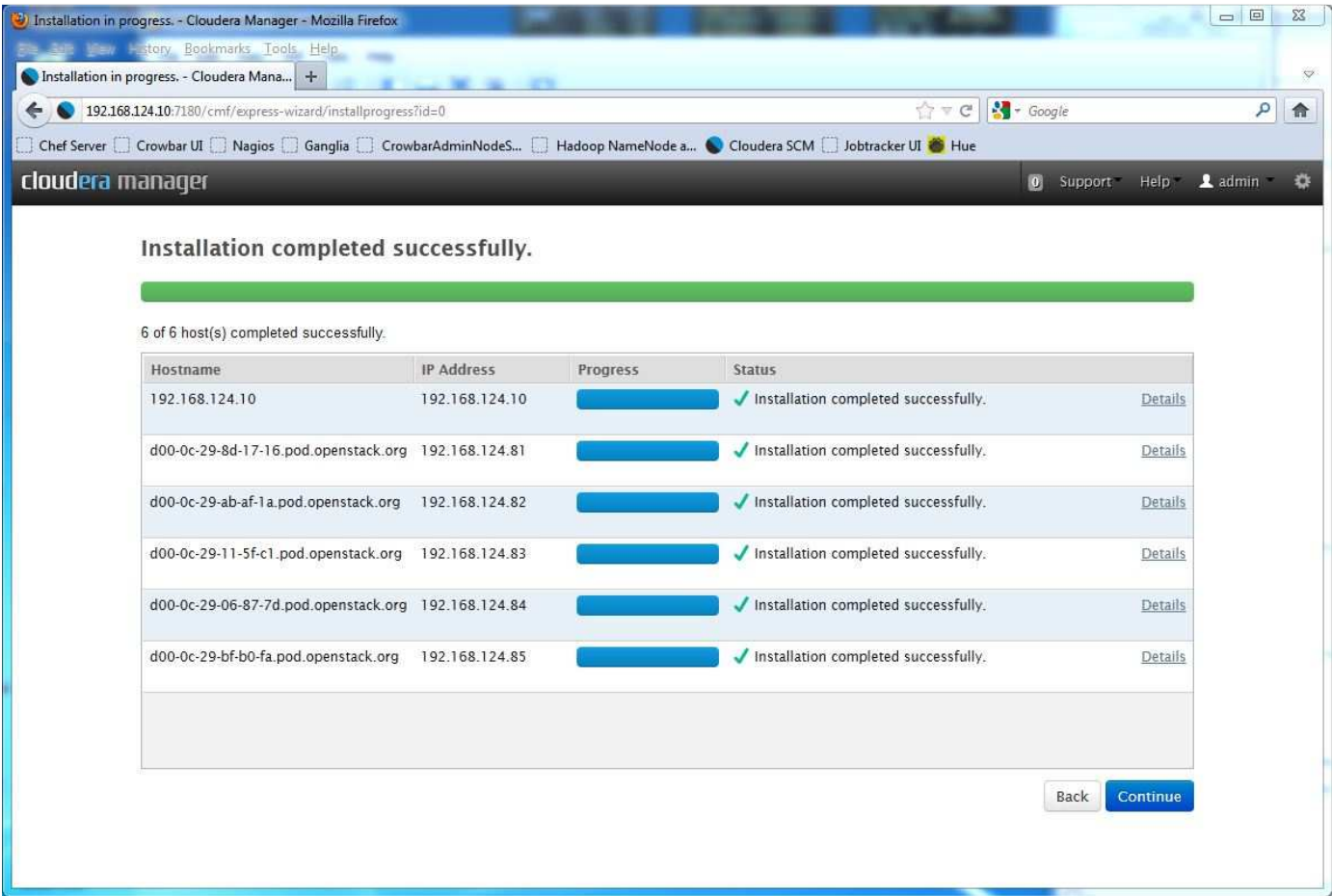
Confirm Password:

Back Start Installation

Package Install Screen

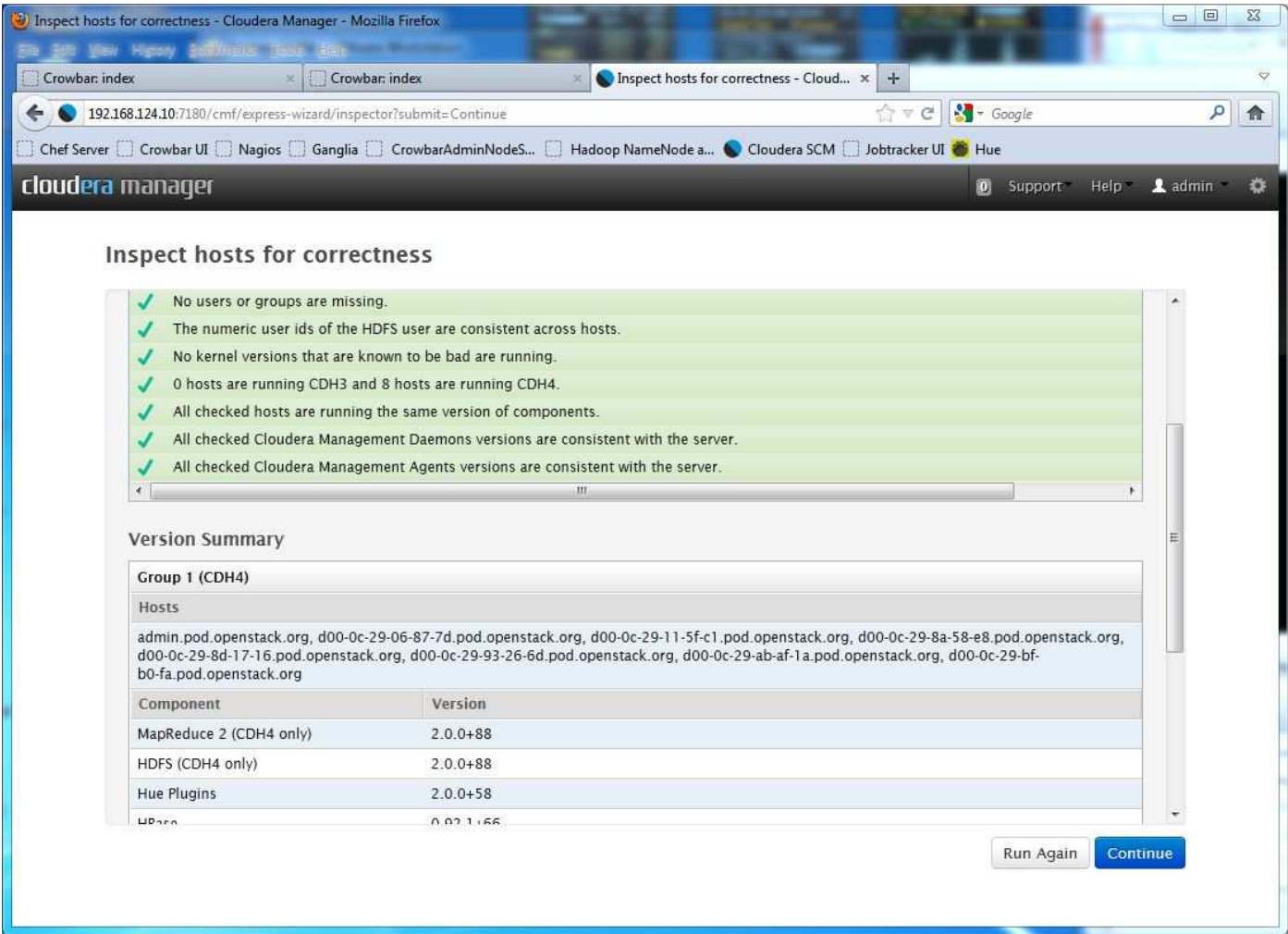
- You will see bar graphs next to each node and the name of the package it is installing.
- Wait for the installation process to complete.
- Click the **Continue** button to proceed.

