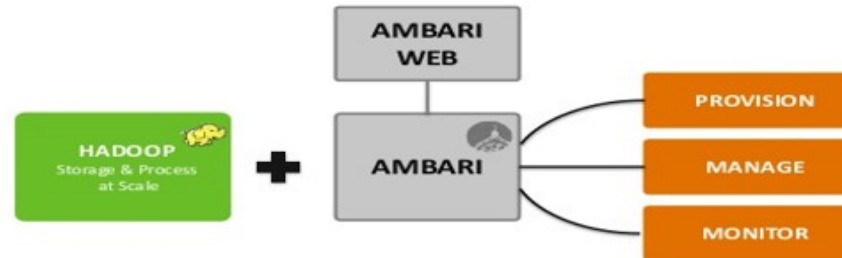# AMBARI

- ✓ What is Ambari?

- ✓ What is Hadoop?

- ✓ Types of managing tools

- ✓ Architecture of Ambari

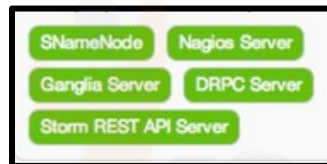- ✓ How to install Ambari?

- ✓ Setting up of Hadoop cluster

What is Ambari?

Apache Ambari is a tool for provisioning, managing, and monitoring Apache Hadoop clusters. Ambari consists of a set of RESTful APIs and a browser-based management interface.
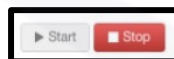
# Ambari operations:



Provision:



- Virtual, physical and cloud Environments.
- Deploy 10s, 100s, 1000s of Hadoop servers

Manage:-

- Advance configuration & host Controls.
- Single point for Host controls.

Monitor:-

- ✓ Pre-configuration metrics and alerts.
- ✓ Single pane of glass for Hadoop & system status.

# What is Hadoop?

- Hadoop is a large-scale and distributed data storage and processing infrastructure using clusters of commodity hosts networked together.

- Monitoring and managing such complex distributed systems is a non-trivial task.

- To help you manage the complexity, Apache Ambari collects a wide range of information from the cluster's nodes and services and presents it to you in an easy-to-read and use, centralized web interface, Ambari Web.

# What is Hadoop?

> Ambari Web displays information such as service-specific summaries, graphs, and alerts.

> You use Ambari Web to create and manage your HDP cluster and to perform basic operational tasks such as starting and stopping services, adding hosts to your cluster, and updating service configurations.

> You also can use Ambari Web to perform administrative tasks for your cluster such as enabling Kerberos security and performing Stack upgrades.

What are the types of managing tools?

# Types of Managing Tools:

> Web Based data colloction → Nutch, Solr, Gora, Hbase

> Mapreduce Programming→ Fair and Capacity schedulers, Oozie

> Moving data→ Hadoop Commands, sqoop, flume, storm

> Monitoring→ Hue, Nagios ,Ganglia

> Analysis with sql → Impala, hive, spark

> ETL→ Pentaho, Talend

> Reporting--> Splunk, Talend

# One liner management tool explanation:-

## Nutch:-

Apache Nutch is a highly extensible and scalable open source web crawler software project.

## Gora:-

The Apache Gora open source framework provides an in-memory data model and persistence for big data. Gora supports persisting to column stores, key value stores, document stores and RDBMSs, and analyzing the data with extensive Apache Hadoop MapReduce.

## Solr:-

Apache solr is a standalone full text search server with Apache Lucene at the backend, User for Web application for text search, A Wrapper around Apache Lucene, Written at Cnet, now at Apache.

## Hbase:-

HBase is a NoSQL databases which experienced a tremendous increase in popularity during recent years.

# One liner management tool explanation:-

**Oozi:-**

Provides workflow management and coordination of those workflows, Manages Directed Acyclic graph of Actions.

**Sqoop:-**

sqoop is a tool to transfer data between Hadoop and relational database, Transform data in Hadoop with Mapreduce or hive, Export data back in to RDB.

**Nagios:-**

This is a system and network monitoring tool.

