# Apache Hadoop with Cloudera Tutorial 1

**Venkatesh Vinayakarao**

[venkateshv@cmi.ac.in](mailto:venkateshv@cmi.ac.in)
[http://vvtesh.co.in](http://vvtesh.co.in)
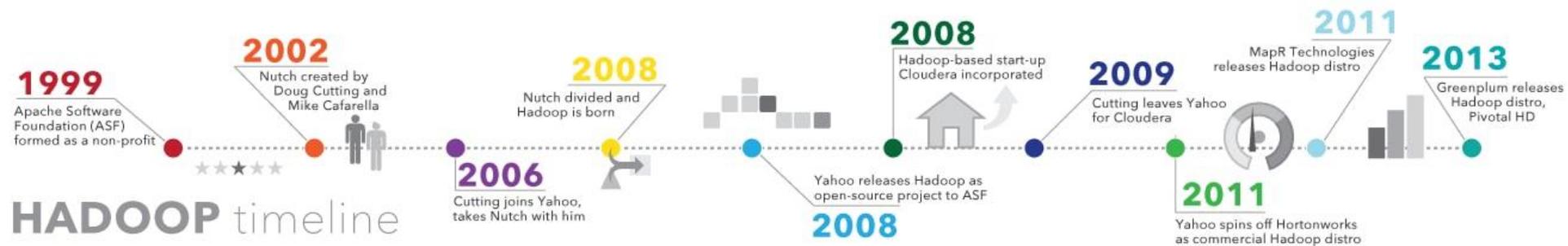
Chennai Mathematical Institute

Based on slides by Suchitra Jayaprakash

# Apache Hadoop

- Hadoop is open-source software framework used for processing data on distributed commodity computing environment.

- A Java based open-source software managed by Apache Software Foundation.

- Founded by Doug Cutting & Mike Cafarella.

- Based on Google's white paper on Google File System & map-reduce.
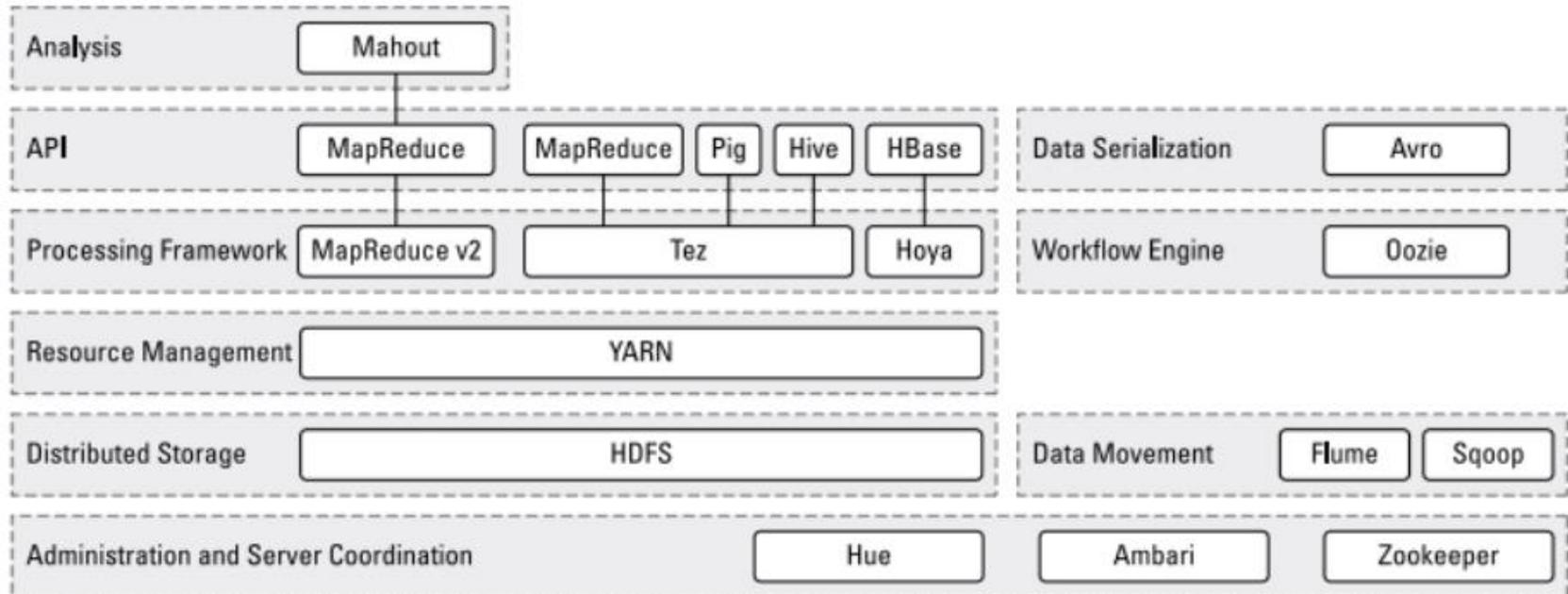
Find Hadoop at https://hadoop.apache.org/

# Hadoop - History



**1999** Apache Software Foundation (ASF) formed as a non-profit

**2002** Nutch created by Doug Cutting and Mike Cafarella

**2006** Cutting joins Yahoo, takes Nutch with him

**2008** Nutch divided and Hadoop is born

**2008** Yahoo releases Hadoop as open-source project to ASF

**2008** Hadoop-based start-up Cloudera incorporated

**2009** Cutting leaves Yahoo for Cloudera

**2011** Yahoo spins off Hortonworks as commercial Hadoop distro

**2011** MapR Technologies releases Hadoop distro

**2013** Greenplum releases Hadoop distro, Pivotal HD

**HADOOP** timeline

Hadoop History - Image Source: https://www.sas.com/en_in/insights/big-data/hadoop.html

