

The Principle of Hadoop :

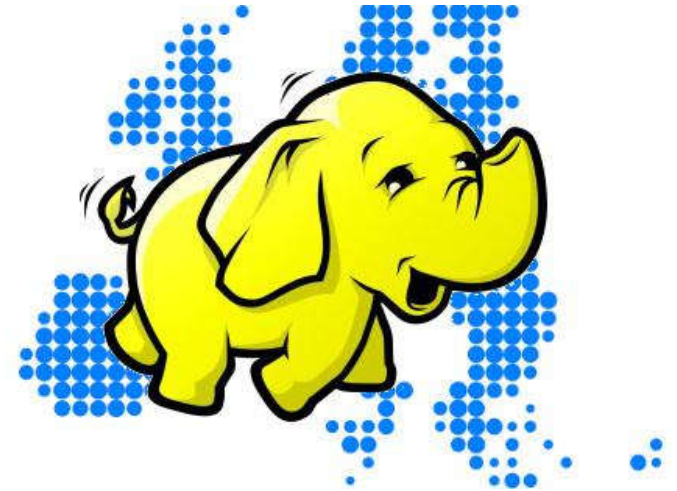
Name: อวชัย ภิรมย์รัตน์

Tel. : 086-813-5354

e-mail : p.Auoychai@gmail.com

Big Data

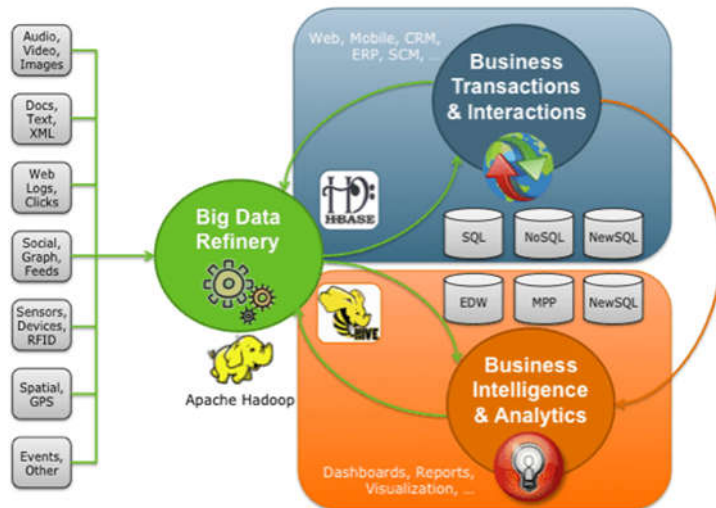
Hadoop World: Why will be Hadoop ?



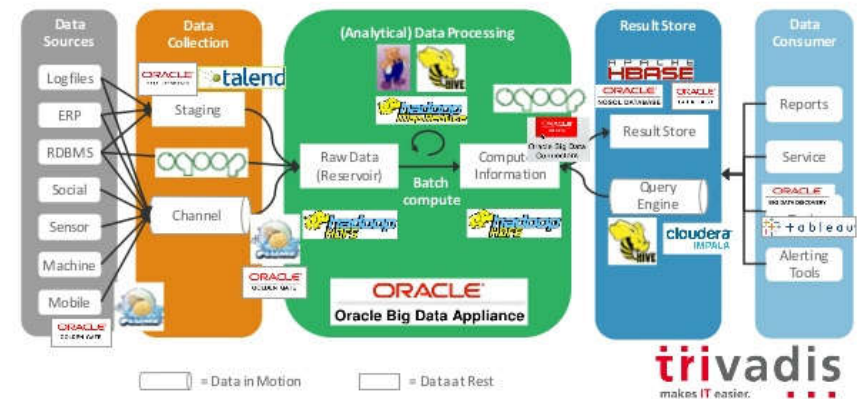
Hadoop was created by Doug Cutting and Mike Cafarella in 2005. It was originally developed to support distribution for the [Nutch search engine](#) project. Doug, who was working at Yahoo! at the time and is now Chief Architect of [Cloudera](#), named the project after his son's toy elephant. Cutting's son was 2 years old at the time and just beginning to talk. He called his beloved stuffed yellow elephant "Hadoop"

Hadoop World: Why will be Hadoop ?

Next-Generation Data Architecture



“Hadoop Ecosystem” Technology Mapping

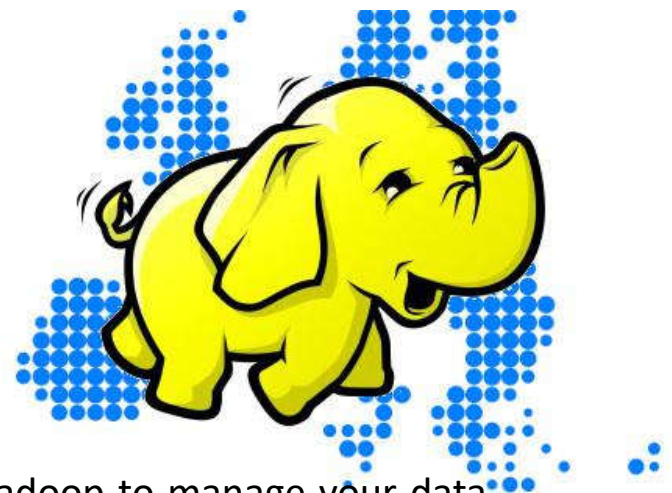


What is Hadoop :

Hadoop is software framework with optimized for distributed processing of very large datasets.

Its has two main features are the Hadoop Distributed File System(HDFS) for storing files and MapReduce for processes the stored information

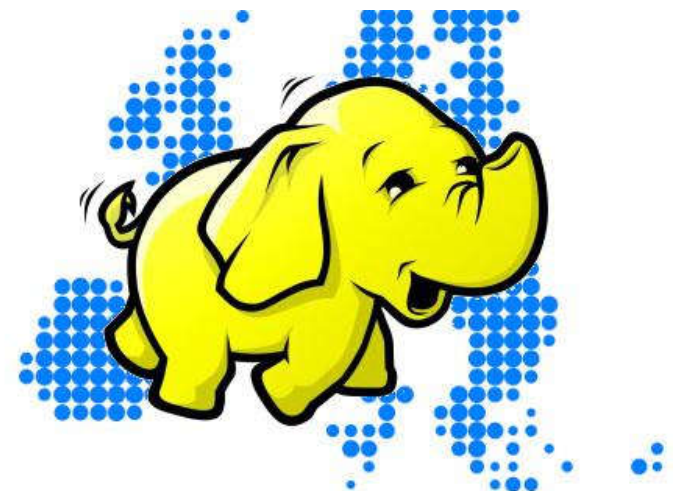
Advantage: Store terabytes of data using any number of inexpensive commodity servers.



What is Hadoop :

Key features that answer – Why Hadoop?

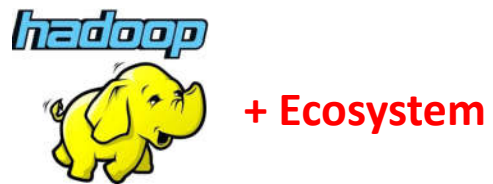
1. Flexible:
2. Scalable/fault tolerant
3. Building more efficient data economy:
4. Robust Ecosystem: MapReduce, Hive, HBase, Zookeeper, HCatalog,...
5. Hadoop is getting more “Real-Time”!
6. Cost Effective:
7. Upcoming Technologies using Hadoop:
8. Hadoop is getting cloudy!



