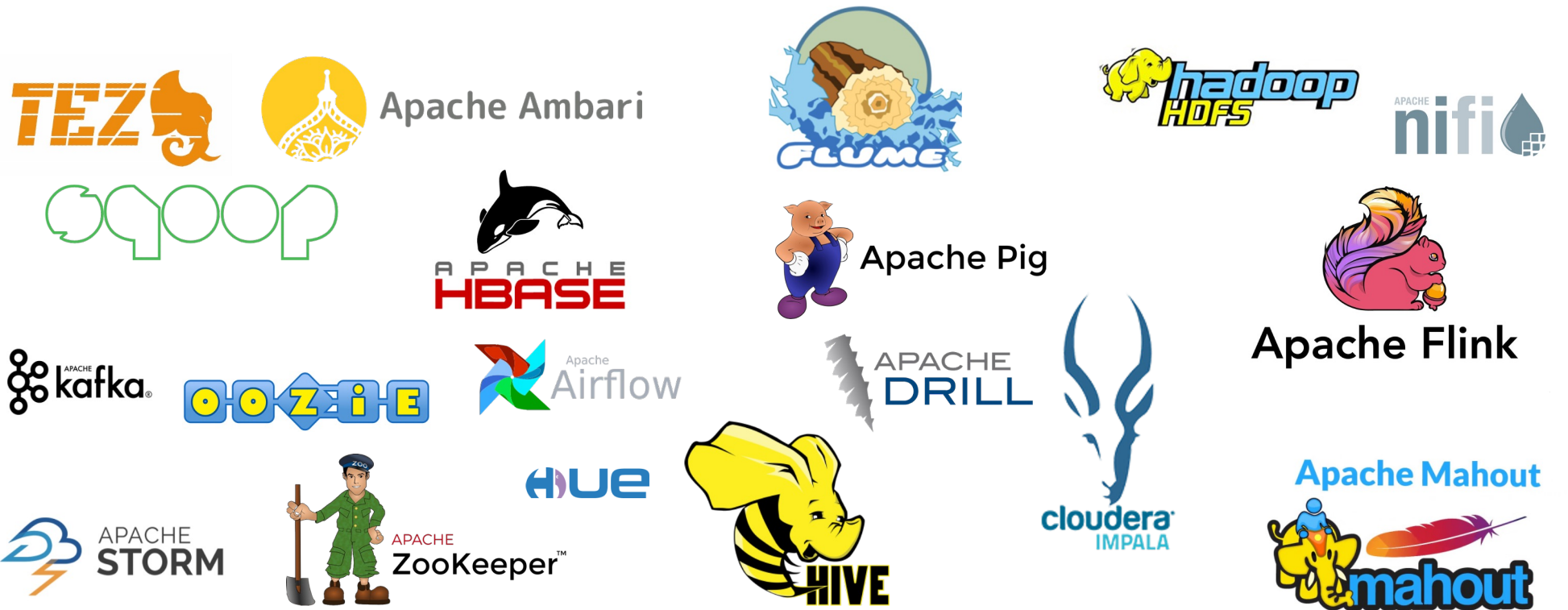


# Ecosistemas Big Data

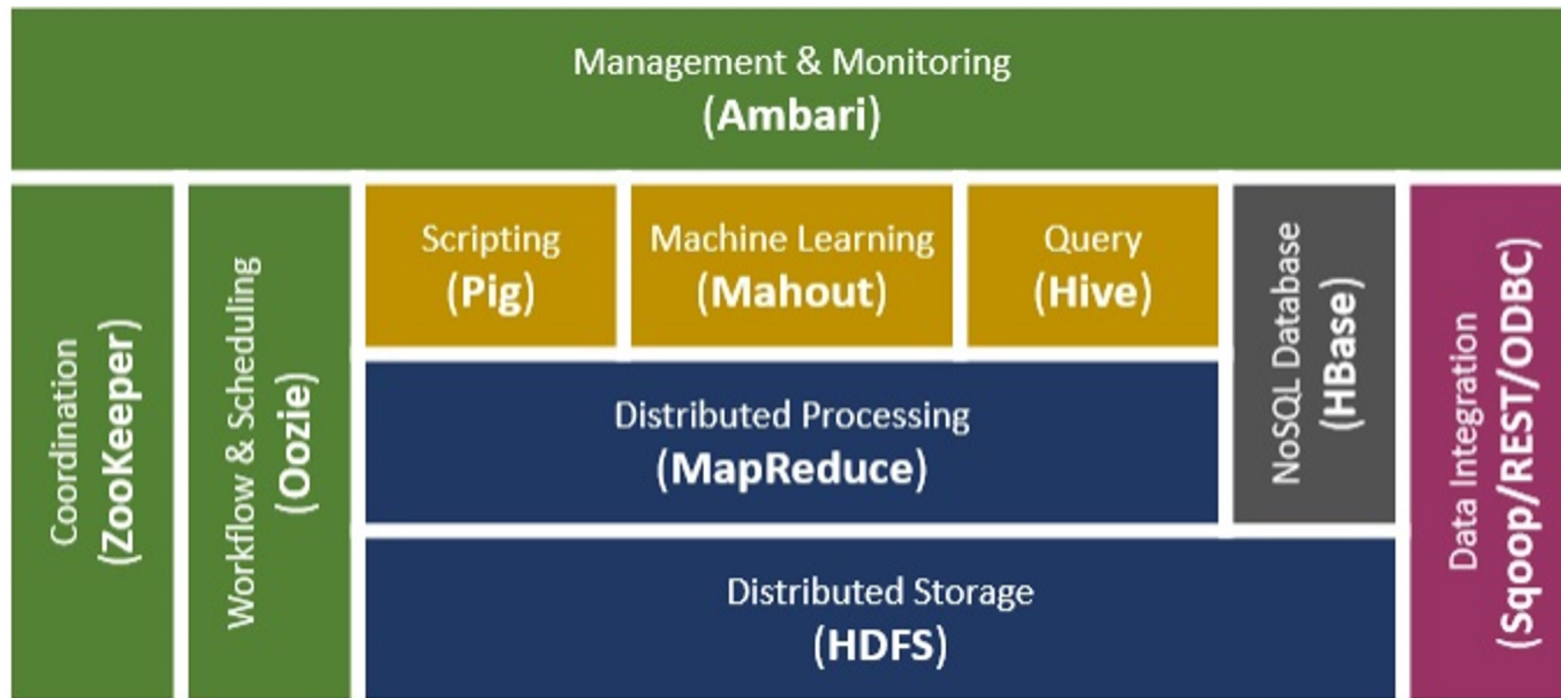
## Hadoop & otros frameworks

# Zoologico de Hadoop

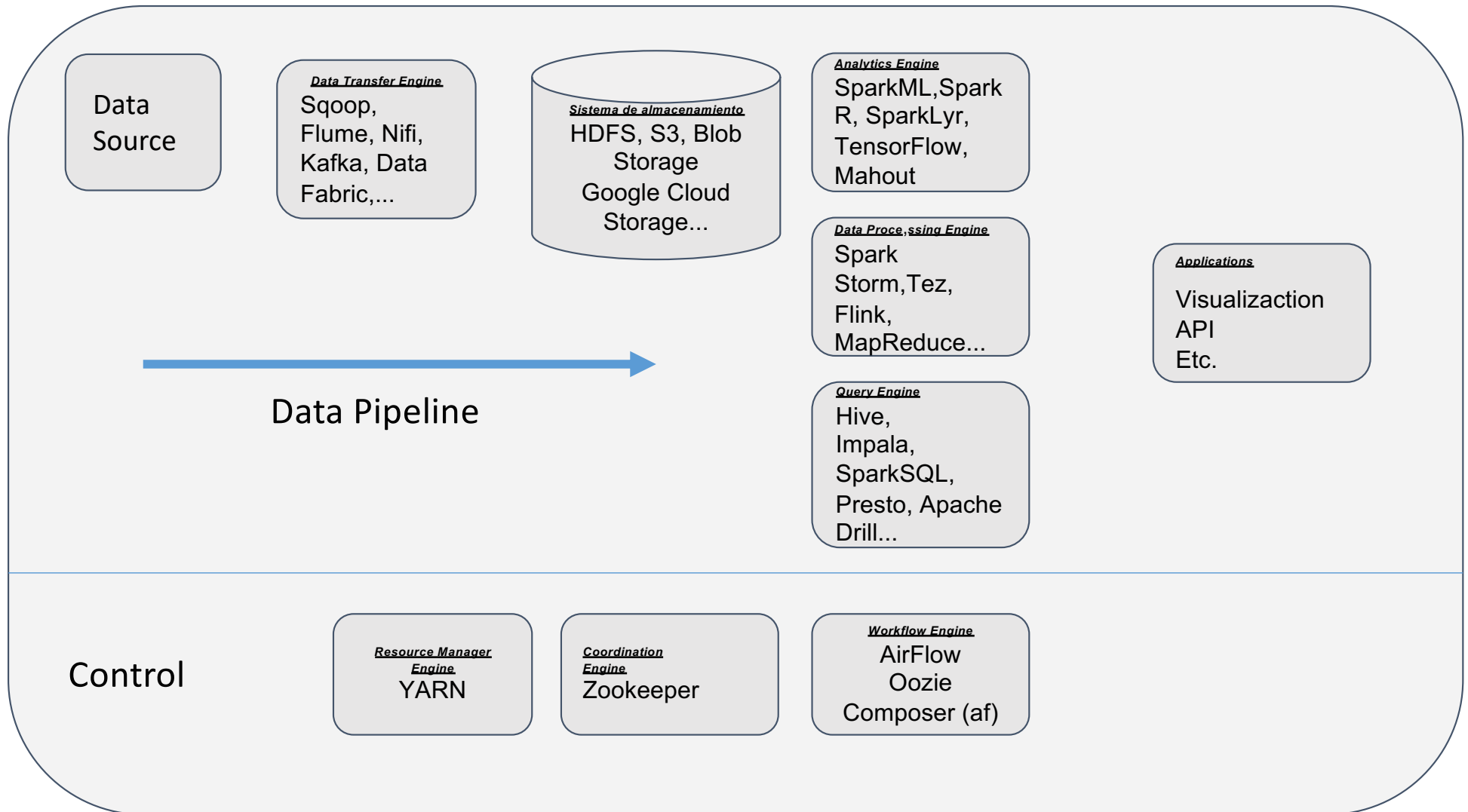


www.educba.com

# Ecosistema Hadoop (estándar)



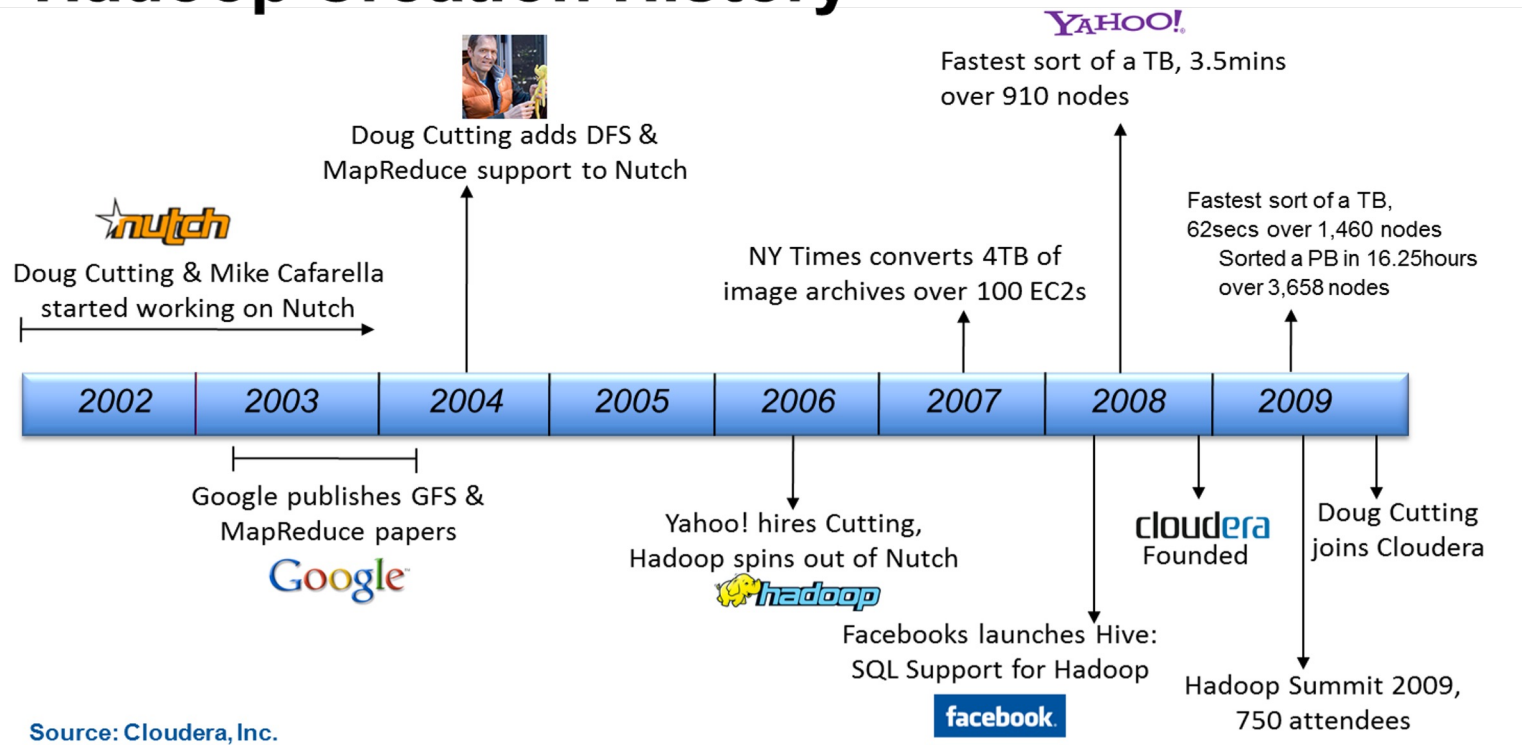
# Ecosistema Big Data



# Ecosistema Hadoop

[https://en.wikipedia.org/wiki/Apache\\_Hadoop](https://en.wikipedia.org/wiki/Apache_Hadoop)