Figure 9 Package Install Screen



Host Inspector Screen

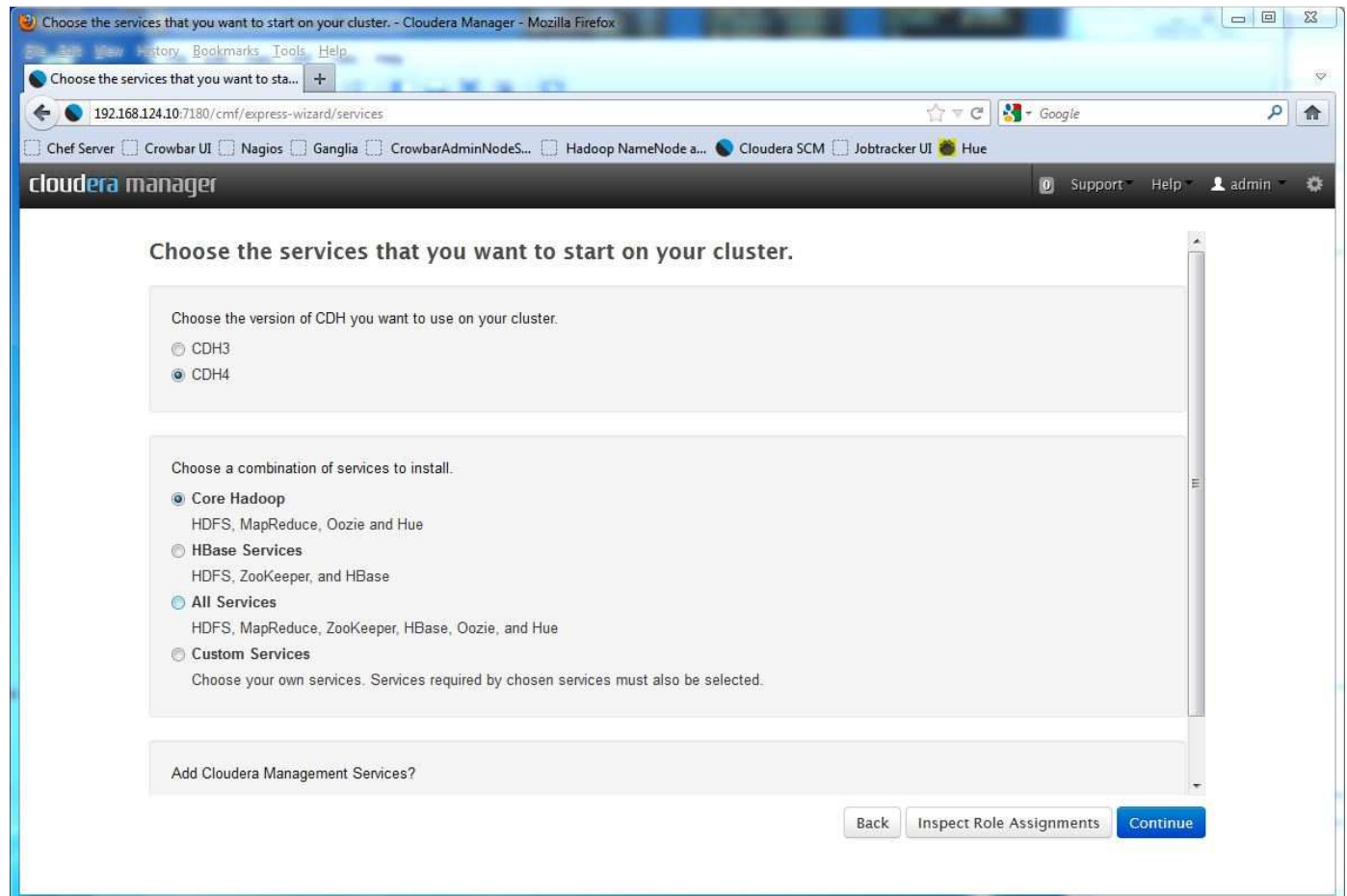
- The Cloudera manager **Host Inspector** runs during this part of the installation process in order to verify proper cluster configuration.
- Wait for this process to complete.
- Click the **Continue** button to proceed.



