



HADOOP Cluster Setup

-

setting up a HADOOP cluster
using Ambari and CentOS 7

By Henrik B, Sørensen @ Beech Grove
Updated November 7, 2016

Version 1.0

Table of Contents

What is Hadoop?.....	3
Purpose.....	3
About the author.....	3
Disclaimer.....	3
Prerequisites.....	4
Server Environment.....	4
Hadoop overview.....	6
Ambari Installation.....	10
Adding Spark.....	22
Test.....	26
Where to go from here?.....	28

What is Hadoop?

The short form :

Hadoop is an ecosystem for cluster computing. Cluster computing is the term used for connecting multiple computers together in a huge "computer" in order to parallel process computational tasks. E.g. bank sometimes connect there computers in order to perform calculations. Hadoop and assistant frameworks offers a way of collecting, importing and processing data including failover.

The official version is :

The Apache Hadoop software library is a framework that allows for the distributed processing of large data sets across clusters of computers using simple programming models. It is designed to scale up from single servers to thousands of machines, each offering local computation and storage. Rather than rely on hardware to deliver high-availability, the library itself is designed to detect and handle failures at the application layer, so delivering a highly-available service on top of a cluster of computers, each of which may be prone to failures.

Purpose

The purpose of this document is to describe the process of installing a Hadoop Cluster setup using Ambari. The software used in this process will only be Open Source and / or free of charge and / or own developed software / scripts.

About the author



Henrik B. Sørensen holds a Bachelors Degree in Electronics Engineering and has worked for the last 16 years in developing software including distributed computing systems for financial institutes.

His work in data measurement and processing has lead him over to Hadoop.

He can be contacted at info@beech-grove.eu for comments regarding this tutorial. This email is not to be used for asking for support of any kind.

Disclaimer

This document is provided as is. No guarantees of any kind is given. Also, it is highly recommended to consult documentation / best practices for e.g. handling security. The document is meant as a way of getting started with Hadoop on a cluster setup and not necessarily for production.

Copying of information / passages from this tutorial is allowed, but please quote this tutorial.

Prerequisites

In order to install Hadoop some things are necessary.

First of all, you need somewhere to install the servers. It may be physical or virtual machines. In this document, I'll use a number of virtual machines located on 2 physical servers running Xen Server. Given the tutorial of the previous document ("HADOOP Cluster Setup - setting up a HADOOP cluster Ambari and CentOS 6") and a multitude of tutorials online, I will save space and not go through the creation of the virtual machines. One note, though. I use XenCenter to control my virtual servers. For some reason, the XenCenter and CentOS 7 graphical installer do not play well together. I found a workaround : Hit Tab when the installation begins and swap `quiet` with `inst.text`. This will give a console / text installation process.

Server Environment

My environment¹ is as follows :

Name	Address	Description
Utilities (utilities)	192.168.1.115	A file server holding a number of script files and demos
Ambari (ambari.cluster)	192.168.1.119	Server hosting the Ambari web site.
Master A (mastera.cluster)	192.168.1.120	Primary NameNode for the cluster
Master B (masterb.cluster)	192.168.1.130	Secondary NameNode for the cluster
Clients A 01-03 (clienta01-3.cluster)	192.168.1.121-123	3 clients in a group
Clients B 01-03 (clientb01-3.cluster)	192.168.1.131-133	3 clients in a group

The naming convention is based on which physical server, the virtual machines are located on.

Each of the machines need to have Java installed (I use the Java-1.8.0-openjdk) and have the `JAVA_HOME` environment variable set to the folder containing the Java Runtime Engine.

```
/usr/lib/jvm/jre-1.8.0-openjdk
```

Logging on to the Ambari server, we need to setup a few things. In order to enable the EPEL repository, setup the SSH key and install the Parallel Distributed Shell (`pdsh` – for remote shell commands), we need to run the following script (`SetupServerEnv.sh`) :

```
#!/bin/bash
wget http://dl.fedoraproject.org/pub/epel/7/x86_64/e/epel-release-7-8.noarch.rpm
rpm -ivh epel-release-7-8.noarch.rpm
yum --enablerepo=epel -y install sshpass

ssh-keygen -t rsa

yum -y install pdsh-rcmd-ssh
```

The script also installs a `sshpass` utility which is used for passing the password to an ssh session.

¹ Please note, that the names and addresses used in this document are dummies and not actual addresses / names.

Setting up the remote logins, I use a little script (CopyHostFiles.sh):

```
#!/bin/bash
cat $1 | while read line
do
  host=`echo $line | gawk '{print $2}'`

  #Skip localhost
  if [ $host != "localhost" ]
  then
    #Skip utilities server
    if [ $host != "utilities" ]
    then
      if [ $host != $HOSTNAME ]
      then
        echo "Processing $host"
        ssh-keyscan $host >> ~/.ssh/known_hosts
        sshpass -p $2 ssh-copy-id -i /root/.ssh/id_rsa.pub $host
        scp /etc/hosts root@$host:/etc/hosts
      fi
    fi
  fi
done
```

This script can of course be expanded to incorporate the Java installation and setup and / or to accommodate other requirements.

Calling this using the following syntax will run through the host file on the server (omitting localhost, current file and the utilities server), add the remote system to known hosts, enable passwordless login and copy the host file to the system using the given password².

```
./CopyHostFiles.sh /etc/hosts [Password]
```

² In this case, the same root password is being used across the entire cluster. This is NOT best practices, but for a demo purpose, this will be ok on a closed environment.

Hadoop overview

Normally, Hadoop is used in huge server systems consisting of many computers. A typical system could look like this :

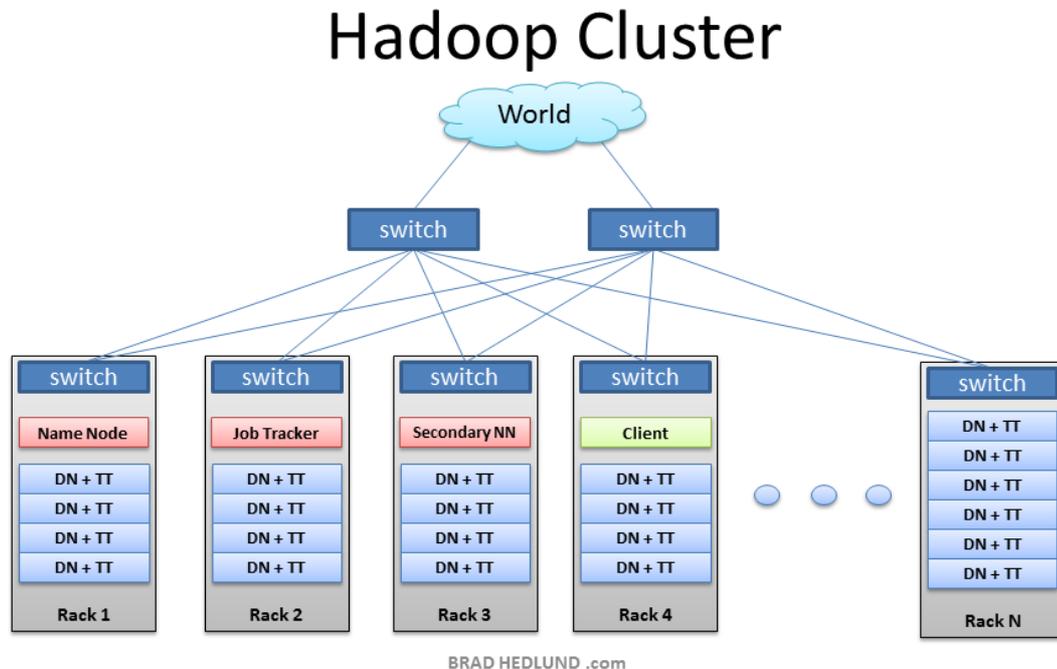


Illustration 1: Hadoop Cluster overview

Hadoop has its own filesystem called HDFS (Hadoop Distributed File System). It also is divided into 4 parts. The Master server contains 2 services called NameNode and JobTracker and is described in the official Hadoop Wiki as :

The **NameNode** is the centerpiece of an HDFS file system. It keeps the directory tree of all files in the file system, and tracks where across the cluster the file data is kept. It does not store the DN data of these files itself.

Client applications talk to the NameNode whenever they wish to locate a file, or when they want to add/copy/move/delete a file. The NameNode responds the successful requests by returning a list of relevant DataNode servers where the data lives.

The NameNode is a Single Point of Failure for the HDFS Cluster. HDFS is not currently a High Availability system. When the NameNode goes down, the file system goes offline. There is an optional SecondaryNameNode that can be hosted on a separate machine. It only creates checkpoints of the namespace by merging the edits file into the fsimage file and does not provide any real redundancy. Hadoop 0.21+ has a BackupNameNode that is part of a plan to have an HA name service, but it needs active contributions from the people who want it (i.e. you) to make it Highly Available.

Client applications talk to the NameNode whenever they wish to locate a file, or when they want to add/copy/move/delete a file. The NameNode responds the successful requests by returning a list of relevant DataNode servers where the data lives.

The NameNode is a Single Point of Failure for the HDFS Cluster. HDFS is not currently a High Availability system. When the NameNode goes down, the file system goes offline. There is an optional SecondaryNameNode that can be hosted on a separate machine. It only creates checkpoints of the namespace by merging the edits file into the fsimage file and does not provide any real redundancy. Hadoop 0.21+ has a BackupNameNode that is part of a plan to have an HA name service, but it needs active contributions from the people who want it (i.e. you) to make it Highly Available.

The **JobTracker** is the service within Hadoop that farms out MapReduce tasks to specific nodes in the cluster, ideally the nodes that have the data, or at least are in the same rack.

Client applications submit jobs to the Job tracker. The JobTracker talks to the NameNode to determine the location of the data. The JobTracker locates TaskTracker nodes with available slots at or near the data. The JobTracker submits the work to the chosen TaskTracker nodes. The TaskTracker nodes are monitored. If they do not submit heartbeat signals often enough, they are deemed to have failed and the work is scheduled on a different TaskTracker.

A TaskTracker will notify the JobTracker when a task fails. The JobTracker decides what to do then: it may resubmit the job elsewhere, it may mark that specific record as something to avoid, and it may even blacklist the TaskTracker as unreliable. When the work is completed, the JobTracker updates its status.

Each of the Clients holds the DataNode and the TaskTracker :

A **DataNode** stores data in the HDFS. A functional filesystem has more than one DataNode, with data replicated across them.

On startup, a DataNode connects to the NameNode; spinning until that service comes up. It then responds to requests from the NameNode for filesystem operations.

Client applications can talk directly to a DataNode, once the NameNode has provided the location of the data. Similarly, MapReduce operations farmed out to TaskTracker instances near a DataNode, talk directly to the DataNode to access the files. TaskTracker instances can, indeed should, be deployed on the same servers that host DataNode instances, so that MapReduce operations are performed close to the data.

DataNode instances can talk to each other, which is what they do when they are replicating data.

A **TaskTracker** is a node in the cluster that accepts tasks - Map, Reduce and Shuffle operations - from a JobTracker.

Every TaskTracker is configured with a set of slots, these indicate the number of tasks that it can accept. When the JobTracker tries to find somewhere to schedule a task within the MapReduce operations, it first looks for an empty slot on the same server that hosts the DataNode containing the data, and if not, it looks for an empty slot on a machine in the same rack.

The TaskTracker spawns a separate JVM processes to do the actual work; this is to ensure that process failure does not take down the task tracker. The TaskTracker monitors these spawned processes, capturing the output and exit codes. When the process finishes, successfully or not, the tracker notifies the

JobTracker. The TaskTrackers also send out heartbeat messages to the JobTracker, usually every few minutes, to reassure the JobTracker that it is still alive. These message also inform the JobTracker of the number of available slots, so the JobTracker can stay up to date with where in the cluster work can be delegated.

