



# Распределенные Вычисления. Apache Hadoop

[NEWPROLAB.COM](http://NEWPROLAB.COM)

# Что такое распределенная система?

- распределенная система это набор компьютеров, представляющихся пользователю единой системой
- распределенная система это такая система, в которой взаимодействие и синхронизация программных компонентов, выполняемых на независимых компьютерах, осуществляется посредством передачи сообщений
- распределенная система - набор независимых компьютеров, не имеющих совместно используемой памяти и общего единого времени и взаимодействующих посредством передачи сообщений

# Основные признаки распределенных систем

- отсутствие единого времени
- отсутствие общей памяти - каждый узел оперирует только своей физической оперативной памятью (возможны варианты виртуальной адресации)
- географическое распределение вычислительных узлов
- независимость и гетерогенность - вычислительные узлы могут иметь разную производительность и разные конфигурации

# Зачем нужны распределенные системы

- увеличение производительности вычислений
- географическое распределение вычислений
- совместное использование ресурсов
- отказоустойчивость

# Восемь заблуждений относительно распределенных вычислений

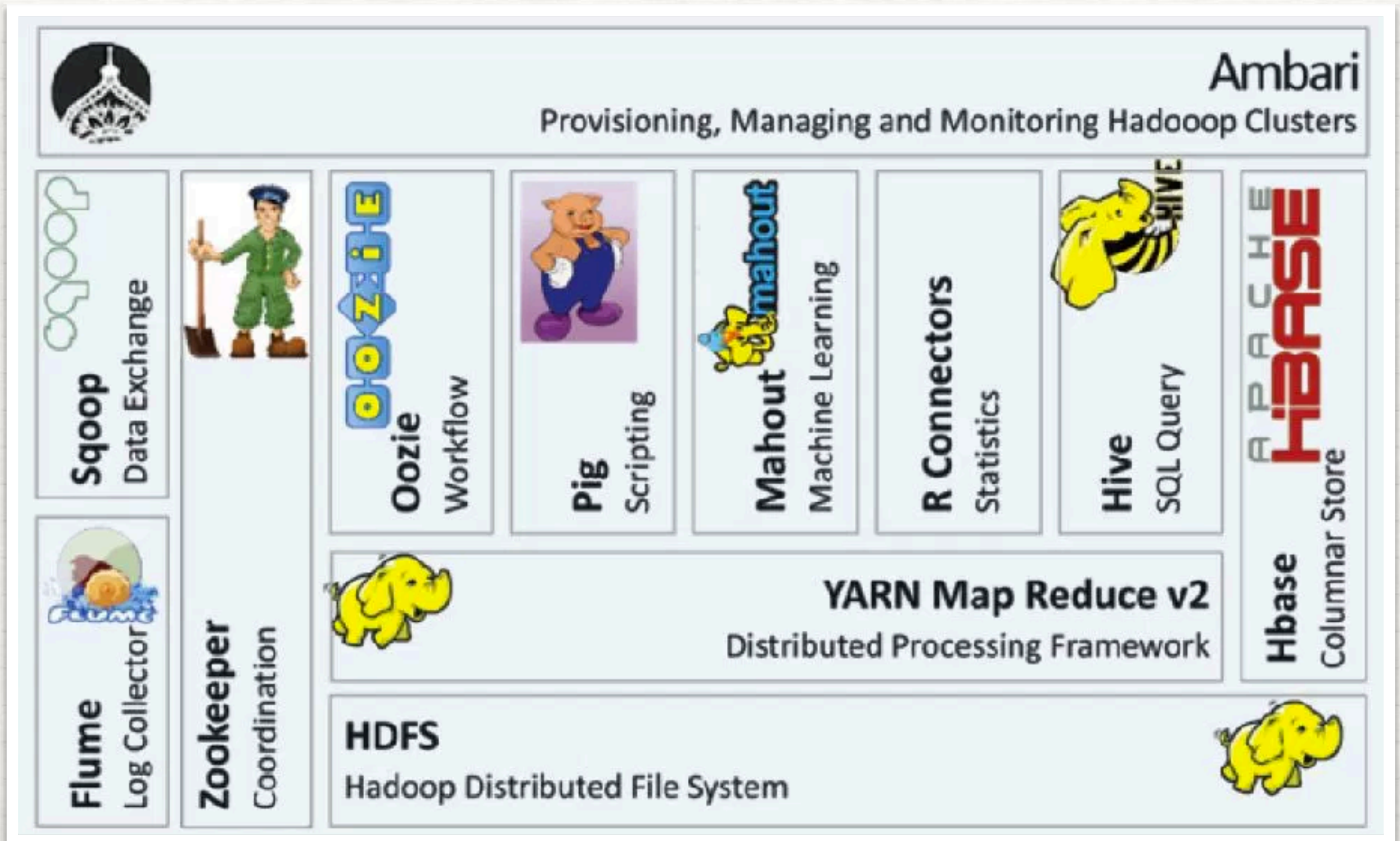
1. Сеть является надежной
2. Задержки передачи сообщения равны нулю
3. Полоса пропускания неограниченно
4. Сеть является безопасной
5. Сетевая топология неизменна
6. Система обслуживается одним администратором
7. Издержки транспортной инфраструктуры равны нулю
8. Сеть является однородной

# HADOOP HISTORY

- 2002 - запуск проекта Nutch
- 2003 - публикация с описанием GFS
- 2004 - создание NDfs (Nutch Distributed File System)
- 2004 - публикация Google и MapReduce
- 2005 - реализация MR в Nutch
- 2006 - выделение проекта Hadoop
- 2008 - выход Hadoop в лидеры ASF (Apache Software Foundation)
- 2008 - Apache Pig
- 2009 - Apache Scoop
- 2010 - выделение проекта HBase (ноябрь - Facebook messages)
- 2010 - Apache Hive
- 2011 - Apache Kafka
- 2012 - Apache Impala
- 2012 - Apache Accumulo
- 2013 - Apache YARN
- 2013 - Apache Spark (начало работы в 2009)



# HADOOP WORLD 1.0



# HADOOP WORLD 2017

## Hadoop Ecosystem

### Data Visualization

SAS Visual Analytics   Tableau   Qlik   SAP Lumine   R   DBLIS   iCharts   Timeline JS   Apache Zeppelin

### System Deployment

Apache Ambari   Apache Mesos   Marathon   Hortonworks HOYA   Apache Bigtop   Daploop   Apache Eagle

Cloudera HUE   Myriad   Brooklyn   Apache Helix   Buildloop   SequenceIQ Cloudbreak