Service Selection Screen

- Select the services that you want to install.
- You can install **All Services** now or **Core Services** and add others later as needed.
- Click on the **Inspect Role Assignments** button to configure the Hadoop cluster services. Do not select **Continue** as this will give you the default role assignments which are probably not what you really want.

Figure 10 Service Selection Screen



Inspect Role Assignments Screen # 1

- Select the role assignments for Hadoop cluster deployment (HDFS/Map Reduce).
- Recommended settings for the Dell Reference Architecture:
 - CM DataNode role – Crowbar nodes which contains the [clouderamanager-slavenode](#) role.
 - CM NameNode – Crowbar node which contains the [clouderamanager-masternamenode](#) role.
 - CM Secondary NameNode – Crowbar node which contains the [clouderamanager-secondarynamenode](#) role.
 - CM TaskTracker roles – Crowbar nodes which contains the [clouderamanager-slavenode](#) role.
 - CM JobTracker role– Crowbar node which contains the [clouderamanager-masternamenode](#) role.
- Please refer to the next diagram (Figure 8 – Screen #2) before clicking the **Continue** button.

Note: The Cloudera Node Inventory page you printed /exported from within the Cloudera Manager barclamp page in Crowbar is very useful for this step to ensure the roles selected in Cloudera Manager are assigned to nodes which have been provisioned (RAID, BIOS, etc ...) specifically for that purpose.

Figure 11 Inspect Role Assignments Screen # 1

Inspect role assignments

You can customize the role assignments for your new cluster here, but note that if assignments are made incorrectly, such as assigning too many roles to a single host, this can significantly impact the performance of your services. Cloudera does not recommend altering assignments unless you have specific requirements, such as having pre-selected a specific host for a specific role.

The host list presented here is prefiltered to remove hosts which are not valid candidates; these include hosts that are: unhealthy, members of other clusters, and/or which have an incompatible version of CDH installed on them.

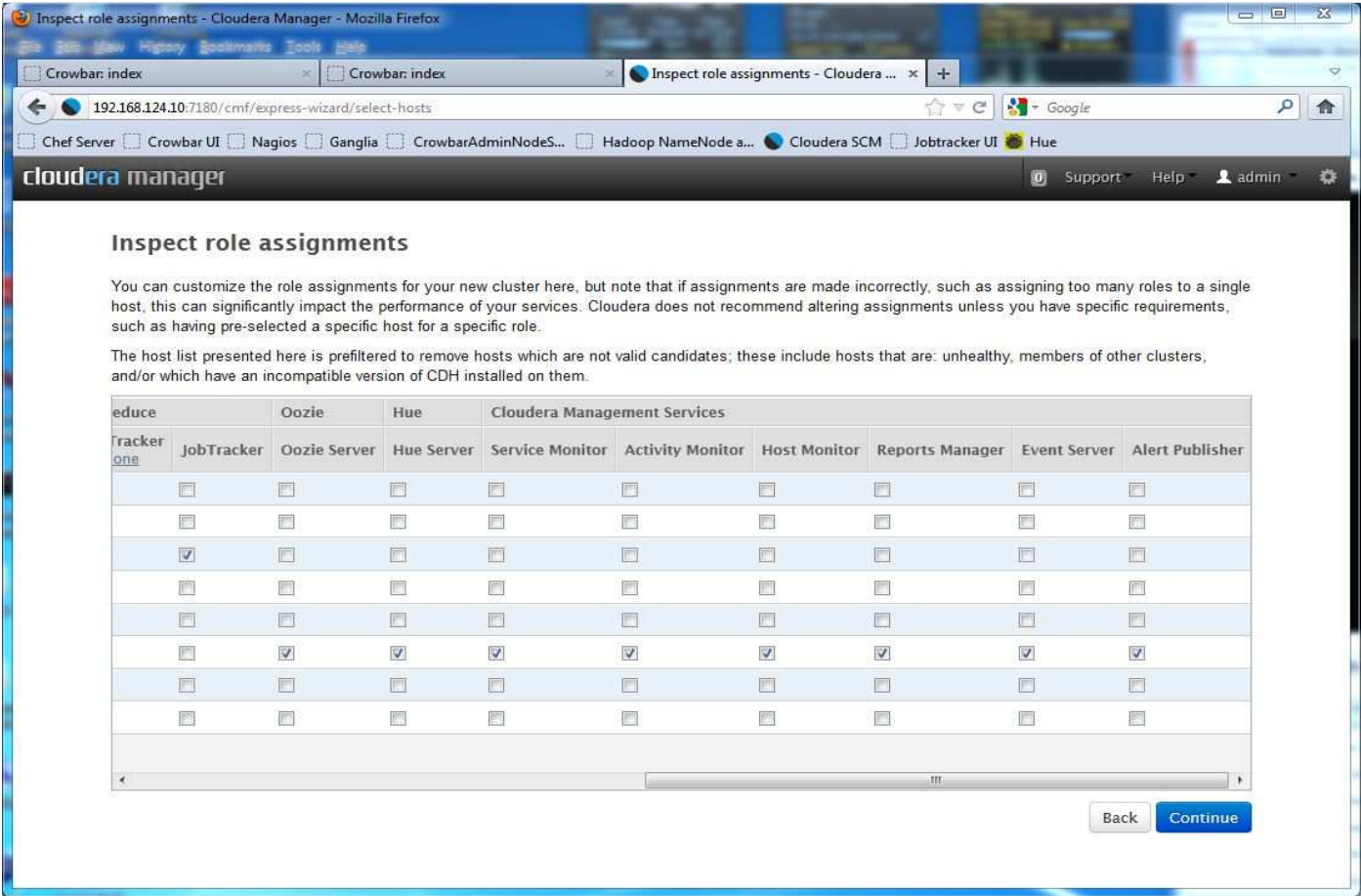
	HDFS				MapReduce			Oozie	Hue
	DataNode All None	NameNode	SecondaryNameNode	HttpFS All None	TaskTracker All None	JobTracker	Oozie Server	Hue Serv	
admin.pod.openstack.org	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	
d00-0c-29-06-87-7d.pod.openstack.org	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	
d00-0c-29-11-5f-e1.pod.openstack.org	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	
d00-0c-29-8a-58-e8.pod.openstack.org	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	
d00-0c-29-8d-17-16.pod.openstack.org	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	
d00-0c-29-93-26-6d.pod.openstack.org	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	
d00-0c-29-ab-af-1a.pod.openstack.org	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	
d00-0c-29-bf-b0-fa.pod.openstack.org	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	

Back Continue

Inspect Role Assignments Screen # 2

- If you entered the Cloudera Manager License key, you will see this additional screen.
- Select the role assignments for Hadoop add-ons services and monitoring services (Activity Monitor, Service Monitor Resource Manager).
- Click the **Continue** button to proceed.