The whole relationship between the various part can be illustrated this way :

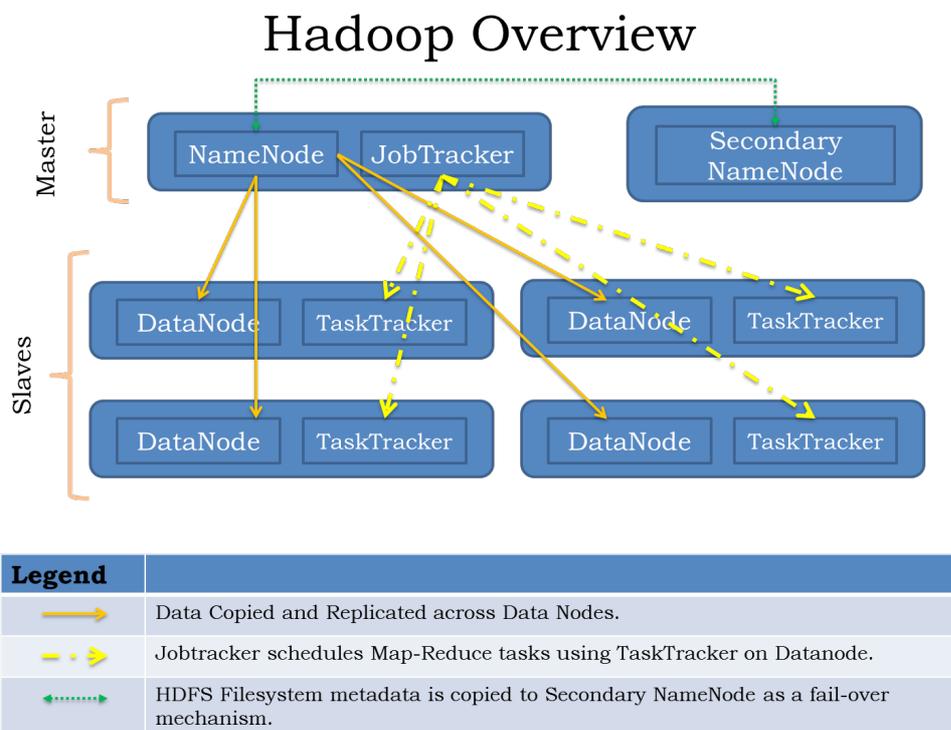


Illustration 2: Hadoop services relation

Normally, a failover strategy will be implemented having a Secondary NameNode, but we will not go into details about this.

The Hadoop eco system includes according to official wiki :

Ambari™	A web-based tool for provisioning, managing, and monitoring Apache Hadoop clusters which includes support for Hadoop HDFS, Hadoop MapReduce, Hive, HCatalog, HBase, ZooKeeper, Oozie, Pig and Sqoop. Ambari also provides a dashboard for viewing cluster health such as heatmaps and ability to view MapReduce, Pig and Hive applications visually alongwith features to diagnose their performance characteristics in a user-friendly manner.
Avro™	A data serialization system.
Cassandra™	A scalable multi-master database with no single points of failure.
Chukwa™	A data collection system for managing large distributed systems.
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.
Mahout™	A Scalable machine learning and data mining library.

Pig™	A high-level data-flow language and execution framework for parallel computation.
Spark™	A fast and general compute engine for Hadoop data. Spark provides a simple and expressive programming model that supports a wide range of applications, including ETL, machine learning, stream processing, and graph computation.
Tez™	A generalized data-flow programming framework, built on Hadoop YARN, which provides a powerful and flexible engine to execute an arbitrary DAG of tasks to process data for both batch and interactive use-cases. Tez is being adopted by Hive™, Pig™ and other frameworks in the Hadoop ecosystem, and also by other commercial software (e.g. ETL tools), to replace Hadoop™ MapReduce as the underlying execution engine.
ZooKeeper™	A high-performance coordination service for distributed applications.

For the setup, we'll use Ambari, Spark and ZooKeeper and setup 2 Master Nodes (1 primary and 1 secondary) and 6 Slave Nodes.

Ambari Installation

In our case, we only need the system for developing solutions for running on Hadoop clusters and also for doing some computations meaning that we will setup up a smaller version. So, we'll set up 2 Master Servers (Name Node) and 6 Workers (Data Nodes) plus a single machine for running the Ambari server.

Using the virtual machines, we now need to setup the systems. Provided we have the 9 virtual images as mentioned in the Server Environment section, we continue to install Ambari on them. PDSH enables us to run commands on other hosts. Hence, we run the InstallAmbari.sh script providing the host file as an argument :

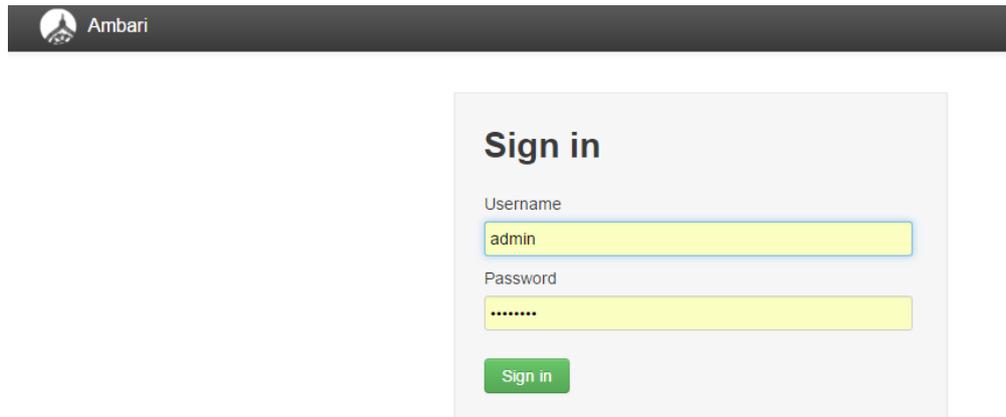
```
#!/bin/bash
cat $1 | while read line
do
    host=`echo $line | gawk '{print $2}'`

    #Skip localhost
    if [ $host != "localhost" ]
    then
        #Skip utilities server
        if [ $host != "utilities" ]
        then
            if [ $host != $HOSTNAME ]
            then
                echo "Processing $host"
                #Retrieve the Repository files
                pdsh -w $host "yum -y install wget"
                pdsh -w $host "wget http://public-repo-
1.hortonworks.com/ambari/centos7/2.x/updates/2.1.2.1/ambari.repo -O
/etc/yum.repos.d/ambari.repo" | sort
                #Install the Ambari Agent
                pdsh -w $host "yum -y install ambari-agent" | sort
                #Setup the Init file to point towards the ambari.cluster server
                pdsh -w $host "sed -i 's/hostname=localhost/hostname=ambari.cluster/g'
/etc/ambari-agent/conf/ambari-agent.ini"
                pdsh -w $host "sed -i 's/hostname=localhost/hostname=ambari.cluster/g'
/etc/ambari-agent/conf/ambari-agent.ini"
                pdsh -w $host "chkconfig ambari-agent on"
                pdsh -w $host "service ambari-agent restart"
                pdsh -w $host "systemctl stop firewalld"
                pdsh -w $host "systemctl disable firewalld"
            fi
        fi
    fi
done

wget http://public-repo-1.hortonworks.com/ambari/centos7/2.x/updates/2.1.2.1/ambari.repo
-O /etc/yum.repos.d/ambari.repo
yum -y install ambari-agent
yum -y install ambari-server
ambari-server setup -j /usr/lib/jvm/jre-1.8.0-openjdk
ambari-server start
systemctl stop firewalld
systemctl disable firewalld
```

Executing the script will take a while – a number of files must be downloaded.

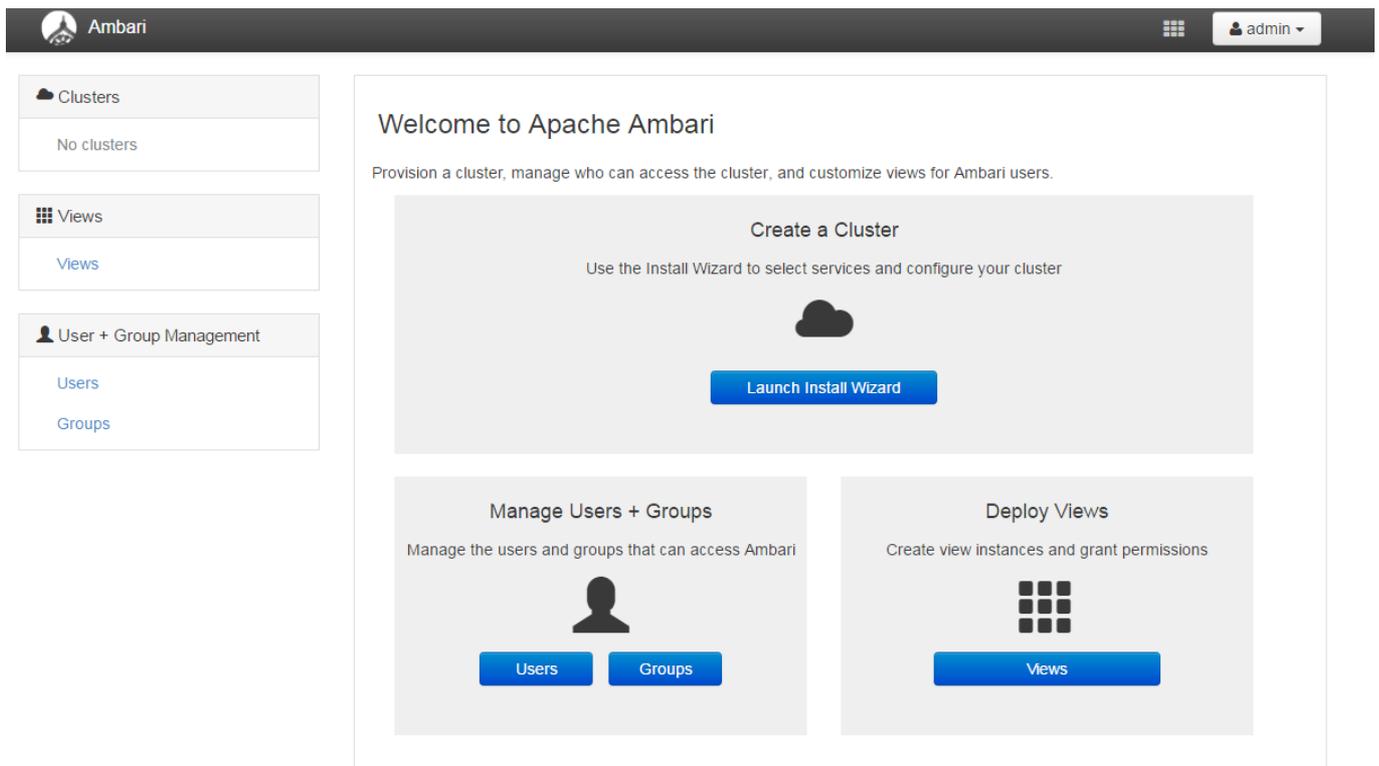
Having waited for the commands to finish the Ambari website should now be available in a browser on port 8080 :



The image shows the Ambari sign-in page. At the top left, there is a dark header with the Ambari logo and the word "Ambari". The main content area is a light gray box titled "Sign in". It contains two input fields: "Username" with the text "admin" and "Password" with masked characters ".....". Below the fields is a green "Sign in" button.

