# OS for Big Data – Hadoop Architecture

Prof Bharati Bhole

# Unit II Syllabus

- Data Product, Building Data Products at Scale with Hadoop, Data Science Pipeline and Hadoop Ecosystem

  Ref: (Chapter 1 - Data Analytics with Hadoop By Benjamin Bengfort & Jenny Kim )

- Operating System for Big Data: Concepts, Hadoop Architecture, Working with Distributed file system, Working with Distributed Computation

  Ref: (Chapter 2 - Data Analytics with Hadoop By Benjamin Bengfort & Jenny Kim )

- Framework for Python and Hadoop Streaming, Hadoop Streaming, MapReduce with Python, Advanced MapReduce.

  Ref: (Chapter 3 - Data Analytics with Hadoop By Benjamin Bengfort & Jenny Kim )

- In-Memory Computing with Spark, Spark Basics, Interactive Spark with PySpark, Writing Spark Applications

  Ref: (Chapter 4 - Data Analytics with Hadoop By Benjamin Bengfort & Jenny Kim )

# Today's Topics

**Operating System for Big Data**

- Hadoop Architecture

- Working with Distributed file system

- Working with Distributed Computation

Ref: (Chapter 2 - Data Analytics with Hadoop By Benjamin Bengfort & Jenny Kim )

# Components of Hadoop

Prof Bharati  Bhole

# Hadoop HDFS

- [Data](Data) is stored in a distributed manner in HDFS. There are two components of HDFS - name node and data node. While there is only one name node, there can be multiple data nodes.

- Hadoop enables you to use commodity machines as your data nodes. This way, you don't have to spend millions of dollars just on your data nodes. However, the name node is always an enterprise server.
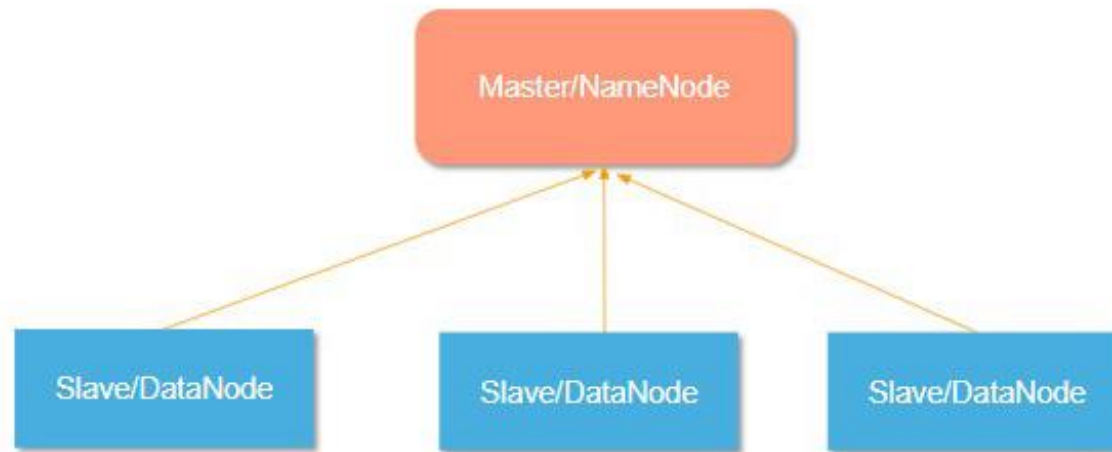
# Hadoop HDFS...

Features of HDFS

- Provides distributed storage

- Can be implemented on commodity hardware

- Provides data security

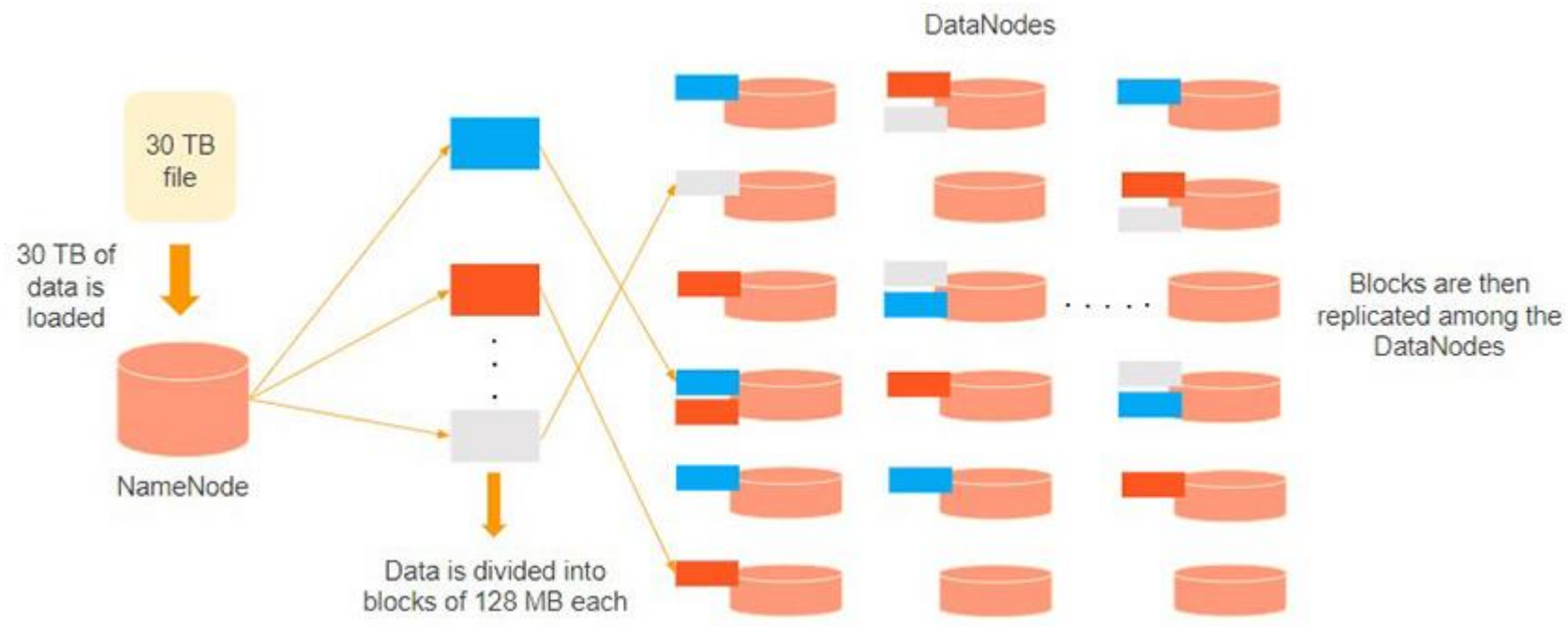- Highly fault-tolerant - If one machine goes down, the data from that machine goes to the next machine

