

# Sistemas BD: procesamiento batch

Sistemas de Big Data

Ricardo García Ródenas  
Ricardo.Garcia@uclm.es



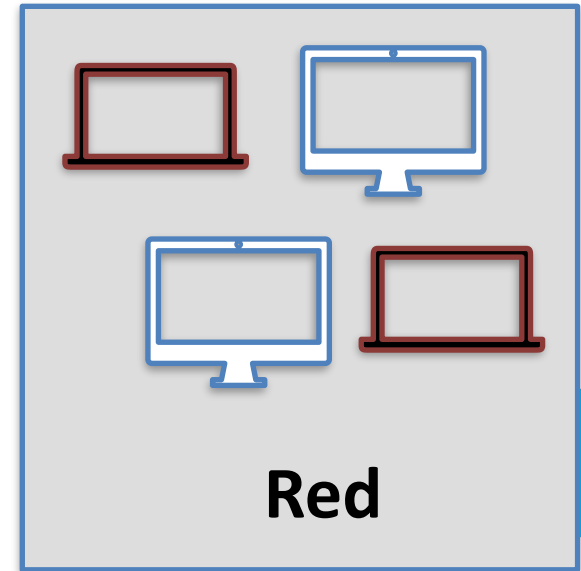
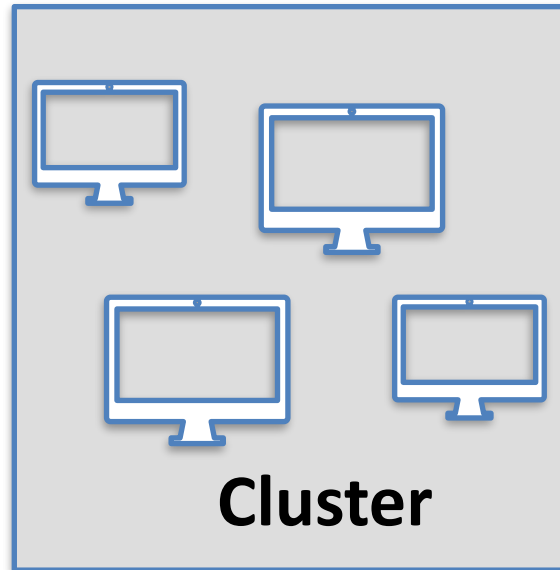
Sistemas de  
Big\_Data



Curso Especialización  
Inteligencia\_Artificial y  
Big\_Data

# Procesamiento batch

- Recolección y almacenamiento de datos **antes** de procesamiento
- Computación distribuida: **MapReduce**