# Hadoop Creation History



[https://en.wikipedia.org/wiki/Apache\\_Hadoop](https://en.wikipedia.org/wiki/Apache_Hadoop)

# Hadoop

- Apache Hadoop es un Framework de Big Data, que provee un modelo de computación distribuida que usa “commodity hardware” de forma flexible y que puede ser escalable.
- Es una arquitectura Master/Slave
- Sigue el paradigma de procesamiento Map/Reduce

# Hadoop

- Emplea 3 componentes básicos dentro de su arquitectura:
  - Un sistema de archivos distribuido (HDFS)
  - Un motor para el procesamiento de grandes volúmenes de datos (MapReduce)
  - Un gestor de recursos (YARN) / planificador (scheduler)



# Motores de almacenamiento



# Hadoop - HDFS

Infraestructura para almacenamiento de grandes volúmenes de datos

# Papers base de HDFS

- Ghemawat, S.; Gobioff, H.; Leung, S. T. (2003). "The Google file system". [\*Proceedings of the nineteenth ACM Symposium on Operating Systems Principles - SOSP '03\*](#) (PDF).  
p. 29. [CiteSeerX 10.1.1.125.789](#). [doi:10.1145/945445.945450](#). [ISBN 1581137575](#).
- Konstantin Shvachko, Hairong Kuang, Sanjay Radia, Robert Chansler. The Hadoop Distributed File System, Yahoo!, Sunnyvale, California USA  
{Shv, Hairong, SRadia, Chansler}@Yahoo-Inc.com .