Illustration 3: Ambari, Initial login

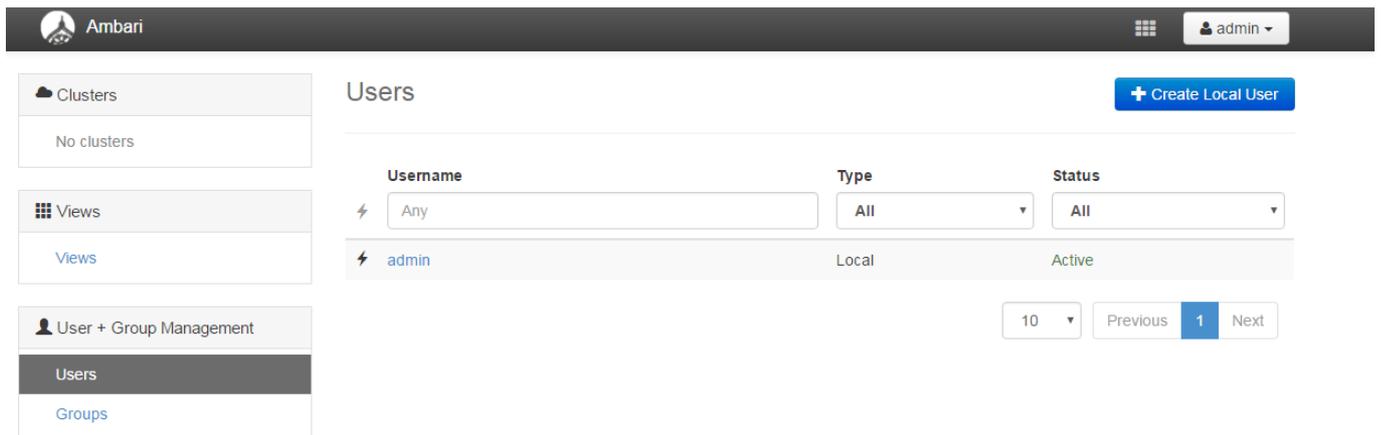
Login using "admin" as both password and username and change it before proceeding :



The image shows the Ambari welcome screen. At the top, there is a dark header with the Ambari logo, the word "Ambari", a grid icon, and a user dropdown menu showing "admin". On the left, there is a sidebar with three main sections: "Clusters" (No clusters), "Views" (Views), and "User + Group Management" (Users, Groups). The main content area is titled "Welcome to Apache Ambari" and contains the text "Provision a cluster, manage who can access the cluster, and customize views for Ambari users." Below this, there are three main action cards: "Create a Cluster" (with a "Launch Install Wizard" button), "Manage Users + Groups" (with "Users" and "Groups" buttons), and "Deploy Views" (with a "Views" button).

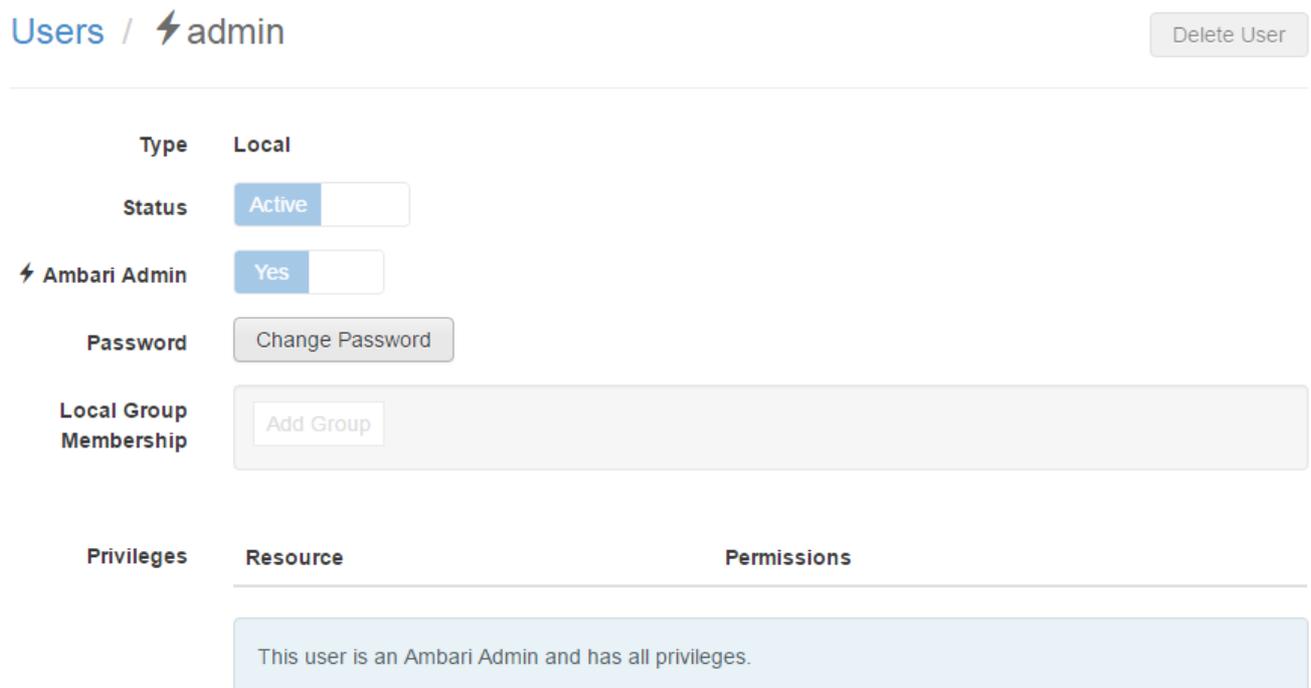
Illustration 4: Ambari, Initial welcome screen

Go to users and click the "admin" user :



The screenshot shows the Ambari interface with the 'Users' section selected in the left-hand navigation menu. The main content area displays a table of users. At the top right, there is a '+ Create Local User' button. The table has columns for Username, Type, and Status. The first row shows the 'admin' user with Type 'Local' and Status 'Active'. Below the table, there are pagination controls showing '10' items per page, with 'Previous', '1', and 'Next' buttons.

Illustration 5: Ambari, User management



The screenshot shows the 'Users / admin' editing page. At the top right, there is a 'Delete User' button. The page is divided into several sections:

- Type:** Local
- Status:** Active (with a dropdown menu)
- Ambari Admin:** Yes (with a dropdown menu)
- Password:** Change Password (button)
- Local Group Membership:** Add Group (button)
- Privileges:** A table with columns 'Resource' and 'Permissions'. A light blue box below the table contains the text: "This user is an Ambari Admin and has all privileges."

Illustration 6: Ambari, User editing

Go back to the frontpage by clicking the Ambari logo.

Now, we need to create a cluster. We click the "Launch Install Wizard" button in order to be guided through the cluster creation :

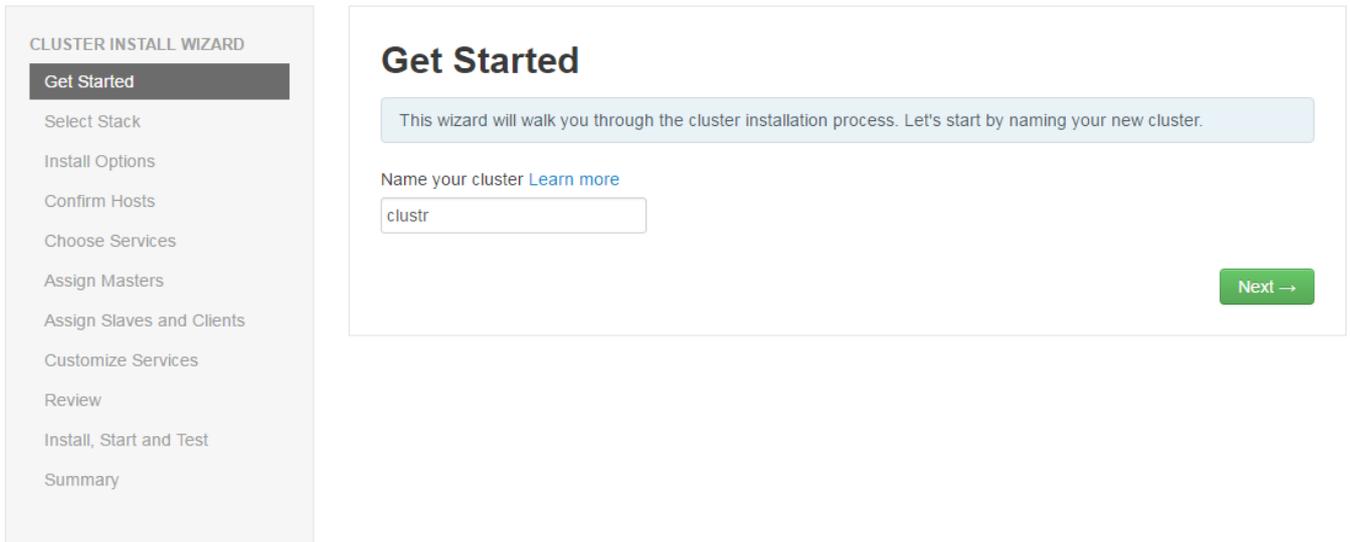


Illustration 7: Ambari Wizard, Name Cluster

First, we name the cluster (in this case "clustr" - not very inspiring).

Next, we select the Hadoop stack (version) to use. The current installation supports up to Hadoop 2.2, so we'll select that.

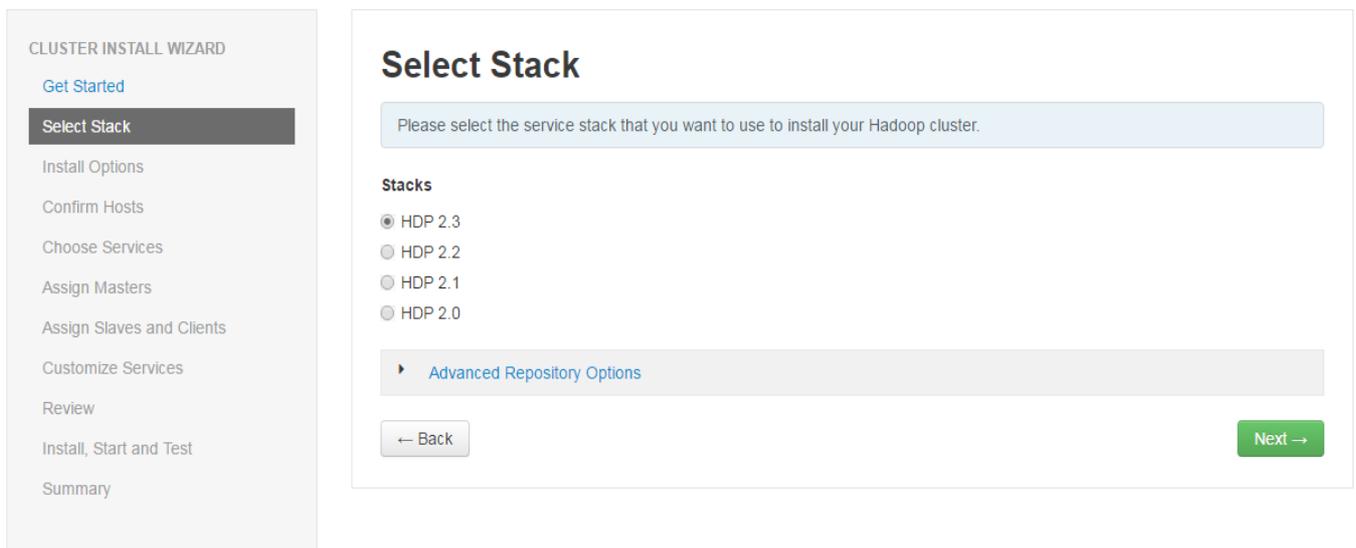
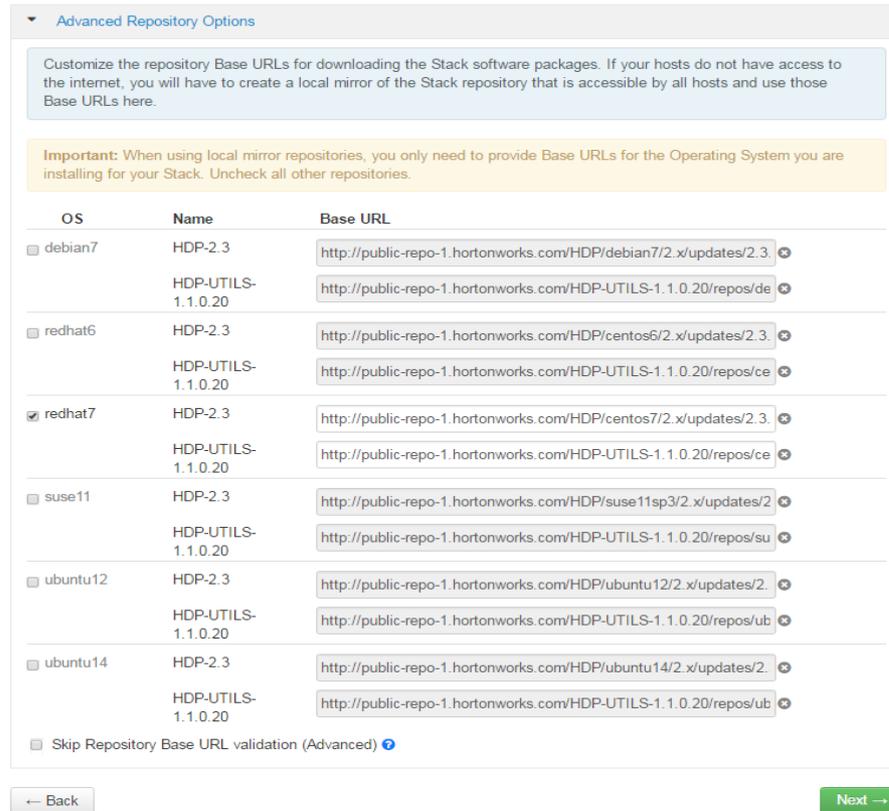


Illustration 8: Ambari Wizard, Stack selection.

