



# Universidad **Ricardo Palma**

RECTORADO

PROGRAMA DE ESPECIALIZACIÓN EN CIENCIA DE DATOS

*Formamos seres humanos para una cultura de paz*

## BIG DATA APLICADO

SESIÓN 2

**Expositores:**

**David Narváez**

**Eder Pineda**

**[bigdataplicado@gmail.com](mailto:bigdataplicado@gmail.com)**

# AGENDA

- INTRODUCCIÓN:
  - QUE ES APACHE HADOOP?
  - HDFS
  - MAPREDUCE
  - YARN
  - ECOSISTEMA HADOOP

# QUE ES APACHE HADOOP?

*El proyecto **Apache™ Hadoop®** desarrolla software de código abierto para computo distribuido, confiable y escalable.*

*La librería **Apache Hadoop** es un framework que permite procesamiento distribuido de grandes set de datos a través de un cluster de servidores, utilizando modelos de programación simples.*

*Está diseñado para escalar de un único servidor a miles de servidores, cada uno ofreciendo computo local y almacenamiento.*

*Está diseñado para la detección y manejo de fallas en su capa de aplicaciones, de esta manera brinda servicio de alta disponibilidad (HA) en un cluster de servidores, los cuales pueden ser propensos a fallas.*





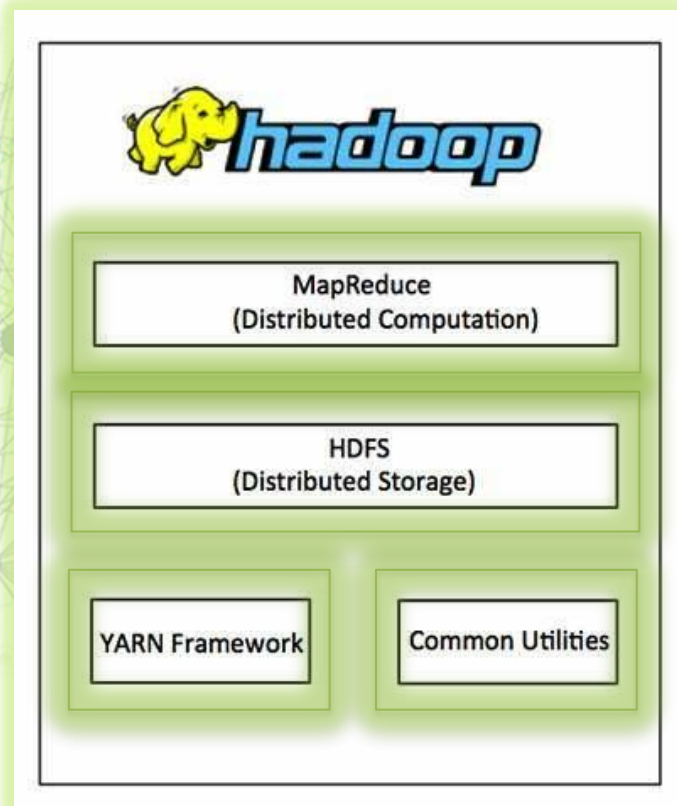
# QUE ES APACHE HADOOP?

## HDFS (Distributed Computation)

Es un sistema de archivos distribuido que proporciona acceso de alto rendimiento a los datos de las aplicaciones.

## Yarn Framework

Es un Framework para la programación de Jobs y administrador de los recursos del cluster



## MapReduce (Distributed Computation)

Un sistema basado en Java necesario para procesamiento paralelo de grandes cantidades de datos. Estas librerías proporcionan un sistema de archivos y los niveles de abstracción del sistema operativo. Contiene los archivos Java necesarios y scripts requeridos para inicializar Hadoop

# HDFS – DISTRIBUTED FILE SYSTEM

- *HDFS es una aplicación escrita en Java que simula un filesystem*
- *Hadoop fue inspirado en los documentos de Google de Map Reduce y Google File System (GFS).*
- *Proporciona almacenamiento redundante para grandes cantidades de datos*
- *HDFS proporciona una shell como cualquier otro sistema de archivos y dispone de una lista de comandos para interactuar con el file system*



