

# Hadoop

Rodolfo Campos  
@camposer

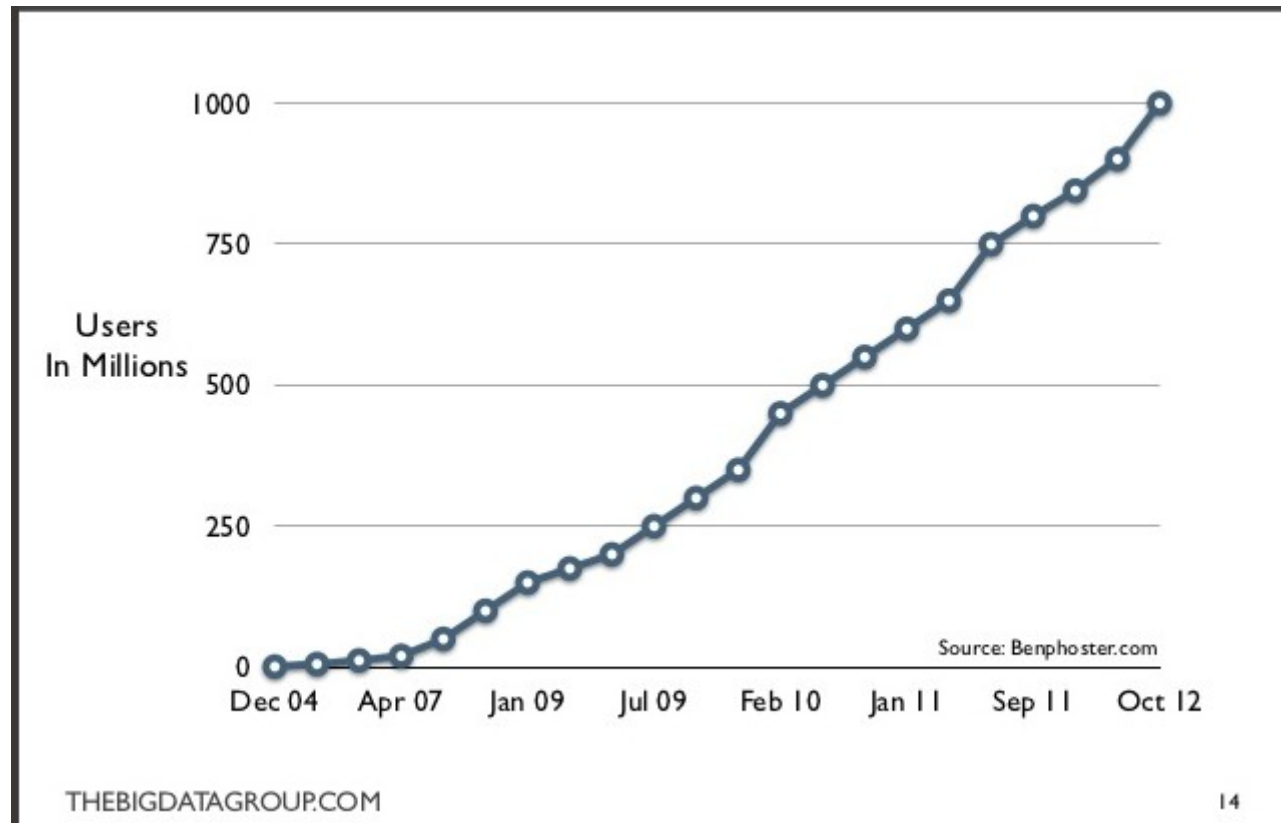
Madrid, Abril de 2014

# Agenda

1. Introducción a Big Data.
2. NoSQL
3. Hadoop

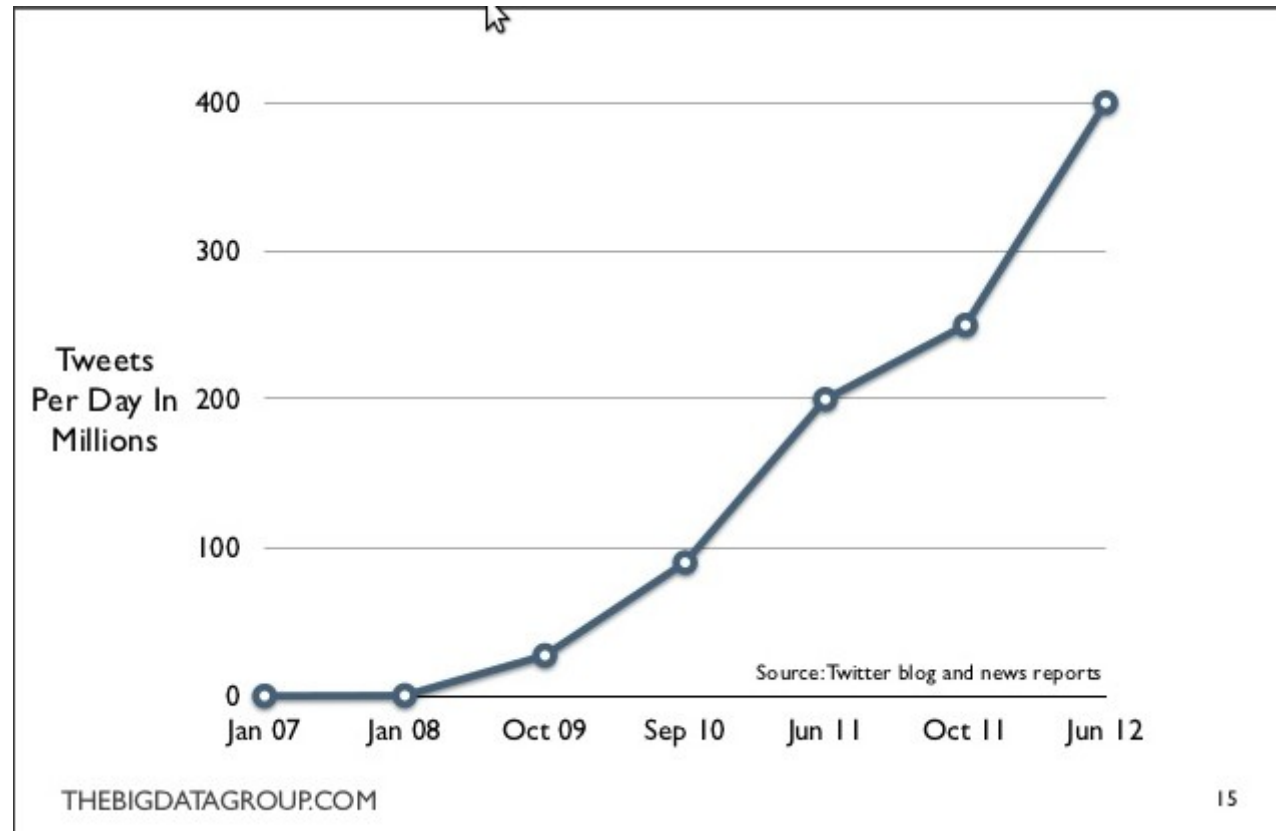
# Justificación

Tendencia: Usuarios de Facebook



# Justificación

Tendencia: Twitter. Tweets por día

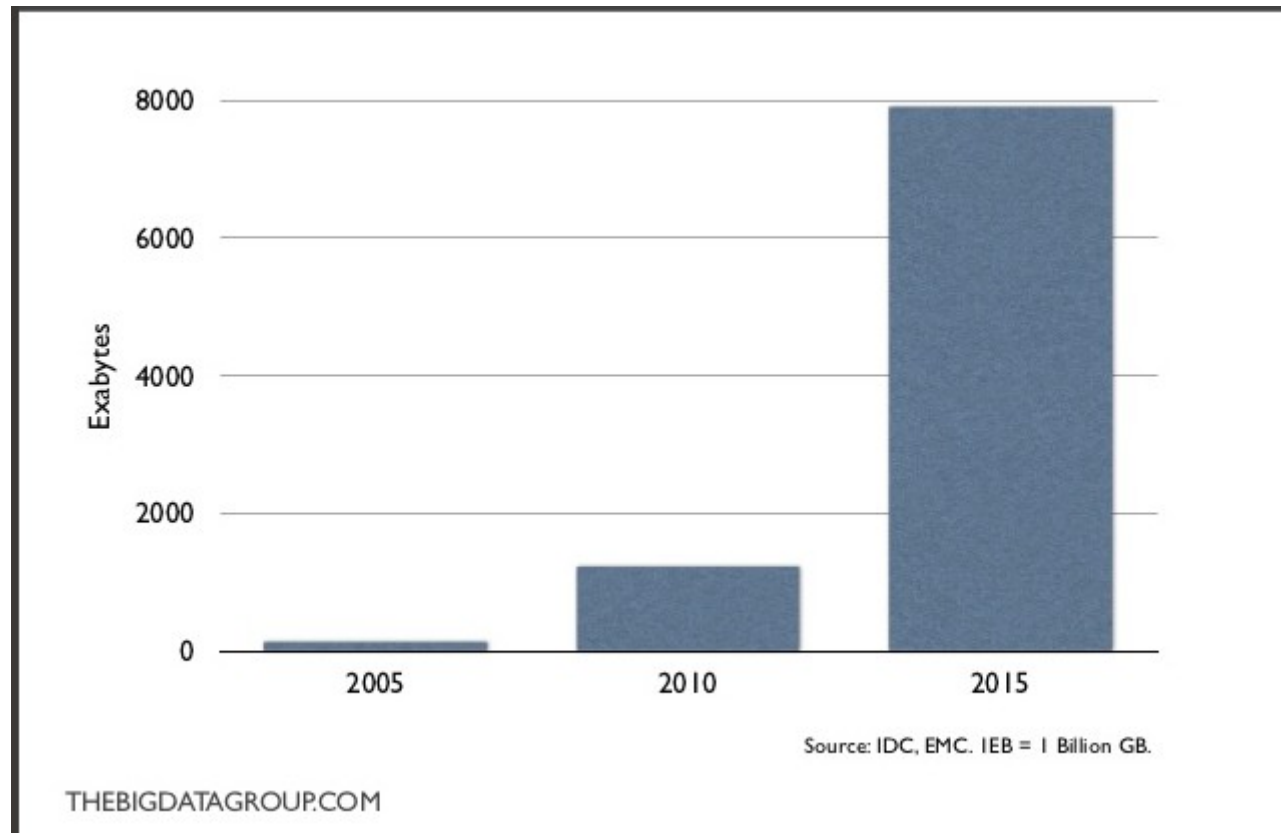


# Justificación

- Walmart maneja más de 1 millón de transacciones por hora
- Google procesa aprox. 24 PB de data al día
- AT&T transfiere aprox. 30 PB de data al día
- Todos los días son enviados aprox. 90 mil millones de correos electrónicos
- World of Warcraft utiliza 1.3 PB de almacenamiento

# Justificación

Tendencia: Crecimiento de la data en el mundo



# Definición

Según Wikipedia en Español, Big Data es:

- Big Data es en el sector de tecnologías de la información y la comunicación una referencia a los sistemas que **manipulan grandes conjuntos de datos** (o *data sets*).