In the Advanced Repository Options, I've disabled all but the RedHat7 repository, since we're using CentOS 7 as our operating system :



Advanced Repository Options

Customize the repository Base URLs for downloading the Stack software packages. If your hosts do not have access to the internet, you will have to create a local mirror of the Stack repository that is accessible by all hosts and use those Base URLs here.

Important: When using local mirror repositories, you only need to provide Base URLs for the Operating System you are installing for your Stack. Uncheck all other repositories.

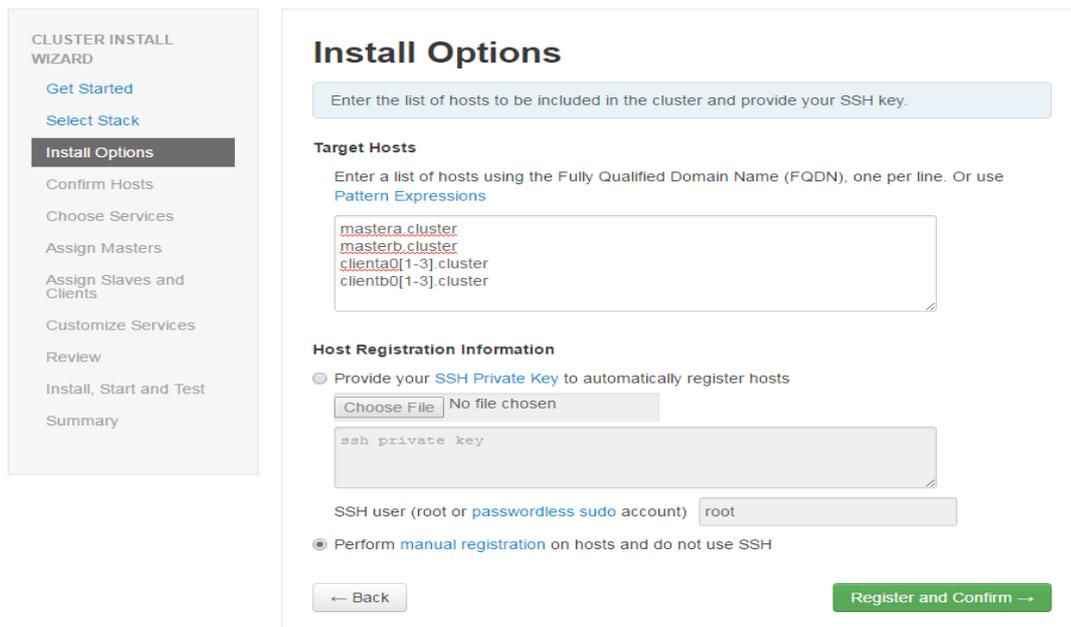
OS	Name	Base URL
<input type="checkbox"/> debian7	HDP-2.3	<input type="text" value="http://public-repo-1.hortonworks.com/HDP/debian7/2.x/updates/2.3."/>
	HDP-UTILS-1.1.0.20	<input type="text" value="http://public-repo-1.hortonworks.com/HDP-UTILS-1.1.0.20/repos/de"/>
<input type="checkbox"/> redhat6	HDP-2.3	<input type="text" value="http://public-repo-1.hortonworks.com/HDP/centos6/2.x/updates/2.3."/>
	HDP-UTILS-1.1.0.20	<input type="text" value="http://public-repo-1.hortonworks.com/HDP-UTILS-1.1.0.20/repos/ce"/>
<input checked="" type="checkbox"/> redhat7	HDP-2.3	<input type="text" value="http://public-repo-1.hortonworks.com/HDP/centos7/2.x/updates/2.3."/>
	HDP-UTILS-1.1.0.20	<input type="text" value="http://public-repo-1.hortonworks.com/HDP-UTILS-1.1.0.20/repos/ce"/>
<input type="checkbox"/> suse11	HDP-2.3	<input type="text" value="http://public-repo-1.hortonworks.com/HDP/suse11sp3/2.x/updates/2."/>
	HDP-UTILS-1.1.0.20	<input type="text" value="http://public-repo-1.hortonworks.com/HDP-UTILS-1.1.0.20/repos/su"/>
<input type="checkbox"/> ubuntu12	HDP-2.3	<input type="text" value="http://public-repo-1.hortonworks.com/HDP/ubuntu12/2.x/updates/2."/>
	HDP-UTILS-1.1.0.20	<input type="text" value="http://public-repo-1.hortonworks.com/HDP-UTILS-1.1.0.20/repos/ub"/>
<input type="checkbox"/> ubuntu14	HDP-2.3	<input type="text" value="http://public-repo-1.hortonworks.com/HDP/ubuntu14/2.x/updates/2."/>
	HDP-UTILS-1.1.0.20	<input type="text" value="http://public-repo-1.hortonworks.com/HDP-UTILS-1.1.0.20/repos/ub"/>

Skip Repository Base URL validation (Advanced)

← Back Next →

Illustration 9: Ambari Wizard, Repositories

Moving on, we have to select the hosts to install the services on.



CLUSTER INSTALL WIZARD

- Get Started
- Select Stack
- Install Options**
- Confirm Hosts
- Choose Services
- Assign Masters
- Assign Slaves and Clients
- Customize Services
- Review
- Install, Start and Test
- Summary

Install Options

Enter the list of hosts to be included in the cluster and provide your SSH key.

Target Hosts

Enter a list of hosts using the Fully Qualified Domain Name (FQDN), one per line. Or use [Pattern Expressions](#)

```

mastera.cluster
masterb.cluster
clienta[1-3].cluster
clientb[1-3].cluster
    
```

Host Registration Information

Provide your **SSH Private Key** to automatically register hosts

Choose File

ssh private key

SSH user (root or [passwordless sudo](#) account)

Perform [manual registration](#) on hosts and do not use SSH

← Back Register and Confirm →

Illustration 10: Ambari Wizard, Install Options

We select all the machines running the agents (besides the one hosting the web site). Please note, that the Fully Qualified Domain Names are used. Due to our setup, we also select manual registration. A verification of the host names is shown in order to validate the patterns added :

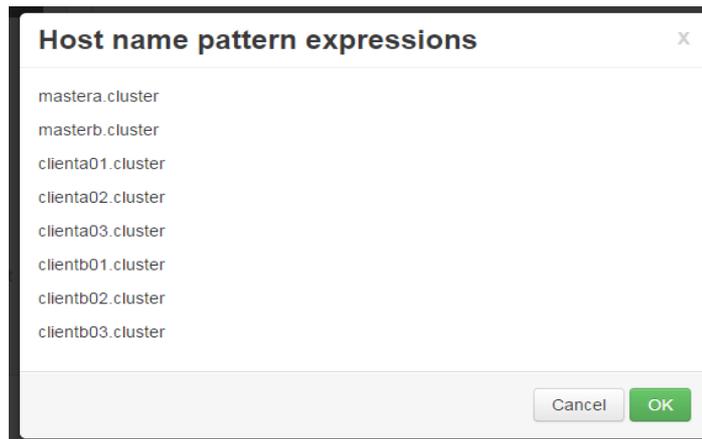


Illustration 11: Ambari Wizard, Pattern validation

Accept the pattern and messages and proceed for the registration process :

CLUSTER INSTALL WIZARD

- Get Started
- Select Stack
- Install Options
- Confirm Hosts**
- Choose Services
- Assign Masters
- Assign Slaves and Clients
- Customize Services
- Review
- Install, Start and Test
- Summary

Confirm Hosts

Registering your hosts.
Please confirm the host list and remove any hosts that you do not want to include in the cluster.

Remove Selected Show: **All (8)** | Installing (0) | Registering (8) | Success (0) | Fail (0)

	Host	Progress	Status	Action
<input type="checkbox"/>	mastera.cluster	<div style="width: 100%; height: 10px; background-color: #00a0e3;"></div>	Registering	Remove
<input type="checkbox"/>	masterb.cluster	<div style="width: 100%; height: 10px; background-color: #00a0e3;"></div>	Registering	Remove
<input type="checkbox"/>	clienta01.cluster	<div style="width: 100%; height: 10px; background-color: #00a0e3;"></div>	Registering	Remove
<input type="checkbox"/>	clienta02.cluster	<div style="width: 100%; height: 10px; background-color: #00a0e3;"></div>	Registering	Remove
<input type="checkbox"/>	clienta03.cluster	<div style="width: 100%; height: 10px; background-color: #00a0e3;"></div>	Registering	Remove
<input type="checkbox"/>	clientb01.cluster	<div style="width: 100%; height: 10px; background-color: #00a0e3;"></div>	Registering	Remove
<input type="checkbox"/>	clientb02.cluster	<div style="width: 100%; height: 10px; background-color: #00a0e3;"></div>	Registering	Remove
<input type="checkbox"/>	clientb03.cluster	<div style="width: 100%; height: 10px; background-color: #00a0e3;"></div>	Registering	Remove

Show: 25 1 - 8 of 8 [←](#) [→](#)

1 Other Registered Hosts

[← Back](#)
[Next →](#)

Illustration 12: Ambari Wizard, Confirm Hosts

Having performed the registration, we proceed to add services

Confirm Hosts

Registering your hosts.
Please confirm the host list and remove any hosts that you do not want to include in the cluster.

Remove Selected
Show: **All (8)** | [Installing \(0\)](#) | [Registering \(0\)](#) | [Success \(8\)](#) | [Fail \(0\)](#)

	Host	Progress	Status	Action
<input type="checkbox"/>	mastera.cluster	<div style="width: 100%; height: 10px; background-color: green;"></div>	Success	Remove
<input type="checkbox"/>	masterb.cluster	<div style="width: 100%; height: 10px; background-color: green;"></div>	Success	Remove
<input type="checkbox"/>	clinta01.cluster	<div style="width: 100%; height: 10px; background-color: green;"></div>	Success	Remove
<input type="checkbox"/>	clinta02.cluster	<div style="width: 100%; height: 10px; background-color: green;"></div>	Success	Remove
<input type="checkbox"/>	clinta03.cluster	<div style="width: 100%; height: 10px; background-color: green;"></div>	Success	Remove
<input type="checkbox"/>	clientb01.cluster	<div style="width: 100%; height: 10px; background-color: green;"></div>	Success	Remove
<input type="checkbox"/>	clientb02.cluster	<div style="width: 100%; height: 10px; background-color: green;"></div>	Success	Remove
<input type="checkbox"/>	clientb03.cluster	<div style="width: 100%; height: 10px; background-color: green;"></div>	Success	Remove

Show: 25 | 1 - 8 of 8

Some warnings were encountered while performing checks against the 8 registered hosts above [Click here to see the warnings.](#)

← Back
Next →

Illustration 13: Ambari Wizard, Registration Completed.

Clicking on the link will display the details. Since, we're using the CentOS 7, which uses the ntpdate service for synchronizing time, we ignore the warning :

Service Issues (1)
⚠

The following services should be up

Service	
ntpd	Not running on 8 hosts

Illustration 14: Service Issues

Proceed to select the services to run on the cluster. In our case, we wish to install a bare minimum Hadoop cluster, also to prevent the installation to timeout. Later, we will go back and add additional services to run on the cluster.

So for now, we will stick to :

- HDFS
- YARN + MapReduce 2
- ZooKeeper
- Ambari Metrics

CLUSTER INSTALL WIZARD

- [Get Started](#)
- [Select Stack](#)
- [Install Options](#)
- [Confirm Hosts](#)
- [Choose Services](#)
- [Assign Masters](#)
- [Assign Slaves and Clients](#)
- [Customize Services](#)
- [Review](#)
- [Install, Start and Test](#)
- [Summary](#)

Choose Services

Choose which services you want to install on your cluster.