# Master and Slave Nodes

- Master and slave nodes form the HDFS cluster. The name node is called the master, and the data nodes are called the slaves.

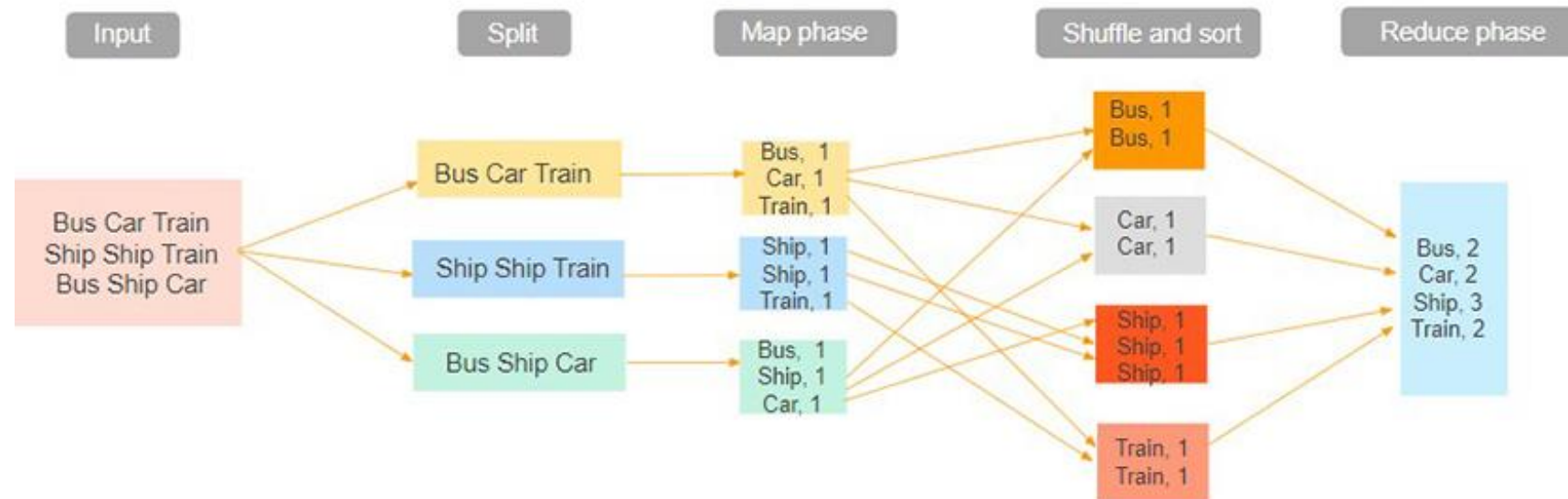- The name node is responsible for the workings of the data nodes. It also stores the metadata.

