

# Aplicaciones del aprendizaje profundo en el análisis de grandes volúmenes de datos

Especialización en Ciencia de Datos

---

Mg. Diego Encinas – Ing. Román Bond



# Agenda-Clase 2

---

Hadoop

- HDFS

Introducción al paradigma MapReduce

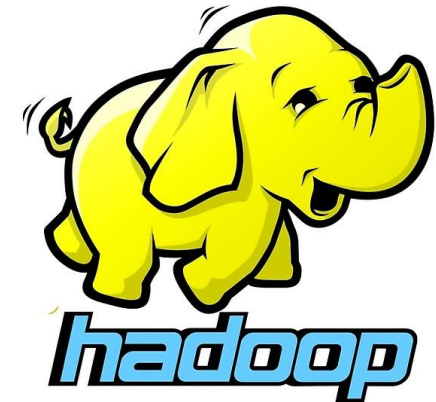
Etapas de un trabajo en MapReduce

- Map
- Shuffle
- Sort
- Reduce

# Hadoop

---

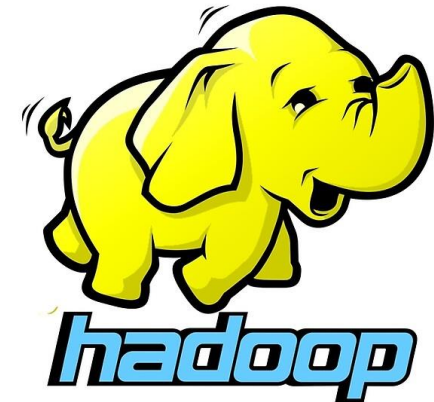
- Para procesar grandes conjuntos de datos, en 2003 Google creó el framework Hadoop capaz de poder procesar grandes volúmenes de datos.
- En 2006, Yahoo continúa con el desarrollo del proyecto Hadoop. Aparece Hadoop MapReduce.
- Actualmente pertenece a Apache
  - Apache Hadoop ([hadoop.apache.org](http://hadoop.apache.org))



# Hadoop

---

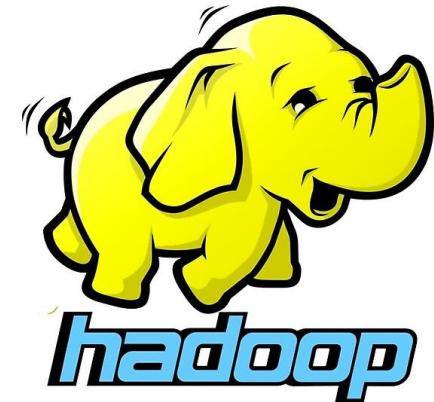
- Es un framework que soporta procesamiento de grandes bases de datos en un ambiente distribuido
- Ejecuta aplicaciones para el tratamiento de grandes volúmenes de datos
- Incluye un sistema de archivos distribuidos (HDFS)
- Tolerante a fallas