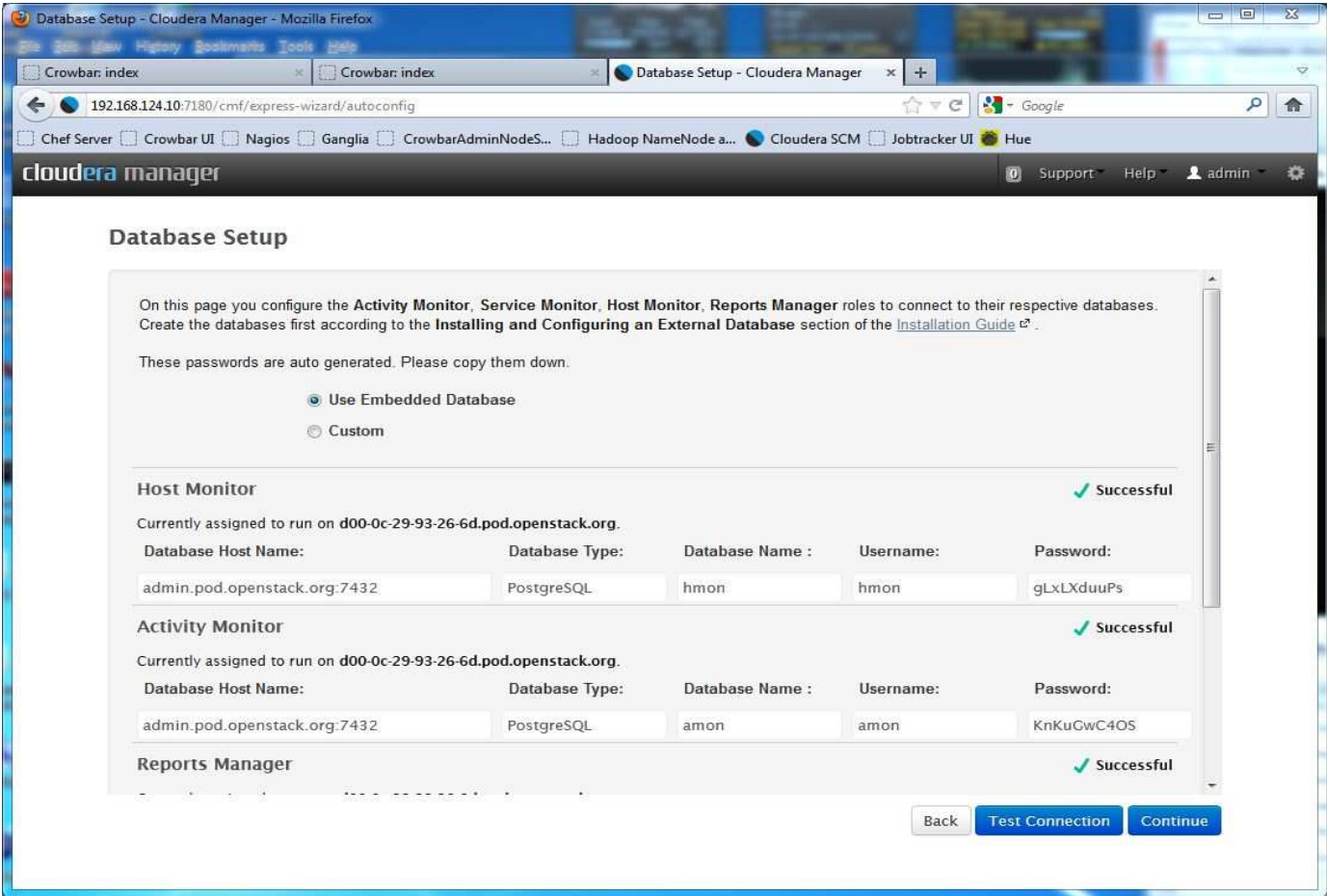
Figure 12 Inspect Role Assignments Screen #2



Monitoring Database Setup Screen

- If you entered the Cloudera Manager License key, you will see this additional screen.
- Select **Use Embedded Database**.
- You can leave the rest of the settings at default values unless you want to change them.
- Click the **Test Connection** button to make sure you can connect to all the databases.
- Click the **Continue** button to proceed.