# Master and Slave Nodes

- The data nodes read, write, process, and replicate the data. They also send signals, known as heartbeats, to the name node. These heartbeats show the status of the data node.
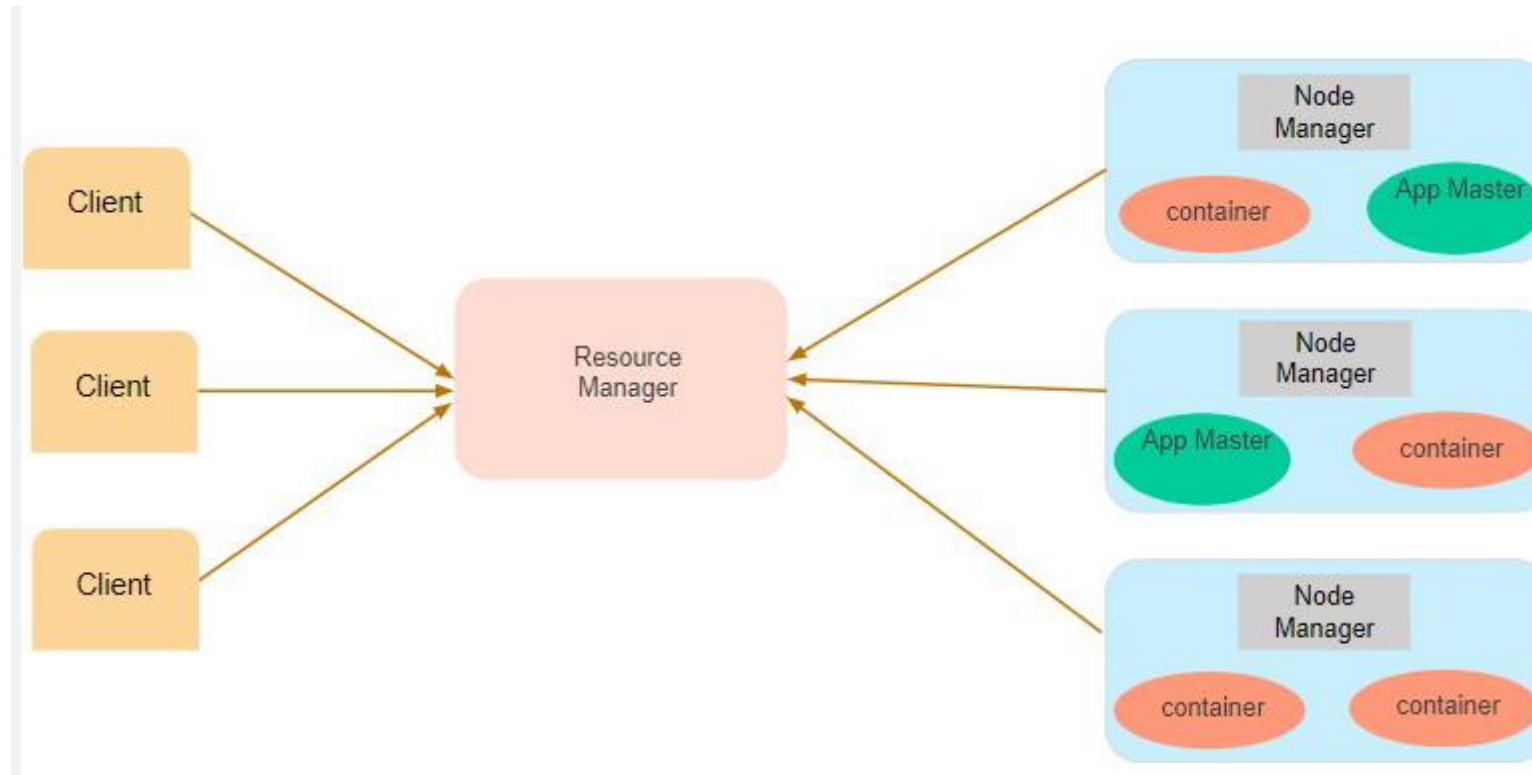
# Hadoop MapReduce

[Hadoop MapReduce](#) is the processing unit of Hadoop. In the MapReduce approach, the processing is done at the slave nodes, and the final result is sent to the master node.

# Hadoop YARN

- [Hadoop YARN](#) stands for Yet Another Resource Negotiator. It is the resource management unit of Hadoop and is available as a component of Hadoop version 2.
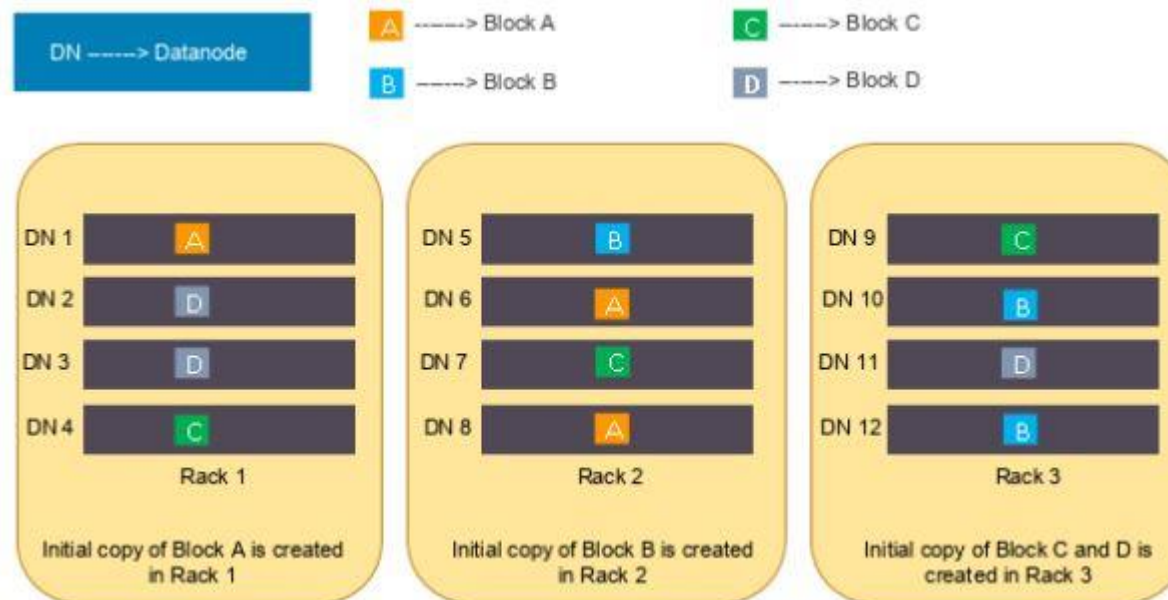
# Architecture of Hadoop

