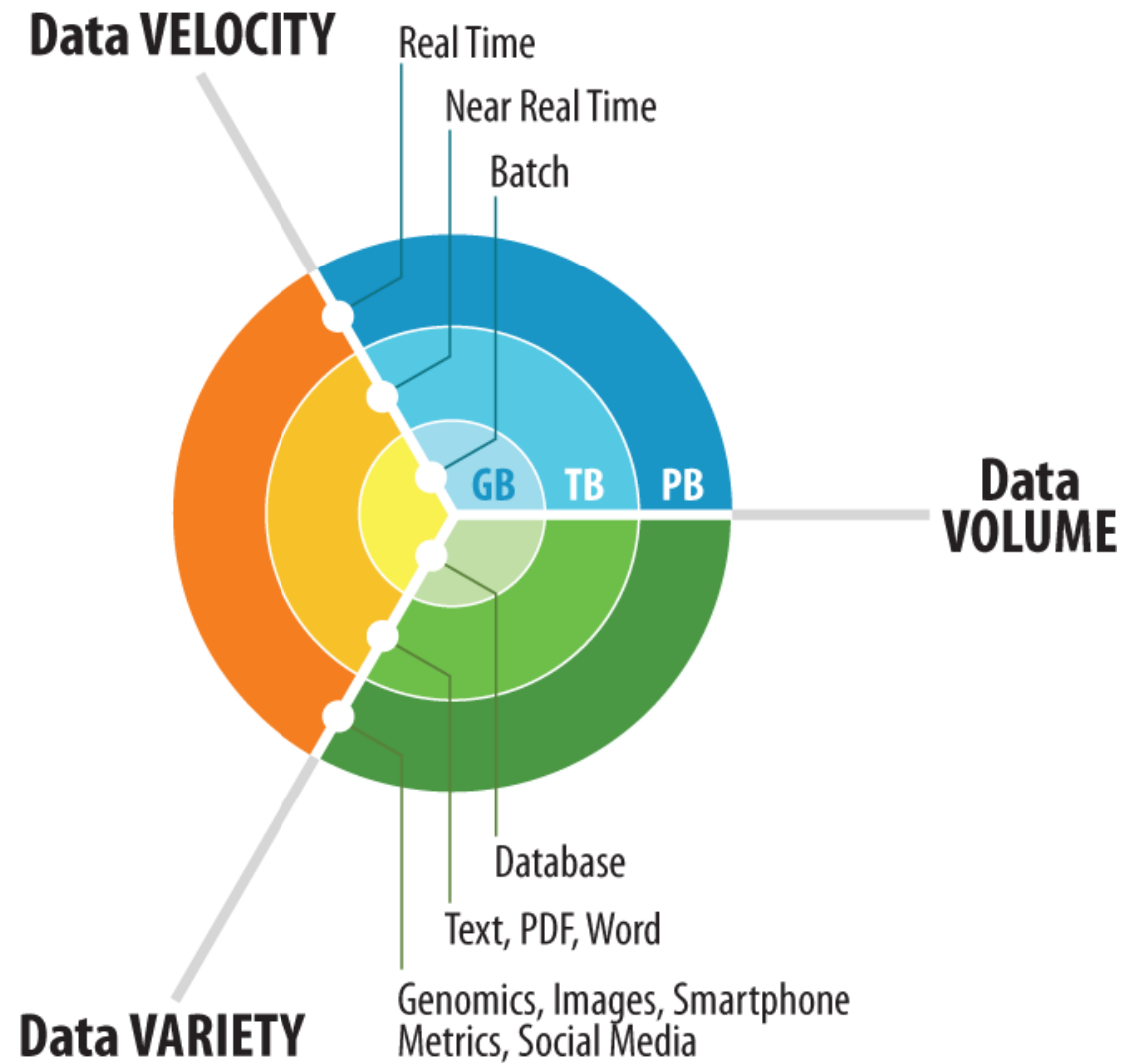


# Big Data

Tianchu.Zhao@uts.edu.au

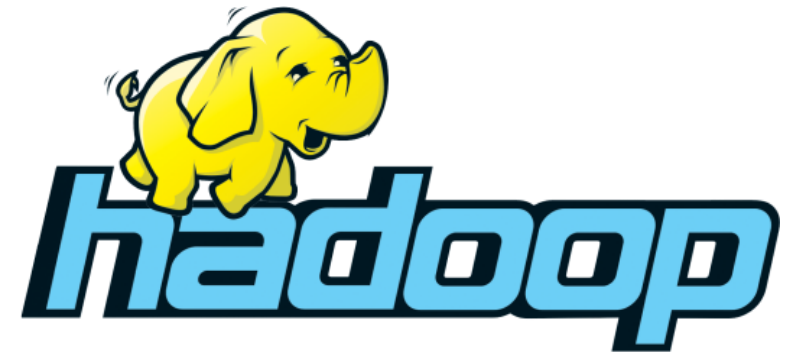
# Big Data?

- **Big Data** is a term for data sets that are so large or complex that traditional data processing applications are inadequate to deal with them.