# MapReduce

Google 2004

- Máquinas procesamiento MapReduce
- Google File System (GFS)
- BigTable



# MapReduce

Apache Software  
Foundation

- Hadoop MapReduce
- Hadoop YARN
- Hadoop Distributed File System (HDFS) y Hbase



# MapReduce

Apache Software  
Foundation

- Entorno programación
- Escalabilidad
- Tolerante a fallos



# MapReduce

## Datos

- Registros **< clave, valor >**

$\langle k, v \rangle$



# MapReduce

## Etapas

- Etapa Map

Función

**map(k,v)**

- Etapa Shuffle

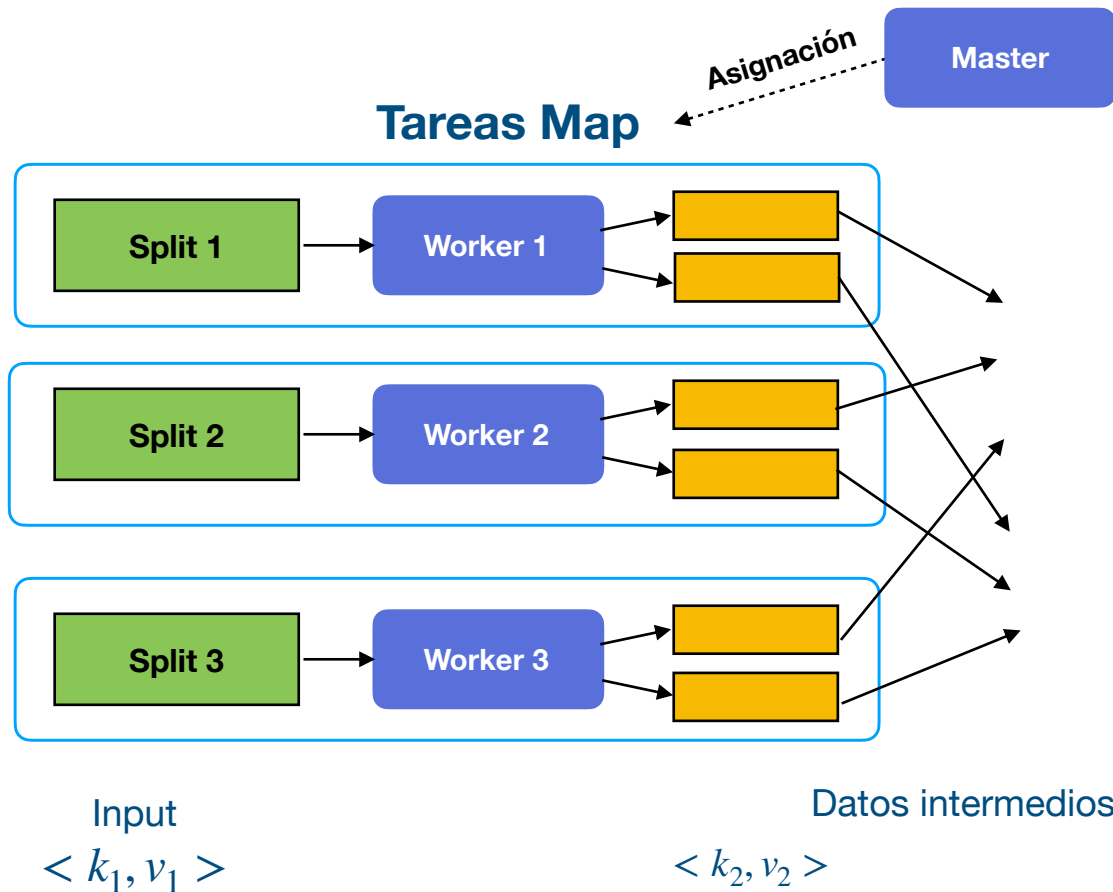
- Etapa Reduce

Función