Service	Version	Description
<input checked="" type="checkbox"/> HDFS	2.7.1.2.3	Apache Hadoop Distributed File System
<input checked="" type="checkbox"/> YARN + MapReduce2	2.7.1.2.3	Apache Hadoop NextGen MapReduce (YARN)
<input type="checkbox"/> Tez	0.7.0.2.3	Tez is the next generation Hadoop Query Processing framework written on top of YARN.
<input type="checkbox"/> Hive	1.2.1.2.3	Data warehouse system for ad-hoc queries & analysis of large datasets and table & storage management service
<input type="checkbox"/> HBase	1.1.1.2.3	A Non-relational distributed database, plus Phoenix, a high performance SQL layer for low latency applications.
<input type="checkbox"/> Pig	0.15.0.2.3	Scripting platform for analyzing large datasets
<input type="checkbox"/> Sqoop	1.4.6.2.3	Tool for transferring bulk data between Apache Hadoop and structured data stores such as relational databases
<input type="checkbox"/> Oozie	4.2.0.2.3	System for workflow coordination and execution of Apache Hadoop jobs. This also includes the installation of the optional Oozie Web Console which relies on and will install the ExtJS Library.
<input checked="" type="checkbox"/> ZooKeeper	3.4.6.2.3	Centralized service which provides highly reliable distributed coordination
<input type="checkbox"/> Falcon	0.6.1.2.3	Data management and processing platform
<input type="checkbox"/> Storm	0.10.0	Apache Hadoop Stream processing framework
<input type="checkbox"/> Flume	1.5.2.2.3	A distributed service for collecting, aggregating, and moving large amounts of streaming data into HDFS
<input type="checkbox"/> Accumulo	1.7.0.2.3	Robust, scalable, high performance distributed key/value store.
<input type="checkbox"/> Ambari Metrics	0.1.0	A system for metrics collection that provides storage and retrieval capability for metrics collected from the cluster
<input type="checkbox"/> Atlas	0.5.0.2.3	Atlas Metadata and Governance platform
<input type="checkbox"/> Kafka	0.8.2.2.3	A high-throughput distributed messaging system
<input type="checkbox"/> Knox	0.6.0.2.3	Provides a single point of authentication and access for Apache Hadoop services in a cluster
<input type="checkbox"/> Mahout	0.9.0.2.3	Project of the Apache Software Foundation to produce free implementations of distributed or otherwise scalable machine learning algorithms focused primarily in the areas of collaborative filtering, clustering and classification
<input type="checkbox"/> Slider	0.80.0.2.3	A framework for deploying, managing and monitoring existing distributed applications on YARN.
<input type="checkbox"/> Spark	1.4.1.2.3	Apache Spark is a fast and general engine for large-scale data processing.

← Back
Next →

Illustration 15: Initial Service Selection

Assigning the Masters of the system. We will attempt to put the load evenly on the 2 masters (mastera and masterb).

CLUSTER INSTALL WIZARD

- [Get Started](#)
- [Select Stack](#)
- [Install Options](#)
- [Confirm Hosts](#)
- [Choose Services](#)
- Assign Masters**
- [Assign Slaves and Clients](#)
- [Customize Services](#)
- [Review](#)
- [Install, Start and Test](#)
- [Summary](#)

Assign Masters

Assign master components to hosts you want to run them on.

SNameNode:

NameNode:

History Server:

App Timeline Server:

ResourceManager:

ZooKeeper Server:

ZooKeeper Server:

Metrics Collector:

mastera.cluster (1.8 GB, 1 cores)

NameNode History Server ResourceManager

ZooKeeper Server

masterb.cluster (1.8 GB, 1 cores)

SNameNode App Timeline Server

ZooKeeper Server Metrics Collector

← Back
Next →

Illustration 16: Ambari Wizard, Assign Masters

Next, we setup the clients :

CLUSTER INSTALL WIZARD

- [Get Started](#)
- [Select Stack](#)
- [Install Options](#)
- [Confirm Hosts](#)
- [Choose Services](#)
- Assign Slaves and Clients**
- [Customize Services](#)
- [Review](#)
- [Install, Start and Test](#)
- [Summary](#)

Assign Slaves and Clients

Assign slave and client components to hosts you want to run them on.
Hosts that are assigned master components are shown with *.
"Client" will install HDFS Client, MapReduce2 Client, YARN Client and ZooKeeper Client.

Host	all none		all none		all none		all none	
mastera.cluster #	<input type="checkbox"/>	DataNode	<input type="checkbox"/>	NFSGateway	<input type="checkbox"/>	NodeManager	<input type="checkbox"/>	Client
masterb.cluster #	<input type="checkbox"/>	DataNode	<input type="checkbox"/>	NFSGateway	<input type="checkbox"/>	NodeManager	<input type="checkbox"/>	Client
clienta01.cluster	<input checked="" type="checkbox"/>	DataNode	<input type="checkbox"/>	NFSGateway	<input checked="" type="checkbox"/>	NodeManager	<input checked="" type="checkbox"/>	Client
clienta02.cluster	<input checked="" type="checkbox"/>	DataNode	<input type="checkbox"/>	NFSGateway	<input checked="" type="checkbox"/>	NodeManager	<input checked="" type="checkbox"/>	Client
clienta03.cluster	<input checked="" type="checkbox"/>	DataNode	<input type="checkbox"/>	NFSGateway	<input checked="" type="checkbox"/>	NodeManager	<input checked="" type="checkbox"/>	Client
clientb01.cluster	<input checked="" type="checkbox"/>	DataNode	<input type="checkbox"/>	NFSGateway	<input checked="" type="checkbox"/>	NodeManager	<input checked="" type="checkbox"/>	Client
clientb02.cluster	<input checked="" type="checkbox"/>	DataNode	<input type="checkbox"/>	NFSGateway	<input checked="" type="checkbox"/>	NodeManager	<input checked="" type="checkbox"/>	Client
clientb03.cluster	<input checked="" type="checkbox"/>	DataNode	<input type="checkbox"/>	NFSGateway	<input checked="" type="checkbox"/>	NodeManager	<input checked="" type="checkbox"/>	Client

Show:
1 - 8 of 8

← Back
Next →

Illustration 17: Ambari Wizard, Assign Slaves and clients

We select all other machines, leaving out the NFSGateway, since we currently do not use it.

Now, we move on to customization of the services. According to the UI, all configurations have been addressed :

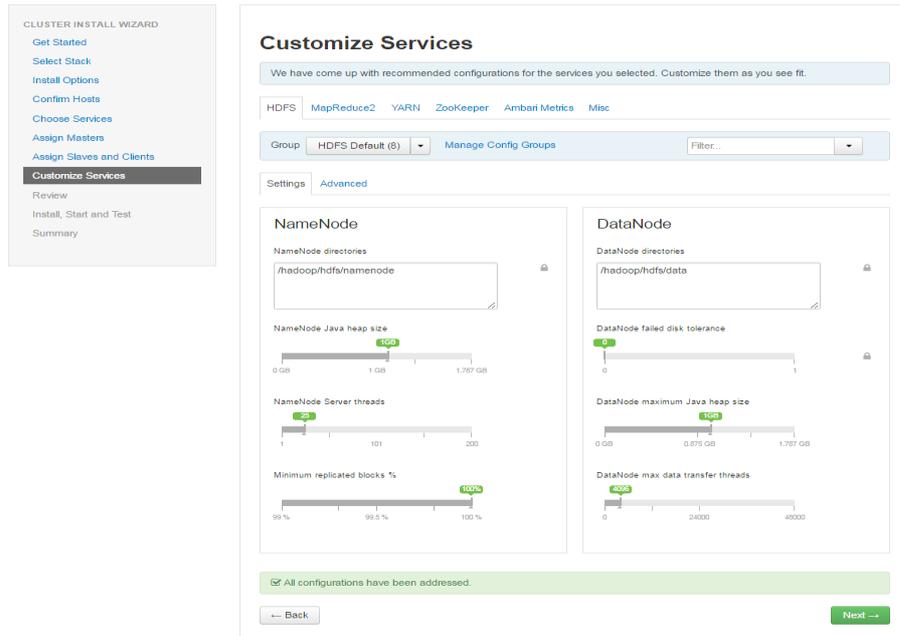


Illustration 18: Ambari Wizard

We can now review our selections and proceed, if we're OK with the settings:

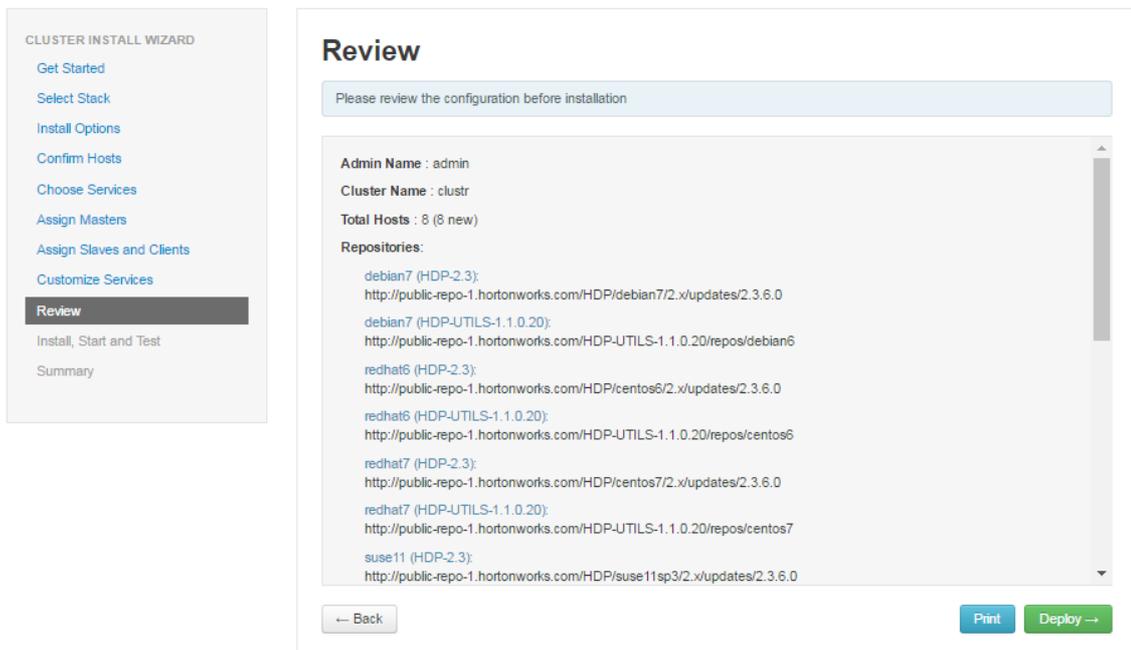


Illustration 19: Installation Review

Upon reviewing the settings, we proceed to start installing the cluster.

Install, Start and Test

Please wait while the selected services are installed and started.

Show: **All (8)** | [In Progress \(8\)](#) | [Warning](#)

Host	Status	Message
clienta01.cluster	<div style="width: 3%;"><div style="width: 3%;"></div></div> 3%	Waiting to install HBase
clienta02.cluster	<div style="width: 3%;"><div style="width: 3%;"></div></div> 3%	Waiting to install HBase
clienta03.cluster	<div style="width: 3%;"><div style="width: 3%;"></div></div> 3%	Waiting to install HBase
clientb01.cluster	<div style="width: 3%;"><div style="width: 3%;"></div></div> 3%	Waiting to install HBase
clientb02.cluster	<div style="width: 3%;"><div style="width: 3%;"></div></div> 3%	Waiting to install HBase
clientb03.cluster	<div style="width: 3%;"><div style="width: 3%;"></div></div> 3%	Waiting to install HBase
mastera.cluster	<div style="width: 3%;"><div style="width: 3%;"></div></div> 3%	Waiting to install Tez Cli
masterb.cluster	<div style="width: 3%;"><div style="width: 3%;"></div></div> 3%	Waiting to install HBase

8 of 8 hosts showing - [Show All](#) Show: 25 1

Illustration 20: Ambari Wizard, Installation

This will take some time based on the internet connection, hardware etc. Having ended the installation and starting of the services (this may fail due to timeouts. But don't worry. Just restart a couple of times until all services startup).