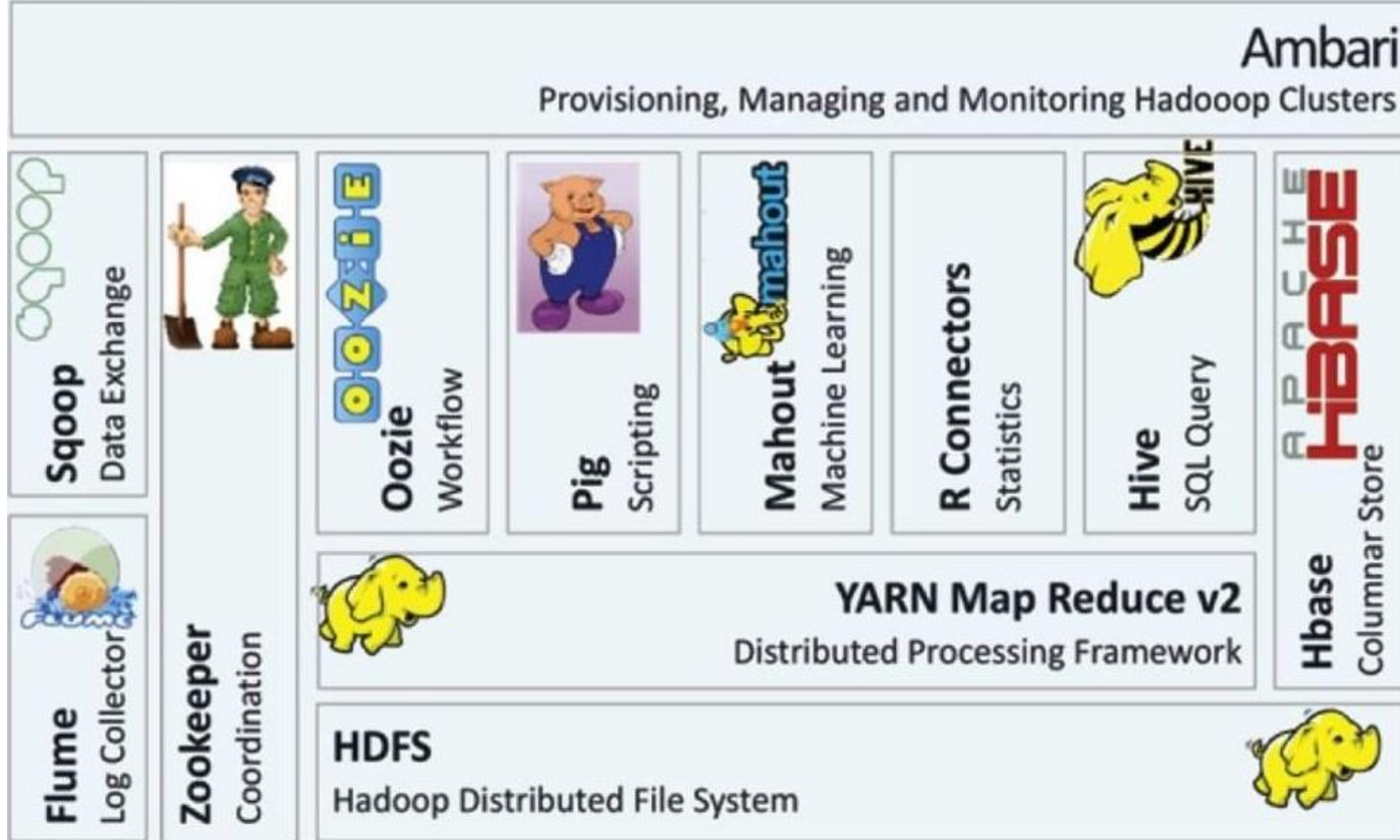


# Hadoop

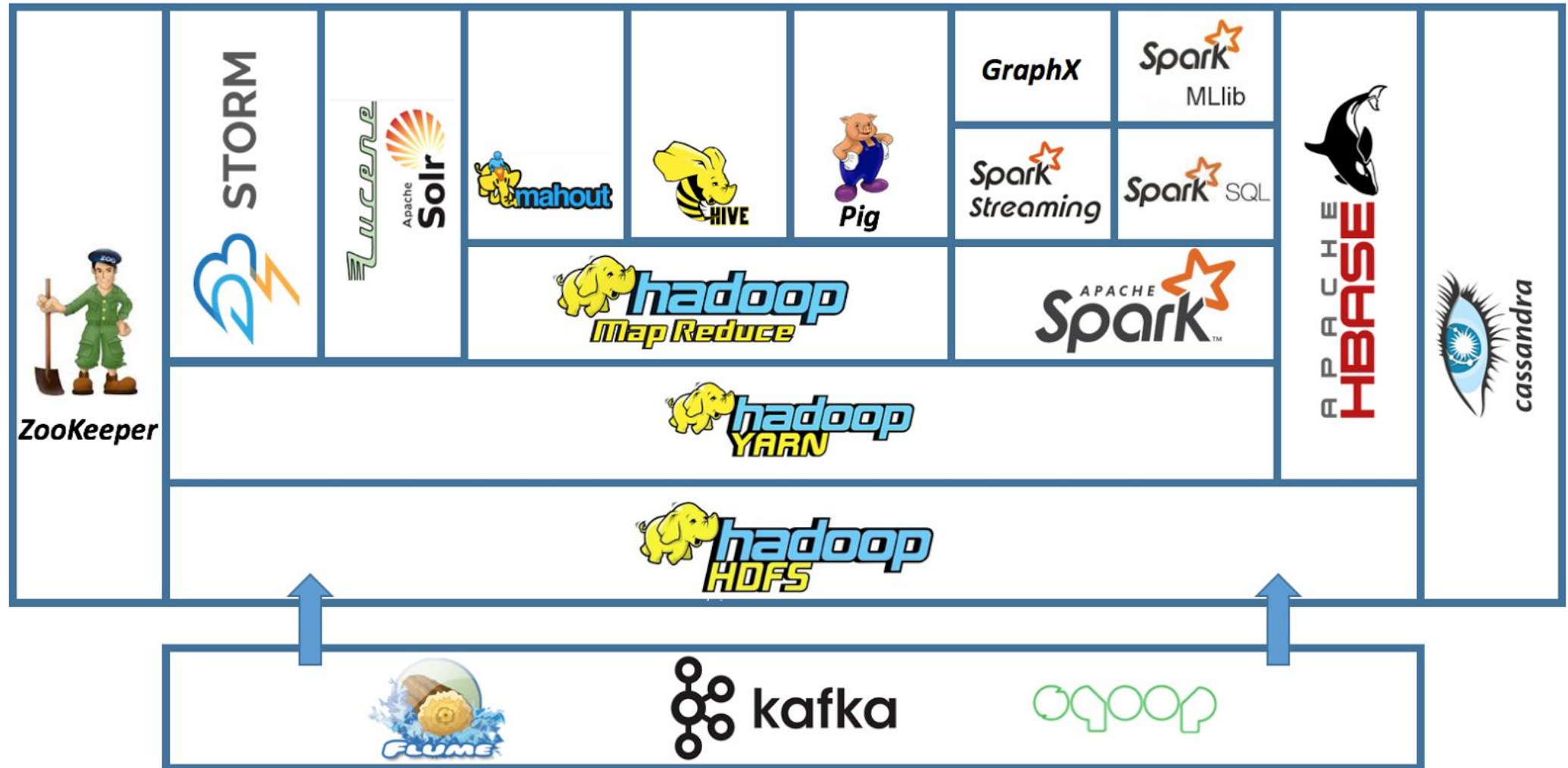
- (Wikipedia) Apache Hadoop is a collection of open-source software utilities that facilitate using a network of many computers to solve problems involving massive amounts of data and computation.



# Hadoop ecosystem



# Hadoop ecosystem



# NoSQL Database

- Key Value
  - Dynamo DB
  - Cassandra
- Column Based
  - HBase
- Document Database
  - MongoDB
  - CouchDB
- Graph Database
  - Polyglot
  - Neo4J

# HDFS

- The Hadoop Distributed File System (HDFS)
- a distributed file system designed to run on commodity hardware
- fault-tolerant
- low-cost hardware