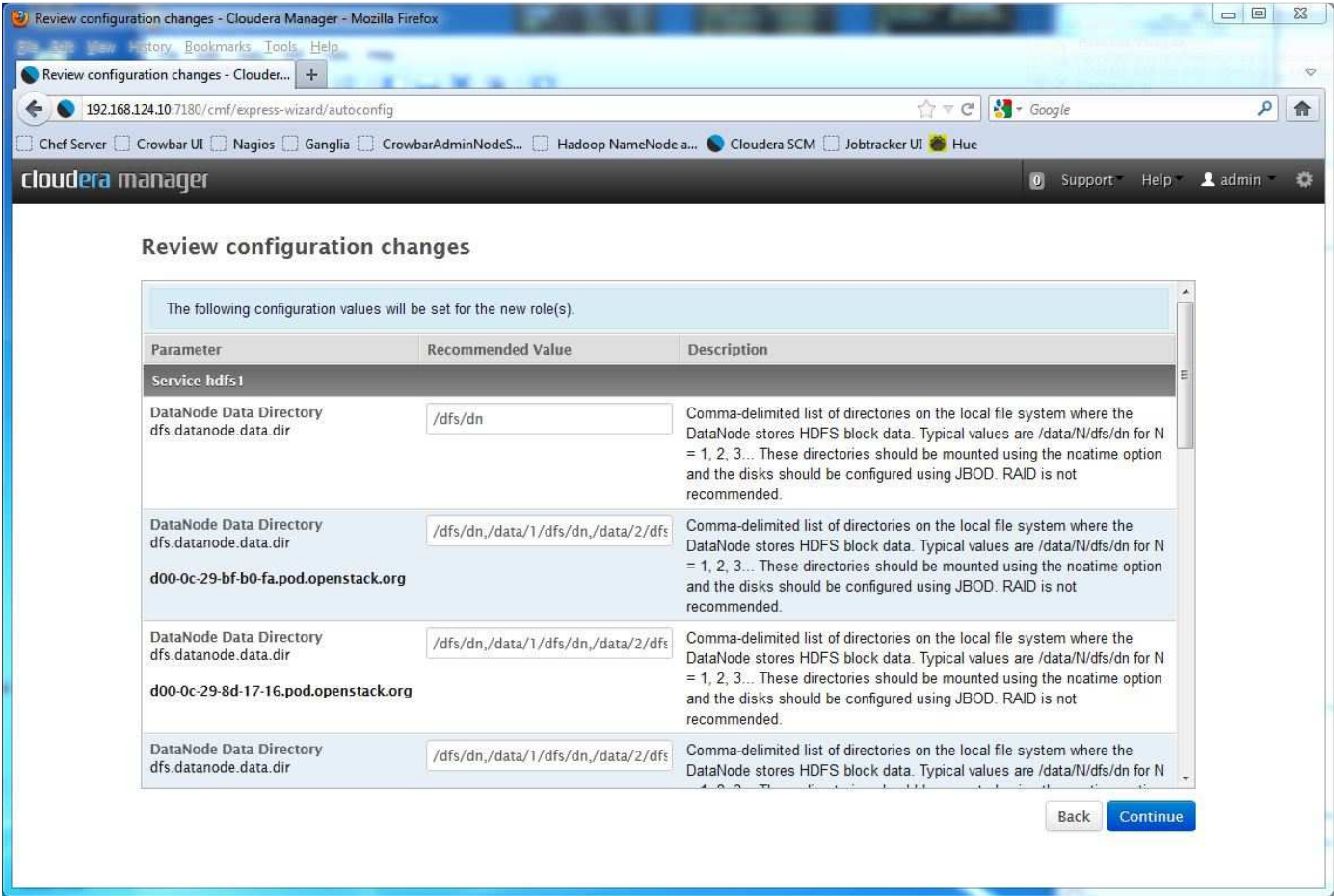
Figure 13 Monitoring Database Setup Screen



Review Configuration Changes Screen

- If you entered the Cloudera Manager License key, you will see this additional screen.
- Set the mail server hostname for alerts (localhost).
- Set the mail server message recipients for alerts.
- Click the **Continue** button to proceed.

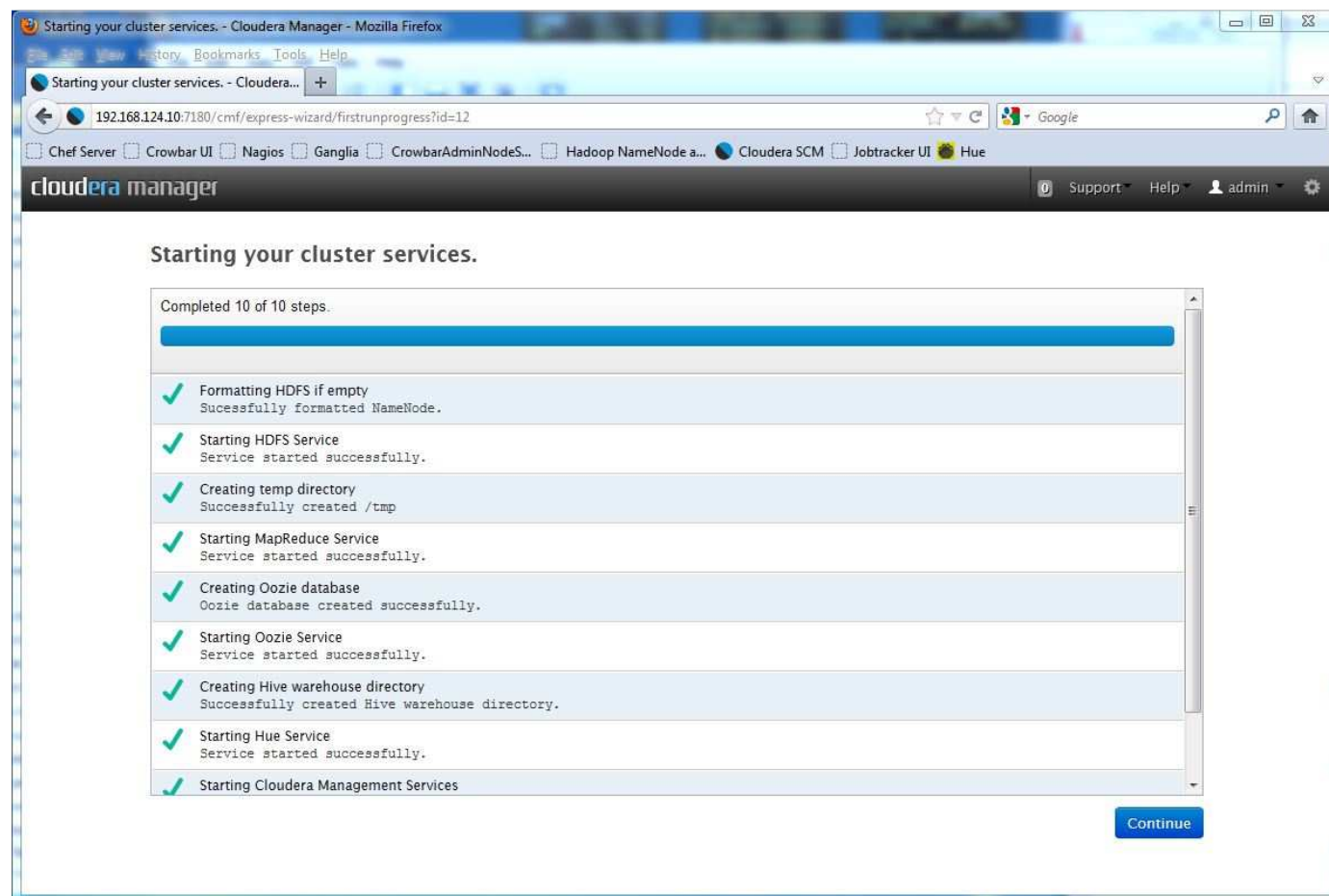
Figure 14 Review Configuration Changes Screen



Cluster Services Initialization Screen

- Wait for the Hadoop cluster configuration process to complete.
- Click the **Continue** button to proceed.