- Hadoop Distributed File System.
- Sistema de Archivos Distribuido.
  - Escalable
  - Soporta Replicación entre nodos. Sin necesidad de RAID
  - Bloque de 64 MB

Hadoop Distributed File System. El Hadoop Distributed File System (**HDFS**) es un sistema de archivos distribuido, escalable y portátil escrito en Java para el framework Hadoop.

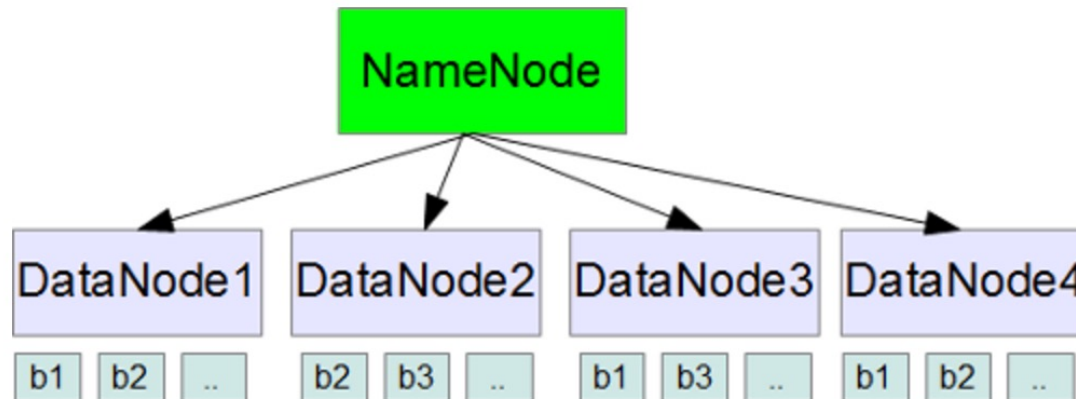
[Hadoop - Wikipedia, la enciclopedia libre](https://es.wikipedia.org/wiki/Hadoop)  
<https://es.wikipedia.org/wiki/Hadoop>

# HDFS - Introducción

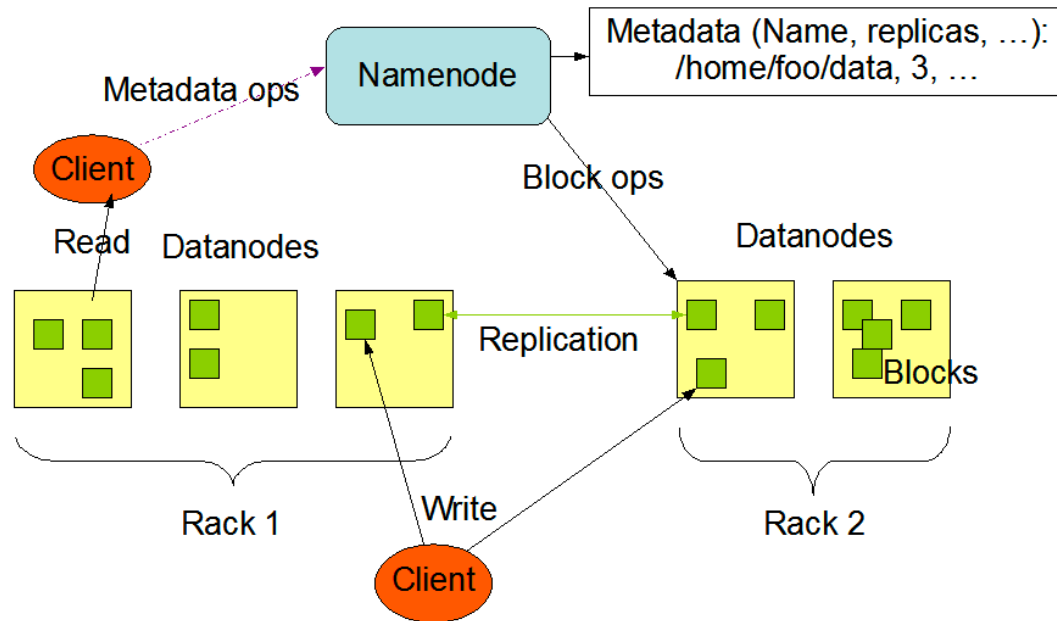
- Es un Sistema de Archivos Distribuidos
- Corre sobre HW commodity
- Altamente tolerable a fallos
- Alto desempeño
- Tamaño de archivos MUY grande y Data Sets grandes
- No compatible con POSIX
- Falla en hw es la norma
- Clusters de cientos o miles de nodos para Almacenar y Procesar.
- Implementado en Java