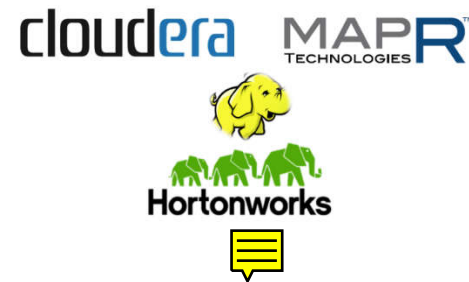
<https://www.edureka.co/blog/why-hadoop/>

Hadoop Classification :

Open Source



Distribution



TERADATA®

ORACLE®
BIG DATA

IBM®

Appliance

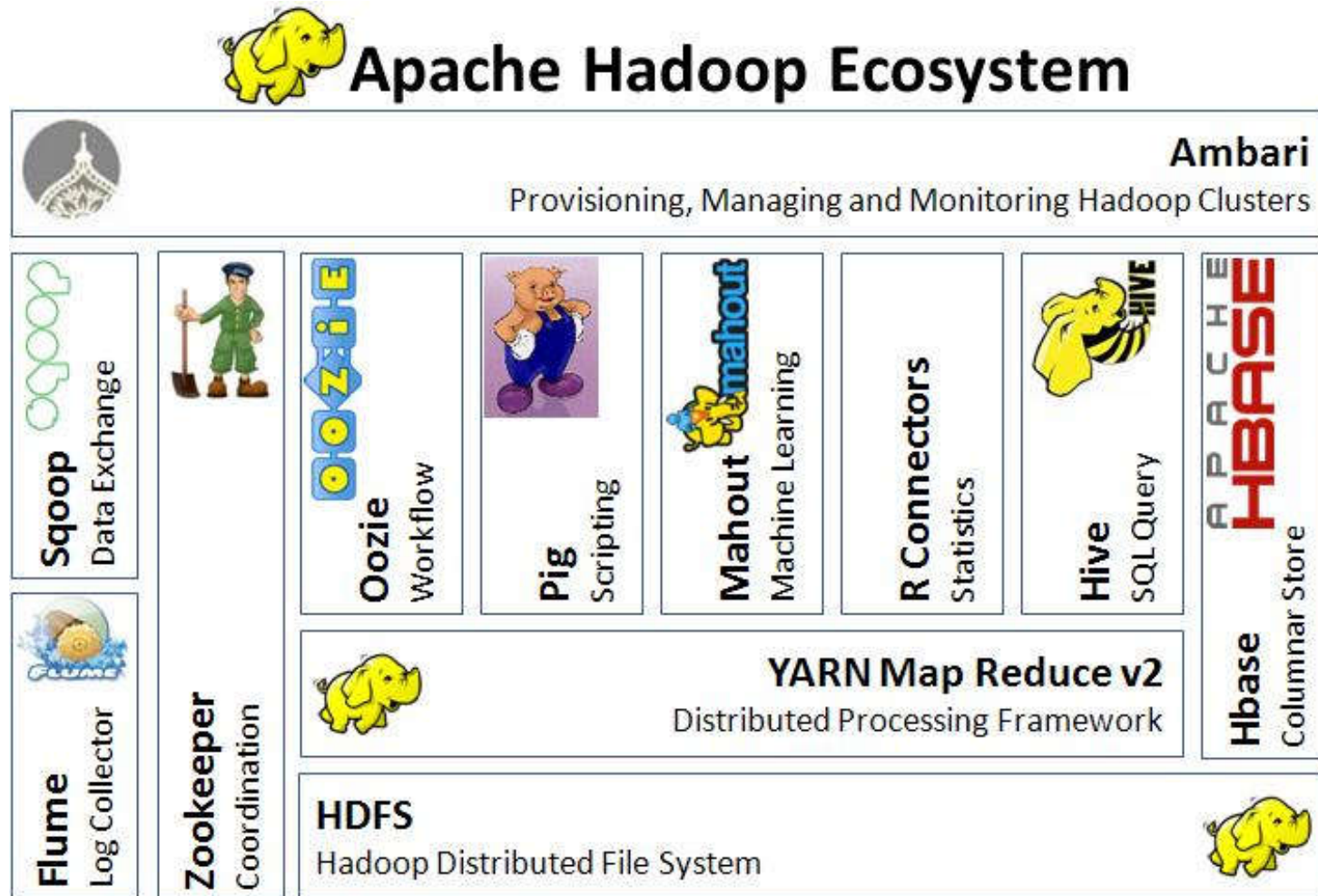
Cloud

Hadoop Component :

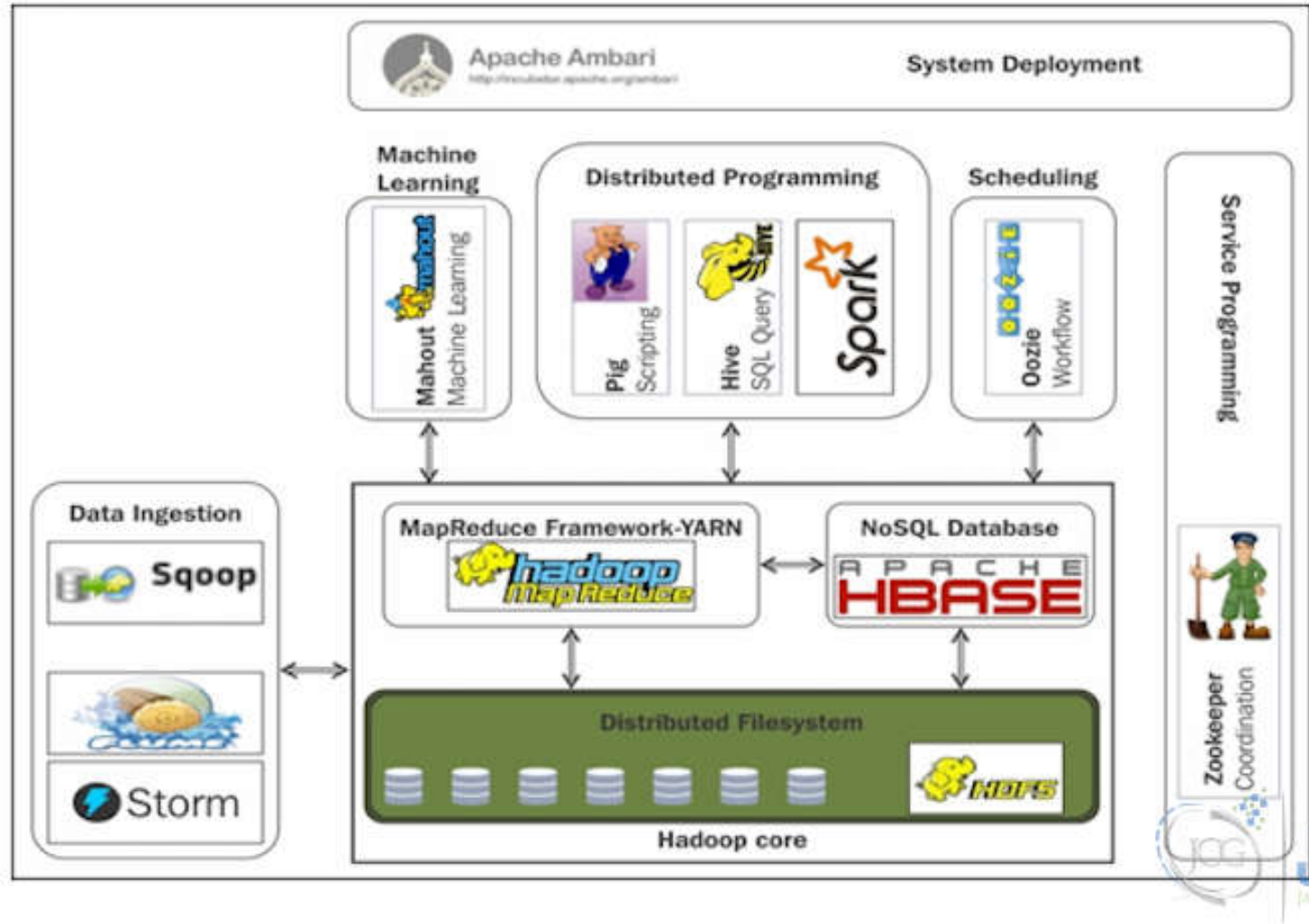
- **Hadoop Distributed File System (HDFS™):** A distributed file system that provides high-throughput access to application data.
- **Hadoop YARN:** A framework for job scheduling and cluster resource management.
- **Hadoop MapReduce:** A YARN-based system for parallel processing of large data sets.