**reduce(k,list(v))**

# MapReduce

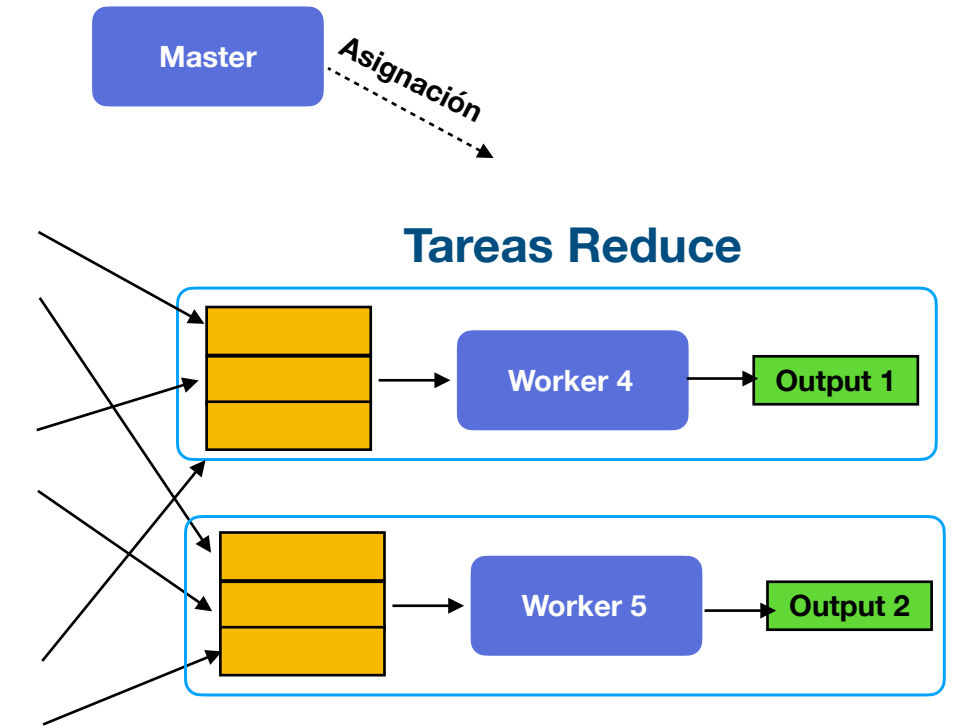
## Etapa Map





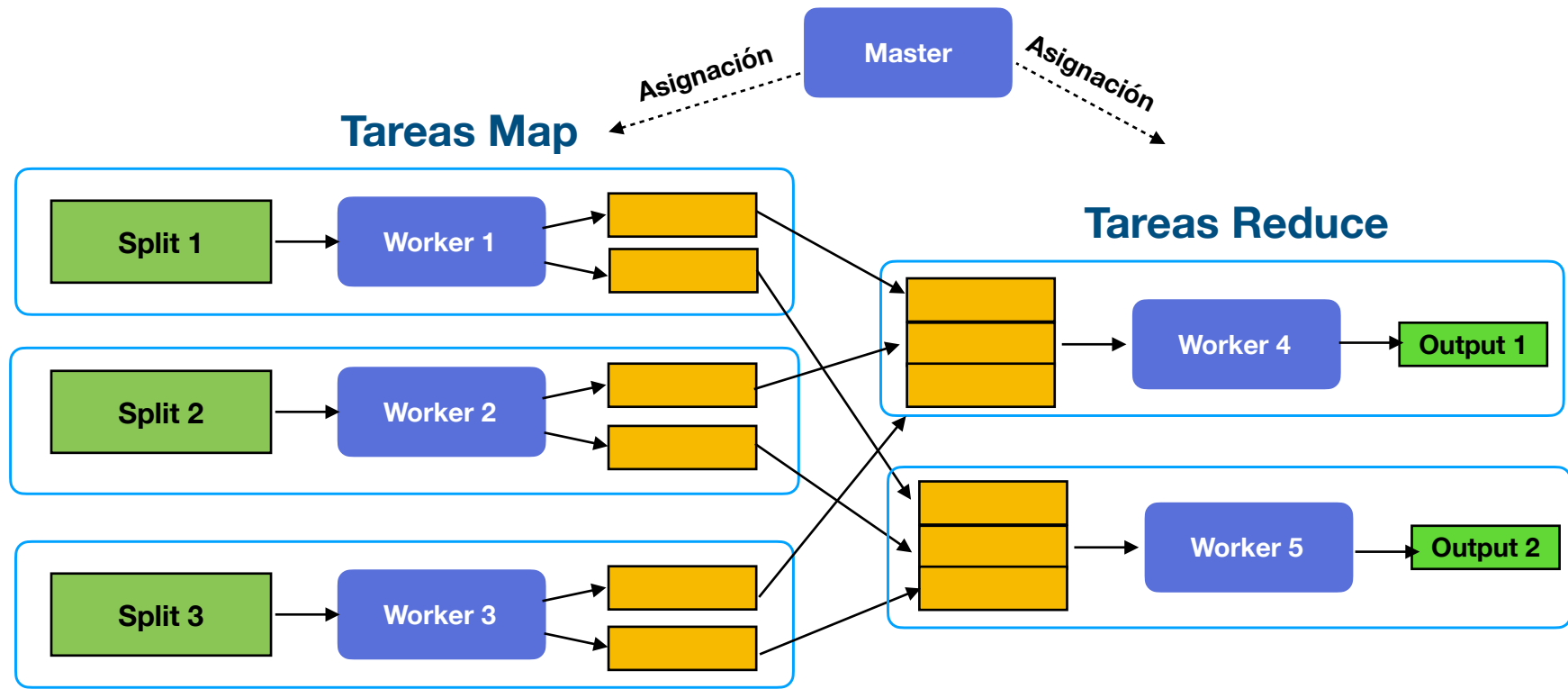
# MapReduce

## Etapa Reduce



Datos intermedios  
 $\langle k_2, [v_2] \rangle$

Output  
 $\langle k_3, v_3 \rangle$



Input  
 $\langle k_1, v_1 \rangle$

Datos intermedios  
 $\langle k_2, v_2 \rangle$        $\langle k_2, [v_2] \rangle$

Output  
 $\langle k_3, v_3 \rangle$

# Sistemas BD procesamiento batch

## MapReduce

Sistemas de Big Data

Ricardo García Ródenas  
Ricardo.Garcia@uclm.es



Sistemas de  
Big\_Data



Curso Especialización  
Inteligencia\_Artificial y  
Big\_Data