# Hadoop HDFS

Hadoop is a framework permitting the storage of large volumes of data on node systems. The Hadoop architecture  allows parallel processing of data using several components:

- Hadoop HDFS to store data across slave machines
- Hadoop YARN for resource management in the Hadoop cluster
- Hadoop MapReduce to process data in a distributed fashion
- Zookeeper to ensure synchronization across a cluster

# Hadoop HDFS...

- DFS in Hadoop Architecture divides large data into different blocks. Replicated three times by default, each block contains 128 MB of data. Replications operate under two rules:

1. Two identical blocks cannot be placed on the same DataNode

2. When a cluster is rack aware, all the replicas of a block cannot be placed on the same rack
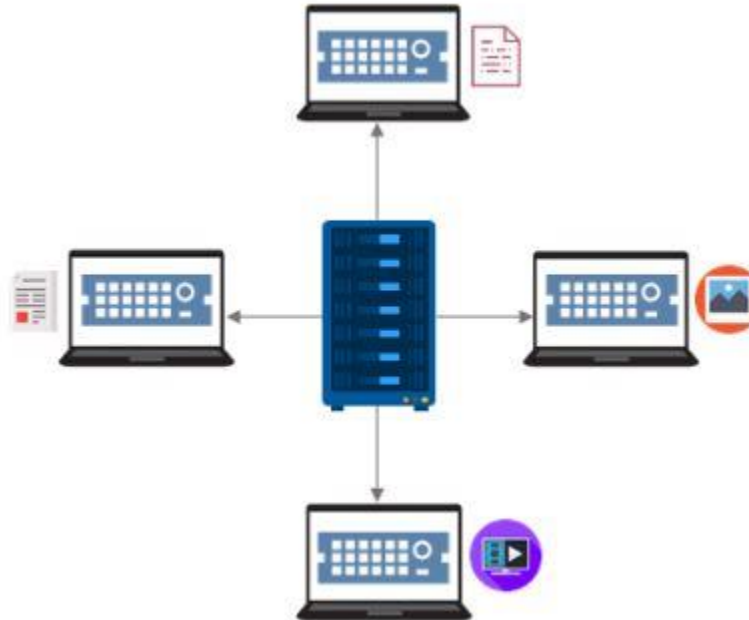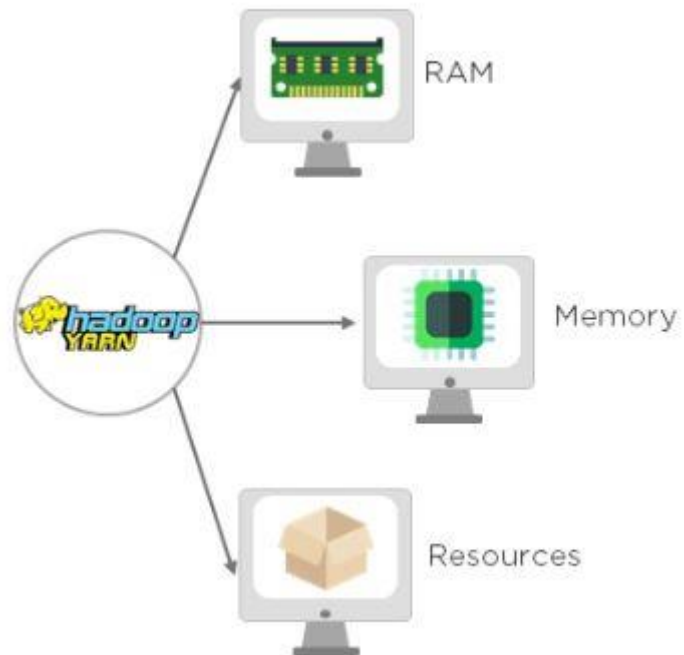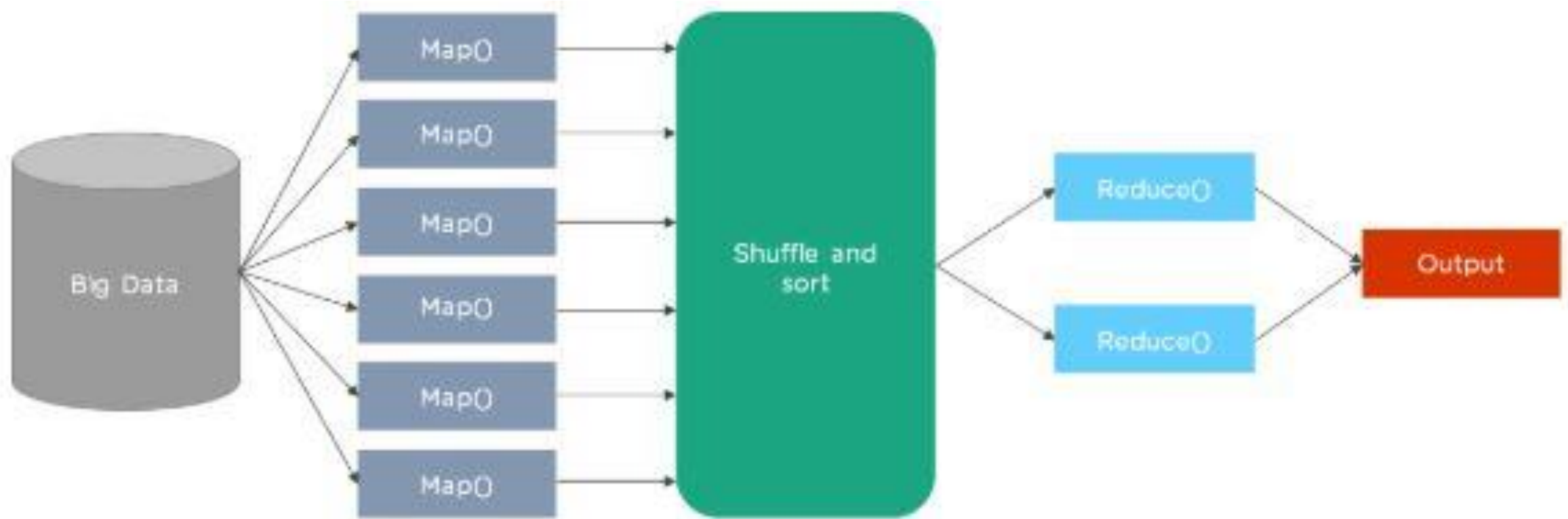
# Hadoop Ecosystem