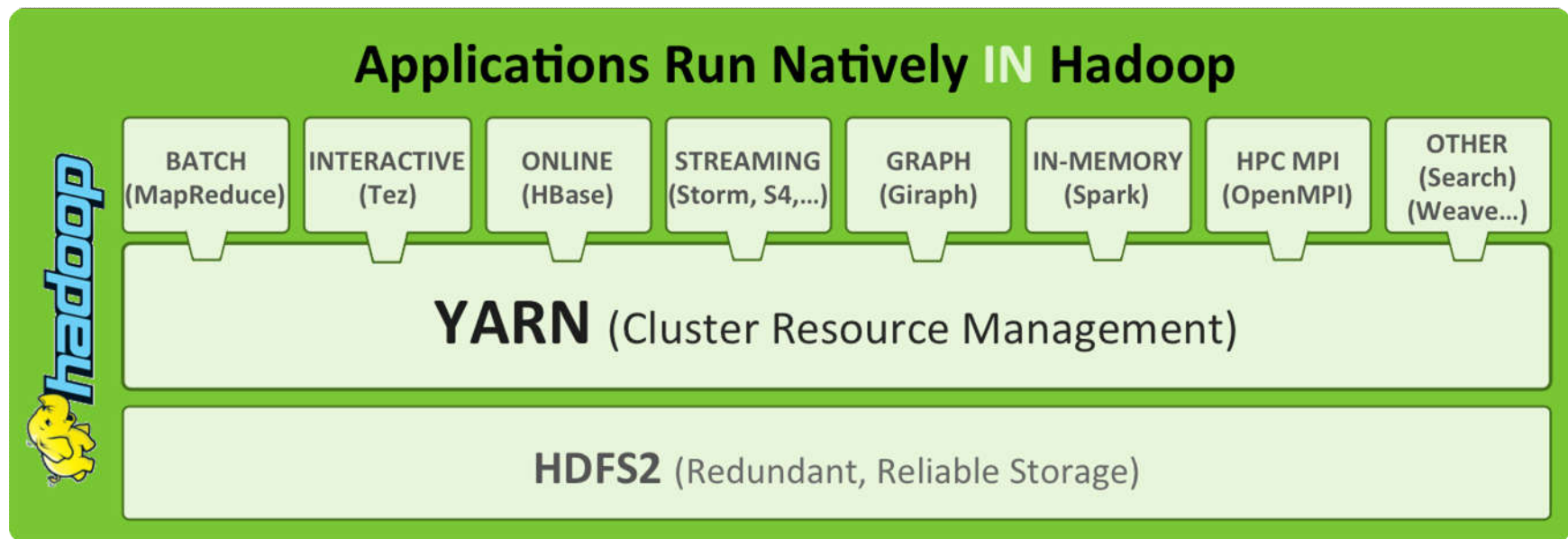
Hadoop Ecosystem :



Hadoop Ecosystem :

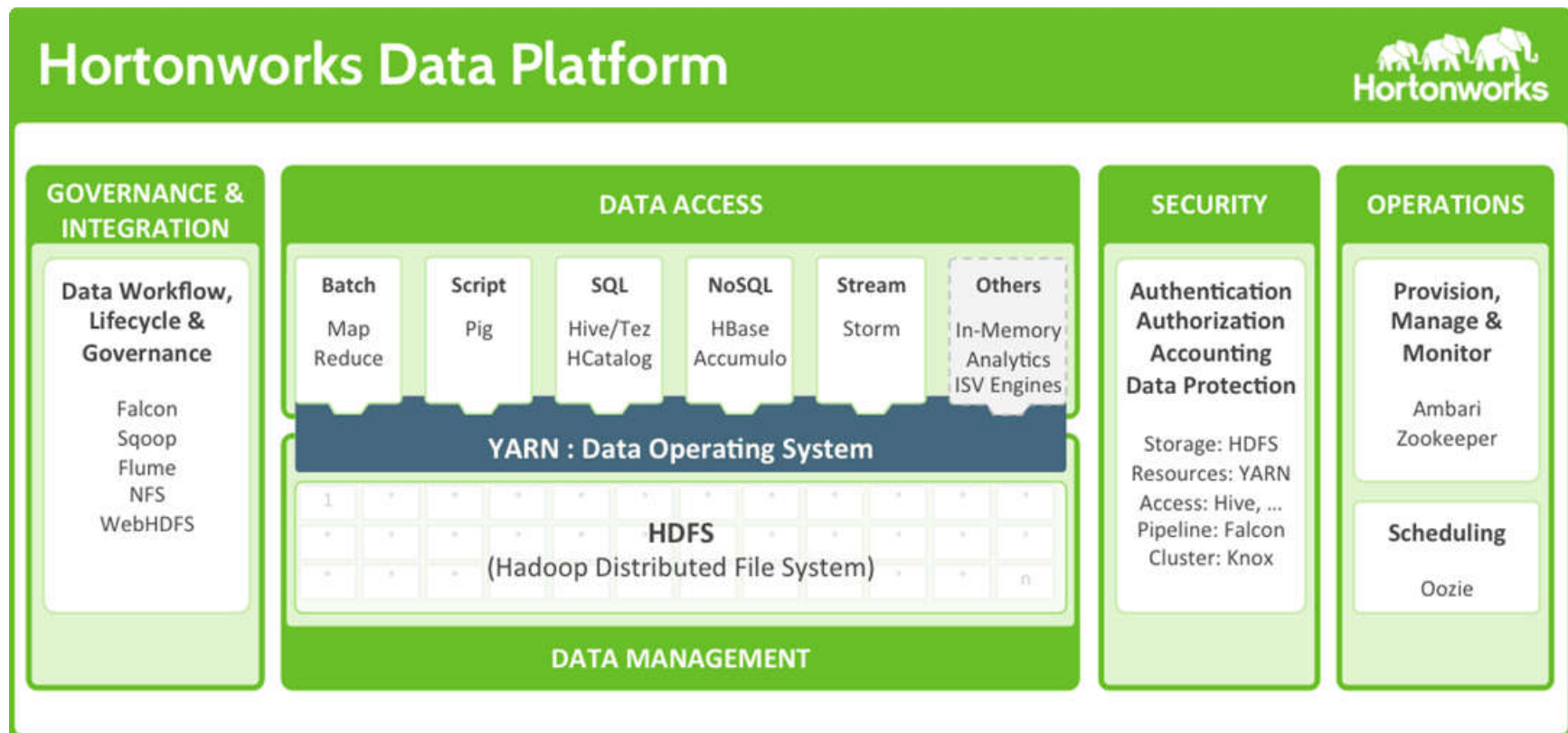


Multi-Purpose Hadoop Component Stack :