# Hadoop

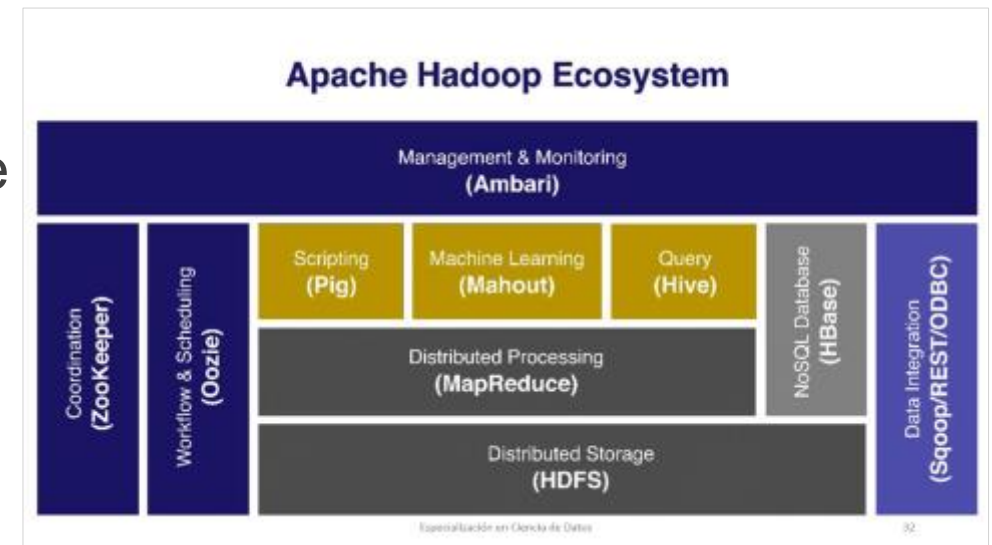
---

- ✓ Diseñado para el procesamiento off-line de los datos (procesamiento en batch)
- ✓ Funciona con la idea de "escriba una sola vez y lea muchas"
- ✗ No permite lectura aleatoria
- ✗ No permite el procesamiento on-line
- Se ejecuta en el "lugar" donde se encuentran los datos



# Componentes Hadoop

- **Common** (I/O, serialización, RPC)
- **HDFS** (file system distribuido)
- **Zookeeper** (servicio de coordinación de procesos)
- **MapReduce** (modelo de procesamiento de datos)
- **Pig** (lenguaje de scripting sobre MapReduce)
- **Cascading** (framework que simplifica el uso de MapReduce)
- **Hive** (lenguaje basado en SQL)



# Componentes Hadoop

- Almacenamiento: Distributed File System (DFS)
  - Los archivos están distribuidos
  - Ofrece transparencia al usuario permitiendo operar con todos los archivos del cluster a través del file system distribuido.
  - Un mismo archivo podría estar almacenado en varias computadoras.
  - Hadoop tiene su propio filesystem distribuido: el HDFS (Hadoop Distributed FileSystem)

Un **sistema de archivos** o **sistema de ficheros**, en informática, es un elemento que controla cómo se almacenan y recuperan los datos. Sus principales funciones son la asignación de espacio a los archivos, la administración del espacio libre y del acceso a los datos resguardados.

Un **sistema de archivos distribuido** o sistema de archivos de red es un sistema de archivos de computadoras que sirve para compartir archivos, impresoras y otros recursos como un almacenamiento persistente en una red de computadoras.

Diagrama adaptado de Wikipedia

18

# DFS

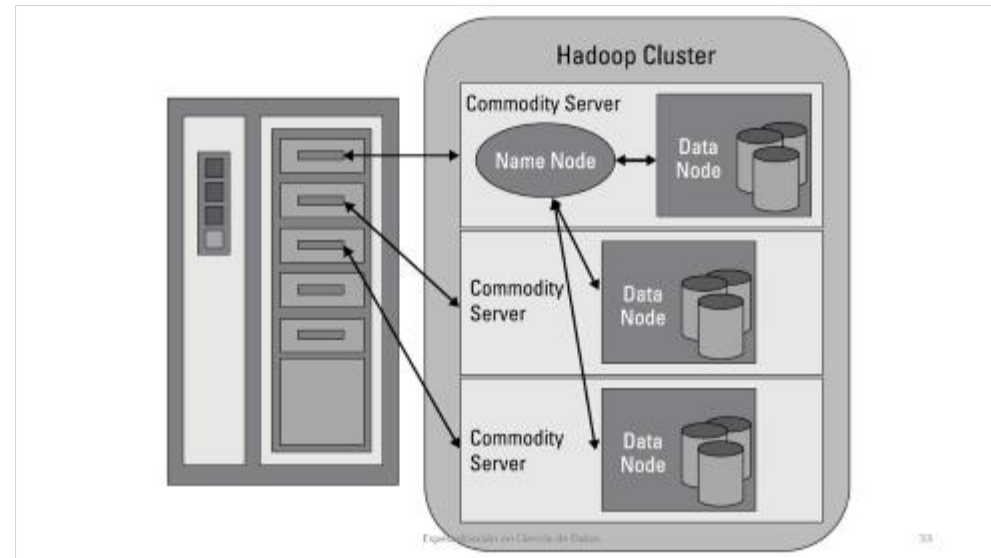
---

- Hay varios sistemas de archivos distribuidos,
  - HDFS
  - HFTP
  - HSFTP
  - HAR
  - FTP
  - S3



# HDFS

- Todos los archivos se dividen en bloques del mismo tamaño (64MB por defecto, aunque es configurable)
- Los bloques pueden estar físicamente en cualquier computadora
- Permite la réplica de bloques para optimización y recupero de fallas



# Procesos del HDFS

---

- Namenode
  - Maneja el árbol del filesystem y los metadatos de cada archivo y carpeta.
  - Conoce para cada bloque del FS que datanode lo maneja.
  - Vínculo con el filesystem del SO
- Datanode
  - Son lo que llevan a cabo la lectura y escritura de los bloques en el filesystem del SO.
  - Lleva a cabo la creación, borrado y replicado de los bloques.
- Secondary namenode: realiza tareas auxiliares al name node.