### Data Ingestion

Apache Flume

Apache Sqoop

Facebook Scribe

Apache Chukwa

Apache Kafka

Netflix Suro

Apache Samza

Cloudera Morphline

IIIO

Apache NIFI

Apache ManifoldCF

### Service Programming

Apache Thrift

Apache Zookeeper

Apache Avro

Apache Curator

Apache Karaf

Twitter Elephant Bird

LinkedIn Norbert

### Scheduling & DR

Apache Dozie

LinkedIn Arkaban

Apache Falcon

Shedoscope

### Security

Apache Sentry

Apache Knox Gateway

Apache Ranger

### Frameworks

Jumbune

Spring XD

Cask Data App Platform

### Metadata

Metascope

Apache Tika

### Machine Learning

Apache Mahout

WEKA

Cloudera Dryx

DeepLearning4j

MADlib

II20

Sparkling Water

Apache SystemML

### Distributed Programming

Apache Ignite

Apache MapReduce

Apache Pig

JAQL

Apache Spark

Apache Storm

Apache Flink

Apache Apex

Netflix PigPen

AMPLAB SIMR

Facebook Corona

Apache REEF

Apache Twill

Damballa Parkour

Apache Hama

Datasalt Pangool

Apache Tez

Apache DataFu

Kangaroo

TinkerPop

Pachyderm MapReduce

Apache Beam

### SQL on Hadoop

Apache Hive

Apache HCatalog

Apache Tealodion

Apache HAWQ

Apache Drill

Cloudera Impala

Facebook Presto

Datasalt Spout SQL

Apache Tajo

Apache Phoenix

Apache MRQL

Kylin

### NoSQL Databases

#### Key-Value

Redis

LinkedIn Voldemort

Base4DB

OpenTSDB

#### Graph

Graph

Neo4j

TitanDB

OrientDB

#### Stream Data Model

EventStore

#### Wide Column

Apache HBase

Apache Cassandra

Hypertable

Apache Accumulo

Apache Kudu

Apache Parquet

#### Document

MongoDB

RethinkDB

ArangoDB

CouchDB

DynamoDB

Gemfire

### NewSQL Databases

TokuDB

HandlerSocket

Akiban Server

Drizzle

Haeinsa

SenseiDB

Sky

BayesDB

InfluxDB

VoltDB

SAP HANA

### Distributed File System

Apache HDFS

Red Hat GlusterFS

Quantcast File System

Ceph File System

Lustre File System

Aluxio

GridGain

XtreemFS



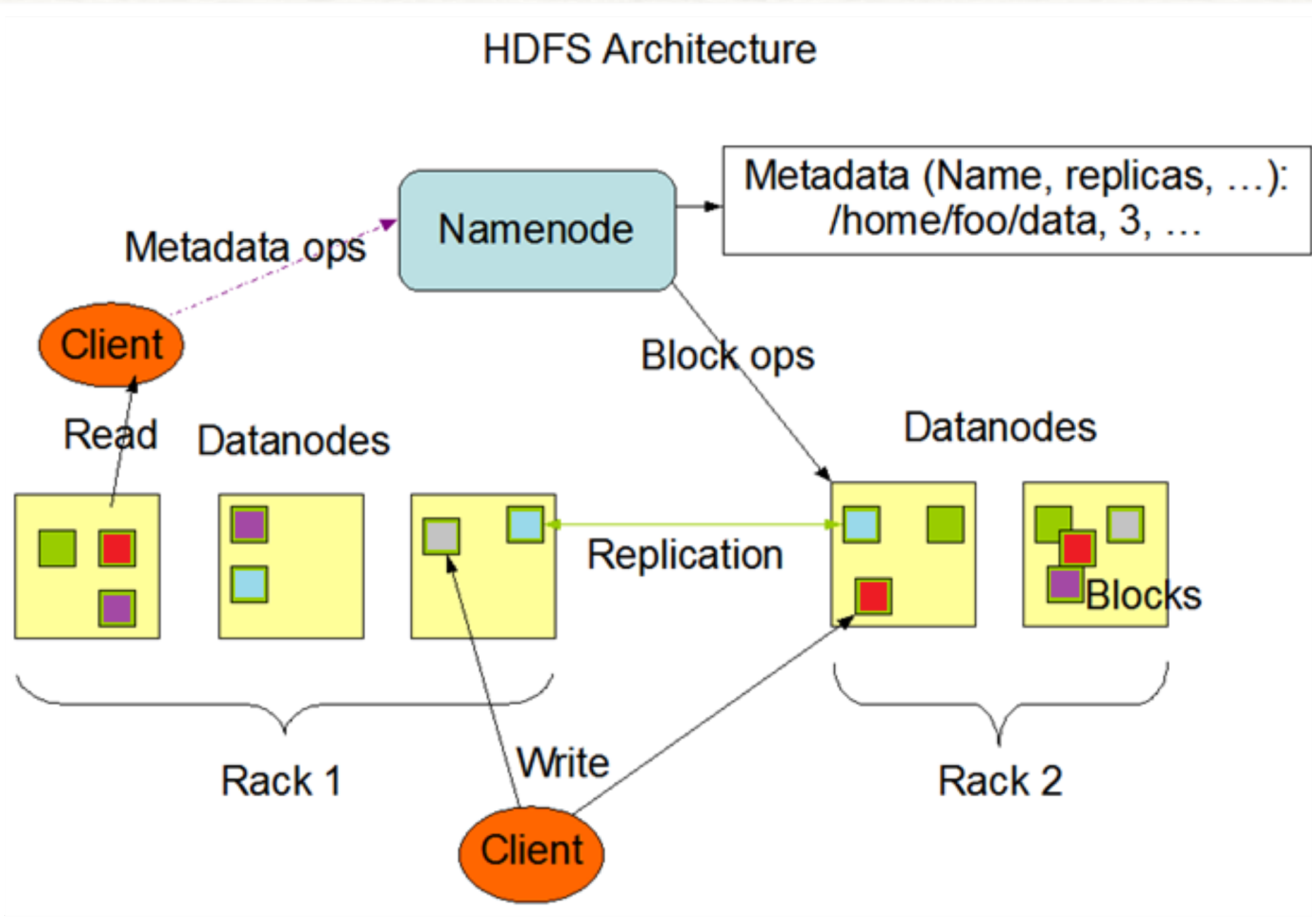
# HADOOP BASICS

1.Hadoop Distributed File System (HDFS):