Prof Bharati  Bhole

# Hadoop Ecosystem

# HDFS

Prof Bharati  Bhole

# YARN

Prof Bharati  Bhole

# MapReduce

# Sqoop



Hadoop data

Relational database and enterprise data warehouse

# Flume

Prof Bharati  Bhole
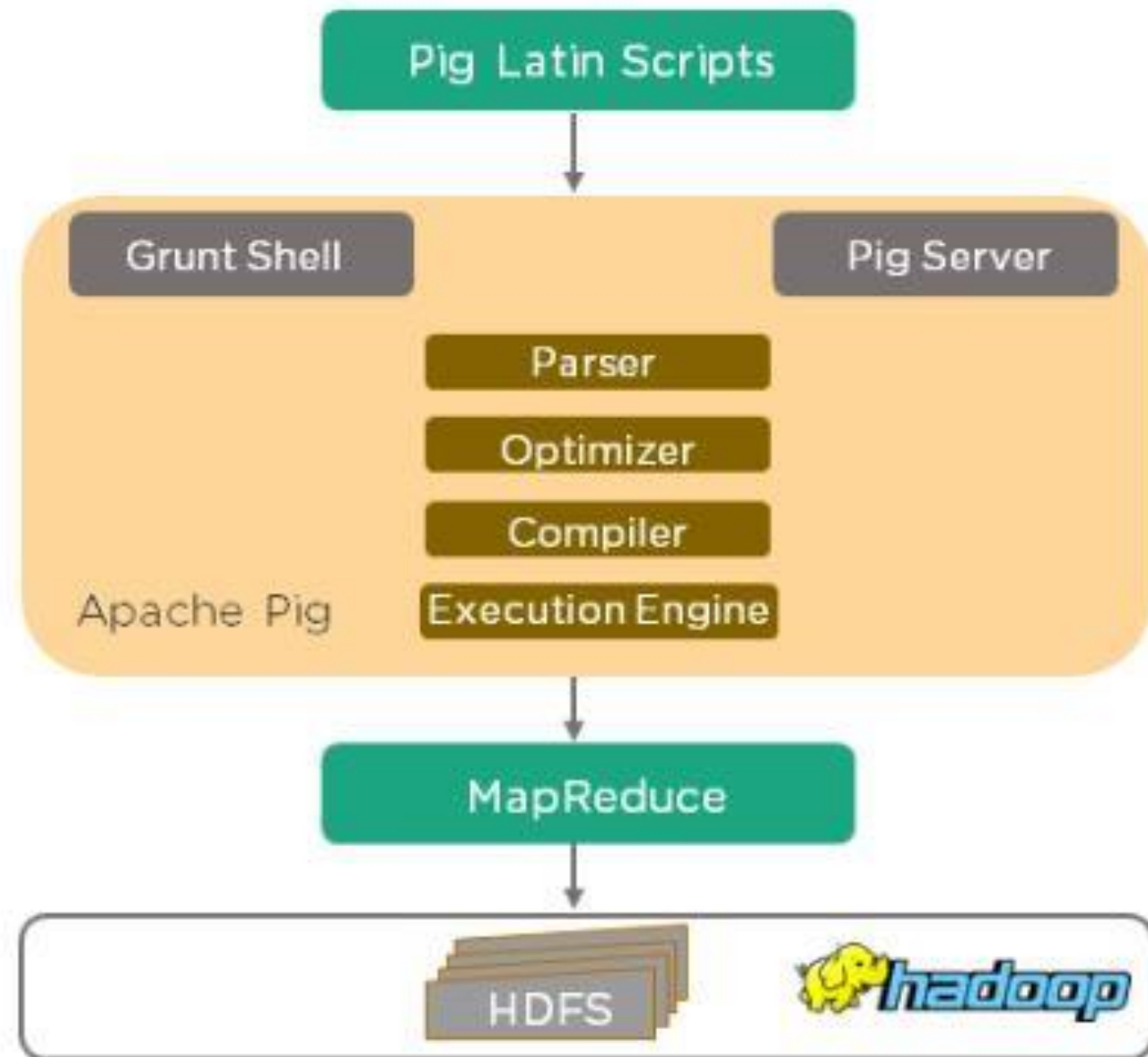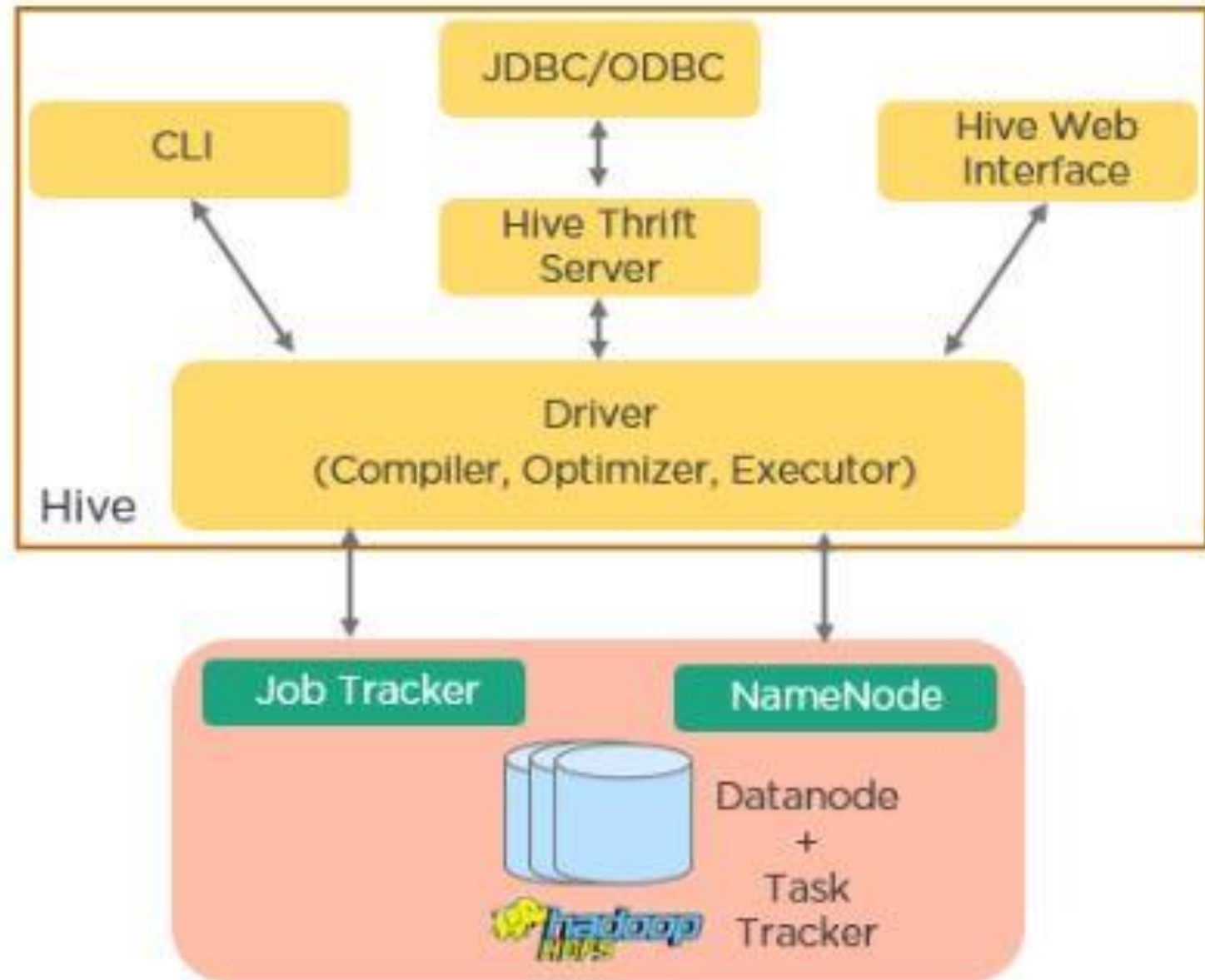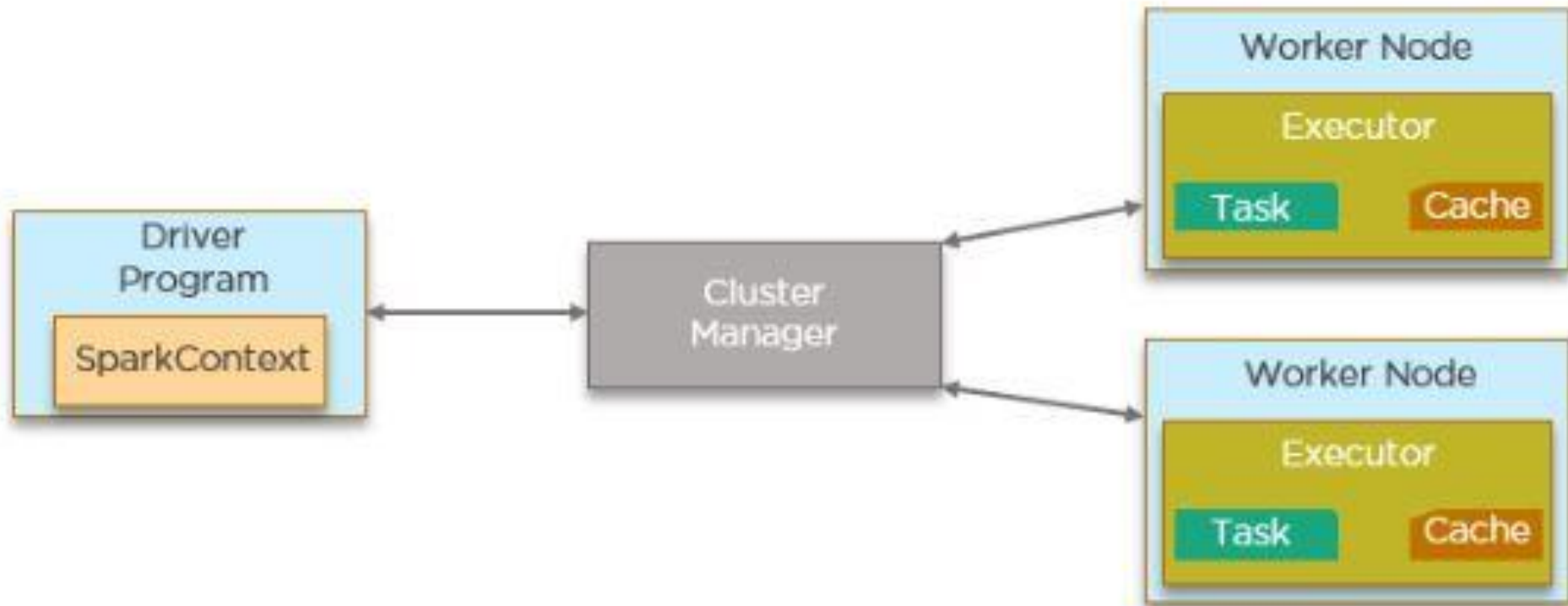
# Pig