- Las dificultades más habituales en estos casos se centran en la **captura, el almacenado, búsqueda, compartición, análisis, y visualización**.
- La tendencia a manipular ingentes cantidades de datos se debe a la necesidad en muchos casos de incluir los datos relacionados del análisis en un gran conjunto de datos relacionado, tal es el ejemplo de los análisis de negocio, los datos de enfermedades infecciosas, o la lucha contra el crimen organizado.

# La búsqueda

- La búsqueda de datos se realiza a través de diferentes herramientas, generalmente: lenguajes SQL, pseudo SQL y específicos.
- Debido a la naturaleza, generalmente distribuida del almacenamiento, algunos motores se valen de entornos de trabajo específicos como MapReduce que buscan operar sobre los diferentes nodos del sistema de forma independiente y generar resultados independientes a partir de estos.



# El almacenado

- El almacenamiento se lleva a cabo generalmente apoyándose en motores de bases de datos NoSQL
- El volumen de datos es un reto, por lo que la escalabilidad y la tolerancia a fallos son de los principales retos a manejar.
- Generalmente el almacenamiento se realiza en Sistemas Distribuidos en lugar de Centralizados

# El análisis

- El análisis de los datos se realiza de forma inline y offline, dependiendo del problema a resolver y del motor de base de datos utilizado.
- La mayoría de motores de bases de datos (en la mayoría NoSQL) proveen conectores para trabajar con los lenguajes de programación más populares.
- Si el problema a resolver no requiere de un mayor procesamiento, el análisis de los datos se puede elaborar de forma inline, en caso contrario, se suelen utilizar técnicas offline (Ej. trabajos por lote)

# NoSQL

Según la Wikipedia en Español, NoSQL se refiere a:

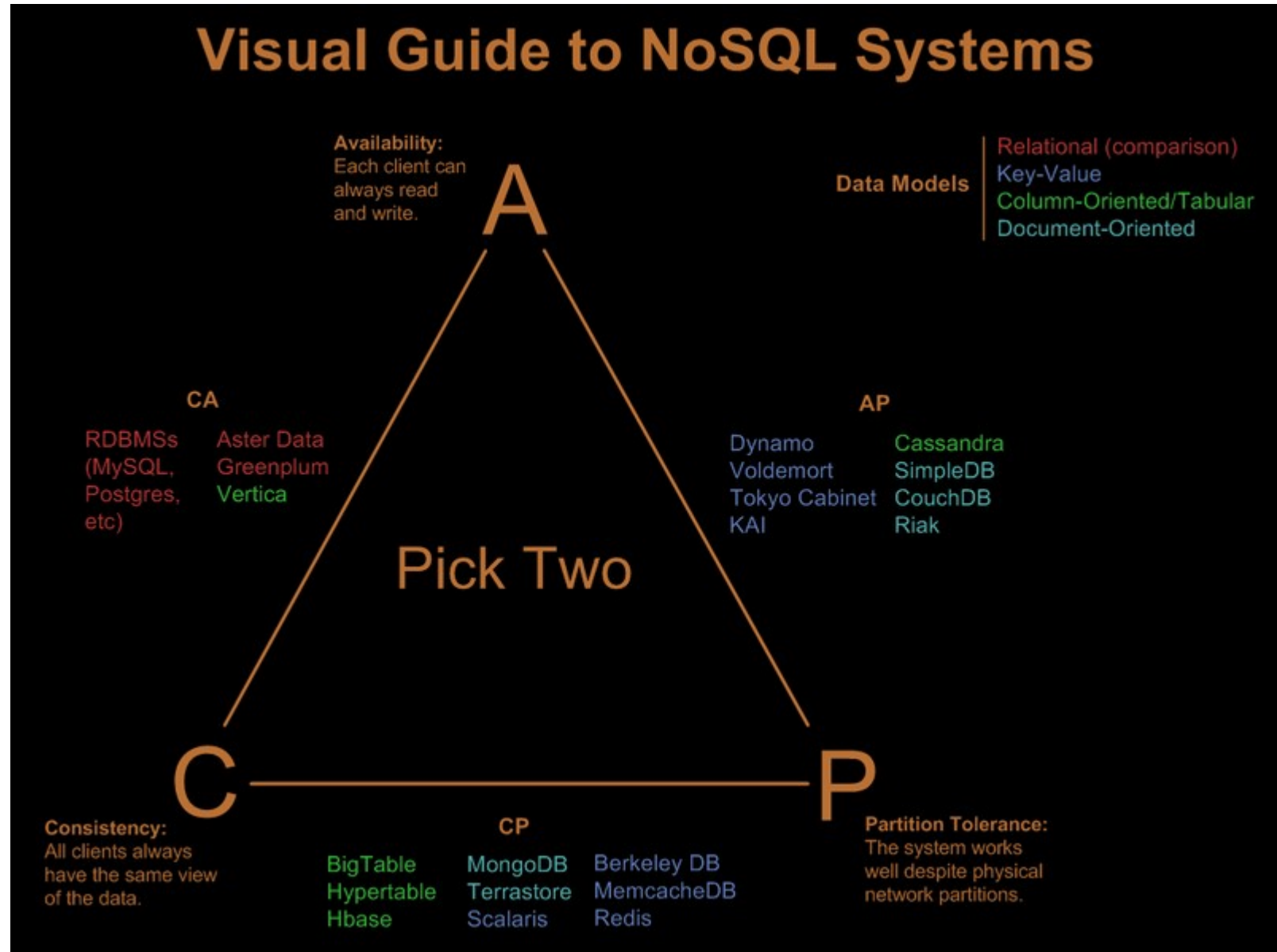
- Una amplia clase de sistemas de gestión de bases de datos que difieren del modelo clásico del sistema de gestión de bases de datos relacionales (RDBMS) en aspectos importantes, el más destacado que no usan SQL como el principal lenguaje de consultas.
- Los datos almacenados no requieren estructuras fijas como tablas, normalmente no soportan operaciones JOIN, ni garantizan completamente ACID (atomicidad, coherencia, aislamiento y durabilidad), y habitualmente escalan bien horizontalmente

# Tipos de motores NoSQL

Los tipos de motores NoSQL más importantes son:

- Tabulares u Orientados a Columnas (Ej. HBase, Cassandra)
- Orientados a documentos (Ej. MongoDB)
- De clave-valor (Ej. Redis o Memcached)
- Orientados a grafos (Ej. Neo4J)