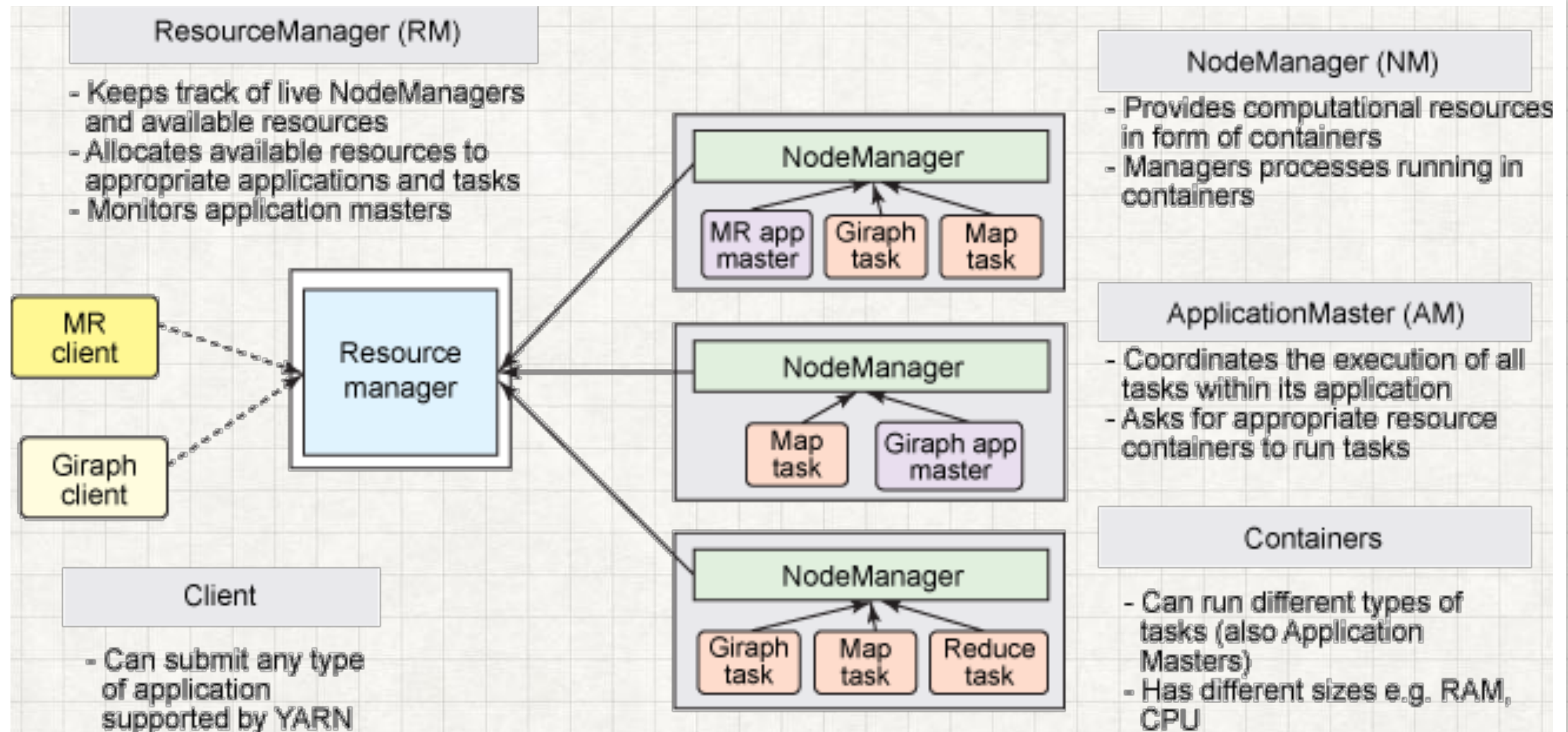
2.Hadoop YARN

3.Hadoop MapReduce

# HDFS BASICS

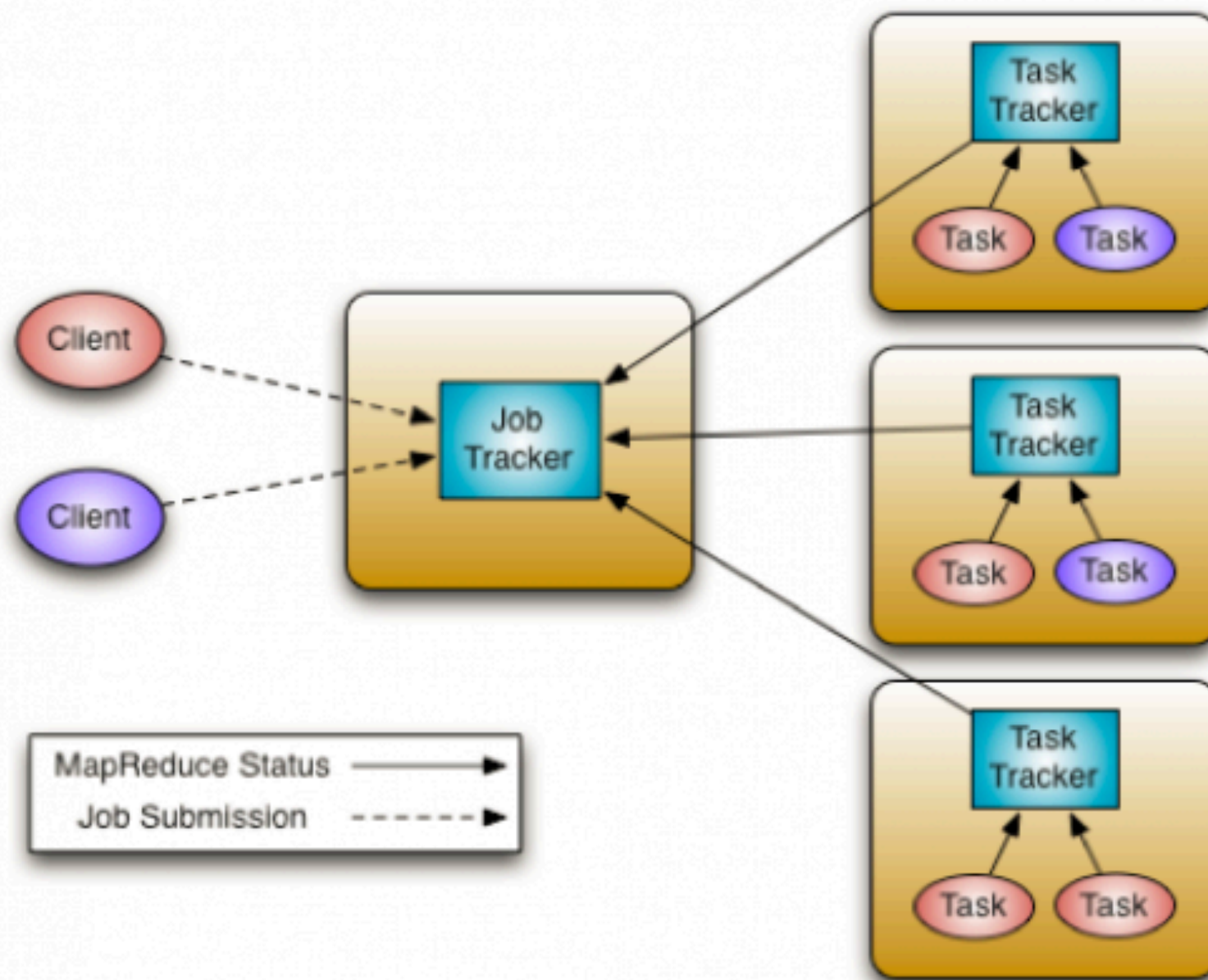