# Arquitectura HDFS

- **NameNode**: Gestiona los metadatos del sistema de archivos.
- **DataNode**: Almacenamiento de archivos en bloques 'b\*' y replicado entre varios nodos.



## HDFS Architecture



# Componentes

- **NameNode & DataNode**
- Normalmente hay uno (1) o dos (2) **NameNode** por Cluster ( primary + secondary NameNode)
  - Gestiona el espacio de nombres del File System
  - Acceso a archivos por los clients
  - Operaciones del file system: open, close, rename, directories.
  - Mapeo de Bloques a Datanodes.
  - Almacena metadatos del HDFS
- Muchos **Datanode** (almacena datos)
  - Almacena bloques (64 MB default)
  - Almacenamiento por bloques: Un archivo es dividido en BLOCKs
  - Cada bloque es replicado en diferentes Datanodes.



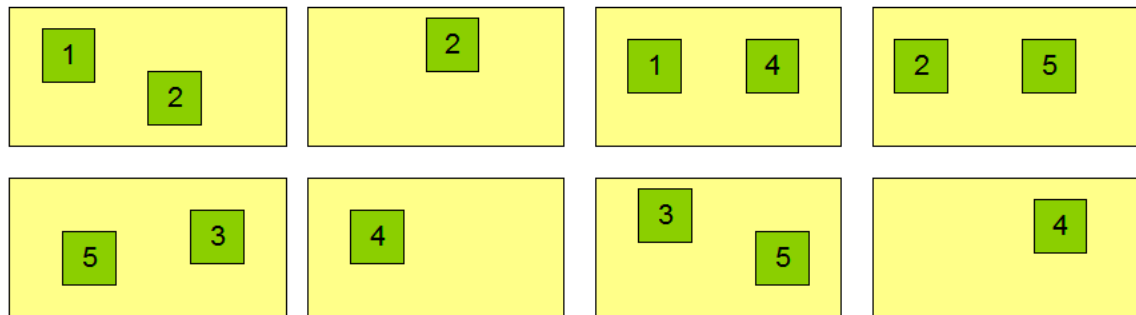
# Replicación de datos

- Los bloques de un archivo son replicados para soportar tolerancia a fallos.
- El factor de replicación -> default 3 (configurable)
- Tamaño de Bloque (64 MB by default, configurable mayor)
- Los archivos en HDFS son WORM
- Namenode monitorea cada Datanode. Protocolo Heartbeat y un Reporte de Bloques son enviados de regreso.

## Block Replication

Namenode (Filename, numReplicas, block-ids, ...)  
/users/sameerp/data/part-0, r:2, {1,3}, ...  
/users/sameerp/data/part-1, r:3, {2,4,5}, ...

## Datanodes



# Como se llevan datos al HDFS?

- Desde Shell o HUE/Ambari de forma interactiva (Lab incluido)
- Desde Bases de datos Relacionales externas (Sqoop)
- Desde fuentes en tiempo real (IoT) Ej: Kafka
- Desde fuentes Streaming (Twitter) Ej: con Flume