# HDFS – DISTRIBUTED FILE SYSTEM

## *Características del HDFS*

- *Alto rendimiento*
- *Tolerante a fallos*
- *Gestión centralizada relativamente simple: arquitectura de maestro-esclavo*
- *Seguridad: Dos niveles entre los cuales elegir (Kerberos - Token)*
- *Optimizado para procesamiento distribuido: localidad de data*
- *Escalabilidad*



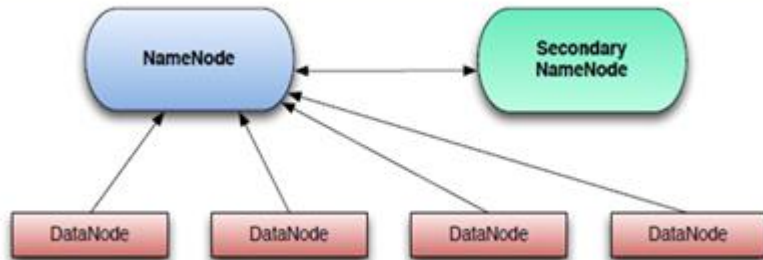




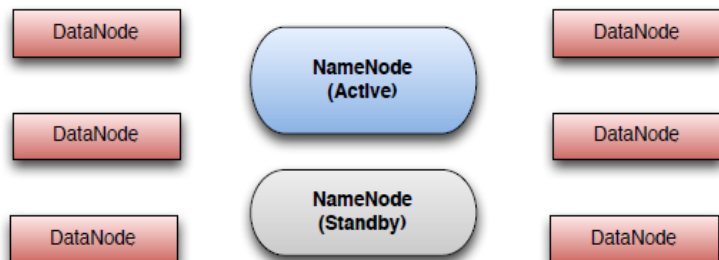
# HDFS – DISTRIBUTED FILE SYSTEM

## Arquitectura HDFS

- *HDFS utiliza una arquitectura maestro/esclavo, donde el maestro consiste en un único **NameNode** que administra la metadata del file system y uno o mas esclavos. Mantiene toda la metada en memoria RAM.*
- ***DataNodes** que almacena los datos.*



*The NameNode is a single point of failure (SPOF)  
The Secondary NameNode is not a failover NameNode!*



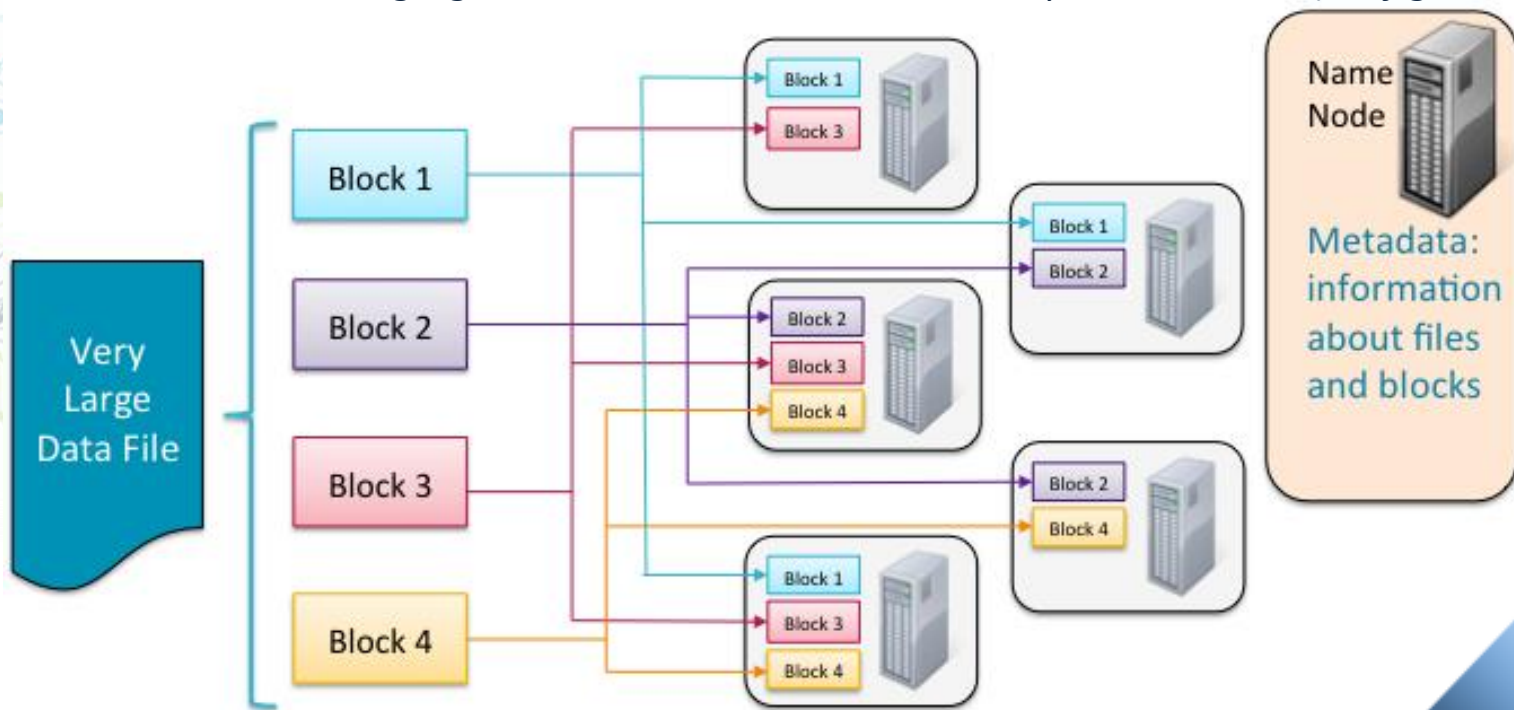
*HDFS With High Availability  
Standby NameNode takes over when active  
NameNode fails*