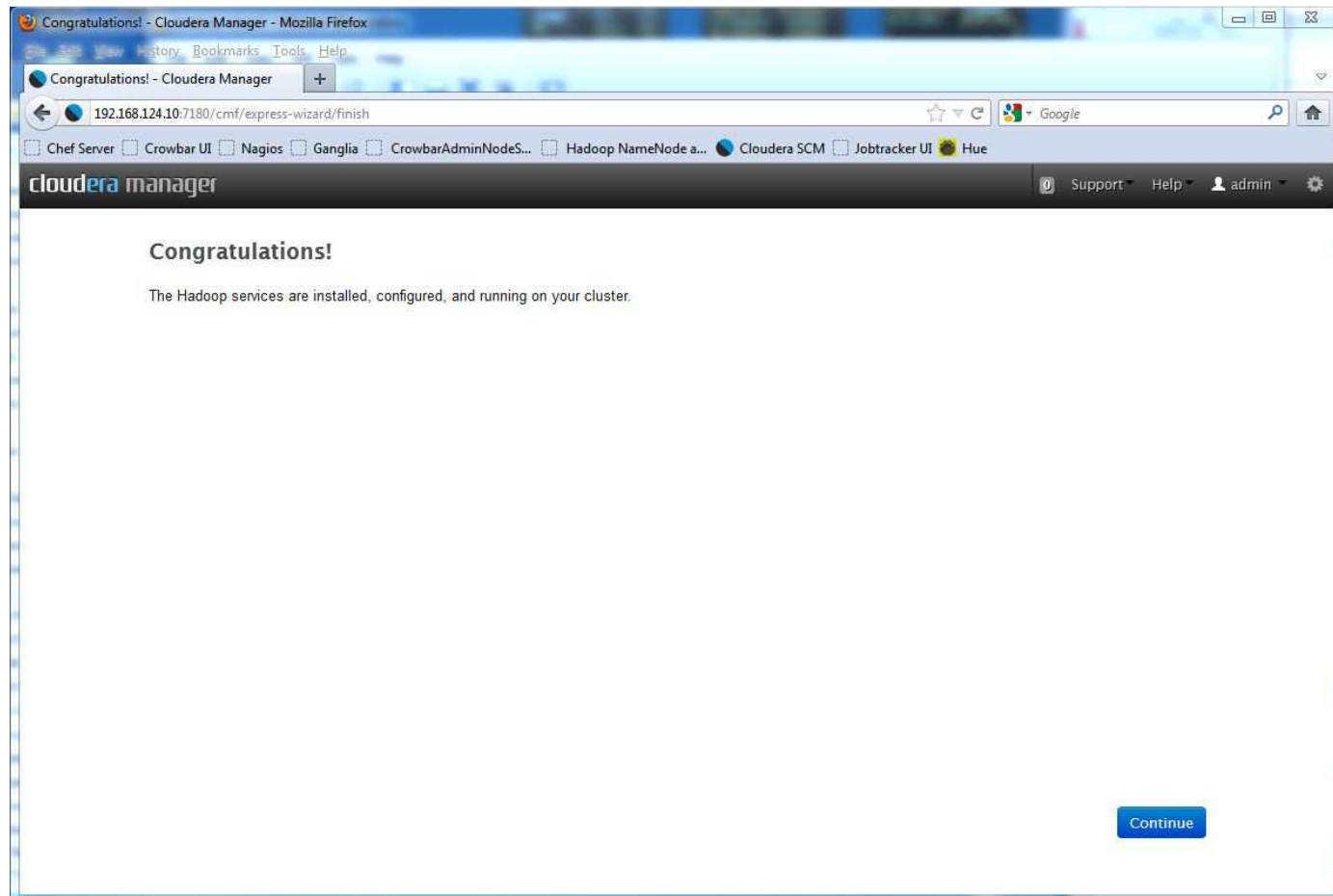
Figure 15 Cluster Services Initialization Screen



Configuration Completion Screen

- If the Hadoop configuration steps complete successfully, you will see the final Cloudera Manager confirmation screen.
- Click the **Continue** button to start using Cloudera Manager.