**Hue:-**

Hue is a lightweight Web server that help you to use hadoop directly from your browser. Hue is just a view on top of any Hadoop distribution and can be installed on any machine.

# One liner management tool explanation:-

Imapala:-

        General purpose sql query Engine, works both for analytical and transactional/Single row workloads.

Talend:-

        Eclipse-base visual programming Editor, generates executable java code. Talend to bring an open source integration tool for easily connecting Apache Hadoop to hundreds of data systems without having to write code.

Now, let us talk about Minimum System Requirements.

# Minimum System Requirements:-

**Hardware Recommendations:-**

      There is no single set of hardware recommendations for installing Hadoop.

**Operating Systems Requirements:-**

      Red Hat Enterprise Linux (RHEL) v5.x or 6.x (64-bit)

        · CentOS v5.x or 6.x (64-bit)

        · SUSE Linux Enterprise Server (SLES) 11, SP1 (64-bit)

Browser Requirements:-

      The Ambari Install Wizard runs as a browser-based Web app. You must have a machine capable of running a graphical browser to use this tool.

The supported browsers are:

· Mac OS X_10.6 or later

      ➢ Firefox latest stable release

# Minimum System Requirements:-

> Safari latest stable release
> Google Chrome latest stable release

- Linux_RHEL, CentOS

  - Firefox latest stable release
  - Google Chrome latest stable release

- Windows Vista, 7

  - Internet Explorer 9.0 and higher

  - Firefox latest stable release

  - Safari latest stable release

  - Google Chrome latest stable release

# Software Requirements:-

On each of your hosts:

- yum
- rpm
- scp
- curl
- wget
- pdsh

- Database Requirements:-

Hive or HCatalog requires a MySQL database for its use. You can choose to use a current instance or Let the Ambari install wizard create one for you.

Let us know about the architecture of Ambari.

# Ambari Architecture:-

Ambari system Architecture:-



Ambari have two components. They are:

· Ambari server –

Master process which communicates with Ambari agents installed on each node participating in the cluster. This has postgres database instance which is used to maintain all cluster related metadata.

# Ambari Architecture:-

**Ambari Agent–**

These are acting agents for Ambari on each node. Every agent periodically sends his own health status along with different metrics, installed services status and many more things. According master decides on next action and conveys back to the agent to act.

**Ambari Installation:**

·   Ambari installation is easy a task of few commands.

·   We will cover Ambari installation and cluster setup.

·   We are assumed of having 4 nodes. Node1, Node2, Node3 and Node4. And we are picking Node1 as our Ambari server.

·   These are installation steps on the RHEL based system, for debian and other systems steps will vary little.

# Installation of Ambari:-

From Ambari server node (Node 1 as we decided)

i.  **Download Ambari public repo**

> sudo wget http://public-repo-1.hortonworks.com/ambari/centos7/2.x/updates/2.1.2/ambari.repo -O /etc/yum.repos.d/ambari.repo

This command will add Hortonworks Ambari repository into yum which is a default package manager for RHEL systems.

ii.  **Install Ambari RPMS**

> sudo yum install -y ambari-server

# Installation of Ambari:-

## iii. Configuring Ambari server

The next thing to do after Ambari installation is to configure Ambari and set it up to provision the cluster.

Following step will take care of this:-

```
sudo ambari-server setup
```

# Installation of Ambari:-



💡 Wish to go with the default options which we do often use? Use-silent option.

# Installation of Ambari:-

## iv. Start the server and Login to web UI

Start the server with



sudo ambari-server start

· Now we can access Ambari web UI (hosted on 8080 port).

· Login into Ambari with default username "admin" and default password "admin".

How to setup Hadoop Cluster?

# Setting up Hadoop cluster:-

---

I.    **Landing page:**



Click on "Launch Install Wizard" to start cluster setup

2.  **Cluster Name:**
·.      Give your cluster a good name.

**Note:** This is just a simple name for cluster, it is not that significant, so don't worry about it and choose any name for it.