# YARN BASICS

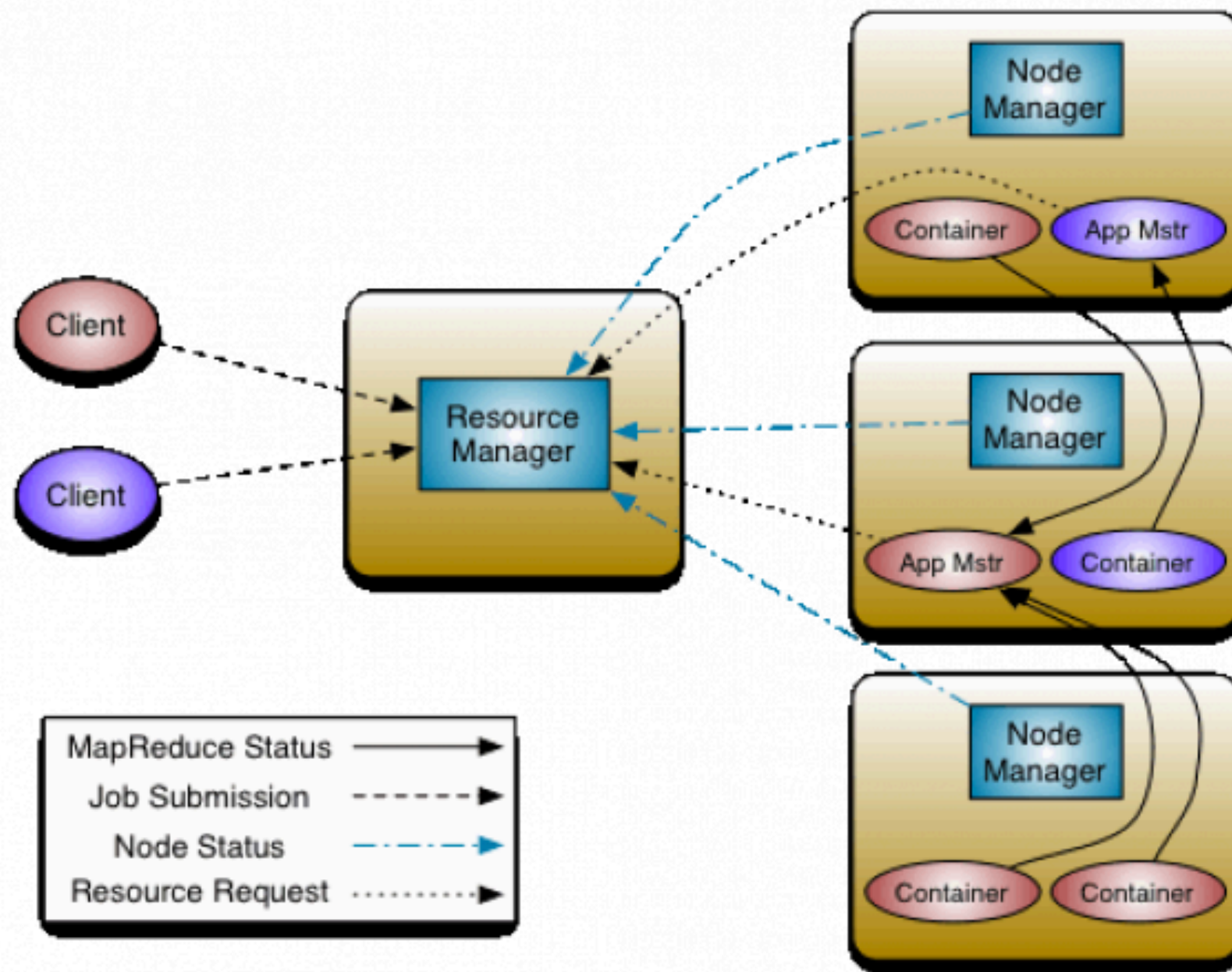




# MR 1

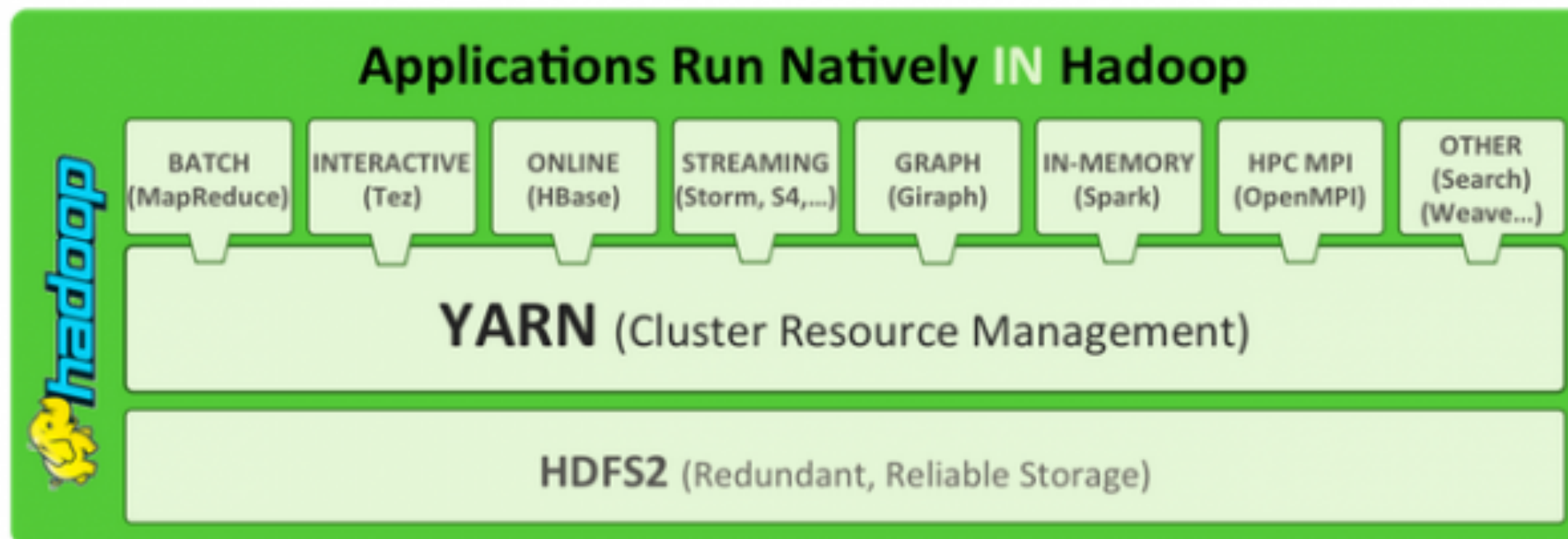
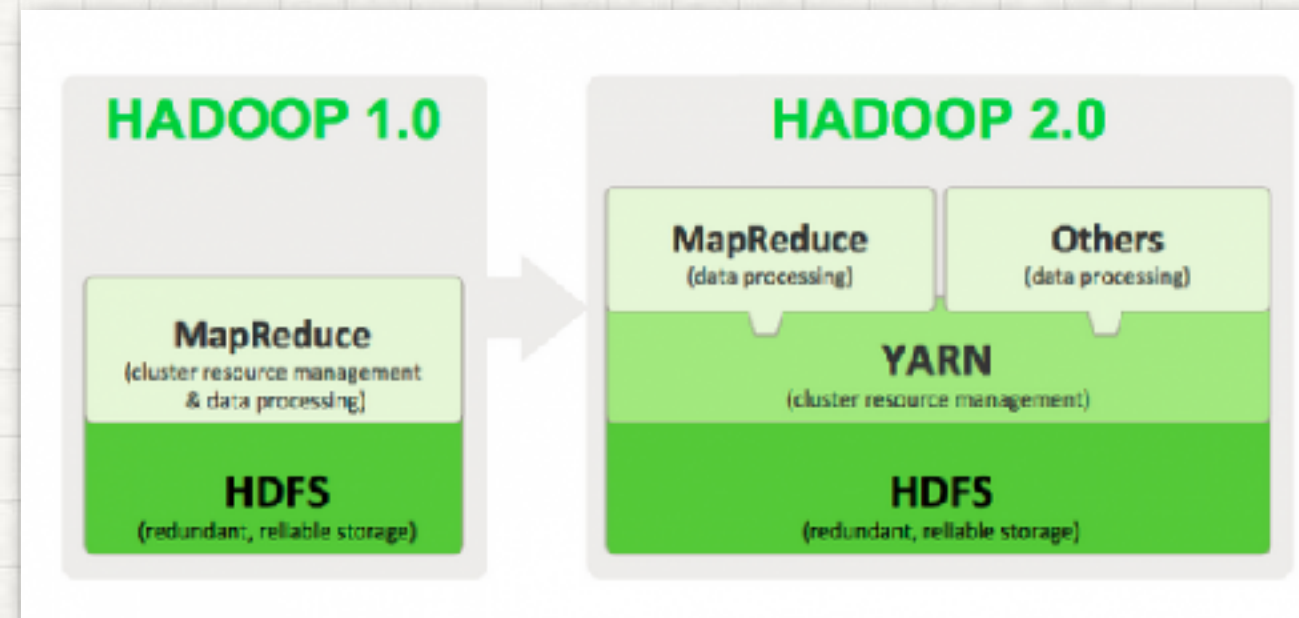