CLUSTER INSTALL WIZARD

- Get Started
- Select Slack
- Install Options
- Confirm Hosts
- Choose Services
- Assign Masters
- Assign Slaves and Clients
- Customize Services
- Review
- Install, Start and Test
- Summary

Install, Start and Test

Please wait while the selected services are installed and started.

100 % overall

Show: **All (8)** | [In Progress \(0\)](#) | [Warning \(0\)](#) | [Success \(8\)](#) | [Fail \(0\)](#)

Host	Status	Message
mastera.cluster	<div style="width: 100%;"><div style="width: 100%;"></div></div> 100%	Success
masterb.cluster	<div style="width: 100%;"><div style="width: 100%;"></div></div> 100%	Success
clienta01.cluster	<div style="width: 100%;"><div style="width: 100%;"></div></div> 100%	Success
clienta02.cluster	<div style="width: 100%;"><div style="width: 100%;"></div></div> 100%	Success
clienta03.cluster	<div style="width: 100%;"><div style="width: 100%;"></div></div> 100%	Success
clientb01.cluster	<div style="width: 100%;"><div style="width: 100%;"></div></div> 100%	Success
clientb02.cluster	<div style="width: 100%;"><div style="width: 100%;"></div></div> 100%	Success
clientb03.cluster	<div style="width: 100%;"><div style="width: 100%;"></div></div> 100%	Success

8 of 8 hosts showing - [Show All](#) Show: 25 1 - 8 of 8

Successfully installed and started the services.

[Next →](#)

Illustration 21: End of Service installation

After the installation process has ended, go on to get the summary of the installation :

Reading the summary, we are informed, that everything worked out fine.

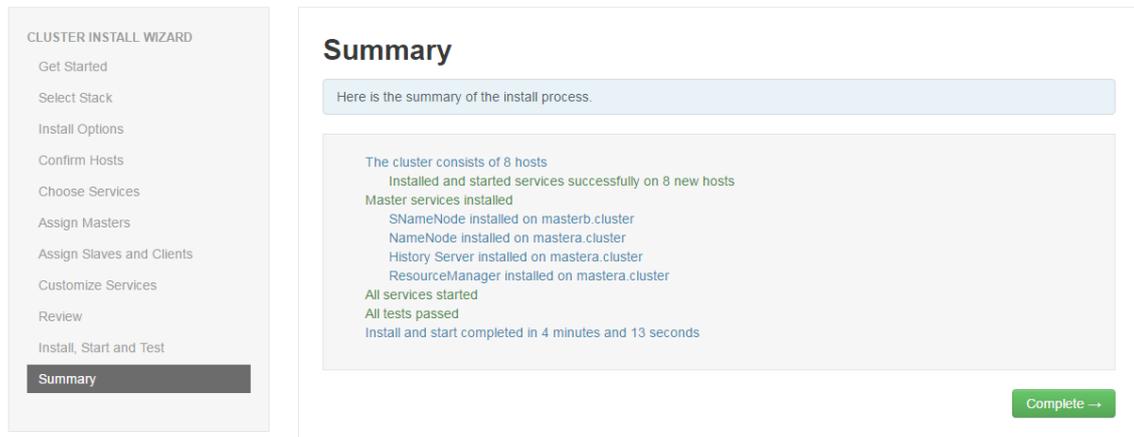


Illustration 22: Installation Summary

Finally, we get an updated frontpage showing us the status of the cluster :

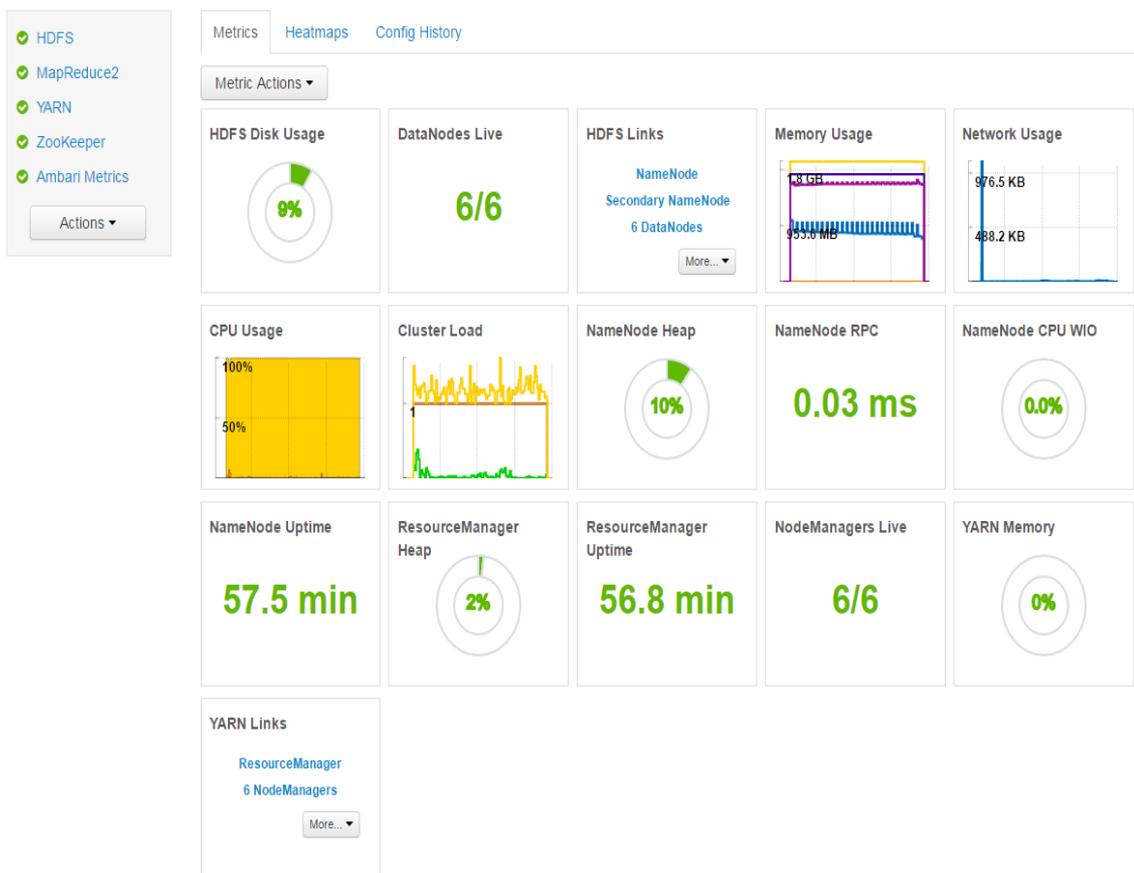


Illustration 23: Cluster control Panel

Now, we should start adding more services to the cluster.

Adding Spark

When installing the cluster, we did not install Spark right away. This would have been the easiest. However, given the internet connection it would result in a time out due to the long time downloading packages³.

On the frontpage, select Add Service in the Actions menu

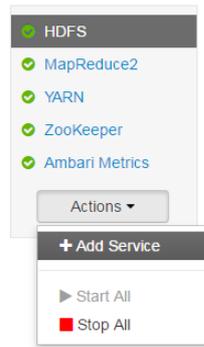


Illustration 24: Add Service

Scroll down and select Spark :



Illustration 25: Select Spark

We assign the mastera to run the Spark History server :

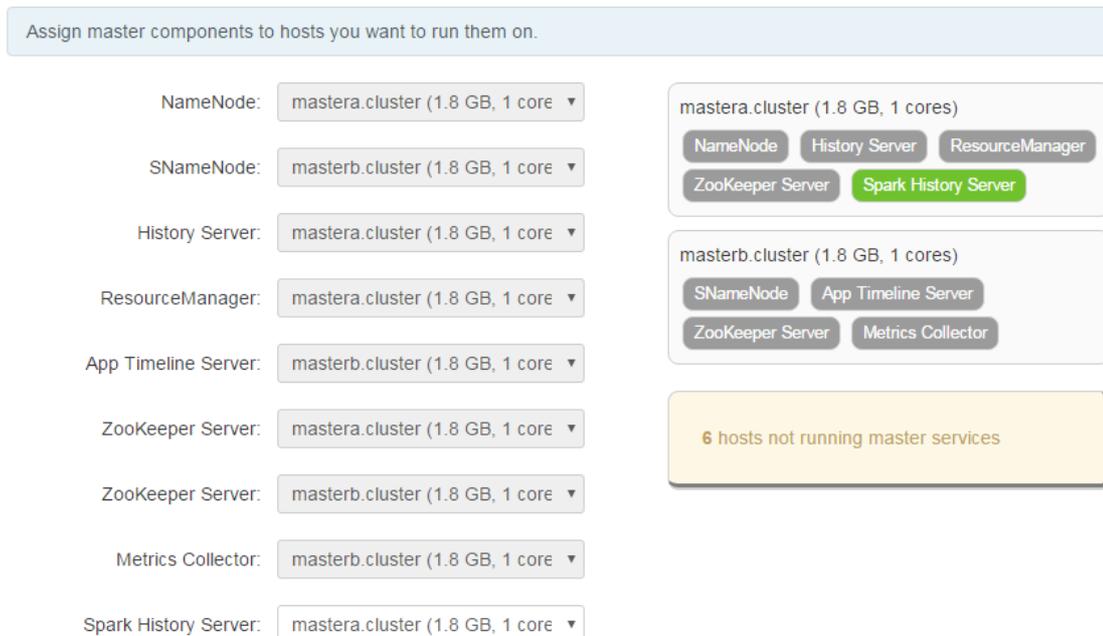


Illustration 26: Select Spark History server

³ There are ways around this. E.g. creating a local repository of packages. This, however, is beyond the scope of this document.

Select all (except masters) as clients :

ADD SERVICE WIZARD

- [Choose Services](#)
- [Assign Masters](#)
- [Assign Slaves and Clients](#)
- [Customize Services](#)
- [Configure Identities](#)
- [Review](#)
- [Install, Start and Test](#)
- [Summary](#)

Assign Slaves and Clients

Assign slave and client components to hosts you want to run them on.
Hosts that are assigned master components are shown with *.
"Client" will install Spark Client

Host	all none	all none	all none	all none
mastera.cluster *	<input type="checkbox"/> DataNode	<input type="checkbox"/> NFSGateway	<input type="checkbox"/> NodeManager	<input type="checkbox"/> Client
masterb.cluster *	<input type="checkbox"/> DataNode	<input type="checkbox"/> NFSGateway	<input type="checkbox"/> NodeManager	<input type="checkbox"/> Client
clienta01.cluster	<input checked="" type="checkbox"/> DataNode	<input type="checkbox"/> NFSGateway	<input checked="" type="checkbox"/> NodeManager	<input checked="" type="checkbox"/> Client
clienta02.cluster	<input checked="" type="checkbox"/> DataNode	<input type="checkbox"/> NFSGateway	<input checked="" type="checkbox"/> NodeManager	<input checked="" type="checkbox"/> Client
clienta03.cluster	<input checked="" type="checkbox"/> DataNode	<input type="checkbox"/> NFSGateway	<input checked="" type="checkbox"/> NodeManager	<input checked="" type="checkbox"/> Client
clientb01.cluster	<input checked="" type="checkbox"/> DataNode	<input type="checkbox"/> NFSGateway	<input checked="" type="checkbox"/> NodeManager	<input checked="" type="checkbox"/> Client
clientb02.cluster	<input checked="" type="checkbox"/> DataNode	<input type="checkbox"/> NFSGateway	<input checked="" type="checkbox"/> NodeManager	<input checked="" type="checkbox"/> Client
clientb03.cluster	<input checked="" type="checkbox"/> DataNode	<input type="checkbox"/> NFSGateway	<input checked="" type="checkbox"/> NodeManager	<input checked="" type="checkbox"/> Client

Show: 25 1 - 8 of 8

← Back
Next →

Illustration 27: Assign Spark Clients

I'll go for the default configuration :

Customize Services

- [Configure Identities](#)
- [Review](#)
- [Install, Start and Test](#)
- [Summary](#)

HDFS MapReduce2 YARN ZooKeeper Ambari Metrics **Spark** Misc

Group: Spark Default (8) Manage Config Groups Filter...