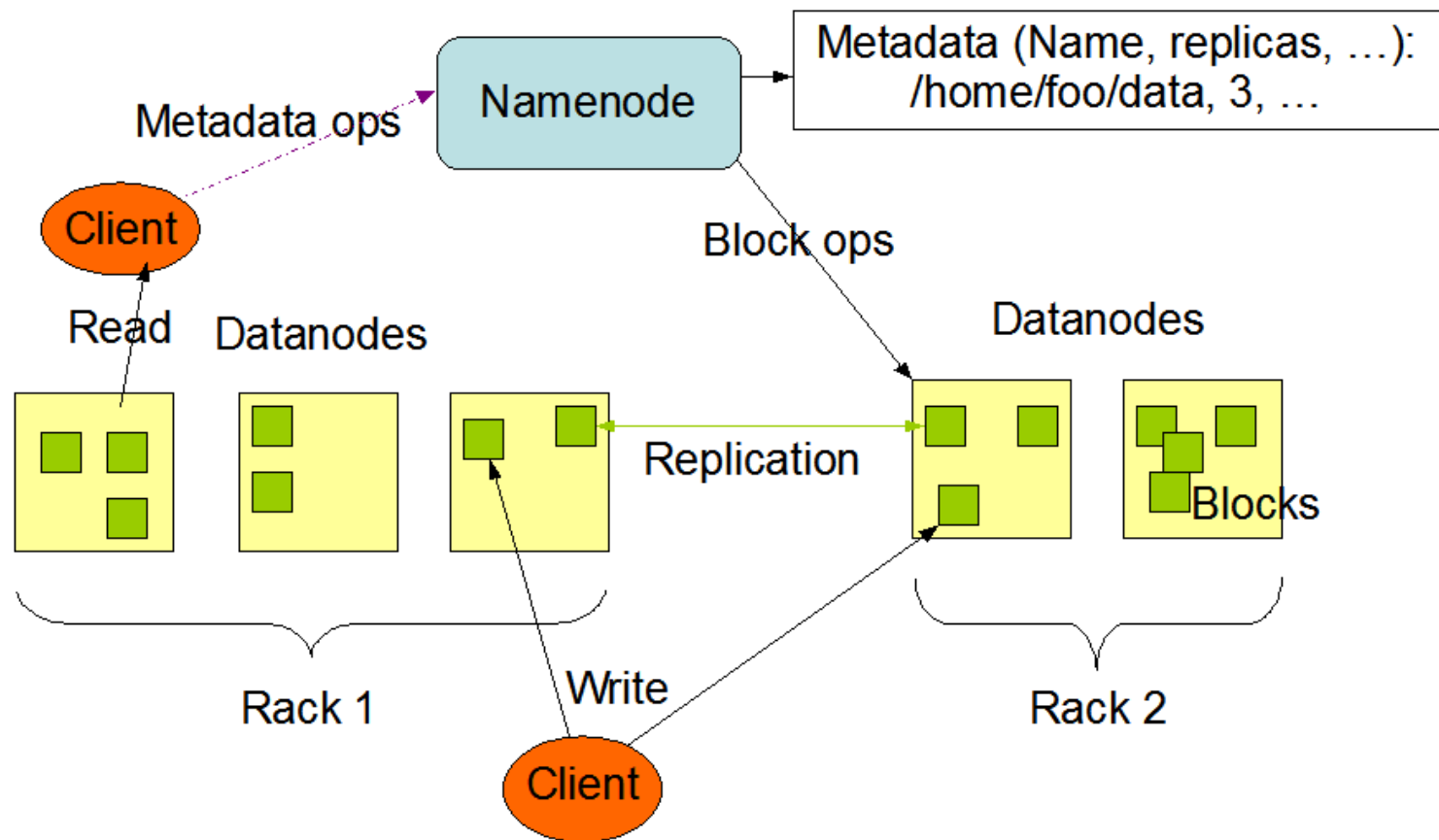


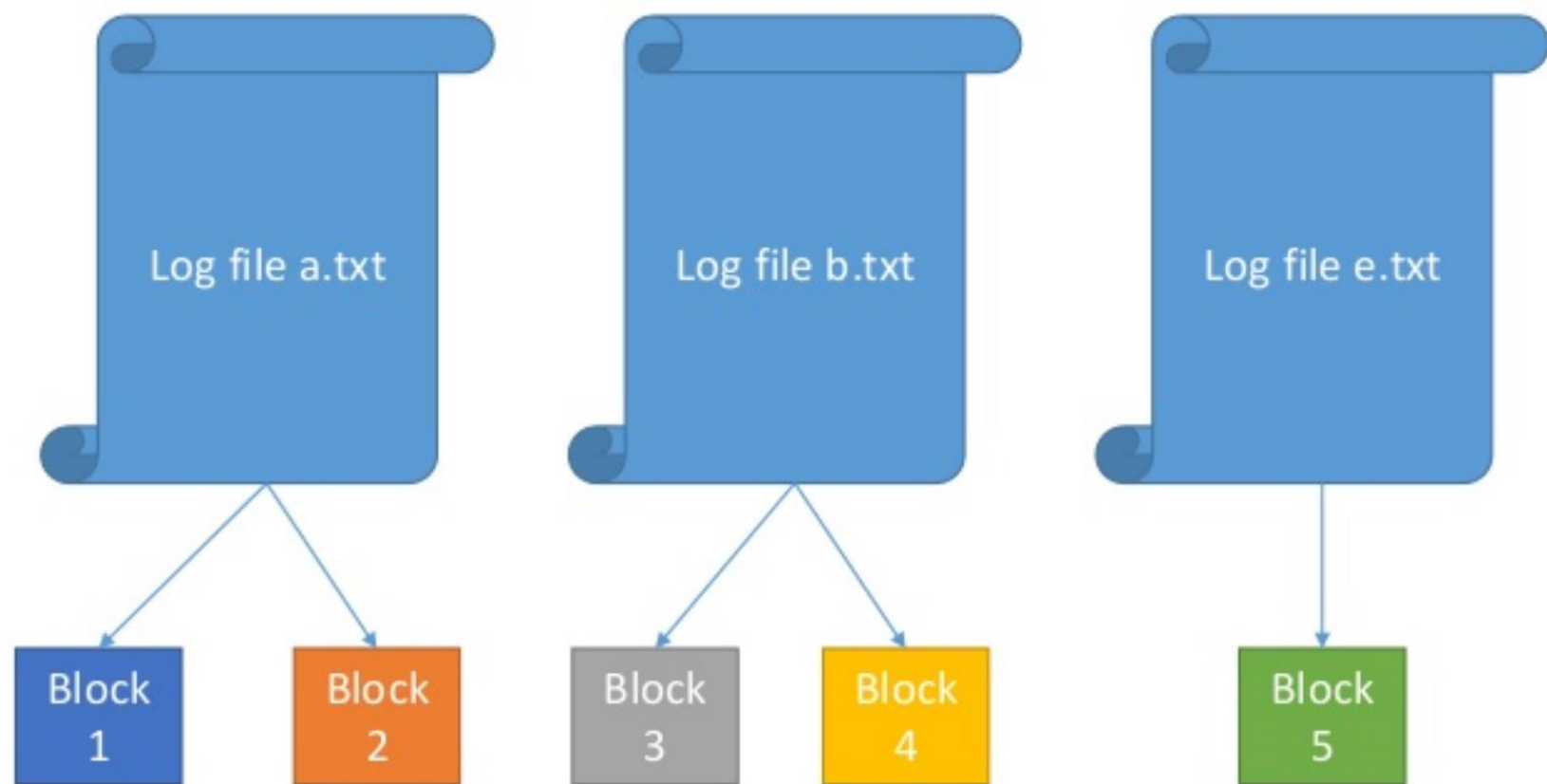


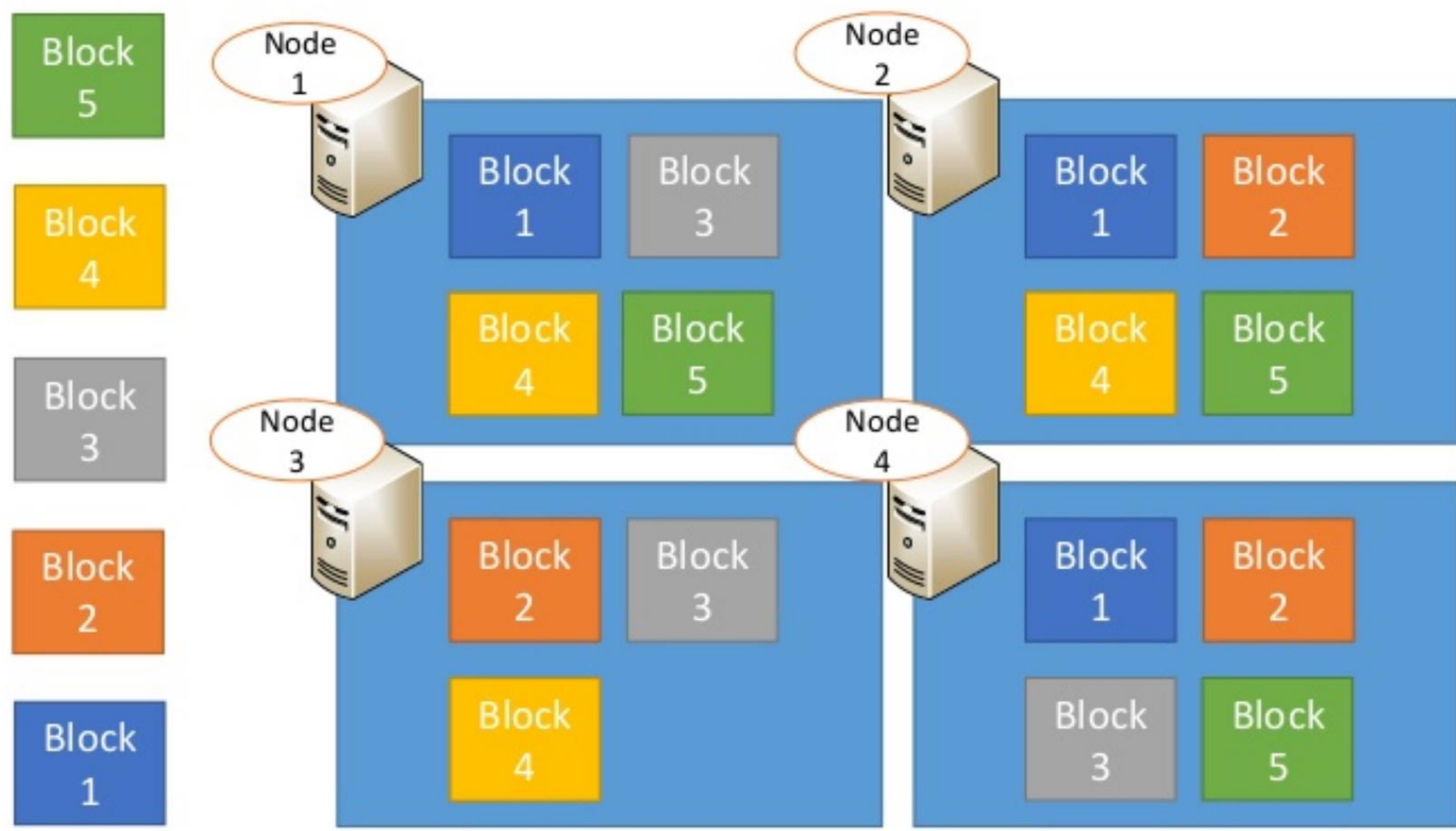
# HDFS Architecture

- HDFS has a master/slave architecture
- It consists of NameNode
  - A master server that manages the file system namespace and regulates access to files by clients
- and DataNode
  - Responsible for serving read and write requests from the file system's clients