# YARN (MR 2)

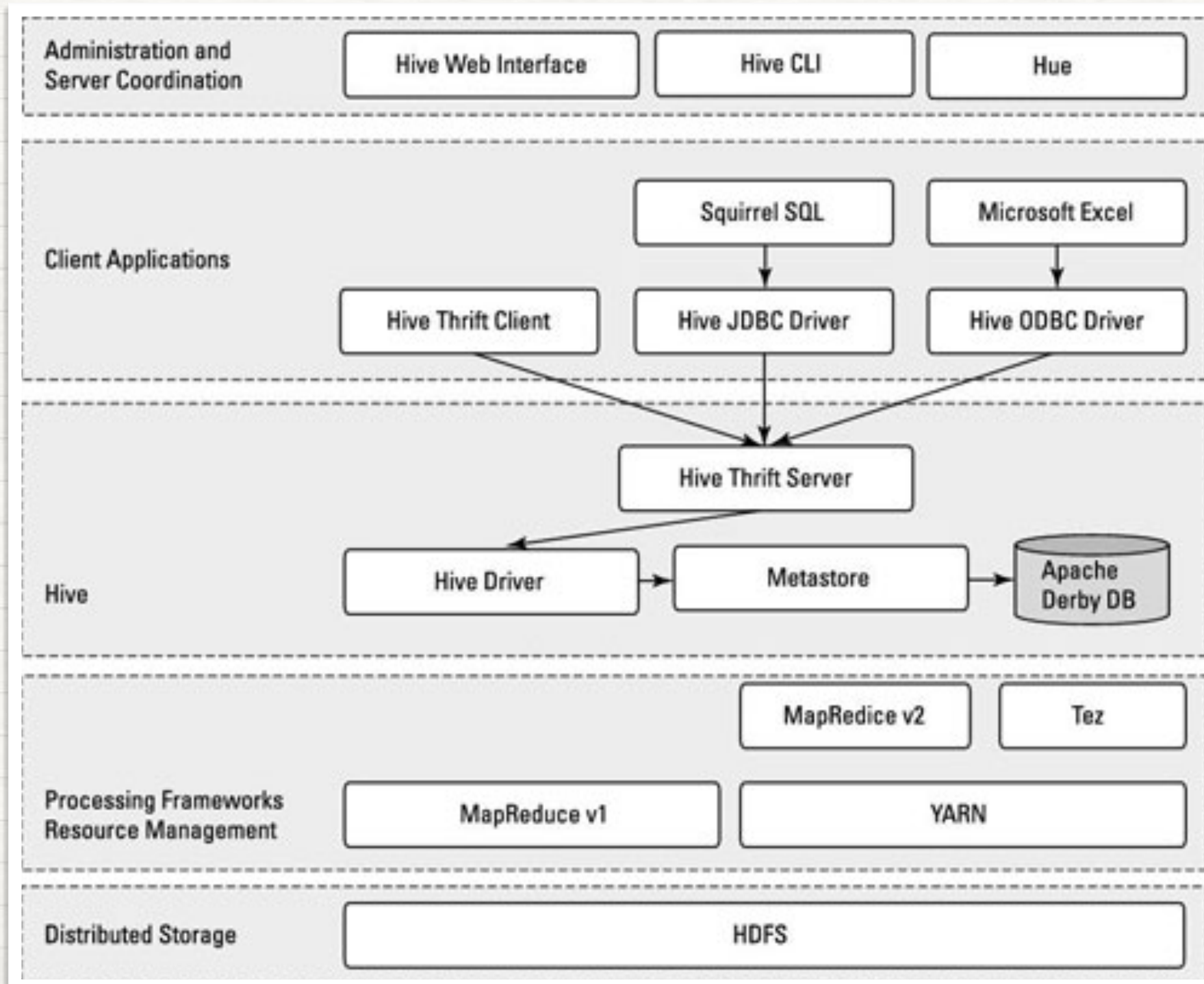