- ▶ [Advanced spark-defaults](#)
- ▶ [Advanced spark-env](#)
- ▶ [Advanced spark-javaopts-properties](#)
- ▶ [Advanced spark-log4j-properties](#)
- ▶ [Advanced spark-metrics-properties](#)
- ▶ [Custom spark-defaults](#)
- ▶ [Custom spark-javaopts-properties](#)
- ▶ [Custom spark-log4j-properties](#)
- ▶ [Custom spark-metrics-properties](#)

All configurations have been addressed.

← Back
Next →

Illustration 28: Customize Spark Configuration

Review the configuration and start the deployment process :

ADD SERVICE WIZARD

- [Choose Services](#)
- [Assign Masters](#)
- [Assign Slaves and Clients](#)
- [Customize Services](#)
- [Configure Identities](#)
- Review**
- [Install, Start and Test](#)
- [Summary](#)

Review

Please review the configuration before installation

Admin Name : admin

Cluster Name : clustr

Total Hosts : 8 (0 new)

Repositories:

```

debian7 (HDP-2.3):
http://public-repo-1.hortonworks.com/HDP/debian7/2.x/updates/2.3.6.0
debian7 (HDP-UTILS-1.1.0.20):
http://public-repo-1.hortonworks.com/HDP-UTILS-1.1.0.20/repos/debian6
redhat6 (HDP-2.3):
http://public-repo-1.hortonworks.com/HDP/centos6/2.x/updates/2.3.6.0
redhat6 (HDP-UTILS-1.1.0.20):
http://public-repo-1.hortonworks.com/HDP-UTILS-1.1.0.20/repos/centos6
redhat7 (HDP-2.3):
http://public-repo-1.hortonworks.com/HDP/centos7/2.x/updates/2.3.6.0
redhat7 (HDP-UTILS-1.1.0.20):
http://public-repo-1.hortonworks.com/HDP-UTILS-1.1.0.20/repos/centos7
suse11 (HDP-2.3):
http://public-repo-1.hortonworks.com/HDP/suse11sp3/2.x/updates/2.3.6.0
                    
```

← Back
Print
Deploy →

Illustration 29: Review Spark Configuration

Wait for the deployment to end :

ADD SERVICE WIZARD

- [Choose Services](#)
- [Assign Masters](#)
- [Assign Slaves and Clients](#)
- [Customize Services](#)
- [Configure Identities](#)
- [Review](#)
- Install, Start and Test**
- [Summary](#)

Install, Start and Test

Please wait while the selected services are installed and started.

7 % overall

Host	Status	Message
clienta01.cluster	<div style="width: 3%; height: 10px; background-color: #00a0e3;"></div> 3%	Waiting to install Spark Client
clienta02.cluster	<div style="width: 3%; height: 10px; background-color: #00a0e3;"></div> 3%	Waiting to install Spark Client
clienta03.cluster	<div style="width: 3%; height: 10px; background-color: #00a0e3;"></div> 3%	Waiting to install Spark Client
clientb01.cluster	<div style="width: 3%; height: 10px; background-color: #00a0e3;"></div> 3%	Waiting to install Spark Client
clientb02.cluster	<div style="width: 3%; height: 10px; background-color: #00a0e3;"></div> 3%	Waiting to install Spark Client
clientb03.cluster	<div style="width: 12%; height: 10px; background-color: #00a0e3;"></div> 12%	Installing Spark Client
mastera.cluster	<div style="width: 3%; height: 10px; background-color: #00a0e3;"></div> 3%	Waiting to install Spark History Server
masterb.cluster	<div style="width: 33%; height: 10px; background-color: #00a0e3;"></div> 33%	Install complete (Waiting to start)

8 of 8 hosts showing - [Show All](#)
Show: 1 - 8 of 8 ⏪ ⏩

Next →

Illustration 30: Spark Installation and Initialization

This will take a while.

ADD SERVICE WIZARD

- Choose Services
- Assign Masters
- Assign Slaves and Clients
- Customize Services
- Configure Identities
- Review
- Install, Start and Test
- Summary

Install, Start and Test

Please wait while the selected services are installed and started.

100 % overall

Host	Status	Message
clienta01.cluster	<div style="width: 100%; height: 10px; background-color: #4CAF50;"></div> 100%	Success
clienta02.cluster	<div style="width: 100%; height: 10px; background-color: #4CAF50;"></div> 100%	Success
clienta03.cluster	<div style="width: 100%; height: 10px; background-color: #4CAF50;"></div> 100%	Success
clientb01.cluster	<div style="width: 100%; height: 10px; background-color: #4CAF50;"></div> 100%	Success
clientb02.cluster	<div style="width: 100%; height: 10px; background-color: #4CAF50;"></div> 100%	Success
clientb03.cluster	<div style="width: 100%; height: 10px; background-color: #4CAF50;"></div> 100%	Success
mastera.cluster	<div style="width: 100%; height: 10px; background-color: #4CAF50;"></div> 100%	Success
masterb.cluster	<div style="width: 100%; height: 10px; background-color: #4CAF50;"></div> 100%	Success

8 of 8 hosts showing - [Show All](#) Show: 25 1 - 8 of 8

Successfully installed and started the services.

Next →

Illustration 31: End of installation

Having completed the installation, the Summary informs us, that certain services need to be restarted

Add Service Wizard

ADD SERVICE WIZARD

- Choose Services
- Assign Masters
- Assign Slaves and Clients
- Customize Services
- Configure Identities
- Review
- Install, Start and Test
- Summary

Summary

Important: You may also need to restart other services for the newly added services to function properly (for example, HDFS and YARN/MapReduce need to be restarted after adding Oozie). After closing this wizard, please restart all services that have the restart indicator ↻ next to the service name.

Here is the summary of the install process.

The cluster consists of 8 hosts
 Installed and started services successfully on 8 new hosts
 Install and start completed in 27 minutes and 33 seconds

Complete →

Illustration 32: Spark Installation Summary

Having installed the Spark service, we are once more presented with the frontpage. The only change is now another service

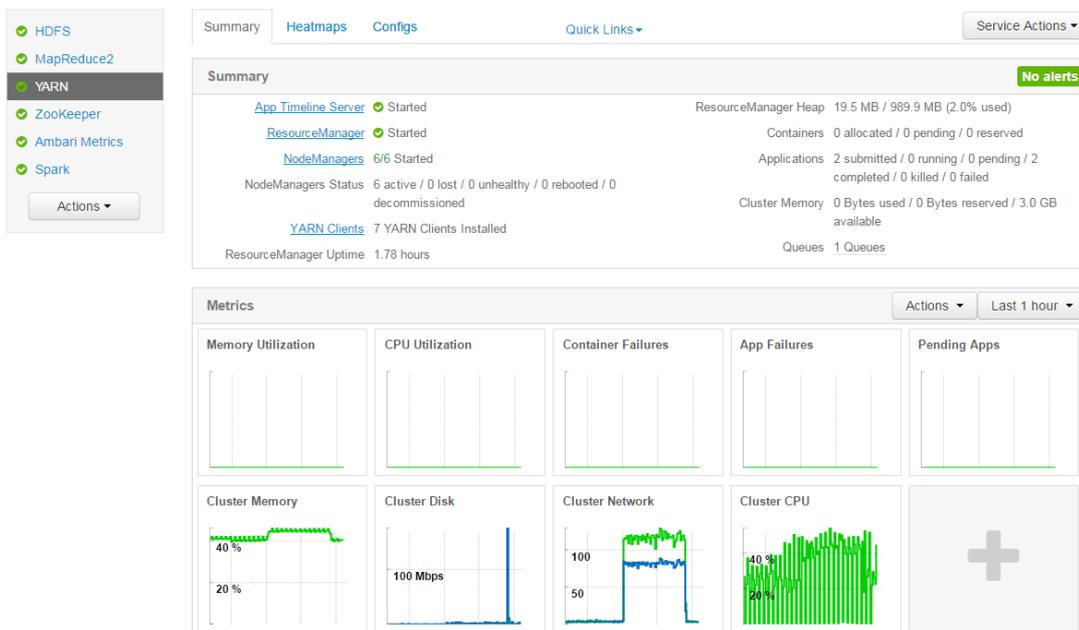


Illustration 33: Spark installation Frontpage

Test

In order to test our Spark installation, we use the demo "Spark Pi" application.

Login into mastera and run the following code :

```
cd /usr/hdp/current/spark-client

su spark

./bin/spark-submit --class org.apache.spark.examples.SparkPi --master yarn-client --num-executors 1 --driver-memory 512m --executor-memory 512m --executor-cores 1 lib/spark-examples*.jar 10
```