# El comando HDFS

---

- HDFS permite crear, borrar, renombrar archivos y carpetas dentro del FS distribuido.
  - Ofrece dos operaciones adicionales
  - Copiar un archivo del FS local al HDFS
- Copiar un archivo del HDFS al FS local

# Componentes Hadoop

---

- En Hadoop la administración de los procesos que se ejecutan en el cluster la lleva a cabo un framework llamado Yarn MapReduce.
- Básicamente Yarn realiza los trabajos usando dos procesos diferentes:
  - Job tracker: maneja todos los trabajos a ser procesados. Tiene en cuenta el mapa del cluster al momento de crear los procesos Task
  - Task tracker: son los encargados de realizar el procesamiento de los datos

# Agenda-Clase 2

---

## Hadoop

- HDFS
- Componentes Hadoop

## Introducción al paradigma MapReduce

### Etapas de un trabajo en MapReduce

- Map
- Shuffle
- Sort
- Reduce

# Paradigma MapReduce

---

- Es un framework para distribuir tareas en múltiples nodos
- El espíritu de MapReduce es "escriba una vez y lea muchas"
- Ventajas
  - Paralelización y distribución de tareas automática
  - Escalable
  - Tolerante a fallos
  - Monitoreo y capacidad de seguridad
  - Flexibilidad de programación (Java, Python, C#, Ruby, C++)
  - Abstracción al programador

# MapReduce

---

- Es, a su vez, un paradigma de programación.
- Hay que pensar como resolver un problema sin tener acceso a todos los datos
- Ejemplo (cálculo del promedio):

```
acum = 0
```

```
for d in datos:
```

```
    acum = acum + d
```

```
prom = acum / len(datos)
```

# MapReduce

---

- El problema del cálculo del promedio se debe "repensar".

*Pedirle a cada nodo que sume y cuente sus datos*

```
acum = 0; n = 0
```

```
for nodo in cluster:
```

```
    acum = acum + nodo.acum
```

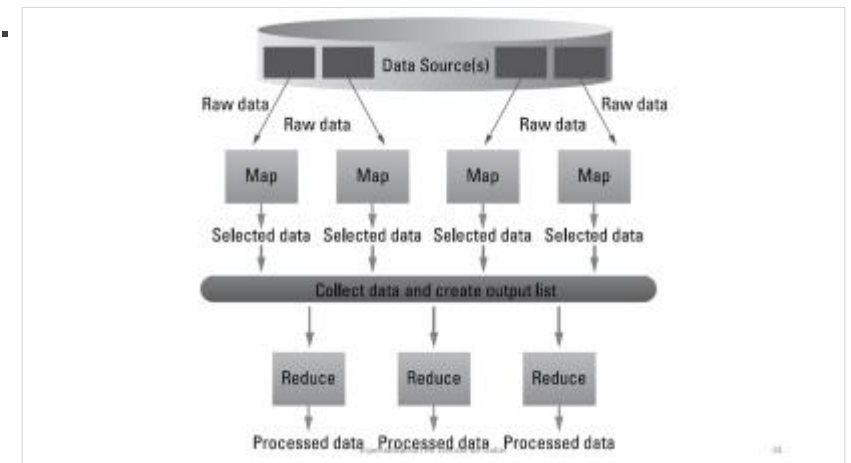
```
    n = n + nodo.n
```

```
promedio = acum / n
```



# MapReduce

- Se pensó en un proceso genérico que permita resolver cualquier problema
  - Paradigma MapReduce
- Toda tarea MapReduce se divide en dos fases:
  - Fase **map**: en la que los datos de entrada son procesados, uno a uno, y transformados en un conjunto intermedio de datos.
  - Fase **reduce**: se reúnen los resultados intermedios y se reducen a un conjunto de datos resumidos, que es el resultado final de la tarea.