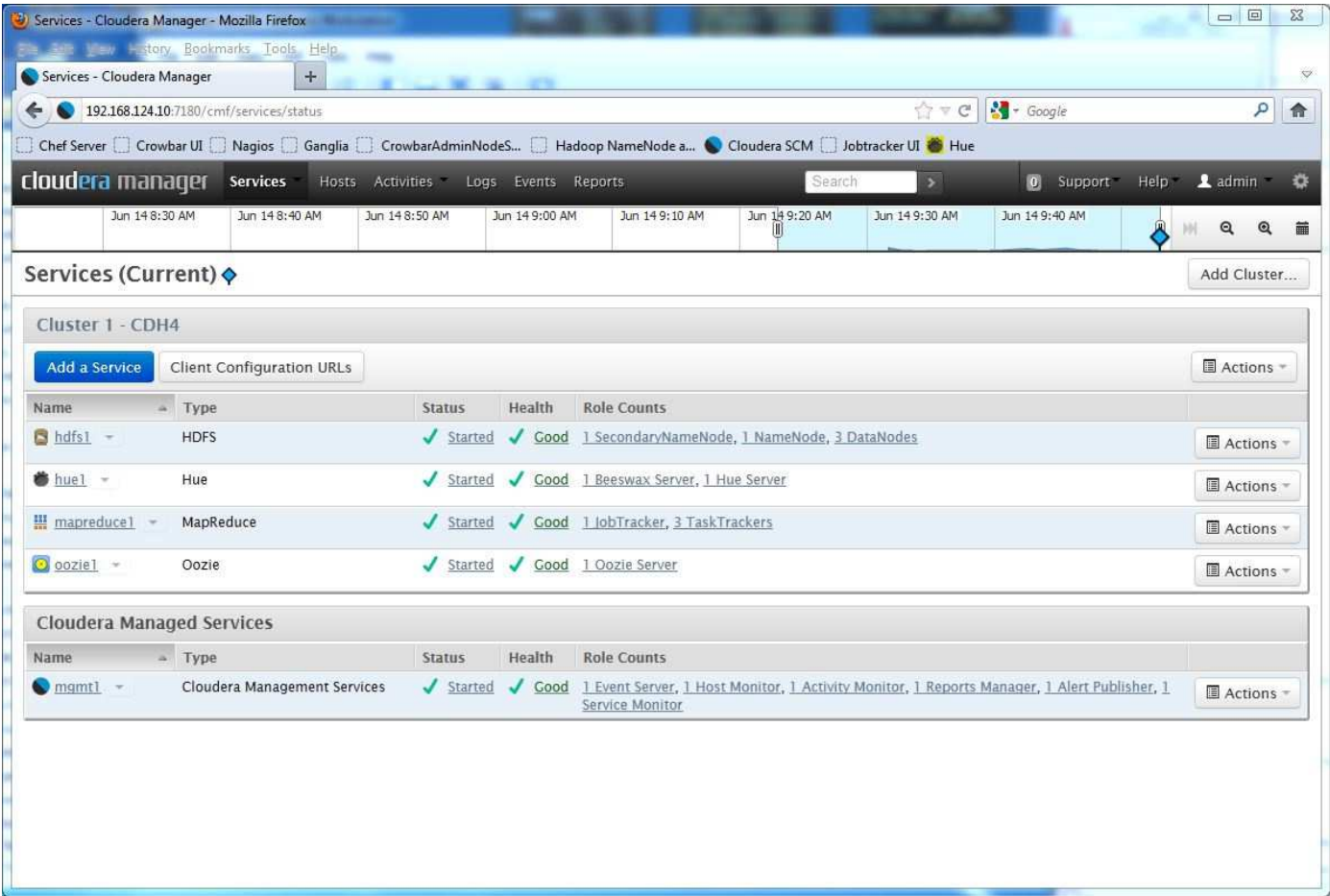
Figure 16 Configuration Completion Screen



Service Display Screen

- This is the normal startup screen after Cloudera Manager has completed the installation steps.
- Please refer to the Cloudera Manager User Guide for additional information on operating Cloudera Manager.

Figure 17 Service Display Screen



Pig Barclamp

Apache Pig is a platform for analyzing large data sets that consists of a high-level language for expressing data analysis programs, coupled with infrastructure for evaluating these programs. The salient property of Pig programs is that their structure is amenable to substantial parallelization, which in turns enables them to handle very large data sets.

Pig's infrastructure layer consists of a compiler that produces sequences of MapReduce programs, for which large-scale parallel implementations already exist (e.g., the Hadoop subproject). Pig's language layer currently consists of a textual language called Pig Latin, which has the following key properties:

- **Ease of programming:** It is trivial to achieve parallel execution of simple, "embarrassingly parallel" data analysis tasks. Complex tasks comprised of multiple interrelated data transformations are explicitly encoded as data flow sequences, making them easy to write, understand, and maintain.
- **Optimization opportunities:** The way in which tasks are encoded permits the system to optimize their execution automatically, allowing the user to focus on semantics rather than efficiency.
- **Extensibility:** Users can create their own functions to do special-purpose processing.

Table 4-34: Pig Barclamp Parameters

Name	Description	Required	Default
java_home	JAVA_HOME environment variable.	true	/usr/java/jdk1.6.0_27/jre
log4jconf	log4jconf log4j configuration file.	true	./conf/log4j.properties
brief	brief logging - no timestamps.	true	false
cluster	Clustername, name of the hadoop jobtracker. If no port is defined port 50020 will be used.	false	
debug_level	Debug level, INFO is default.	true	INFO
file	A file that contains pig script.	false	
jar	Load jarfile, colon separated.	false	
verbose	Verbose print all log messages to screen (default to print only INFO and above to screen).	true	false
exectype	Exectype local or mapreduce - mapreduce is default.	true	mapreduce
ssh_gateway	HOD gateway property.	false	
hod_expect_root	HOD expect root property.	false	
hod_expect_uselatest	HOD use latest root property.	false	
hod_command	HOD command root property.	false	
hod_config_dir	HOD config directory property.	false	
hod_param	HOD param property.	false	
pig_spill_size_threshold	Do not spill temp files smaller than this size (bytes).	true	5000000
pig_spill_gc_activation_size	EXPERIMENT: Activate garbage collection when spilling a file bigger than this size (bytes). This should help reduce the number of files being spilled.	true	40000000
log_file	Log file location.	false	

Hive Barclamp

Hive is a data warehouse system for Hadoop that facilitates easy data summarization, ad-hoc queries, and the analysis of large datasets stored in Hadoop compatible file systems. Hive provides a mechanism to project structure onto this data and query the data using a SQL-like language called HiveQL. This language also allows traditional map/reduce programmers to plug in their custom mappers and reducers when it is inconvenient or inefficient to express this logic in HiveQL.