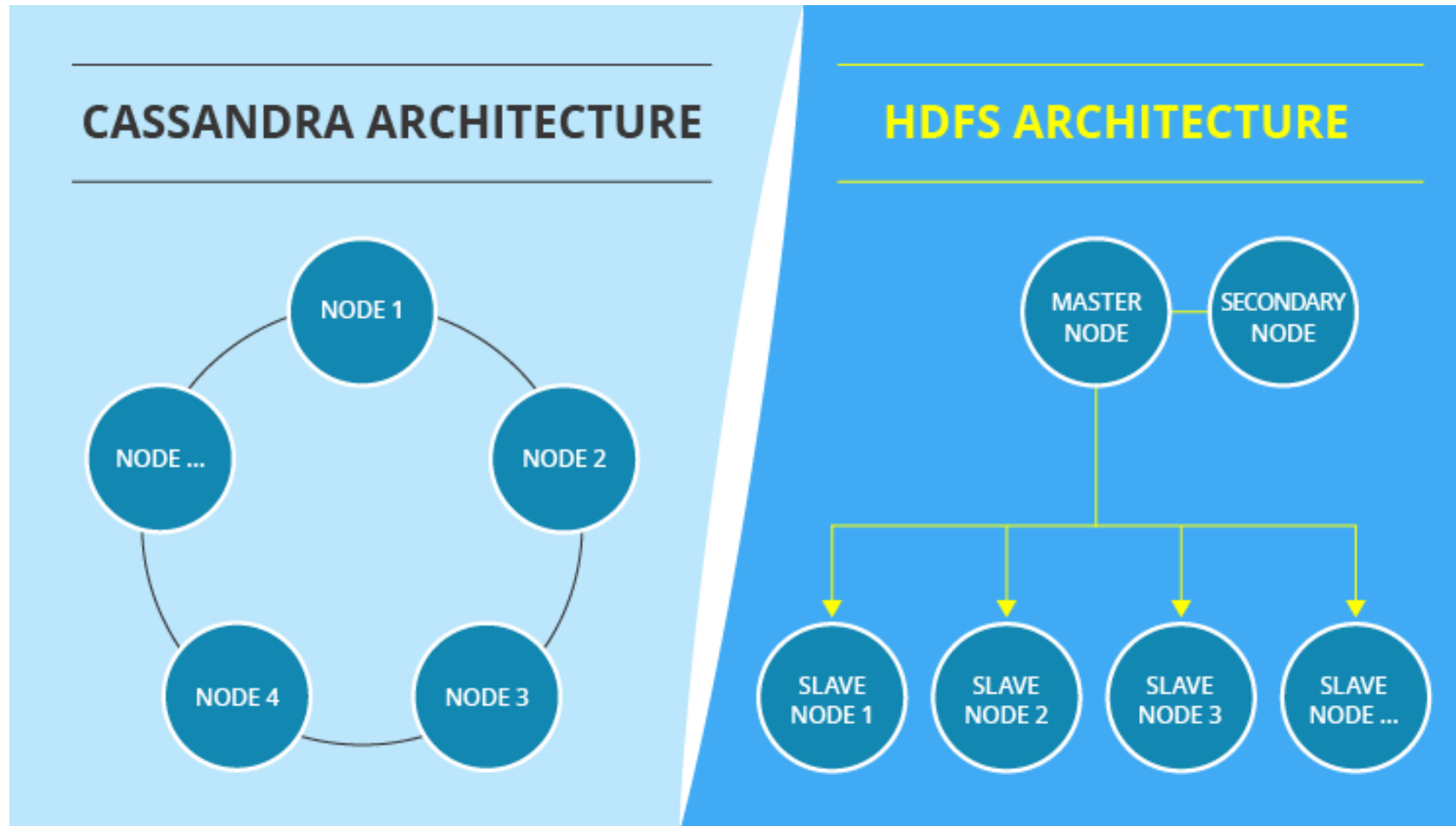
## HDFS Architecture







# Cassandra

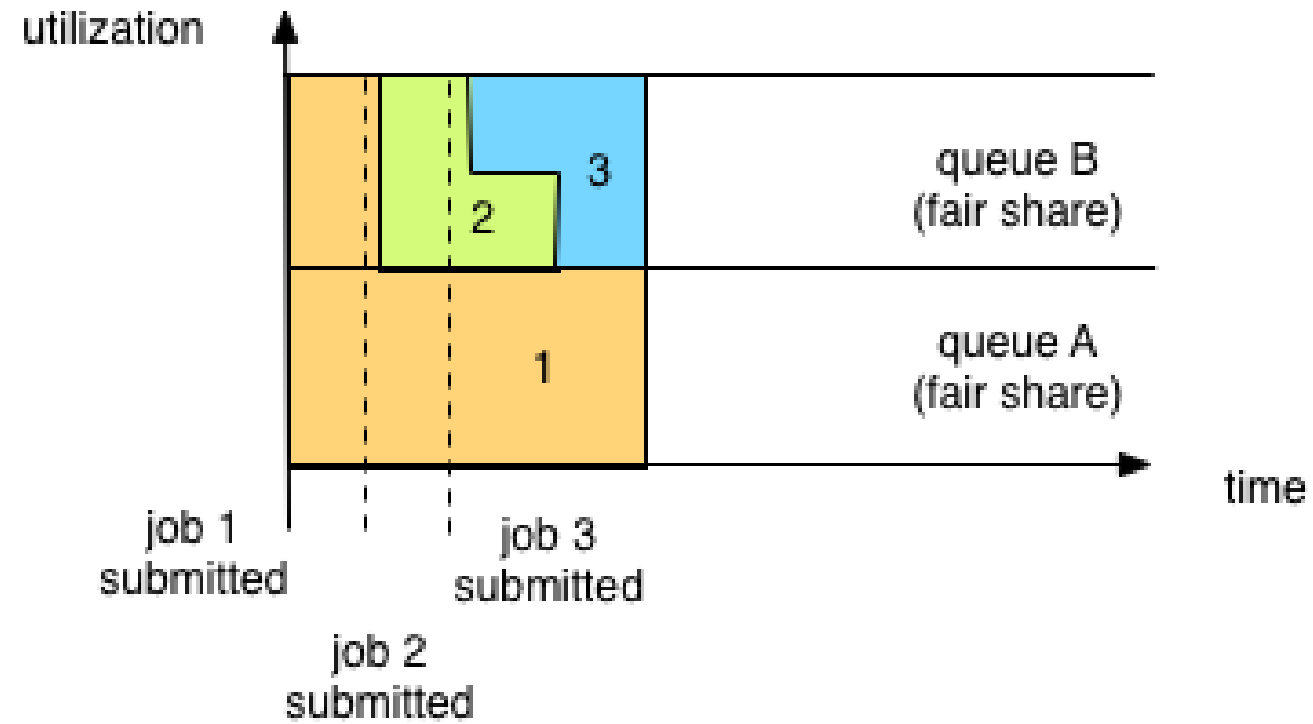


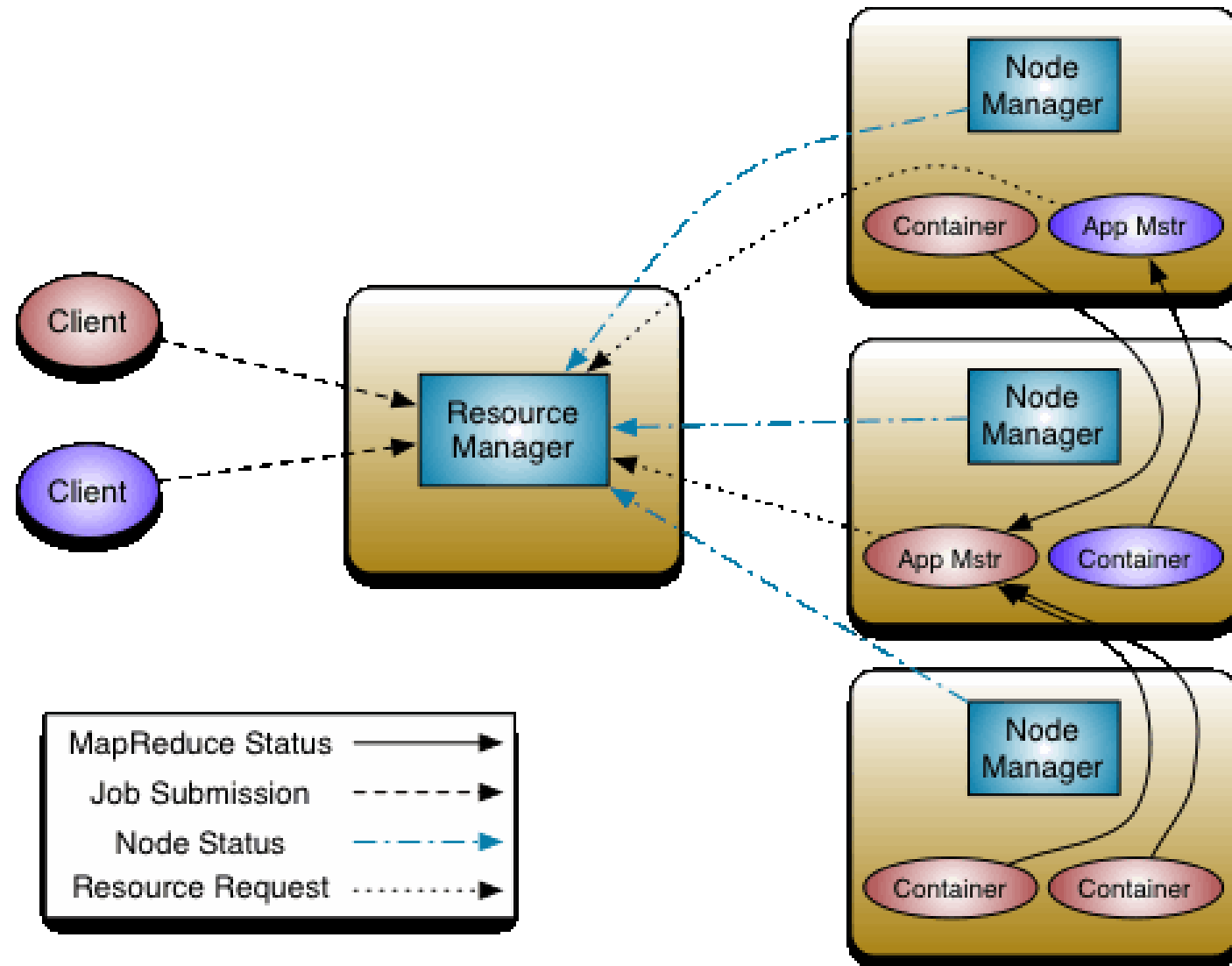
# YARN

- **Hadoop YARN** is a cluster management technology
- The fundamental idea of YARN is to split up the functionalities of resource management and job