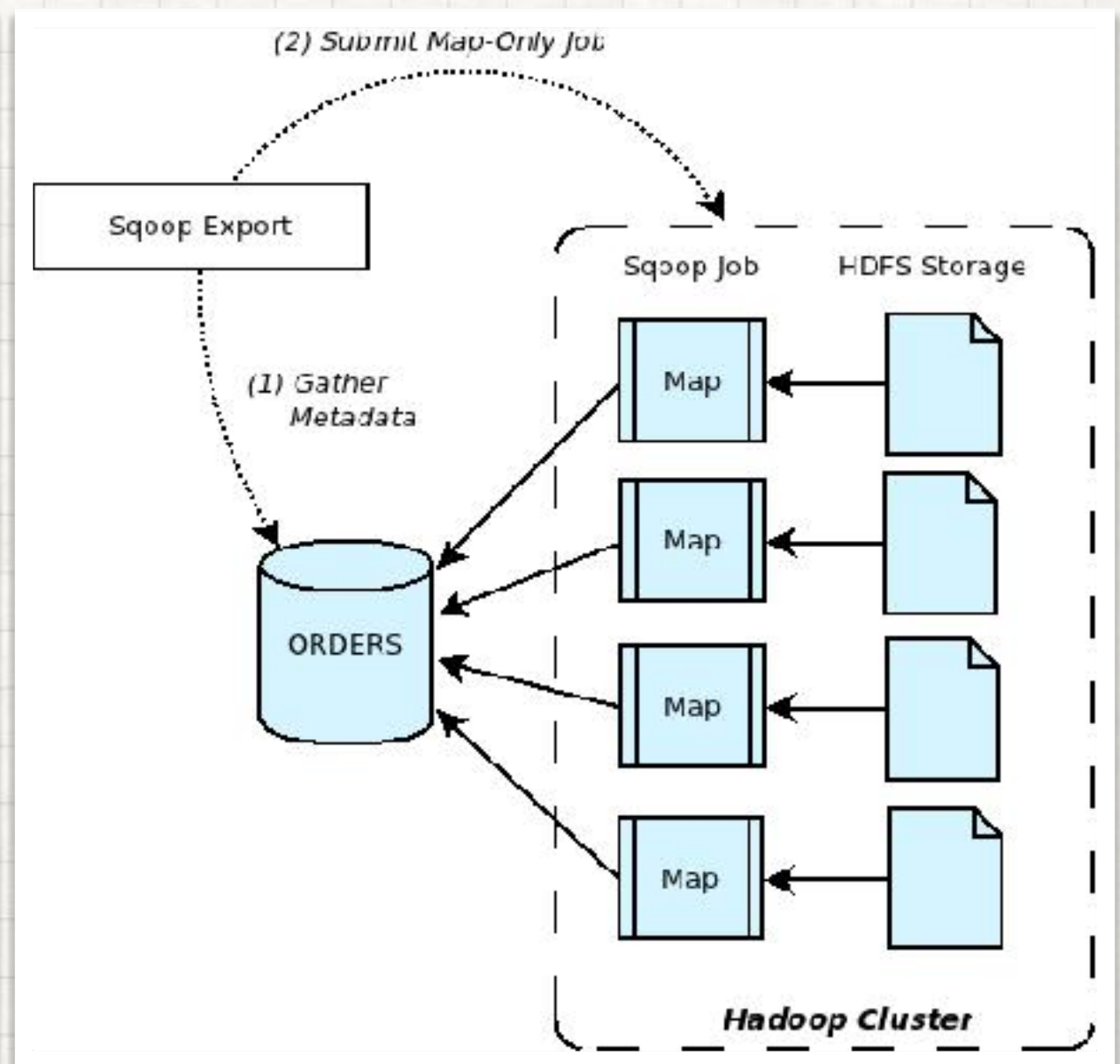
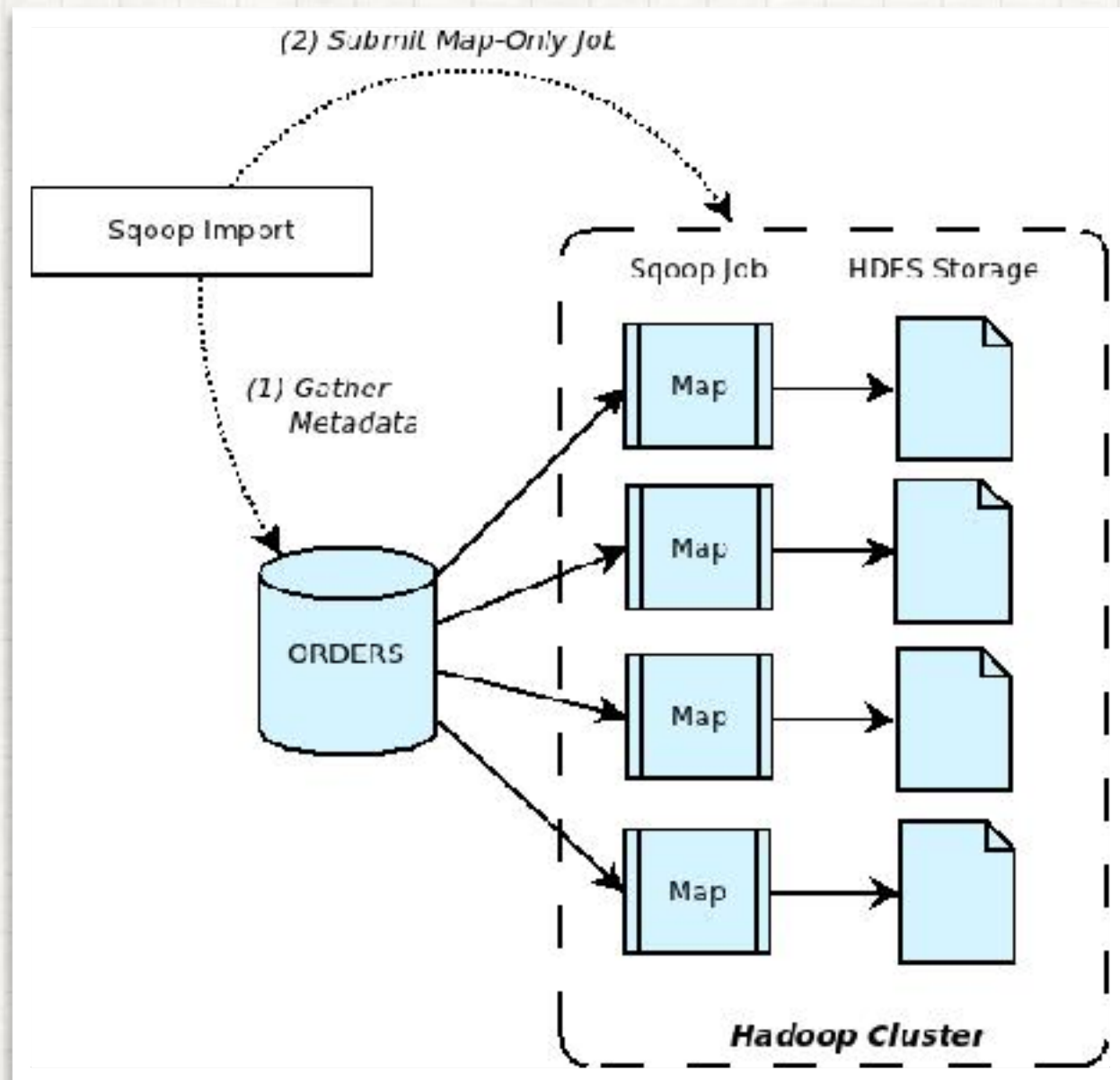
# YARN