# Setting up Hadoop cluster:-

## 3. Stack selection



·     This page will list stacks available to install.
·     Each stack is pre-packaged with Hadoop ecosystem component.
·     These stacks are from Hortonworks. (We can install plain Hadoop too.)

# Setting up Hadoop cluster:-

## 4. Hosts Entry and SSH key entry

Prior moving further this step we should have password less SSH setup for all the participating nodes.



· Add the hostnames of the nodes, single entry on each line. [ Add FQDN which can be obtained by hostname –f command].

· Select private key used while setting up password less SSH and username using which private key was created.

# Setting up Hadoop cluster:-

## 5. Hosts registration status



· You can see some operations being performed; these operations include setting Ambari-agent on each node, creating basic setups on each nodes. Once we see ALL GREEN we are ready to move on. Sometimes this may take time as it installs few packages.

Sometime registering hosts fails, retry twice at least or Install Ambari-agent manually on each node.

# Setting up Hadoop cluster:-

**6. Choose services you wish to install:**

· As per selected stacks in step 3, we have number of services that we can install in the cluster. You can choose one you want.

· Ambari intelligently selects dependent services if you haven't selected it.

· For instance, you selected HBase but not Zookeeper it will prompt same and will add Zookeeper also to the cluster.

# Setting up Hadoop cluster:-

## 6. Choose services you wish to install:

# Setting up Hadoop cluster:-

**7. Master services mapping with Nodes:**

As you are aware of Hadoop ecosystem has tools which are based on master-slave architecture.

✓ In this step we will associate master processes with the node.

# Setting up Hadoop cluster:-

✓ Here make sure you properly balance your cluster.
✓ Also keep in mind primary and secondary services like Namenode and secondary Namenode are not on the same machine.

**Always distribute services in such a way that you have only one master on one node. So that node failure does not disturb multiple services, keep good amount of main memory for optimum performance.**

# Setting up Hadoop cluster:-

**8. Slaves mapping with Nodes:-**

      Similar to masters, map slave services on the nodes. In general, all the nodes will have slave process running at least for Datanodes and Nodemanagers.

# Setting up Hadoop cluster:-



**9. Customize services**

This is very important page for Administrators!

· Here you can configure properties for your cluster to make it most suited to your use cases.

· Also it will have some required properties like Hive metastore password (if hive is selected) etc. These will be pointed with Red error like symbols.

# Setting up Hadoop cluster:-

**10. Review and start provisioning**

Make sure you review the cluster configuration before launch as this will save from unknowingly set wrong configurations.