# HDFS – DISTRIBUTED FILE SYSTEM

## Bloques HDFS

- Cuando un archivo es agregado al HDFS, este es dividido en bloques de 128MB (configurable)



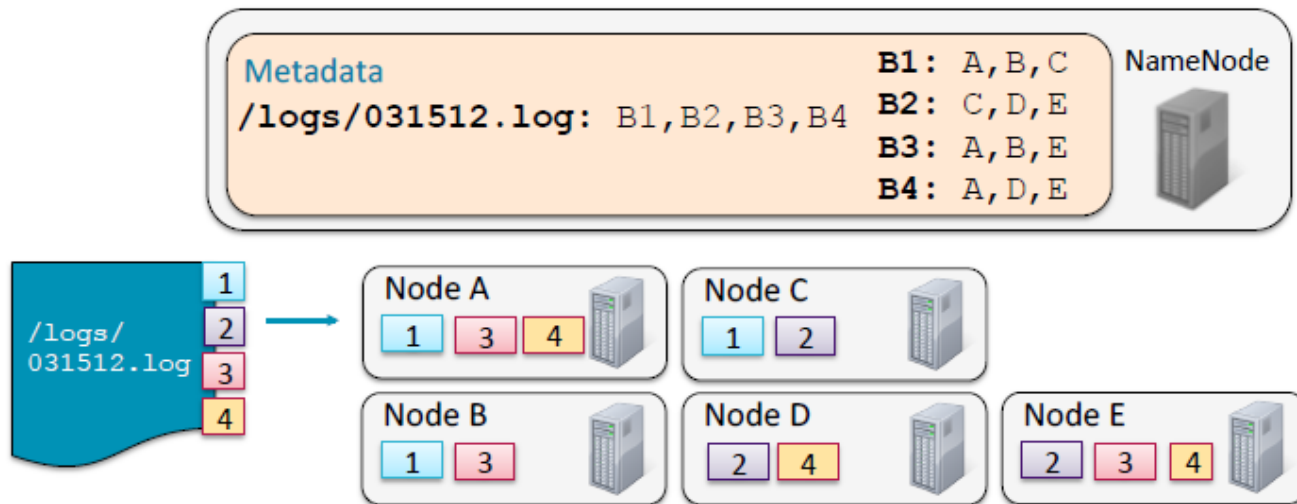




# HDFS – DISTRIBUTED FILE SYSTEM

## Replicación HDFS

- Los bloques son replicados en los nodos del cluster
  - Basados en el factor de replicación (tres por defecto)
- La replicación aumenta la confiabilidad y el rendimiento
  - Confiabilidad: los datos pueden tolerar la pérdida de todas las réplicas excepto una
  - Rendimiento: mas oportunidades por localidad de datos



# MAPREDUCE

- **Hadoop MapReduce** es un framework de software para el desarrollo de aplicaciones que procesan grandes cantidades de datos (multi-terabyte datasets) en paralelo, en clusters grandes (miles de nodes) de commodity hardware de una manera confiable y tolerante a fallos.