# HIVE BASICS

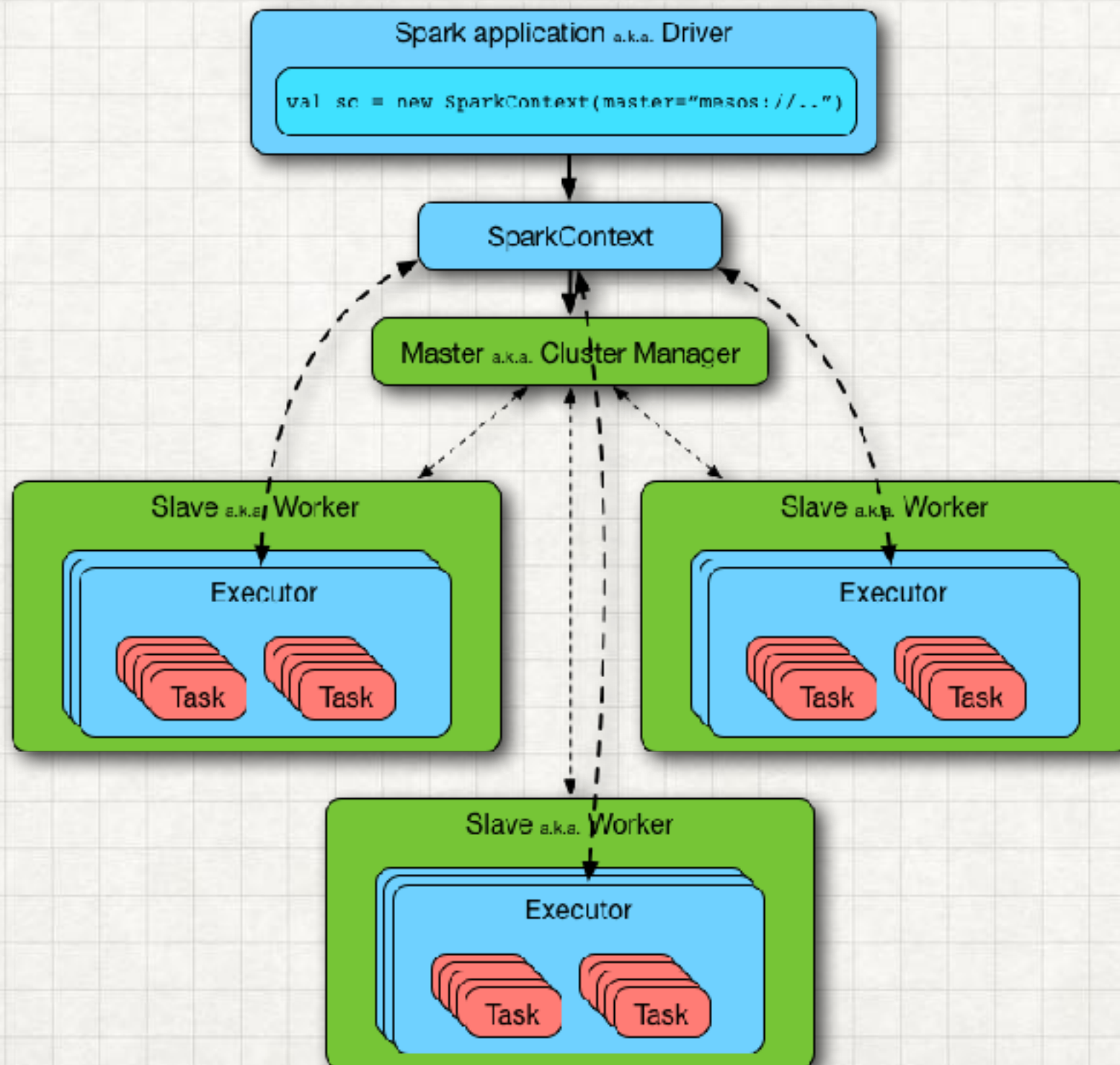


# SQOOP BASICS

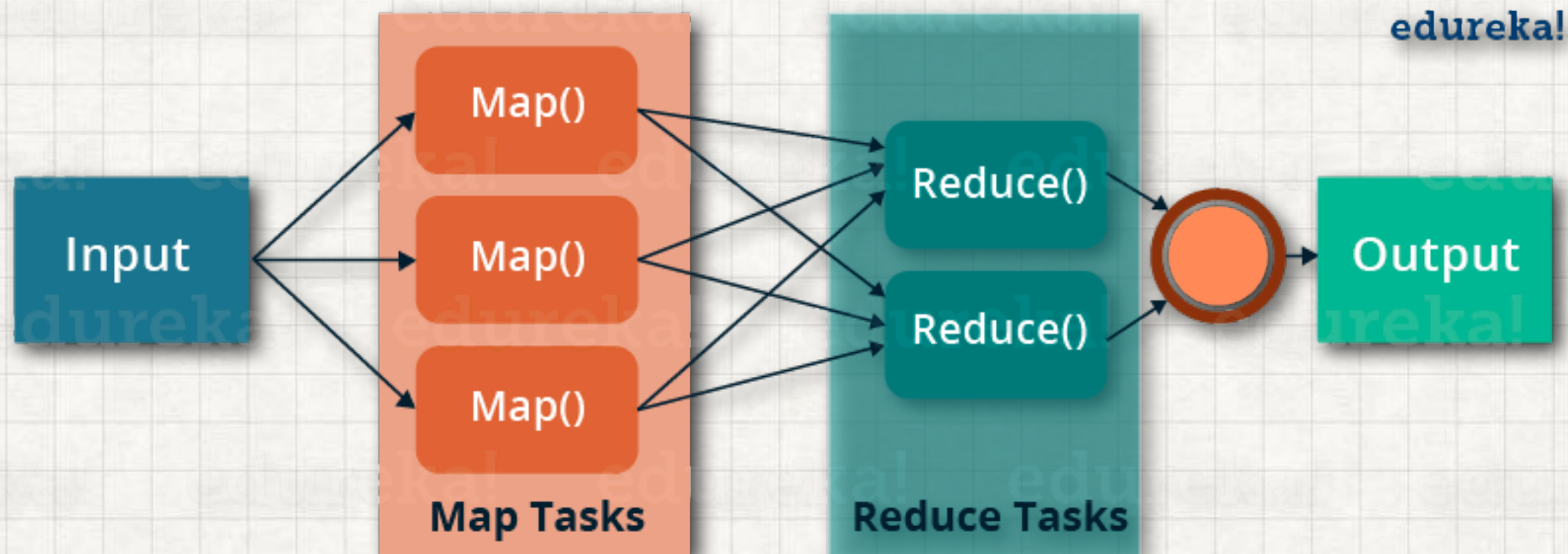




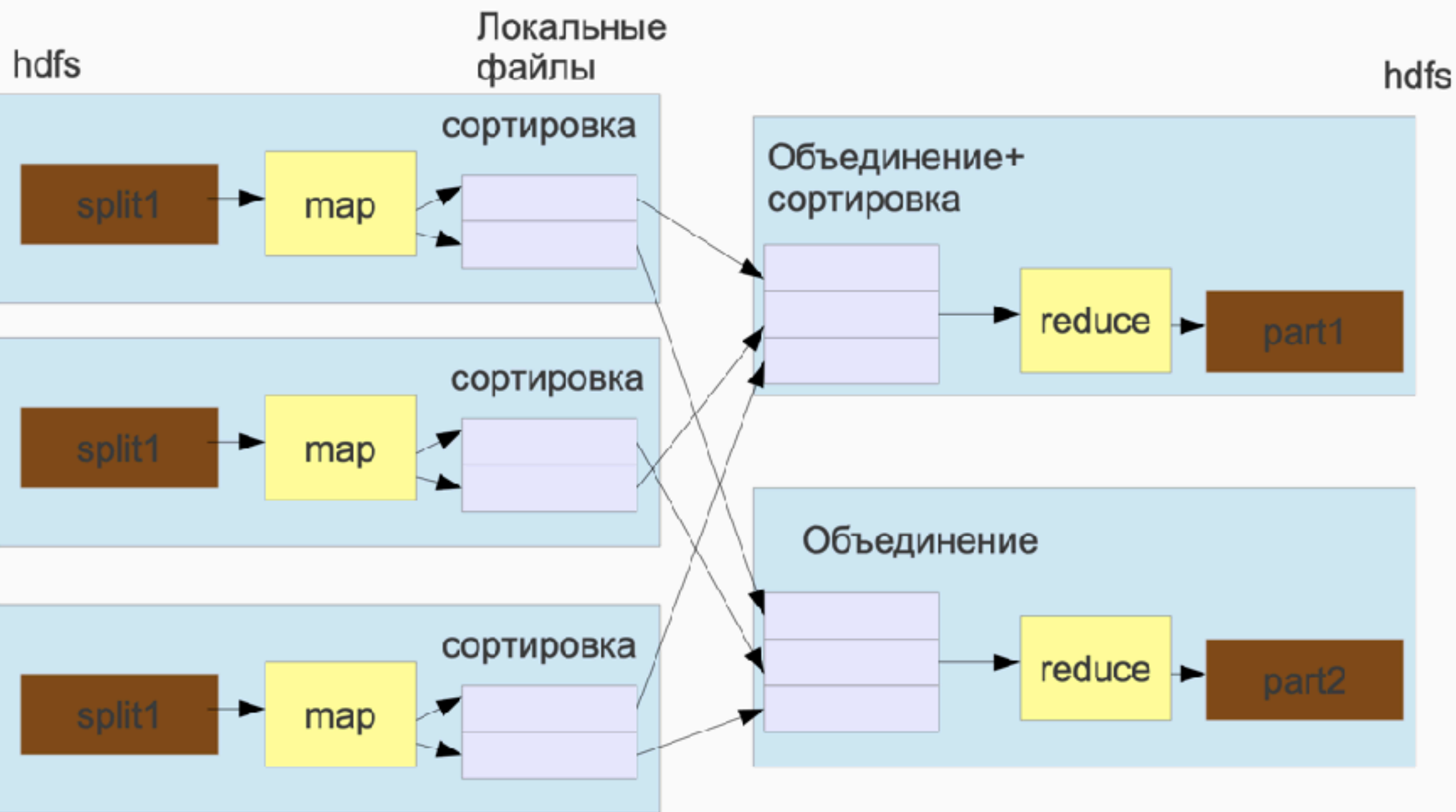
# SPARK BASICS



# MAP-REDUCE

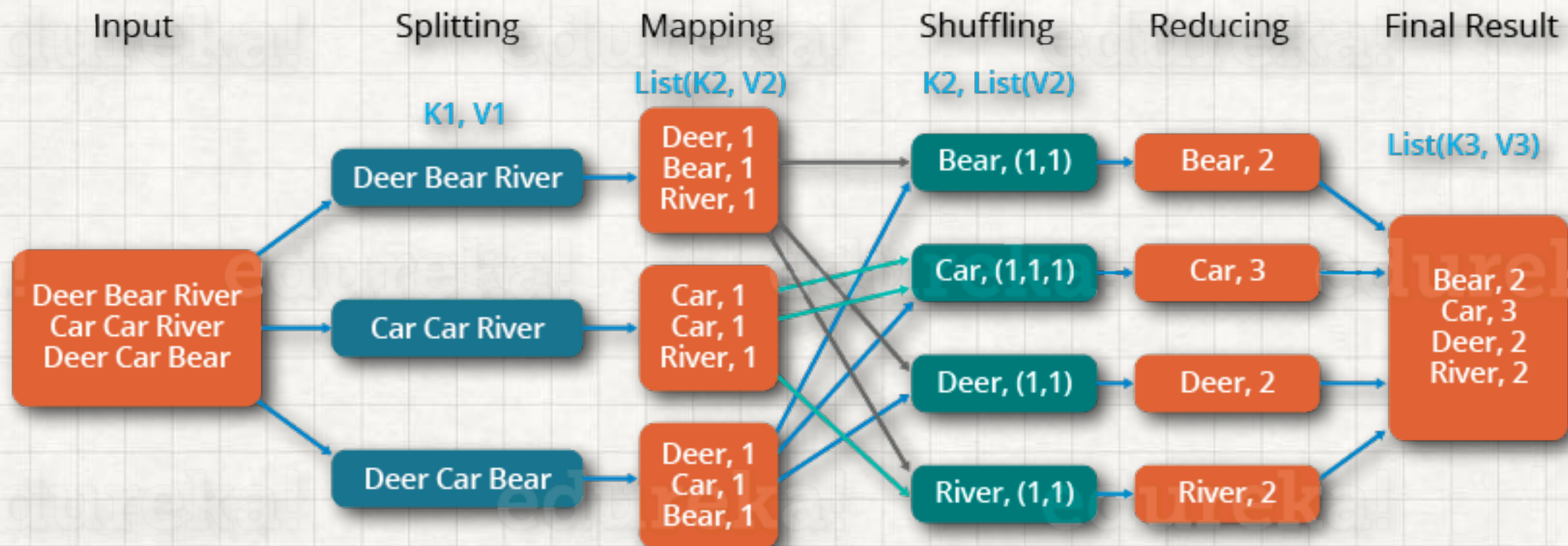






## The Overall MapReduce Word Count Process

edureka!



# INPUTFORMAT

- Основные функции :
  - Разбивает исходные данные на splits
  - Обеспечивает чтение данных из split в процессе работы mapper
  -



# БАЗОВЫЕ КЛАССЫ

```
public abstract class InputFormat {  
  
    public abstract List getSplits(JobContext context)..;  
  
    public abstract RecordReader createRecordReader( InputSplit  
        split, TaskAttemptContext context) ..  
  
}
```

# БАЗОВЫЕ КЛАССЫ

```
public abstract class InputSplit {  
  
    public abstract long getLength() throws IOException,  
        InterruptedException;  
  
    public abstract String[] getLocations() throws ...;  
  
}
```



# СЕРИАЛИЗАЦИЯ

- **Сериализация — преобразование объекта в байтовый поток данных**
- **Десериализация — преобразование потока байтов в объект.**
- **Базовый интерфейс объекта сериализации в Hadoop — `org.apache.hadoop.io.Writable`**
- `public interface Writable { void write(DataOutput out) throws IOException; void readFields(DataInput in) throws IOException; }`

# MAPPER

- **Базовый класс для реализации mapper :**  
`org.apache.hadoop.mapreduce.Mapper`
- **Базовый метод для перегрузки**
- `protected void map(KEYIN key, VALUEIN value, Mapper.Context context)`
- **Для того чтобы добавить результат функции map требуется вызвать `context.write(KEYOUT key, VALUEOUT value)`**

# MAPPER

```
public class WordMapper extends Mapper {  
  
    @Override protected void map(LongWritable key, Text value, Context context)  
  
        throws IOException, InterruptedException {  
  