# Modules

- The project includes these modules:
  - **Hadoop Common**: The common utilities that support the other Hadoop modules.
  - **Hadoop Distributed File System (HDFS™)**: A distributed file system that provides high-throughput access to application data.
  - **Hadoop YARN**: A framework for job scheduling and cluster resource management.
  - **Hadoop MapReduce**: A YARN-based system for parallel processing of large data sets

# Hadoop Ecosystem



(source: Hadoop for Dummies)

# Hadoop Distribution

- Need for commercial distributions
  - Customization as per industry needs
  - Maintenance
  - Ease of installation and use
  - Reliability add-ons etc.

- Several vendors offer Hadoop distribution
  - We shall use cloudera for this tutorial

# Cloudera

- Founded in 2008 by three engineers from Google, Yahoo! and Facebook (Christophe Bisciglia, Amr Awadallah and Jeff Hammerbacher).

- Major code contributor of Apache Hadoop ecosystem.

- First company to develop and distribute Apache Hadoop based software in March 2009.

- Additional feature includes user interface, security, interface for third party application integration.

- Offers customer support for installing , configuring , optimising Cloudera distribution through its enterprise subscription service.

- Provides a proprietary Cloudera Manager for easy installation , monitoring & trouble shooting.

# Cloudera



An illustration of Cloudera's open-source Hadoop distribution (source: cloudera website).

# Tutorial - Objectives

- Understand and Install Cloudera distribution of Hadoop
  - Understand Containers
  - Understand Docker
  - Virtual Machines (VM) Vs. Docker
  - Docker Installation
  - Cloudera Quickstart
- Enjoy Hadoop!

# Role of Containers

Why do we need such containers?

# Docker

- Docker is an open-source tool that uses containers to create, deploy, and manage distributed applications.

- Developers use containers to create packages for applications that include all libraries that are needed to run the application in isolation.

# Docker Installation

- Sign up to [https://docs.docker.com/](https://docs.docker.com/)
- Follow instructions at [https://docs.docker.com/docker-for-windows/install/](https://docs.docker.com/docker-for-windows/install/)
- For this tutorial, we install
  - Docker desktop on Windows
- To check docker installation is proper, open docker terminal (or Command Prompt on Windows) and type: docker run hello-world

# Docker Installation

- Successful Installation!

# How to Run Cloudera Hadoop?

- Pull latest cloudera image
  - docker pull cloudera/quickstart:latest

- List the images we have
  - docker images

- Run Cloudera in a docker
  - docker run --hostname=quickstart.cloudera --privileged=true -t -i -p 8888:8888 -p 8080:8080 -p 8088:8088 -p 7180:7180 -p 50070:50070 cloudera/quickstart /usr/bin/docker-quickstart

# Cloudera Quickstart



To exit, [root@quickstart /]# exit

# Some useful commands

- docker ps



- docker stats [containerid]

- docker inspect [CONTAINER ID]

# We are all set!

NameNode: http://localhost:50070/



Yarn page: http://localhost:8088/

HUE- http://localhost:8888/
Default username / password : cloudera / cloudera

# We could do all that on the cloud!

- Go to Google Cloud
    - https://cloud.google.com/
- Activate Cloud Shell



- You should see this

# Let us install cloudera here!

- Docker is pre-installed.



```
CLOUD SHELL
Terminal        (flash-spot-319715)  ✕   +  ▾
```

```
For more examples and ideas, visit:
 https://docs.docker.com/get-started/

vvtesh@cloudshell:~ (flash-spot-319715)$ docker pull hello-world
Using default tag: latest
latest: Pulling from library/hello-world
Digest: sha256:dcba6daec718f547568c562956fa47e1b03673dd010fe6ee58ca806767031d1c
Status: Image is up to date for hello-world:latest
docker.io/library/hello-world:latest
vvtesh@cloudshell:~ (flash-spot-319715)$ docker run hello-world

Hello from Docker!
This message shows that your installation appears to be working correctly.
```

# Install Cloudera

```
vvtesh@cloudshell:~ (flash-spot-319715)$ docker pull cloudera/quickstart:latest
latest: Pulling from cloudera/quickstart
Image docker.io/cloudera/quickstart:latest uses outdated schema1 manifest format. Please upgrade to a
ormation at https://docs.docker.com/registry/spec/deprecated-schema-v1/
1d00652ce734: Downloading [===>                                        ]  347.9MB/4.444GB
```

```
CLOUD SHELL
Terminal        (flash-spot-319715) ×  + ▾                                                    ✎ Open Editor    ▣  ⚙  ▣  ▢  ⋮    ↕  ◩  ✕

1d00652ce734: Pull complete
Digest: sha256:f91bee4cdfa2c92ea3652929a22f729d4d13fc838b00f120e630f91c941acb63
Status: Downloaded newer image for cloudera/quickstart:latest
docker.io/cloudera/quickstart:latest
vvtesh@cloudshell:~ (flash-spot-319715)$ docker run --hostname=quickstart.cloudera --privileged=true -t -i -p 8888:8888 -p 8080:8080 -p 8088:8088 -p 7180:718
0 -p 50070:50070 cloudera/quickstart /usr/bin/docker-quickstart
Starting mysqld:                                          [  OK  ]
```

# Try the Preview Port

Click here
Change port to 50070

# Namenode Information Opens in a Browser

# Thank You