# CAP



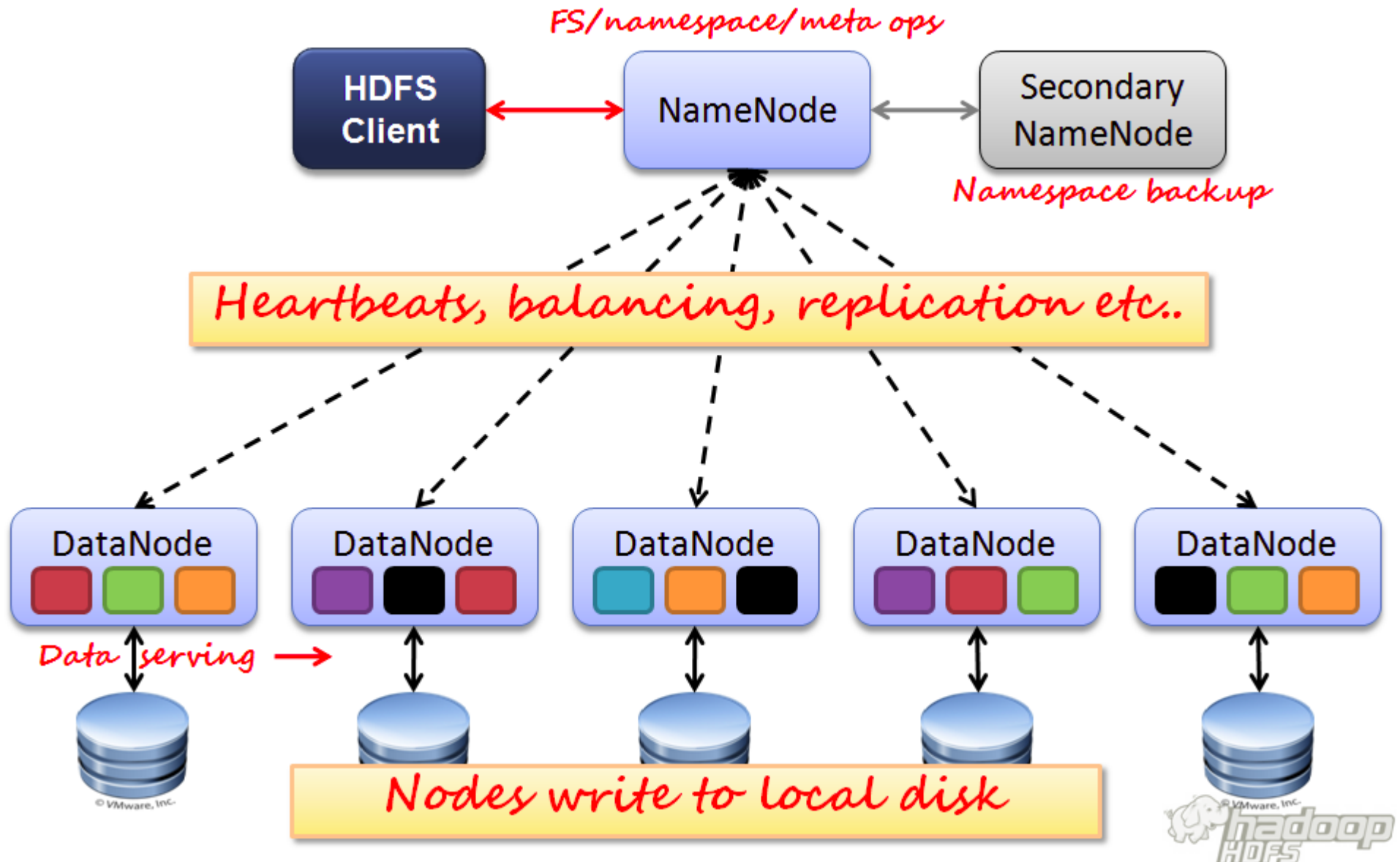
# Hadoop

- Apache Hadoop es un marco de trabajo que facilita la computación distribuida para grandes conjuntos de datos sobre clusters de máquinas utilizando modelos de programación sencillos.
- Fue diseñado para escalar rápidamente y está compuesto por los siguientes módulos:
  - Hadoop Common. Utilidades comunes usadas por todos los módulos.
  - Hadoop Distributed File System (HDFS). Un sistema de archivos distribuido que provee un alto rendimiento en el acceso a datos.
  - Hadoop YARN (Yet Another Resource Negotiator). Un marco de trabajo para planificación de trabajos (jobs) y gestión de los recursos de cluster.
  - Hadoop MapReduce. Un sistema basado en YARN para el procesamiento de grandes volúmenes de datos.