# MapReduce

---

- En el ejemplo del promedio:

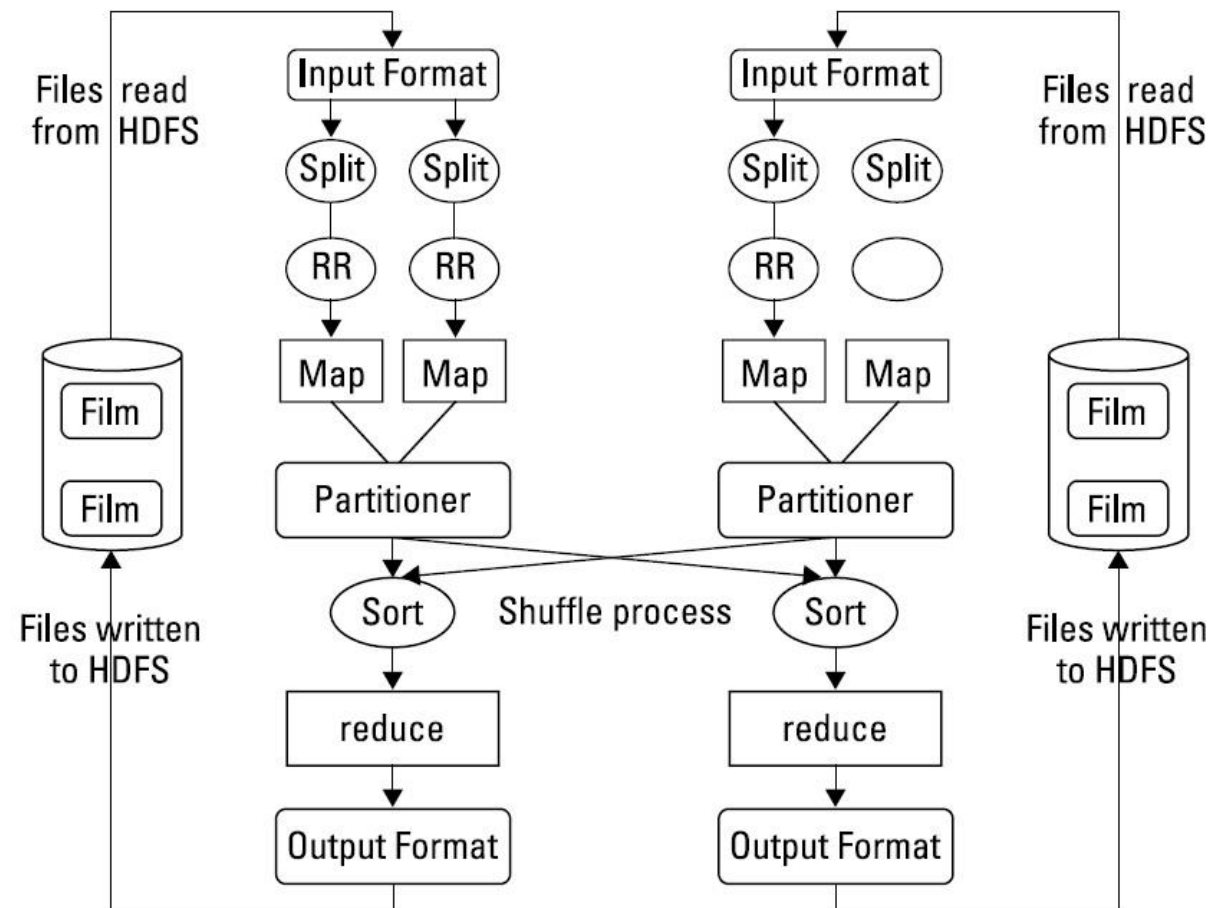
*Pedirle a cada nodo que sume y cuente sus datos*

Map

```
acum = 0; n = 0
for nodo in cluster:
    acum = acum + nodo.acum
    n = n + nodo.n
promedio = acum / n
```