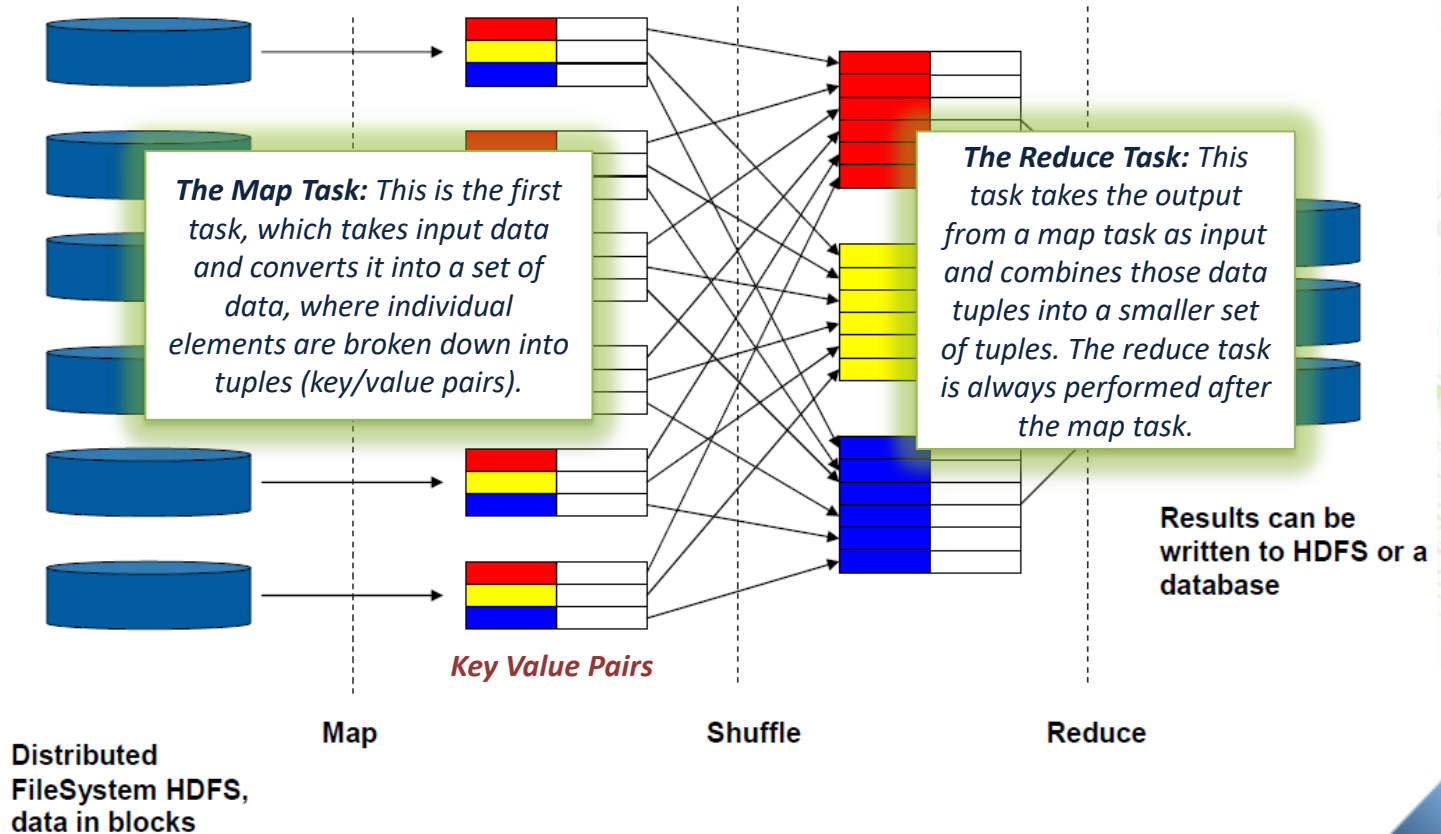
# MAPREDUCE

- *Un Job MapReduce normalmente divide los set de datos de entrada en fracciones independientes que son procesadas por el **proceso de mapeo (Map)** de forma completamente paralela. El framework ordena los outputs del proceso de mapeo, los cuales son input para el **proceso de reducción (Reduce)**. Ambos, el input y el output de los jobs son almacenados en un file-system. El framework se encarga de programar las tareas, monitorearlas y re-ejecutar los procesos caídos.*