# Hadoop

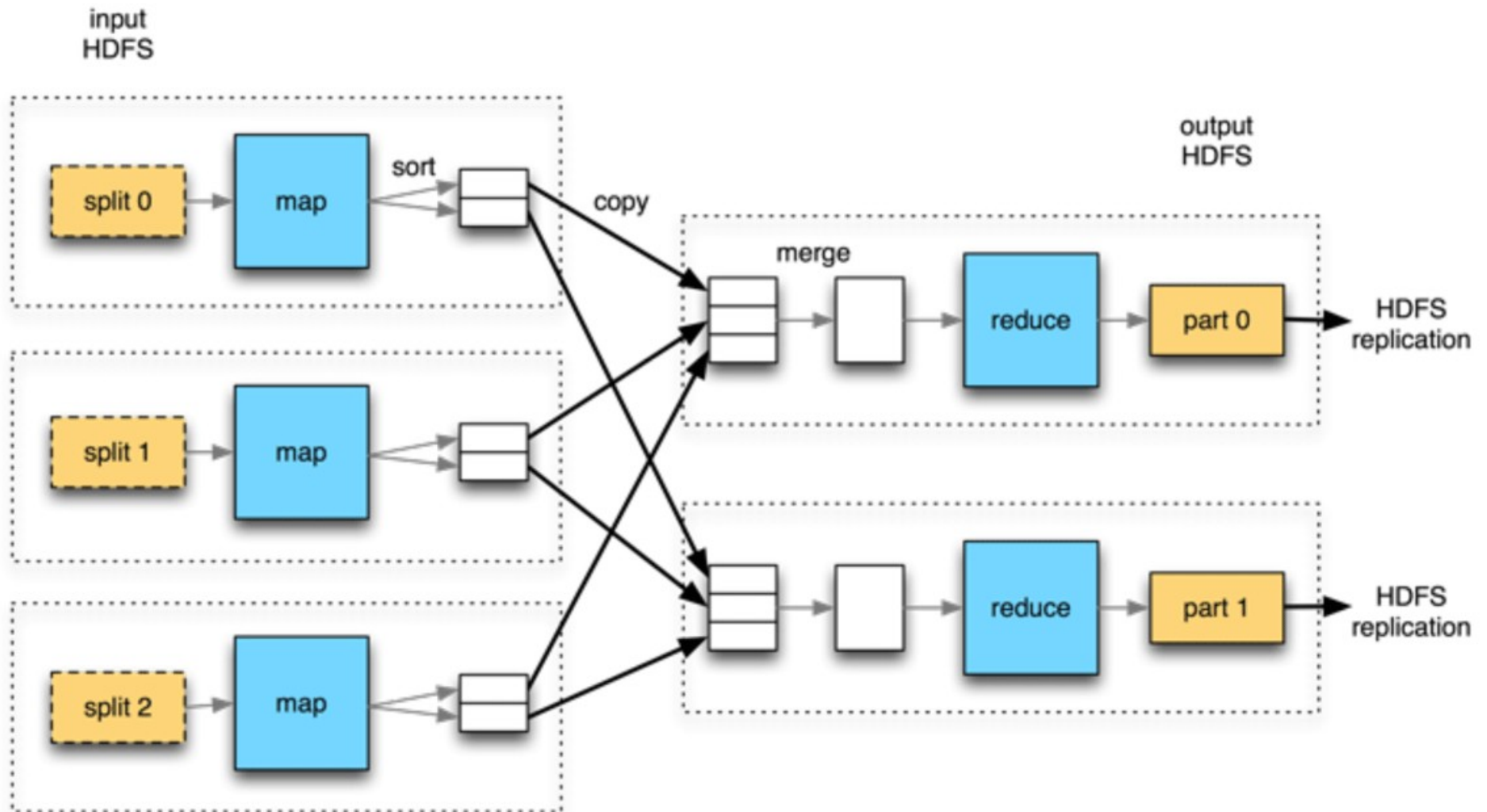
- Aparte de los módulos mencionados anteriormente, hay un conjunto de tecnologías relacionadas al proyecto. Algunas de las más importantes:
  - Pig. Lenguaje de flujo de datos para consultar HDFS o HBase. Las sentencias son traducidas a MapReduce.
  - Hive. Interfaz SQL para consultar HDFS o HBase. Las sentencias son traducidas a MapReduce.
  - HBase. Manejador de Bases de Datos NoSQL tabular que provee acceso aleatorio rápido, sobre HDFS y soporta MapReduce.
  - Sqoop. Permite transferir datos entre RDBMS y Hadoop.
  - Oozie. Automatiza la gestión de trabajos (jobs).
  - Flume. Herramientas de agregación distribuida.
  - Avro. Sistema de serialización de datos flexible basado en JSON.
  - Mahout. Librería scalable de Aprendizaje Automático (*Machine Learning*).