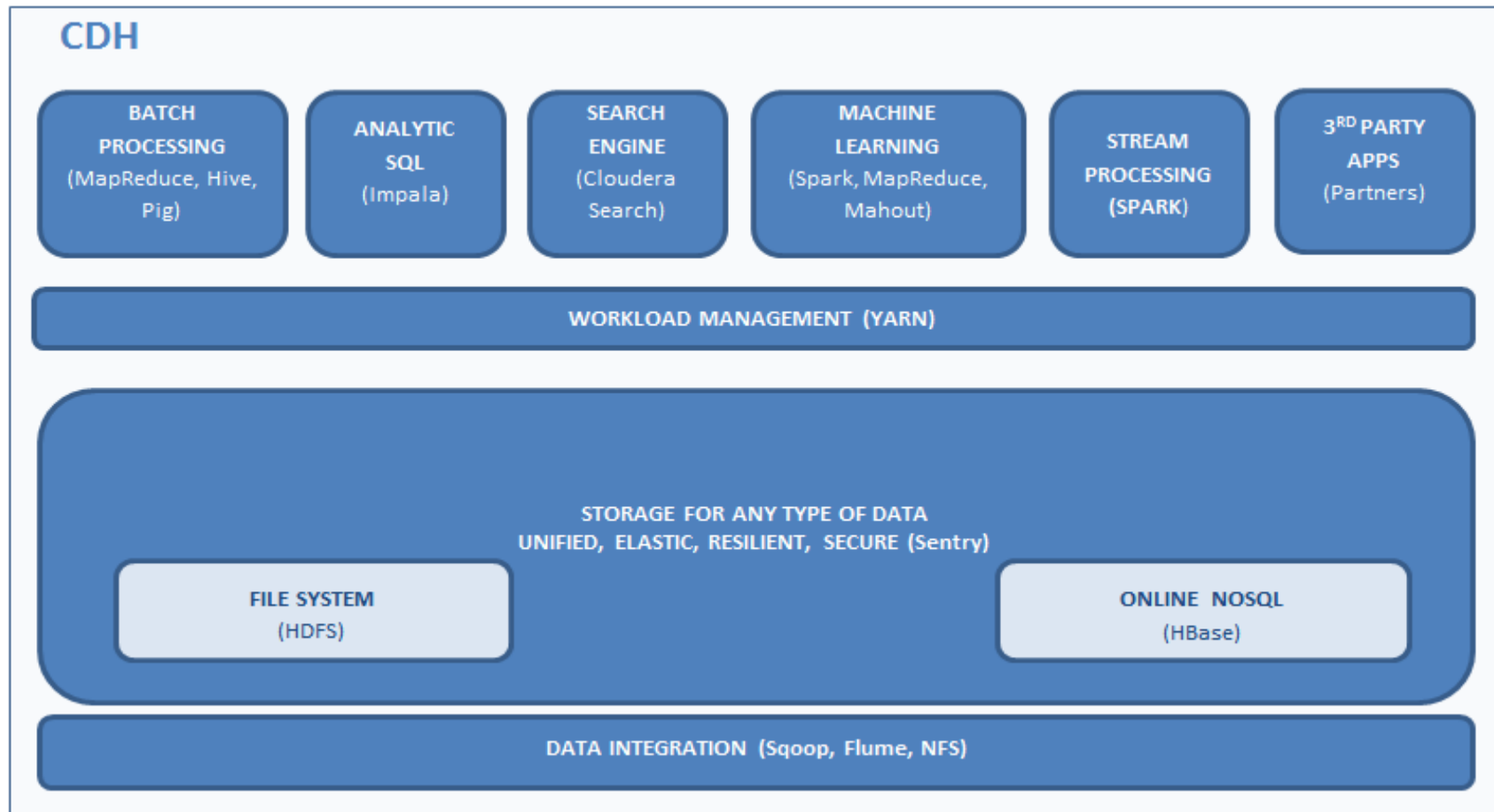


<http://tm.durusau.net/wp-content/uploads/2013/06/YARN2.png>

Enterprise Hadoop Component Stack :



Enterprise Hadoop Component Stack:



Hadoop Component & Terminology :

Master/Slave Architecture:

Master / Secondary NameNode

DataNodes

Data Management:

HDFS

Data Bock

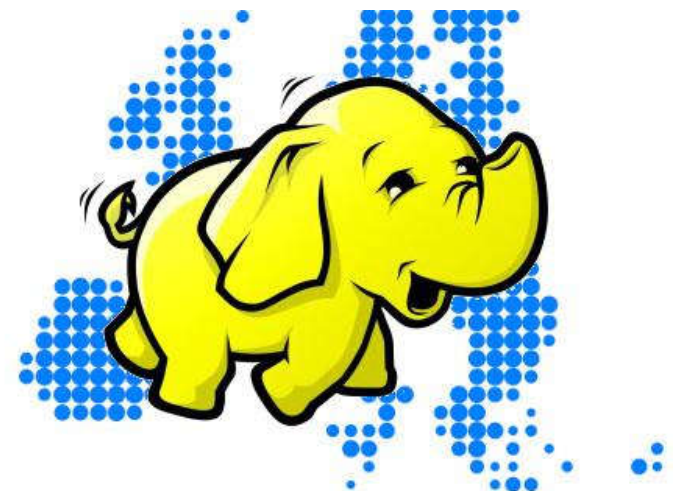
Replication Management

Rack Awareness

HDFS Read/Write – Behind the scenes

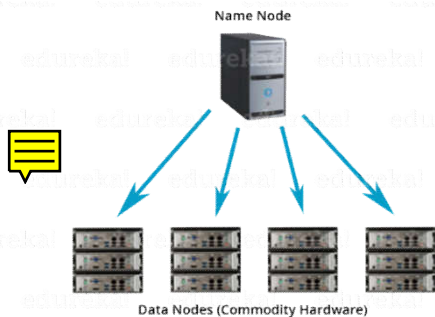
Process & Resource Management :

YARN



Hadoop Component & Terminology :

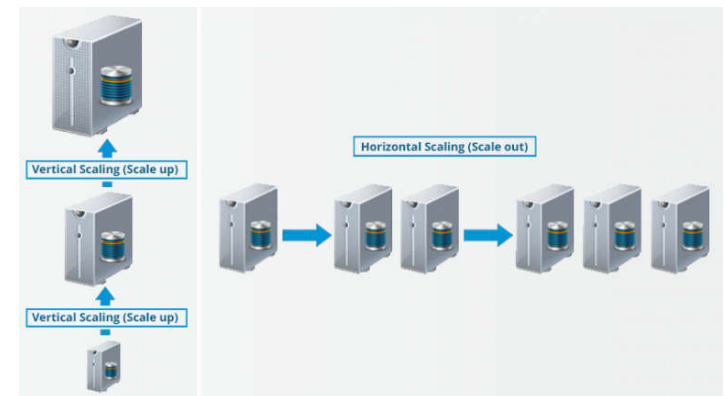
Hadoop is software framework with optimized for distributed processing of very large datasets. The main features are the Hadoop Distributed File System(HDFS) and MapReduce.



1. Distributed Storage:



2. Distributed & Parallel Computation:



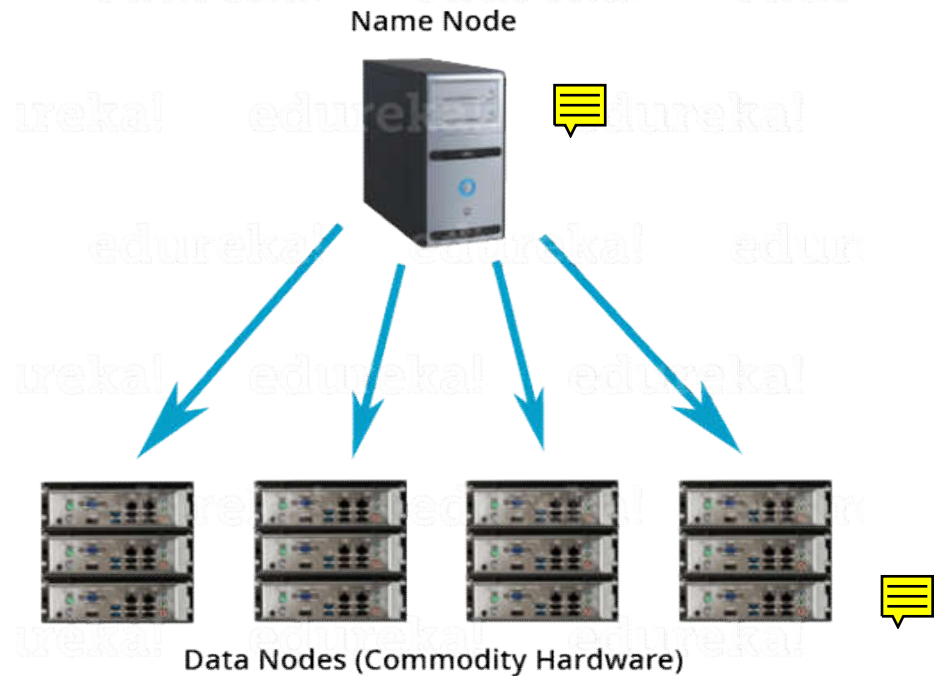
3. Horizontal Scalability:

Hadoop Component & Cluster :

Master/Slave Architecture:

Master / Secondary NameNode

DataNodes



NameNode (Master Node)

- Manage data block on DataNodes with Metadata , block location , file size , permission , hierarchy ,etc.
- Maintain DataNodes on cluster , ensure DNs are live

DataNode (Slave Node)

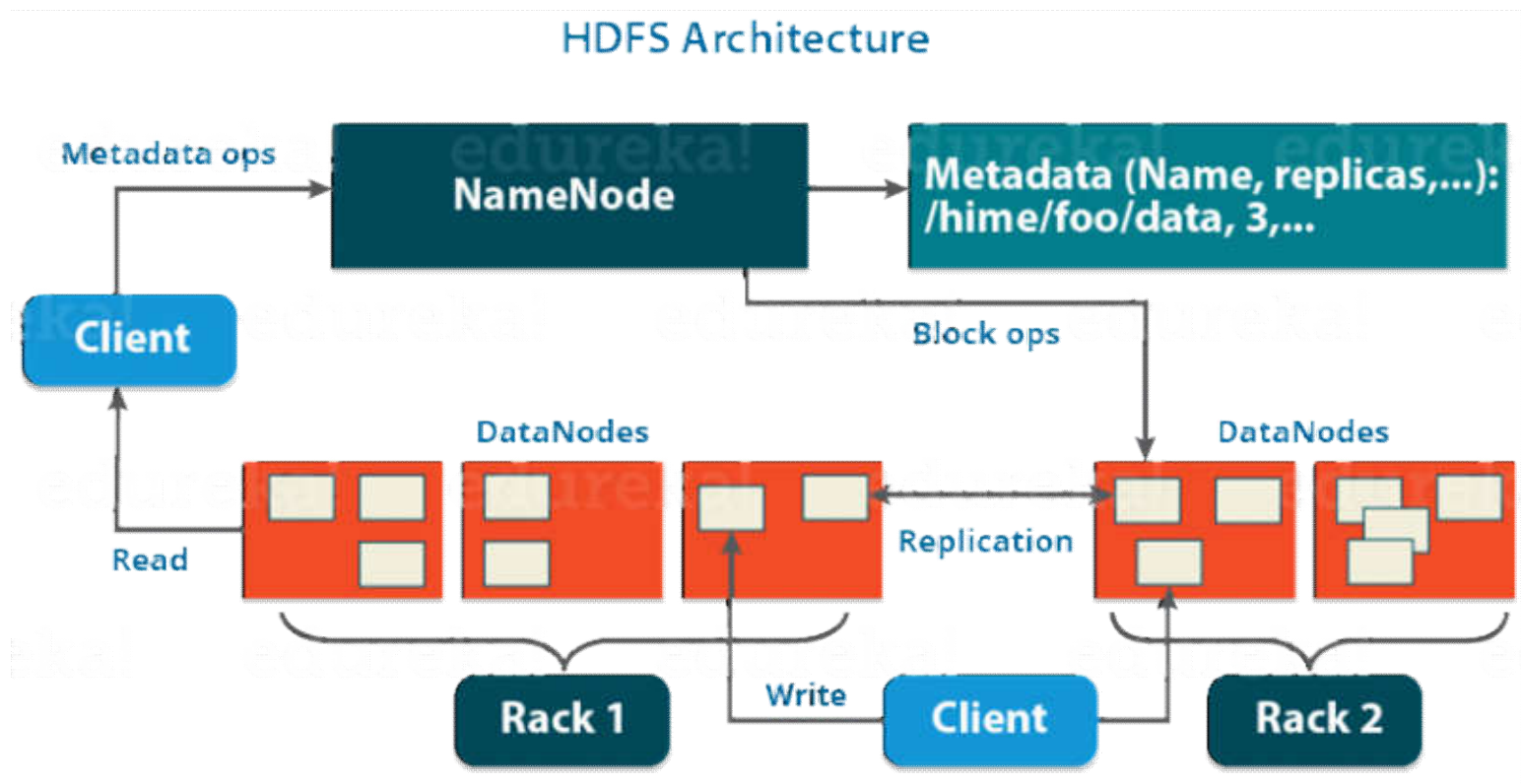
- Actual stored data
- Servicer read and write to Clients.
- Sent Heartbeats to NameNode

HDFS Architecture :

Data Management:

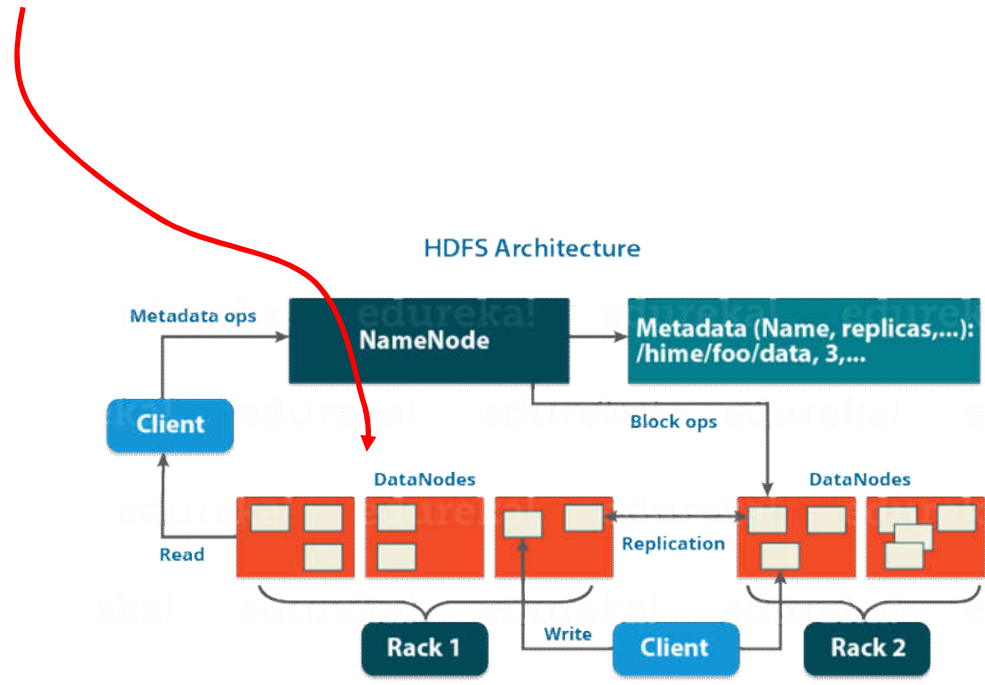
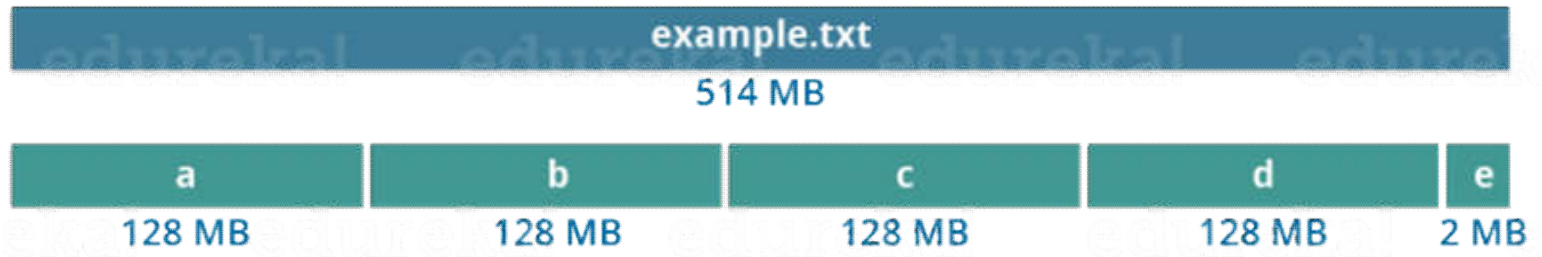
HDFS | Data Block