Reduce

# MapReduce

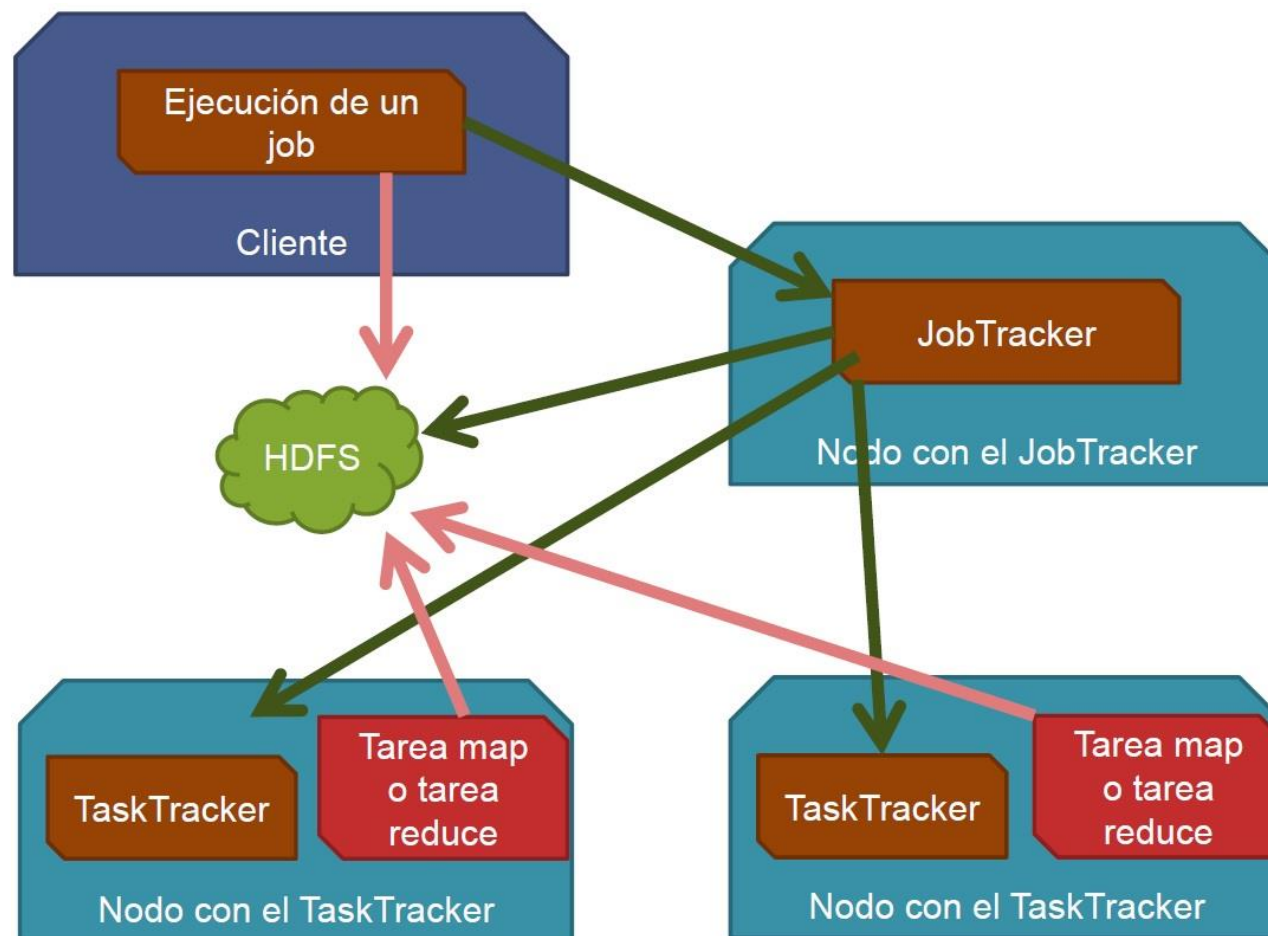


# MapReduce

---

- La unidad de trabajo de MapReduce es un Job
- Un Job se divide en una tarea map y una tarea reduce.
- Los Jobs de MapReduce son controlados por un daemon conocido como JobTracker, el cual reside en el "nodo master"
- Los clientes envían Jobs MapReduce al JobTracker y este distribuye la tarea usando otros nodos del cluster
- Esos nodos se conocen como TaskTracker y son responsables de la ejecución de la tarea asignada y reportar el progreso al JobTracker