        String line = value.toString();  
  
        String[] words = StringTools.split(StringTools.removeAllNonSymbols(line).toLowerCase());  
  
        for (String word : words) {  
  
            context.write(new Text(word), new IntWritable(1));  
  
        }  
  
    }  
  
}
```



# REDUCER

- **Базовый класс для реализации reducer :**  
`org.apache.hadoop.mapreduce.Reducer`
- **Базовый метод для перегрузки**
- `protected void reduce(KEYIN key, Iterable values, Reducer.Context context)`
- **Для того чтобы добавить результат функции map требуется вызвать `context.write(KEYOUT key, VALUEOUT value)`**

# ПРИМЕР

```
public class WordReducer extends Reducer {  
  
    @Override protected void reduce(Text key, Iterable values, Context context) throws IOException,  
        InterruptedException {  
  
        long count=0;  
  
        Iterator iter = values.iterator();  
  
        while(iter.hasNext()) {  
  
            iter.next();  
  
            count++;  
  
        }  
  
        context.write(key, new LongWritable(count)); } }
```



# DRIVER

```
public class WordCountApp {  
    public static void main(String[] args) throws Exception {  
        if (args.length != 2) {  
            System.err.println("Usage: WordCountApp "); System.exit(-1); }  
  
        Job job = Job.getInstance();  
  
        job.setJarByClass(WordCountApp.class);  
  
        job.setJobName("Word count");  
  
        FileInputFormat.addInputPath(job, new Path(args[0]));  
  
        FileOutputFormat.setOutputPath(job, new Path(args[1]));  
  
        job.setMapperClass(WordMapper.class);  
  
        job.setReducerClass(WordReducer.class);  
  
        job.setOutputKeyClass(Text.class);  
  
        job.setOutputValueClass(IntWritable.class);  
    }  
}
```