Replication | Rack Awareness

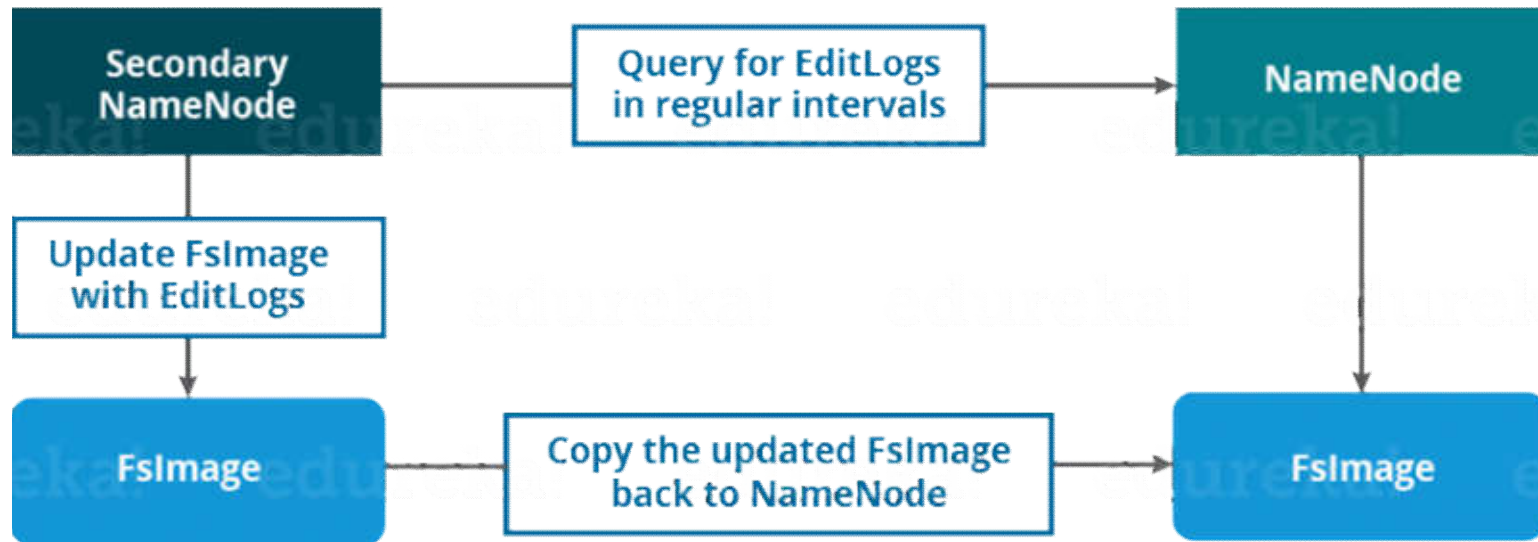


Data Block :

64MB , 128MB / Block



Secondary NameNode :

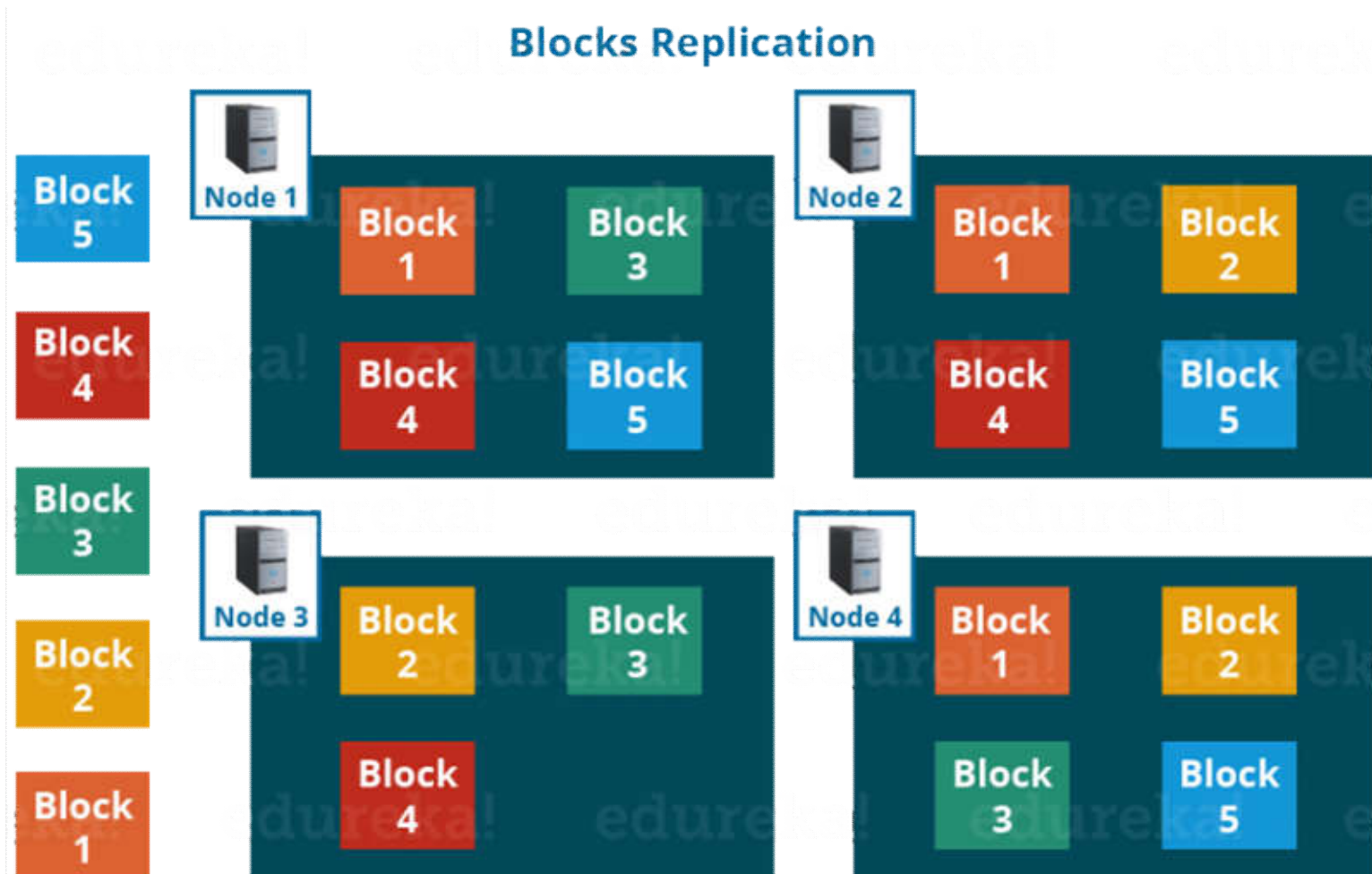


don't be confused about the Secondary NameNode being a **backup NameNode** because it is not.