1.Pig Latin - This is the language for scripting
2.Pig Latin Compiler - This converts Pig Latin code into executable code

Prof Bharati Bhole
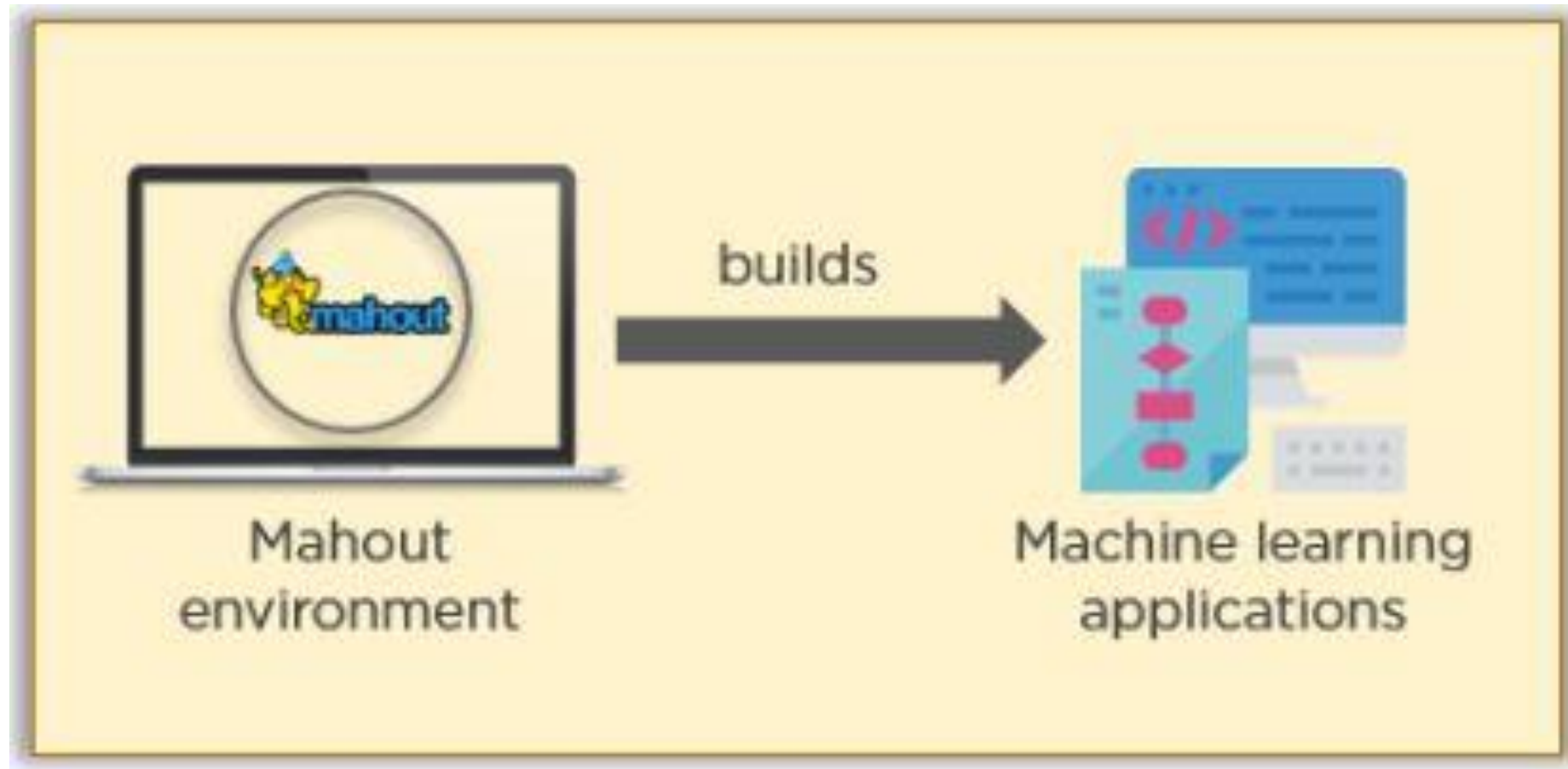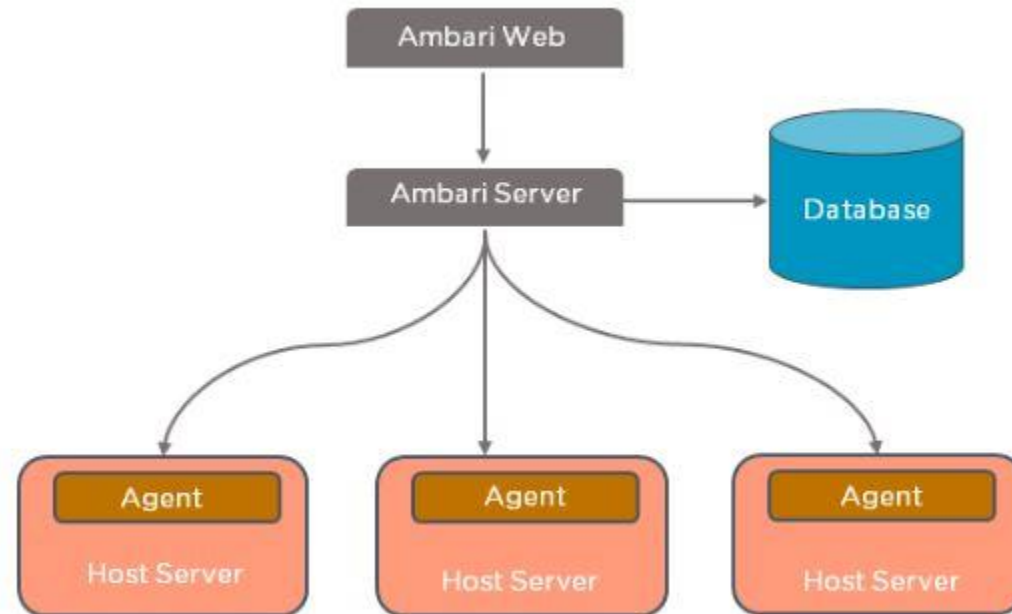
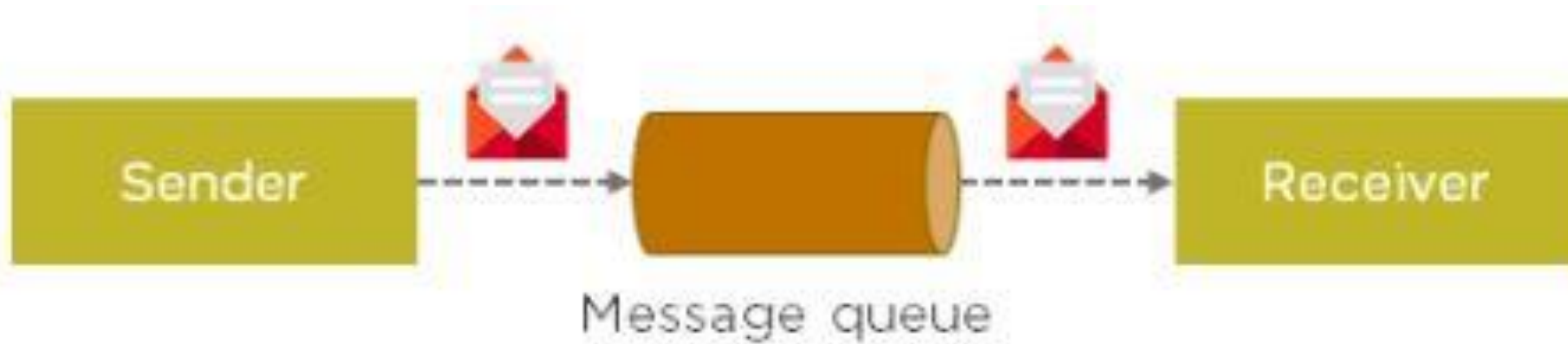# Hive

# Spark

# Mahout

# Ambari

# Kafka

Prof Bharati Bhole

# Storm

# Ranger

Prof Bharati  Bhole

# Knox



Security

Apache Ranger

# Oozie



Workflow system

OOZIE

# Thank You….

# Revise the topics from Syllabus References…

Fill Your Attendance Form….!

Prof Bharati Bhole

# Syllabus References

1. Big Data and Analytics, [Subhashini Chellappan Seema Acharya](), Wiley

2. Data Analytics with Hadoop *An Introduction for Data Scientists,* Benjamin Bengfort and Jenny Kim, O'Reilly

3. Big Data and Hadoop, V.K Jain, Khanna Publishing

   https://books.google.co.in/books?id=i6NODQAAQBAJ&pg=PA122&source=gbs_toc_r&cad=4#v=onepage&q&f=true