Table 4-41: Hive Barclamp Parameters

Name	Description	Required	Default
hive_exec_scratchdir	Scratch space for Hive jobs.	true	/tmp/hive-\${user.name}
hive_metastore_local	Controls whether to connect to remove metastore server or open a new metastore server in Hive Client JVM.	true	true
javax_jdo_option_Connectio nURL	JDBC connect string for a JDBC metastore.	true	jdbc:derby::databaseName=metast ore_db;create=true
javax_jdo_option_Connectio nDriverName	Driver class name for a JDBC metastore.	true	org.apache.derby.jdbc.EmbeddedD river
hive_metastore_metadb_dir	The location of filestore metadata base dir.	true	file:///var/metastore/metadb/
hive_metastore_uris	Comma separated list of URIs of metastore servers. The first server that can be connected to will be used.	true	file:///var/metastore/metadb/
hive_metastore_warehouse_ dir	The location of the default database for the warehouse.	true	/user/hive/warehouse
hive_metastore_connect_ret ries	Number of retries while opening a connection to metastore.	true	5
hive_metastore_rawstore_i mpl	Name of the class that implements org.apache.hadoop.hive.metastore.rawsto re interface. This class is used to store and retrieval of raw metadata objects such as table, database.	true	org.apache.hadoop.hive.metastore. ObjectStore
hive_default_fileformat	Default file format for CREATE TABLE statement. Options are TextFile and SequenceFile.	true	TextFile
hive_map_aggr	Whether to use map-side aggregation in Hive Group By queries.	true	false
hive_join_emit_interval	How many rows in the right-most join operand Hive should buffer before emitting the join result.	true	1000
hive_exec_script_maxerrsize	Maximum number of bytes a script is allowed to emit to standard error (per map-reduce task). This prevents runaway scripts from filling logs partitions to capacity .	true	100000
hive_exec_compress_output	Controls whether the final outputs of a query (to a local/hdfs file or a hive table) is compressed. The compression codec and other options are determined from hadoop config variables mapred.output.compress.	true	false
hive_exec_compress_interm ediate	Controls whether intermediate files produced by hive between multiple map-	true	false

Dell | Cloudera Apache Hadoop Solution

Name	Description	Required	Default
	reduce jobs are compressed. The compression codec and other options are determined from hadoop config variables mapred.output.compress.		

Sqoop Barclamp

Sqoop is an SQL based command-line tool to assist with HDFS data import/export (SQL-to-Hadoop). Sqoop is a tool designed to transfer data between Hadoop and relational databases. You can use Sqoop to import data from a relational database management system (RDBMS) such as MySQL or Oracle into the Hadoop Distributed File System (HDFS), transform the data in Hadoop MapReduce, and then export the data back into an RDBMS.

Sqoop automates most of this process by relying on the database to describe the schema for the data to be imported. Sqoop uses MapReduce to import and export the data, which provides parallel operation as well as fault tolerance.

Table 4-48: Sqoop Barclamp Parameters

Name	Description	Required	Default
sqoop_connection_factories	A comma-delimited list of ManagerFactory implementations which are consulted, in order, to instantiate ConnManager instances used to drive connections to databases.	false	
sqoop_tool_plugins	A comma-delimited list of ToolPlugin implementations which are consulted, in order, to register SqoopTool instances which allow third-party tools to be used.	false	
sqoop_metastore_client_enable_autoconnect	If true, Sqoop will connect to a local metastore for job management when no other metastore arguments are provided.	true	false
sqoop_metastore_client_autoconnect_url	The connect string to use when connecting to a job-management metastore. If unspecified, uses ~/.sqoop/. You can specify a different path here.	false	
sqoop_metastore_client_autoconnect_username	The username to bind to the metastore.	false	
sqoop_metastore_client_autoconnect_password	The password to bind to the metastore.	false	
sqoop_metastore_client_record_password	If true, allow saved passwords in the metastore.	false	
sqoop_metastore_server_location	Path to the shared metastore database files. If this is not set, it will be placed in ~/.sqoop/.	false	
sqoop_metastore_server_port	Port that this metastore should listen on.	false	

Support

Cloudera Support

To obtain support for Hadoop:

- Open a request at Cloudera's support portal: <http://www.cloudera.com/hadoop-support/>

Printed in USA

www.dell.com | support.dell.com

Appendix A: Dell | Hadoop Solution Components

- **Hadoop:** <http://en.wikipedia.org/wiki/Hadoop>
- **Hadoop Distributed File System (HDFS):** A distributed file system that provides high-throughput access to application data (http://en.wikipedia.org/wiki/Hadoop_Distributed_FileSystem#Hadoop_Distributed_File_System).
- **MapReduce:** A software framework for distributed processing of large data sets on compute clusters (<http://en.wikipedia.org/wiki/MapReduce>).
- **HBase:** A scalable, distributed database that supports structured data storage for large tables.
- **Hive:** A data warehouse infrastructure that provides data summarization and ad-hoc querying.
- **ZooKeeper:** A high-performance coordination service for distributed applications.
- **Pig:** A platform for analyzing large data sets that consists of high-level language for expressing data analysis programs, coupled with infrastructure for evaluating these programs.
- **Sqoop:** A tool designed to import data from relational databases into Hadoop. Sqoop uses JDBC to connect to a database.
- **Oozie:** An open-source workflow engine and coordination service to manage data processing jobs within Hadoop.
- **Hue:** A browser based interface for interacting with Hadoop clusters.
- **Crowbar:** A Dell provided, supported, and maintained toolset for system deployment and configuration automation. Crowbar supports the bare-metal bring-up of new hardware and configuration management of existing hardware.

Appendix B: External References

- Cloudera: <http://www.cloudera.com>
- Nagios: <http://www.nagios.org>
- Ganglia: <http://ganglia.sourceforge.net>

To Learn More

For more information on the Dell | Cloudera Apache Hadoop Solution, visit:

www.Dell.com/Hadoop

©2011 Dell Inc. All rights reserved. Trademarks and trade names may be used in this document to refer to either the entities claiming the marks and names or their products. Specifications are correct at date of publication but are subject to availability or change without notice at any time. Dell and its affiliates cannot be responsible for errors or omissions in typography or photography. Dell's Terms and Conditions of Sales and Service apply and are available on request. Dell service offerings do not affect consumer's statutory rights.

Dell, the DELL logo, and the DELL badge, PowerConnect, and PowerVault are trademarks of Dell Inc.