This will soon fill the whole screen with a lot of log messages.

```
16/11/07 00:30:52 INFO YarnClientSchedulerBackend: Application application_1478474972019_0001 has started running.
16/11/07 00:30:52 INFO Utils: Successfully started service 'org.apache.spark.network.netty.NettyBlockTransferService' on port 38176.
16/11/07 00:30:52 INFO NettyBlockTransferService: Server created on 38176
16/11/07 00:30:52 INFO BlockManagerMaster: Trying to register BlockManager
16/11/07 00:30:52 INFO BlockManagerMasterEndpoint: Registering block manager 192.168.1.120:38176 with 267.3 MB RAM, BlockManagerId(driver, 192.168.1.120, 38176)
16/11/07 00:30:52 INFO BlockManagerMaster: Registered BlockManager
16/11/07 00:30:52 INFO YarnClientSchedulerBackend: SchedulerBackend is ready for scheduling beginning after waiting maxRegisteredResourcesWaitingTime: 30000(ms)
16/11/07 00:30:52 INFO YarnHistoryService: Application started: SparkListenerApplicationStart(Spark Pi, Some(application_1478474972019_0001), 147847498657, spark_None, None)
16/11/07 00:30:52 INFO YarnHistoryService: About to POST entity application_1478474972019_0001 with 3 events to timeline service http://masterb.cluster:8188/ws/v1/timeline/
16/11/07 00:30:52 INFO SparkContext: Starting job: reduce at SparkPi.scala:36
16/11/07 00:30:52 INFO DAGScheduler: Got job 0 (reduce at SparkPi.scala:36) with 10 output partitions
16/11/07 00:30:52 INFO DAGScheduler: Final stage: ResultStage 0 (reduce at SparkPi.scala:36)
16/11/07 00:30:52 INFO DAGScheduler: Parents of final stage: List()
16/11/07 00:30:52 INFO DAGScheduler: Missing parents: List()
16/11/07 00:30:52 INFO DAGScheduler: Submitting ResultStage 0 (MapPartitionsRDD[1] at map at SparkPi.scala:32), which has no missing parents
16/11/07 00:30:53 INFO MemoryStore: ensureFreeSpace(1083) called with curMem=0, maxMem=280248975
16/11/07 00:30:53 INFO MemoryStore: Block broadcast_0 stored as value in memory (estimated size 1888.0 B, free 267.3 MB)
16/11/07 00:30:53 INFO MemoryStore: ensureFreeSpace(1202) called with curMem=1888, maxMem=280248975
16/11/07 00:30:53 INFO MemoryStore: Block broadcast_0_piece0 stored as bytes in memory (estimated size 1202.0 B, free 267.3 MB)
16/11/07 00:30:53 INFO BlockManagerInfo: Added broadcast_0_piece0 in memory on 192.168.1.120:38176 (size: 1202.0 B, free: 267.3 MB)
16/11/07 00:30:53 INFO SparkContext: Created broadcast 0 from broadcast at DAGScheduler.scala:161
16/11/07 00:30:53 INFO DAGScheduler: Submitting 10 missing tasks from ResultStage 0 (MapPartitionsRDD[1] at map at SparkPi.scala:32)
16/11/07 00:30:53 INFO YarnScheduler: Adding task set 0.0 with 10 tasks
16/11/07 00:30:57 INFO YarnClientSchedulerBackend: Registered executor: AkkaRpcEndpointRef(Actor[AkkaTcp://sparkExecutor@clientb02.cluster:42895/user/Executor@f883927924] with ID 1)
16/11/07 00:30:57 INFO TaskSetManager: Starting task 0.0 in stage 0.0 (TID 0, clientb02.cluster, PROCESS_LOCAL, 2189 bytes)
16/11/07 00:30:58 INFO BlockManagerMasterEndpoint: Registering block manager clientb02.cluster:37918 with 267.3 MB RAM, BlockManagerId(1, clientb02.cluster, 37918)
16/11/07 00:31:07 INFO BlockManagerInfo: Added broadcast_0_piece0 in memory on clientb02.cluster:37918 (size: 1202.0 B, free: 267.3 MB)
16/11/07 00:31:08 INFO TaskSetManager: Starting task 1.0 in stage 0.0 (TID 1, clientb02.cluster, PROCESS_LOCAL, 2189 bytes)
16/11/07 00:31:08 INFO TaskSetManager: Starting task 2.0 in stage 0.0 (TID 2, clientb02.cluster, PROCESS_LOCAL, 2189 bytes)
16/11/07 00:31:08 INFO TaskSetManager: Starting task 3.0 in stage 0.0 (TID 3, clientb02.cluster, PROCESS_LOCAL, 2189 bytes)
16/11/07 00:31:08 INFO TaskSetManager: Starting task 4.0 in stage 0.0 (TID 4, clientb02.cluster, PROCESS_LOCAL, 2189 bytes)
16/11/07 00:31:08 INFO TaskSetManager: Starting task 5.0 in stage 0.0 (TID 5, clientb02.cluster, PROCESS_LOCAL, 2189 bytes)
16/11/07 00:31:08 INFO TaskSetManager: Starting task 6.0 in stage 0.0 (TID 6, clientb02.cluster, PROCESS_LOCAL, 2189 bytes)
16/11/07 00:31:08 INFO TaskSetManager: Starting task 7.0 in stage 0.0 (TID 7, clientb02.cluster, PROCESS_LOCAL, 2189 bytes)
16/11/07 00:31:08 INFO TaskSetManager: Starting task 8.0 in stage 0.0 (TID 8, clientb02.cluster, PROCESS_LOCAL, 2189 bytes)
16/11/07 00:31:08 INFO TaskSetManager: Starting task 9.0 in stage 0.0 (TID 9, clientb02.cluster, PROCESS_LOCAL, 2189 bytes)
16/11/07 00:31:08 INFO YarnHistoryService: About to POST entity application_1478474972019_0001 with 10 events to timeline service http://masterb.cluster:8188/ws/v1/timeline/
16/11/07 00:31:08 INFO TaskSetManager: Finished task 2.0 in stage 0.0 (TID 2) in 206 ms on clientb02.cluster (1/10)
16/11/07 00:31:08 INFO TaskSetManager: Finished task 4.0 in stage 0.0 (TID 4) in 198 ms on clientb02.cluster (2/10)
16/11/07 00:31:08 INFO TaskSetManager: Finished task 5.0 in stage 0.0 (TID 5) in 188 ms on clientb02.cluster (3/10)
16/11/07 00:31:08 INFO TaskSetManager: Finished task 6.0 in stage 0.0 (TID 6) in 141 ms on clientb02.cluster (4/10)
16/11/07 00:31:08 INFO TaskSetManager: Finished task 7.0 in stage 0.0 (TID 7) in 125 ms on clientb02.cluster (5/10)
16/11/07 00:31:08 INFO TaskSetManager: Finished task 8.0 in stage 0.0 (TID 8) in 102 ms on clientb02.cluster (6/10)
16/11/07 00:31:08 INFO TaskSetManager: Finished task 9.0 in stage 0.0 (TID 9) in 84 ms on clientb02.cluster (7/10)
16/11/07 00:31:08 INFO TaskSetManager: Finished task 1.0 in stage 0.0 (TID 1) in 249 ms on clientb02.cluster (8/10)
16/11/07 00:31:08 INFO TaskSetManager: Finished task 3.0 in stage 0.0 (TID 3) in 199 ms on clientb02.cluster (9/10)
16/11/07 00:31:08 INFO TaskSetManager: Finished task 0.0 in stage 0.0 (TID 0) in 10777 ms on clientb02.cluster (10/10)
16/11/07 00:31:08 INFO DAGScheduler: ResultStage 0 (reduce at SparkPi.scala:36) finished in 15.531 s
16/11/07 00:31:08 INFO YarnScheduler: Removed TaskSet 0.0, whose tasks have all completed, from pool
16/11/07 00:31:08 INFO DAGScheduler: Job 0 finished: reduce at SparkPi.scala:36, took 16.003066 s
```

Illustration 34: Raw output

Going to the Quick Links menu and ResourceManager UI, we can navigate to the overview of running Jobs.

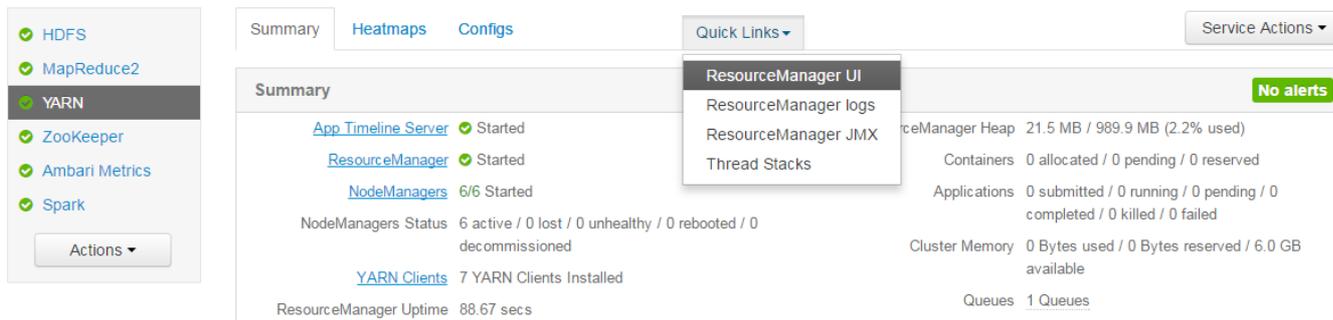
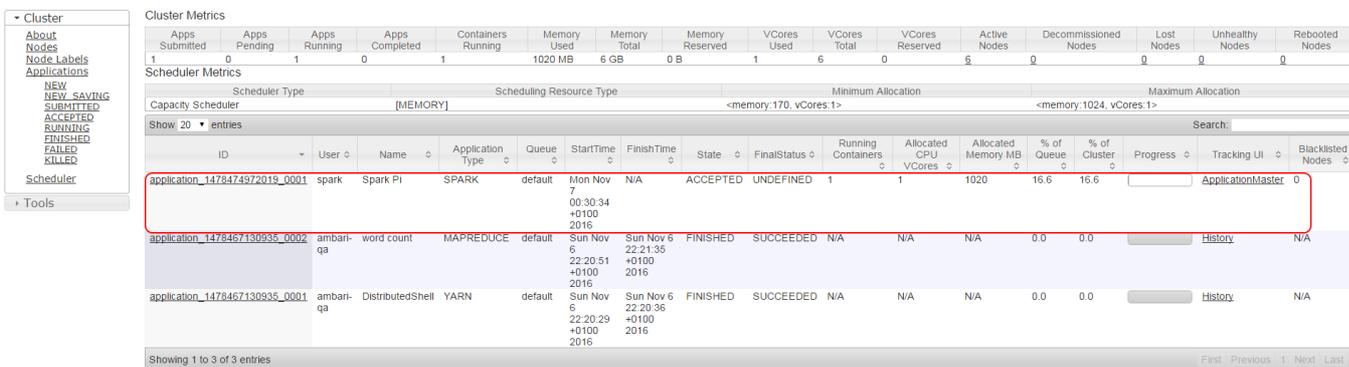


Illustration 35: Find Navigation Links

The application overview shows a list of applications

ID	User	Name	Application Type	Queue	StartTime	FinishTime	State	FinalStatus	Running Containers	Allocated CPU V-Cores	Allocated Memory MB	% of Queue	% of Cluster	Progress	Tracking UI	Blacklisted Nodes
application_1478467180935_0001	spark	Spark Pi	SPARK	default	Mon Nov 7 00:33:34 +0100 2016	N/A	ACCEPTED	UNDEFINED	1	1	1020	16.6	16.6		ApplicationMaster	0
application_1478467180935_0002	ambari-qa	word count	MAPREDUCE	default	Sun Nov 6 22:20:51 +0100 2016	Sun Nov 6 22:21:35 +0100 2016	FINISHED	SUCCEEDED	N/A	N/A	N/A	0.0	0.0		History	N/A
application_1478467130935_0001	ambari-qa	DistributedShell	YARN	default	Sun Nov 6 22:20:29 +0100 2016	Sun Nov 6 22:20:36 +0100 2016	FINISHED	SUCCEEDED	N/A	N/A	N/A	0.0	0.0		History	N/A

Illustration 36: Spark application overview

When the application has run to the end, the result is presented on the screen:

```
16/11/07 00:31:08 INFO DAGScheduler: ResultStage 0 (reduce at SparkPi.scala:36) finished in 15.531 s
16/11/07 00:31:08 INFO YarnScheduler: Removed TaskSet 0.0, whose tasks have all completed, from pool
16/11/07 00:31:08 INFO DAGScheduler: Job 0 finished: reduce at SparkPi.scala:36, took 16.003066 s
Pi is roughly 3.141908
```

The example is not especially usable. However, this is just a test of the installation.

Where to go from here?

Having a Hadoop cluster without using it is no fun. So, what should we use it for?

A Hadoop cluster is normally used for processing data from various sources, compressing and processing it into a simple result. The whole process is shown in the following illustration :

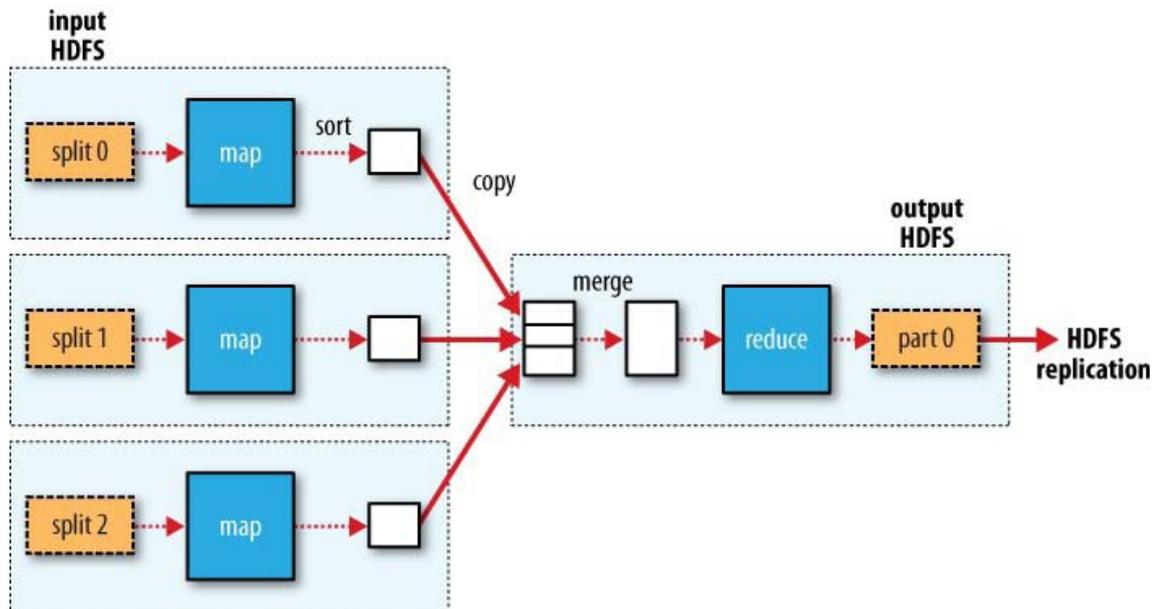


Illustration 37: Hadoop data flow

Having run the various samples on the cluster, verifying the state of the of the cluster, the example really not serve any practical uses.

That being said, the cluster now is ready to start working. In order to get acquainted wick Hadoop and Spark, I suggest you download a large dataset. A listing of these can be found at :

http://hadoopilluminated.com/hadoop_illuminated/Public_Bigdata_Sets.html

The data can be fed into the cluster and processed using various tools such as Spark, R etc.

Also, experiments with other types of services can be performed using the cluster.

So, happy Hadooping, and please feel free to keep an eye out for more intros on Hadoop.

/Henrik