# MapReduce

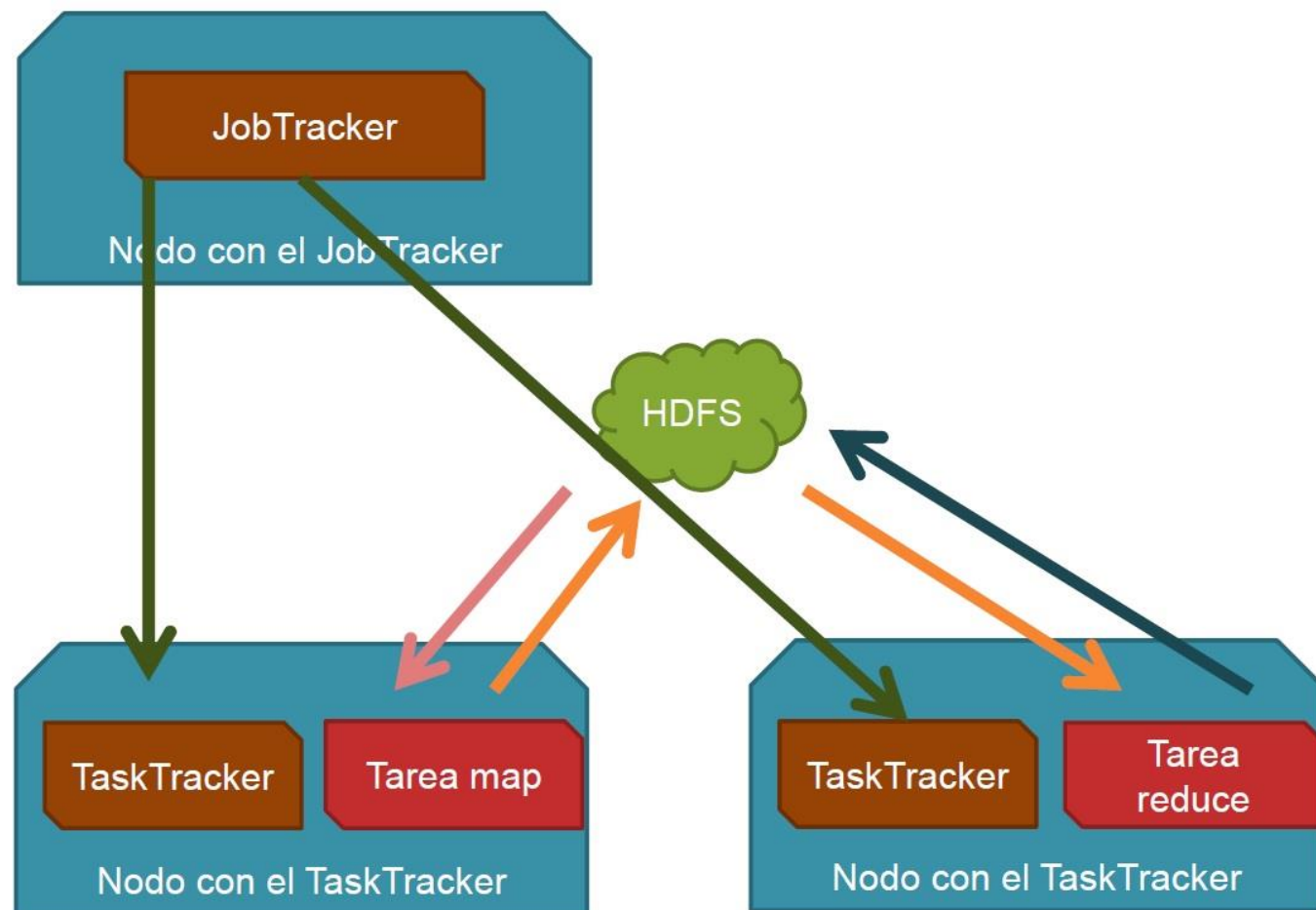


# MapReduce

---

- Un job MapReduce es un proceso que se divide en cuatro fases:
  - Map -> Shuffle -> Sort -> Reduce
- Map y reduce son las tareas que se deben programar para la aplicación.
- Cada TaskTracker ejecuta la tarea encomendada (map o reduce)
- Shuffle y sort son internas en la ejecución del job.

# MapReduce



# Preguntas? O ...

---