scheduling/monitoring into separate daemons. The idea is to have a global ResourceManager (*RM*) and per-application ApplicationMaster (*AM*). An application is either a single job or a DAG of jobs.



- Maximise cluster utilisation





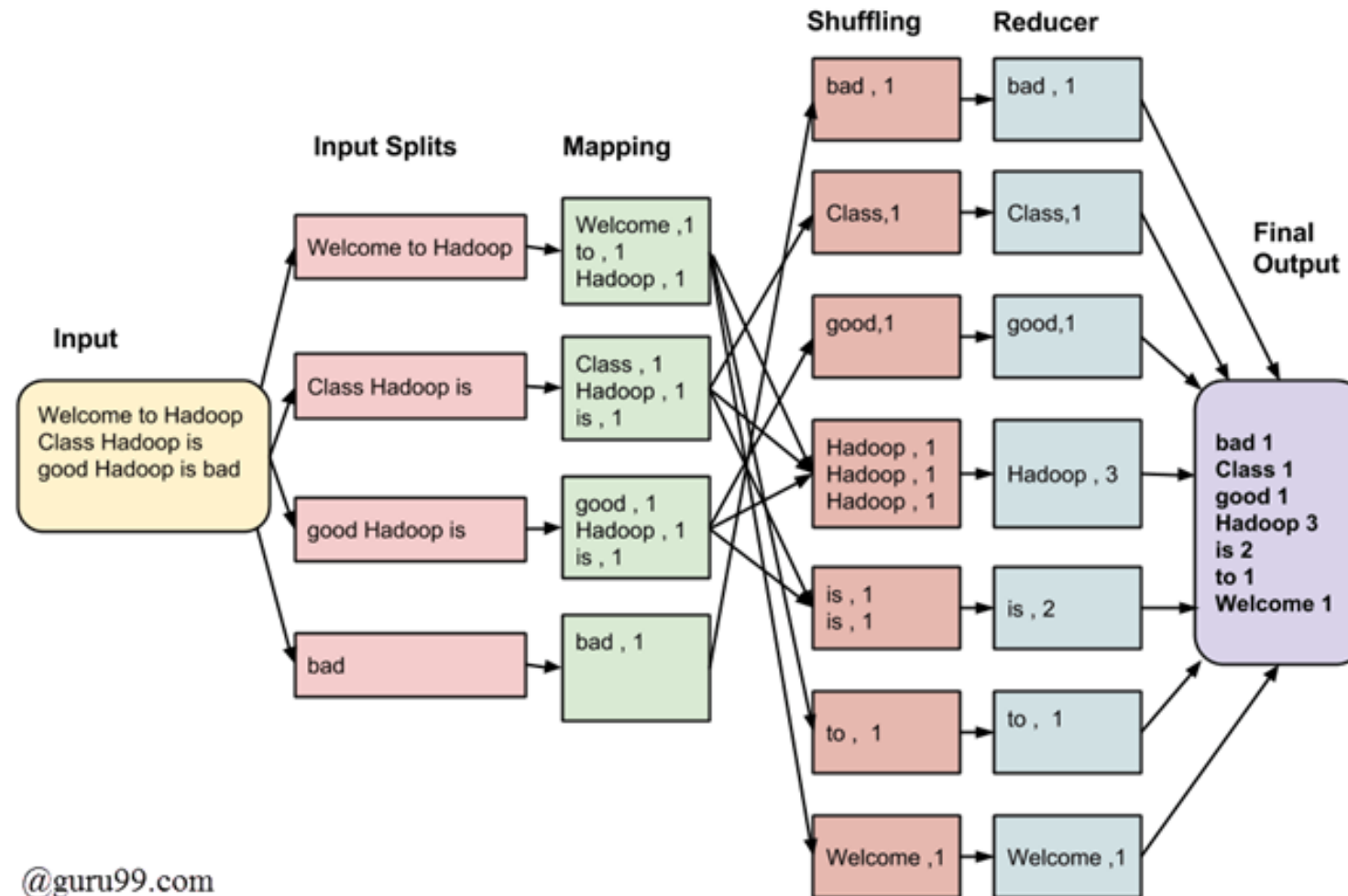


# MapReduce

- MapReduce is a programming model and an associated implementation for processing and generating big data sets with a parallel, distributed algorithm on a cluster.



# MapReduce wordcount



# Spark

- Apache Spark is an open-source distributed general-purpose cluster-computing framework
- data representation: rdd, dataframe
- lazy execution
- In memory