Automatic backup Metadata of Master NameNode , Metadata , FsImage , EditLogs

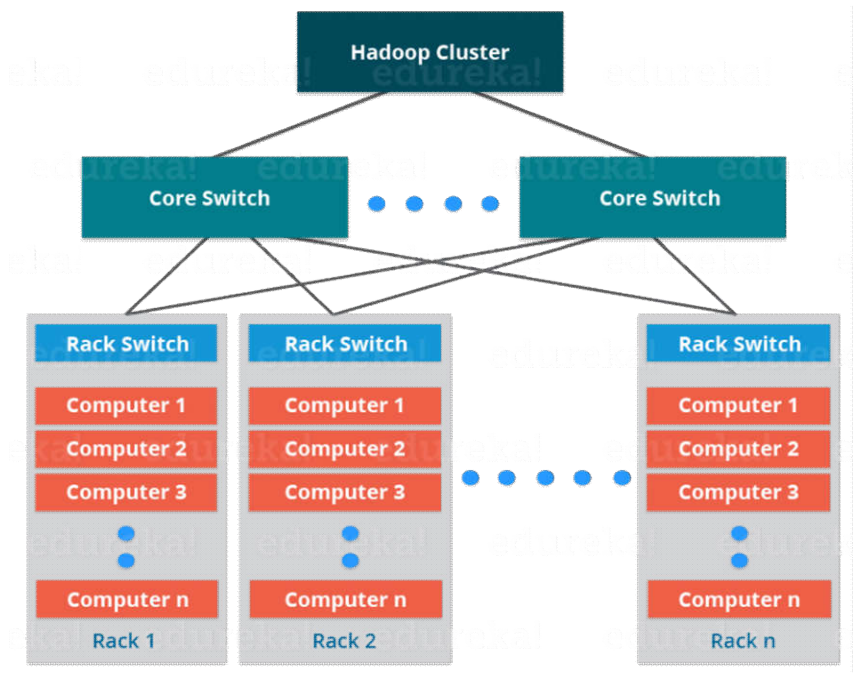
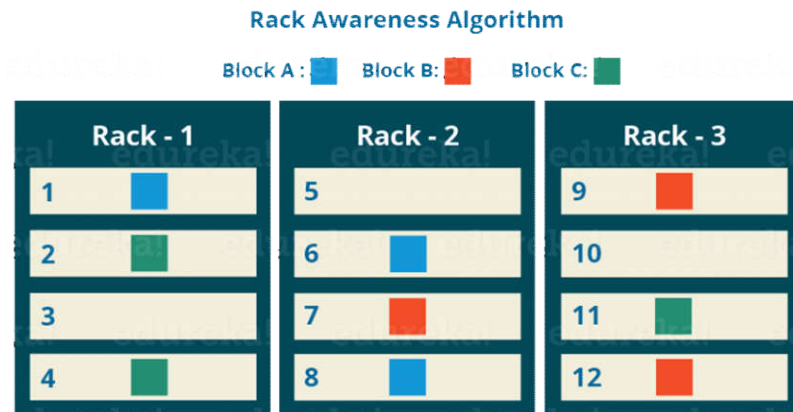
Secondary NameNode act as Checkpoint Node

Replication Management :



Hadoop Topology :

Rack Awareness :

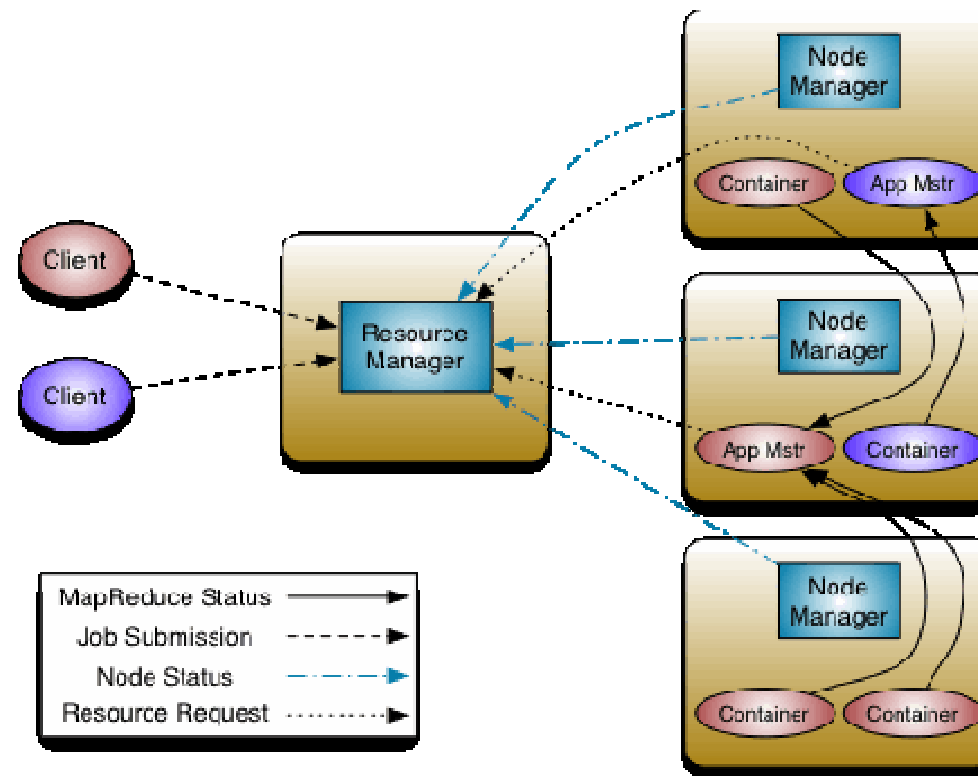


Balance the network traffic for reduce latency

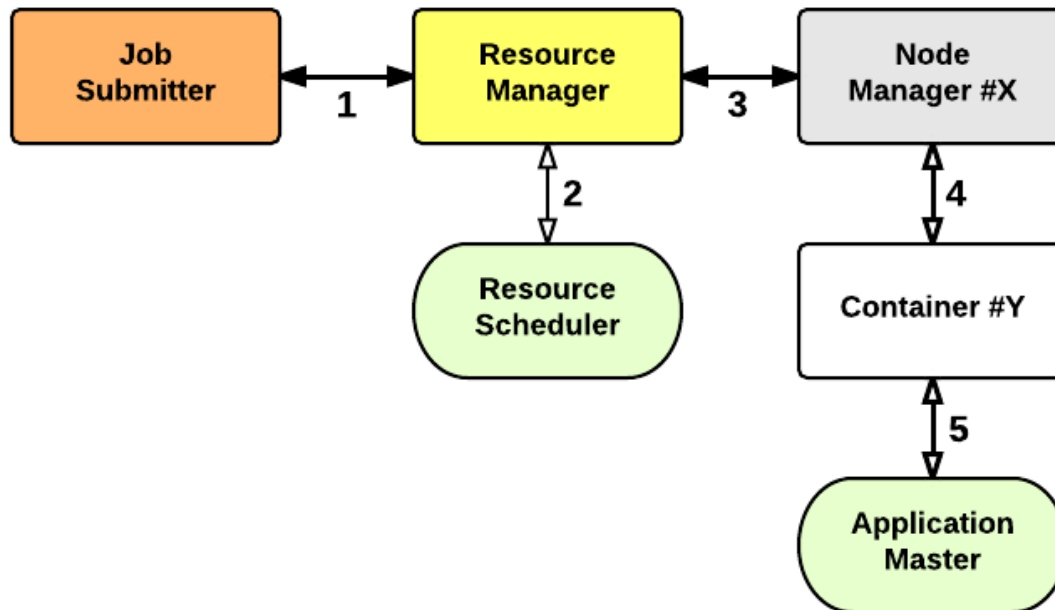
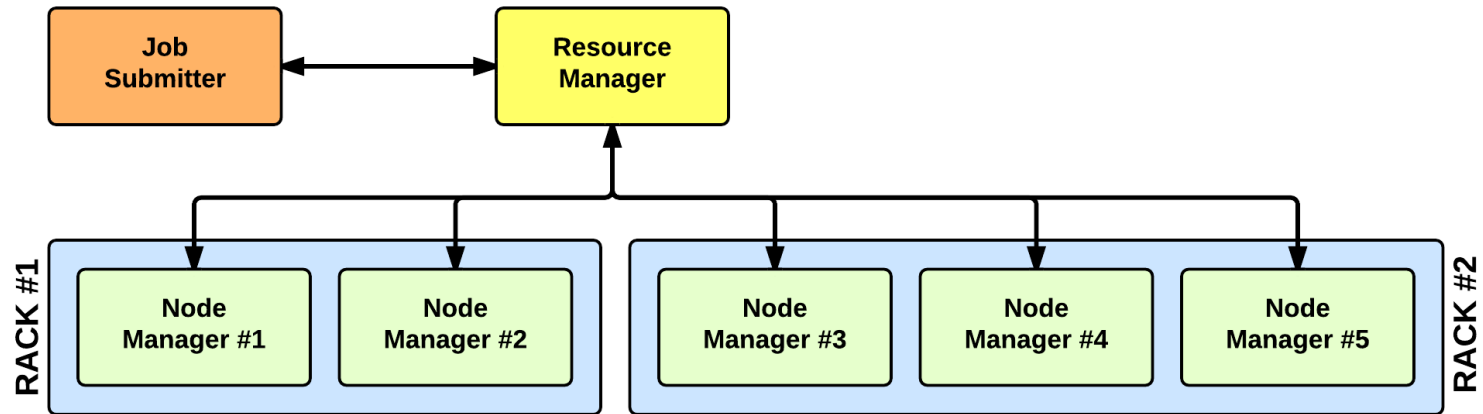
YARN Architecture :