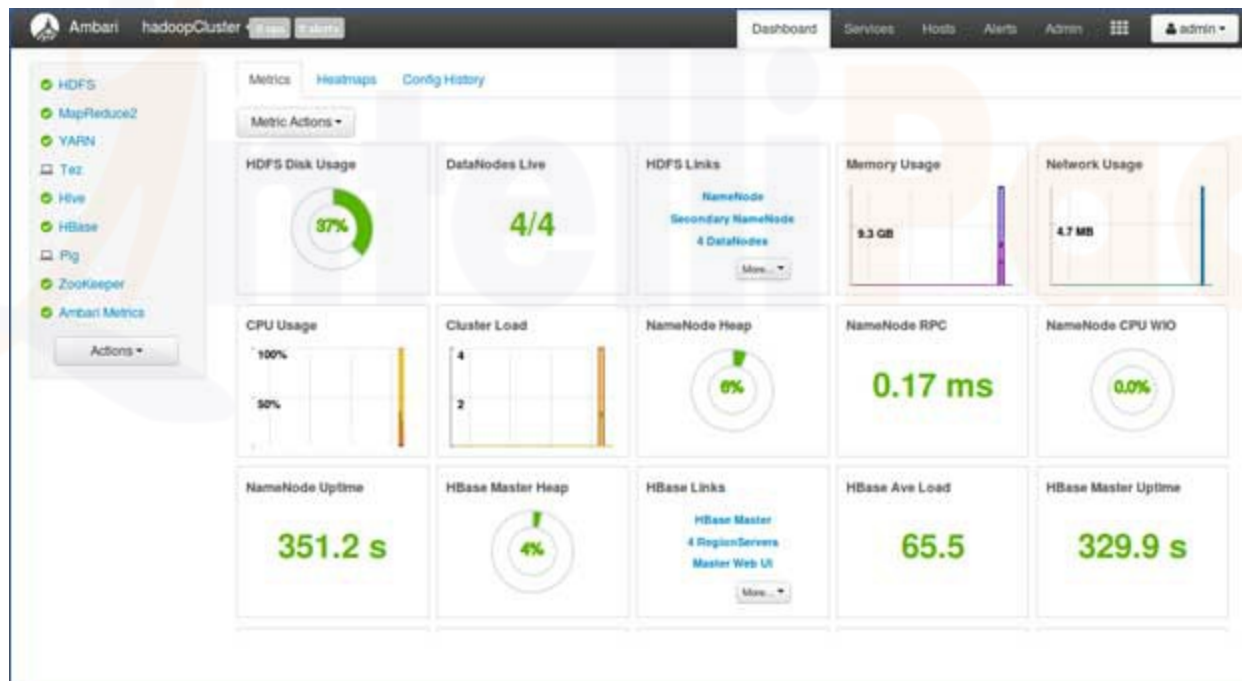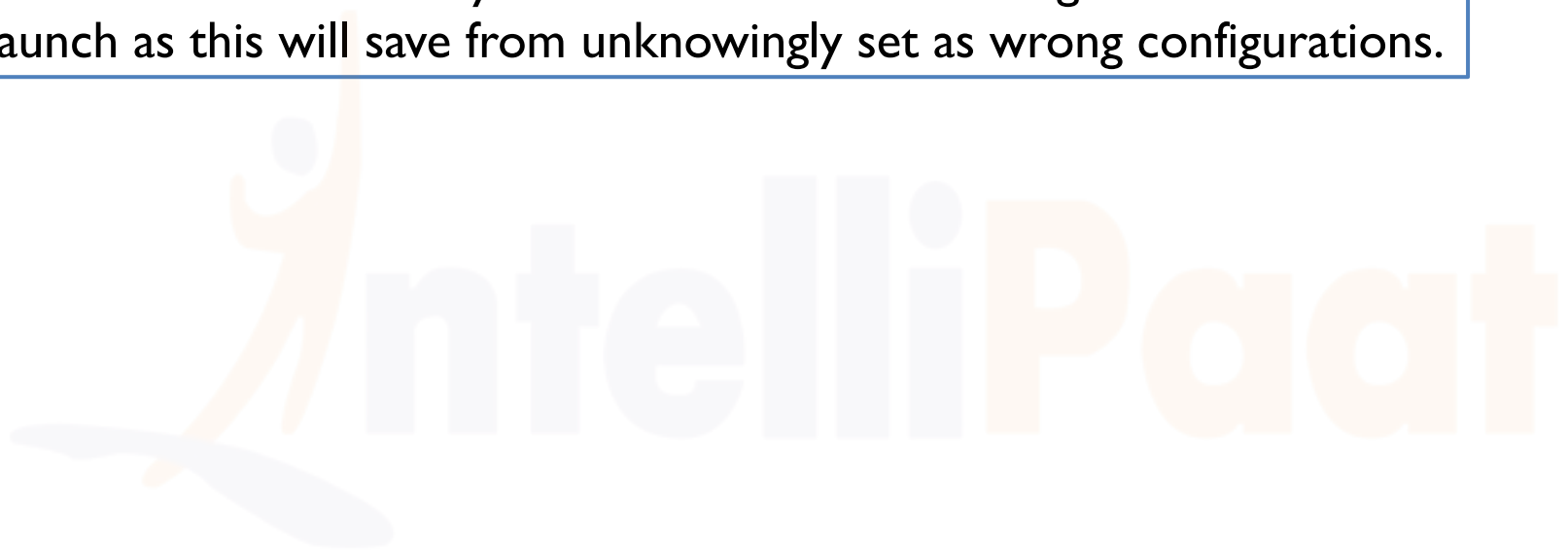**11. Launch and stay back until status becomes GREEN.**
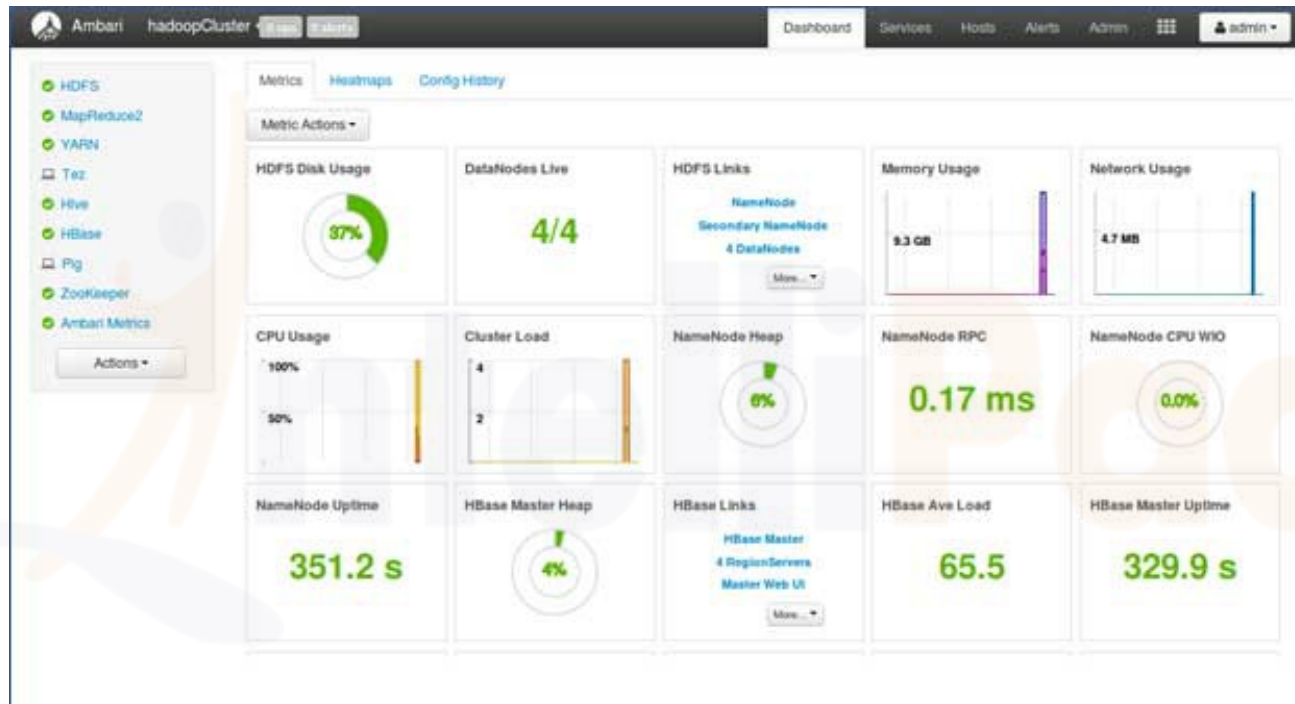
# Setting up Hadoop cluster:-

**10. Review and start provisioning:**
            Make sure you review the cluster configuration before
launch as this will save from unknowingly set as wrong configurations.

# Setting up Hadoop cluster:-

**11. Launch and stay back until status becomes GREEN.**



**Wow!** We have successfully Installed Hadoop and all the components on all the nodes of the cluster. Now we can get ourselves start playing with Hadoop.

# Conclusion:

- We have now learned how to install Hadoop and its components on multi-node cluster
  using a simple web based tool called Apache Ambari.

- Apache Ambari provides us a simpler interface and saves lots of our efforts on installation, monitoring and management which would have be very tedious with so many components and their different installation steps and monitoring controls.

# Thank You

Email us – support@intellipaat.com

Visit us - https://intellipaat.com