# Hadoop - HDFS

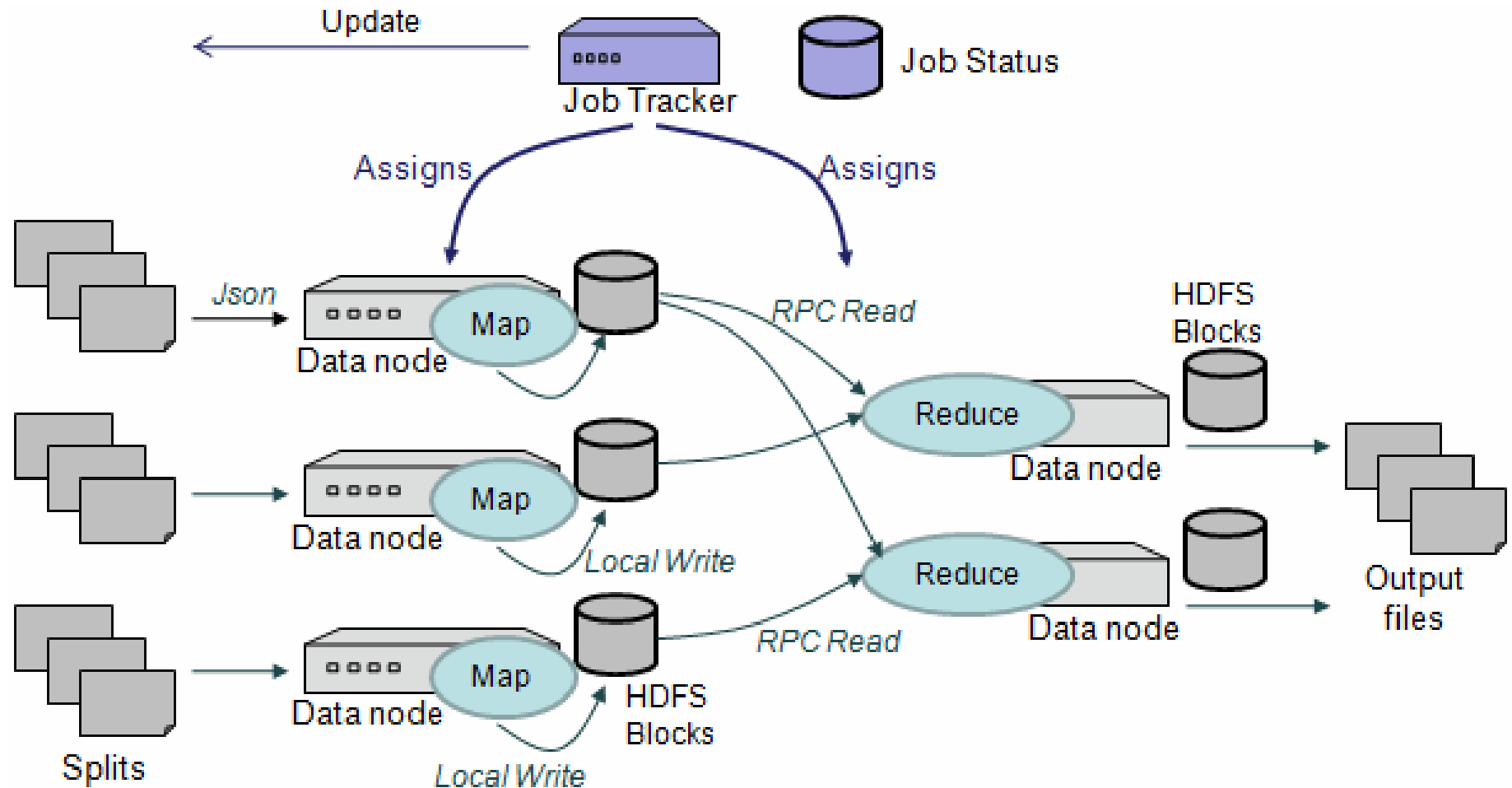




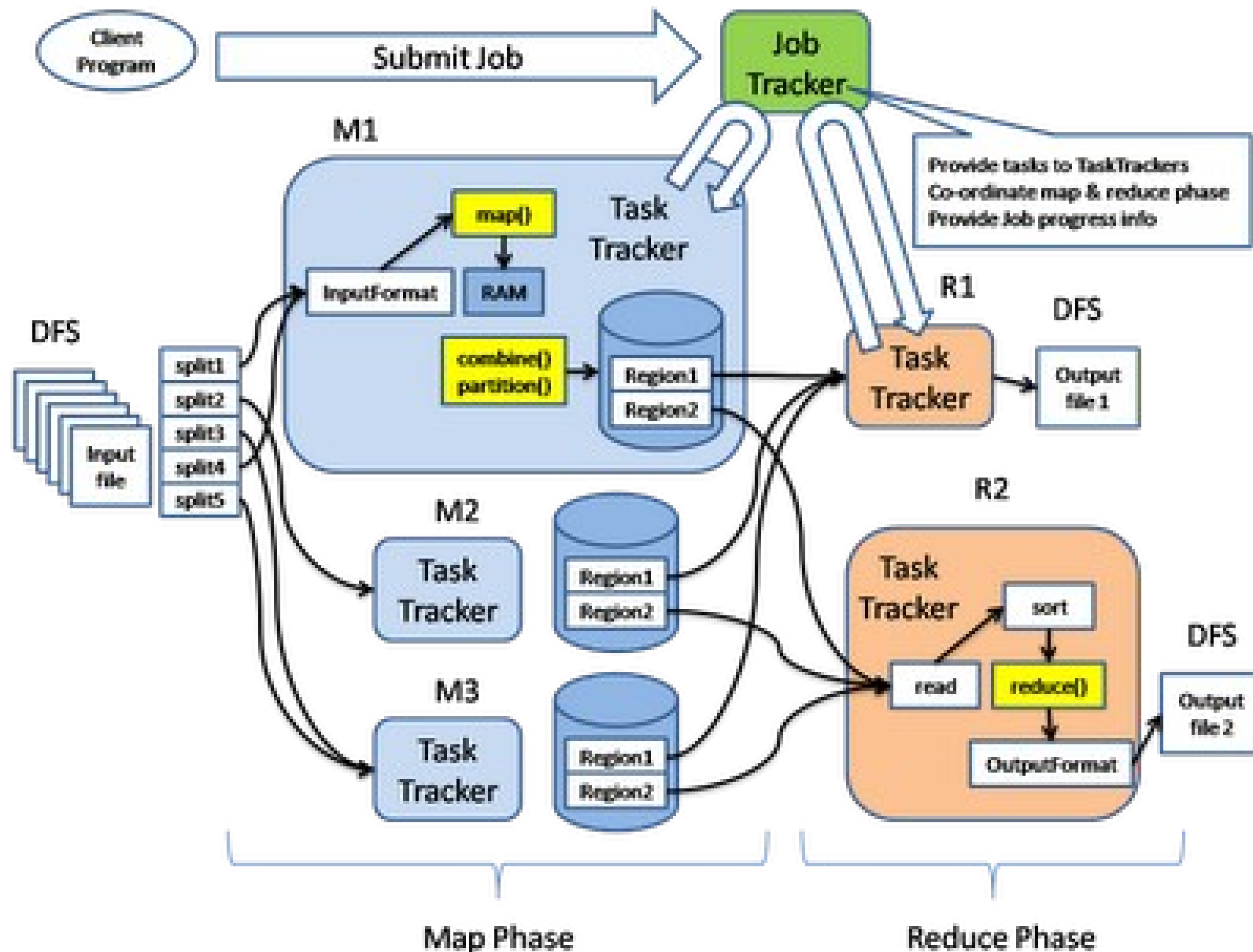
# Hadoop - MapReduce



# Hadoop - MapReduce



# Hadoop - MapReduce



# Pig

- Pig provee una herramienta para el procesamiento de datasets a través de transformaciones. Las diferentes transformaciones son traducidas a operaciones MapReduce.
- Pig tiene dos componentes: el lenguaje (PigLatin) y el entorno de ejecución (sobre Hadoop).
- Pig posee un intérprete de comandos llamado Grunt, aunque también permite scripts y modo embebido (nativo Java)
- El lenguaje permite utilizar estructuras, sentencias, expresiones, tipos, esquemas, funciones predefinidas, macros (procedimientos) y funciones definidas por el usuario.

# Pig

- La documentación completa se puede conseguir en el sitio Web. Luego, la documentación del libro de Hadoop (Capítulo 11) está bastante completa.
- Un ejemplo para discutir:

```
-- max_temp.pig: Finds the maximum temperature by year
records = LOAD 'input/ncdc/micro-tab/sample.txt'
AS (year:chararray, temperature:int, quality:int);
filtered_records = FILTER records BY temperature != 9999 AND
(quality == 0 OR quality == 1 OR quality == 4 OR quality == 5 OR quality == 9);
grouped_records = GROUP filtered_records BY year;
max_temp = FOREACH grouped_records GENERATE group,
MAX(filtered_records.temperature);
DUMP max_temp;
```

# HBase

- Motor NoSQL Orientado a Columnas, específicamente a familia de columnas. Cada fila contiene un conjunto de columnas pertenecientes a una familia.
- Funciona sobre HDFS y ZooKeeper.
- Ofrece acceso aleatorio garantizando alta disponibilidad
- Se pueden agregar dinámicamente nuevas columnas siempre y cuando pertenezcan a familias previamente definidas.
- Las tablas son particionadas horizontalmente automáticamente en regiones.
- Las actualizaciones de filas son atómicas.

# HBase