Process & Resource Management :

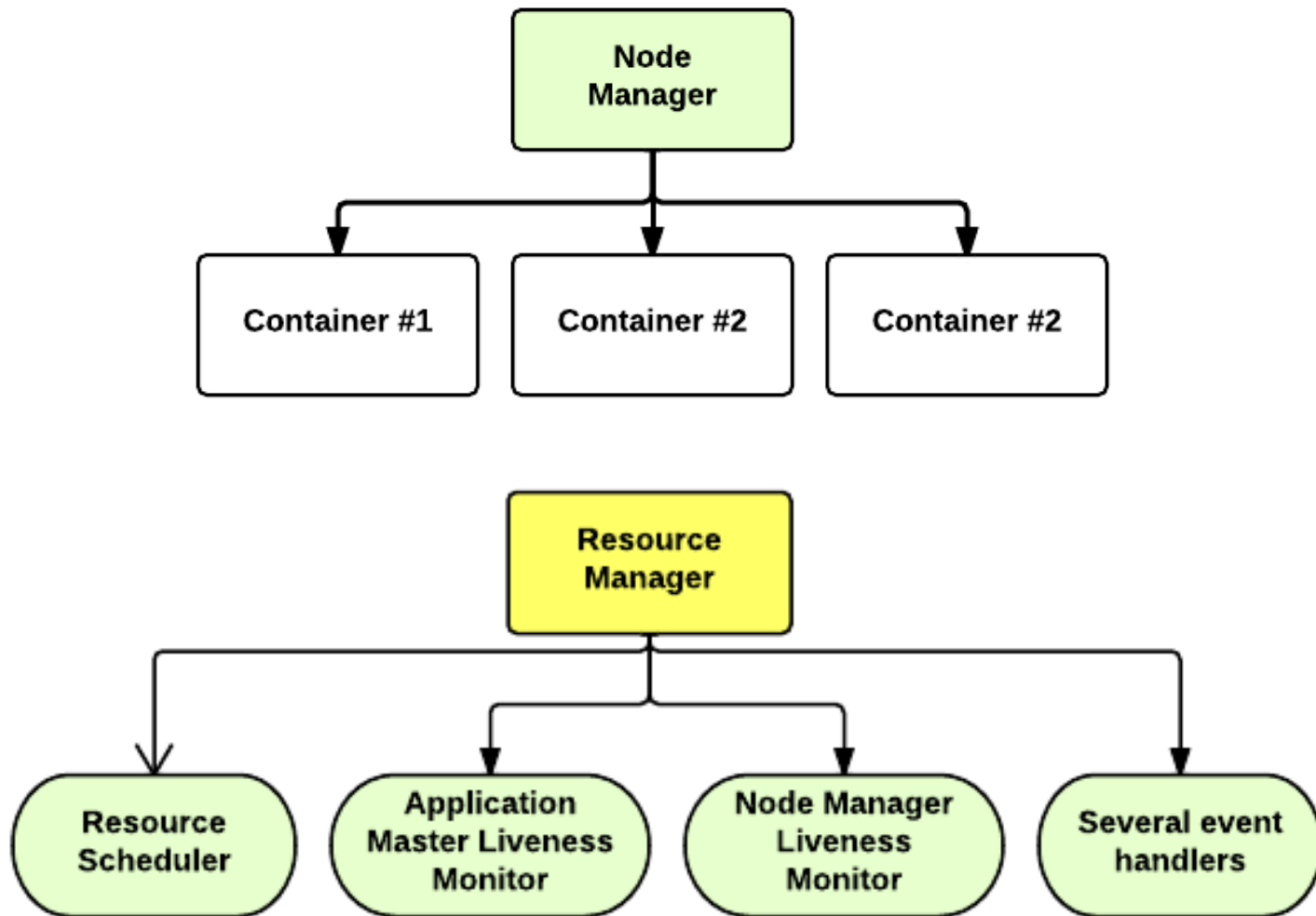
YARN



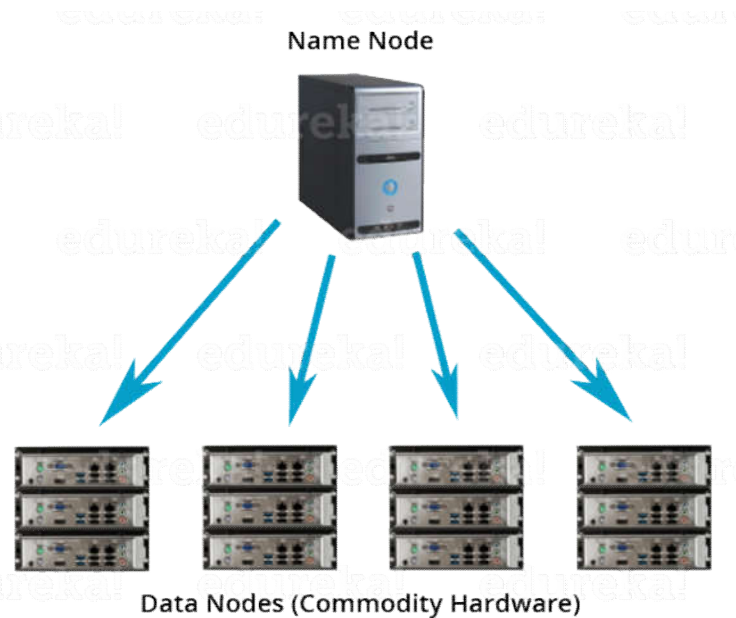
YARN Architecture :



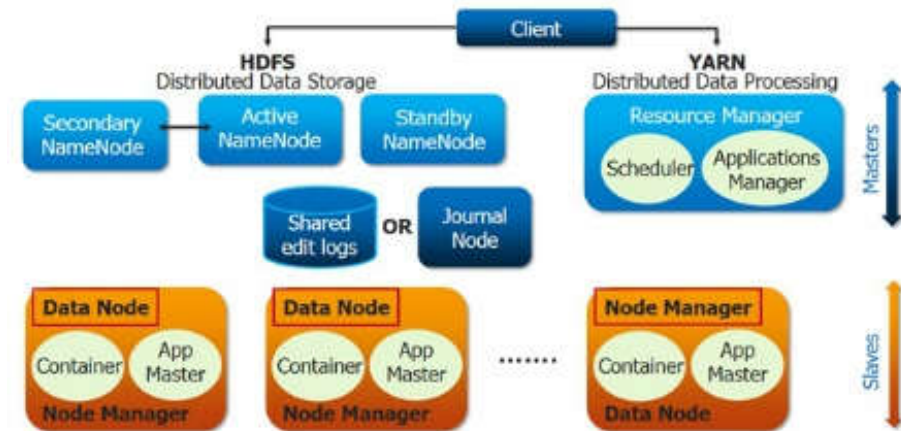
YARN Architecture :



HDFS & YARN Daemon Process :



Apache Hadoop 2.0 and YARN



MasterNode:

NameNode , ResourceManager

SlaveNode:

DataNode , NodeManager

The End

Big
data

Shift