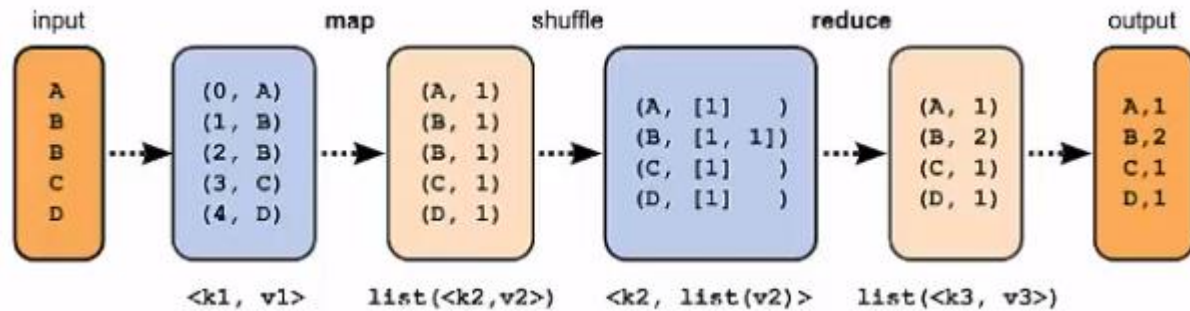


# MAPREDUCE



# MAPREDUCE

*Ejemplo:*



# YARN

- *Yet Another Resource Negotiator (YARN)*
- *YARN permite ejecutar diferentes cargas de trabajo en el mismo clúster Hadoop*
- *YARN permite compartir dinámicamente recursos de memoria y CPU del cluster entre marcos de procesamiento*



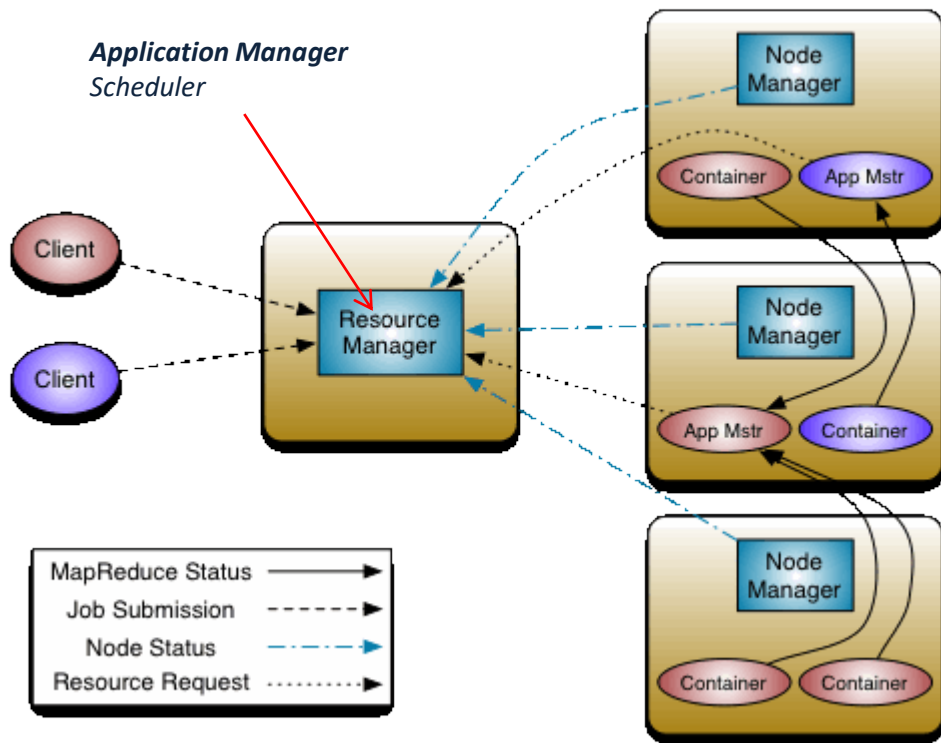




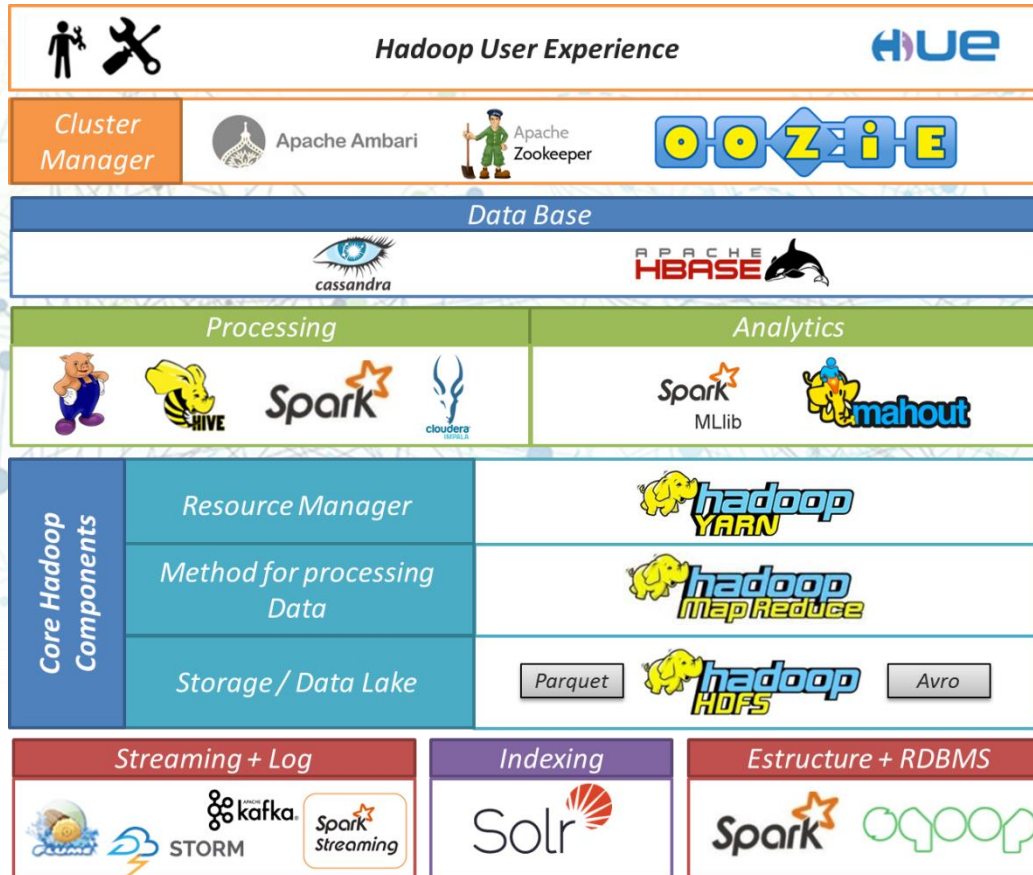
# YARN

*La idea fundamental de YARN es dividir las funcionalidades de administración de recursos y programación/monitoreo de job en demonios separados.*

*La idea es tener un Resource Manager global (RM) y por cada aplicación un ApplicationMaster (AM). Una aplicación es un trabajo único o un DAG de trabajos.*



# ECOSISTEMA HADOOP



# ECOSISTEMA HADOOP

## Data Ingestion

HDFS, Spark, Sqoop, Flume, Kafka, Storm

## Data Processing & Analysis

Hive, Pig, Spark & Impala

## Machine Learning & Analytics

Spark Mllib, Mahout

## Data Base

Cassandra, HBase

## Data Visualization/Discovery

Zeppeling, Tableau\*, Power BI\*



THE  
**APACHE**  
SOFTWARE FOUNDATION



**+tableau**



**Power BI**