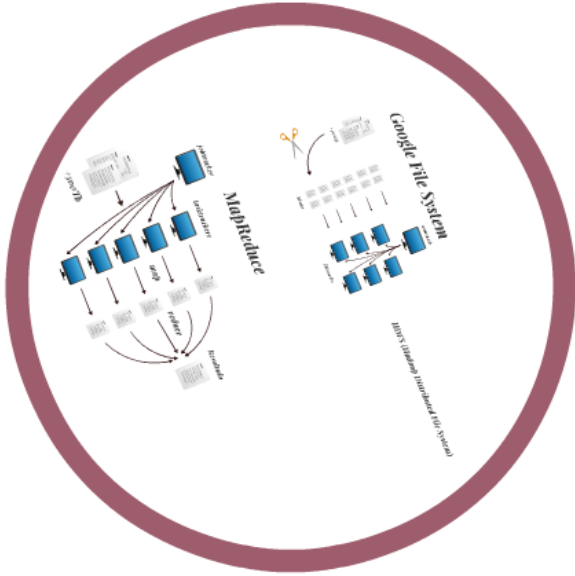
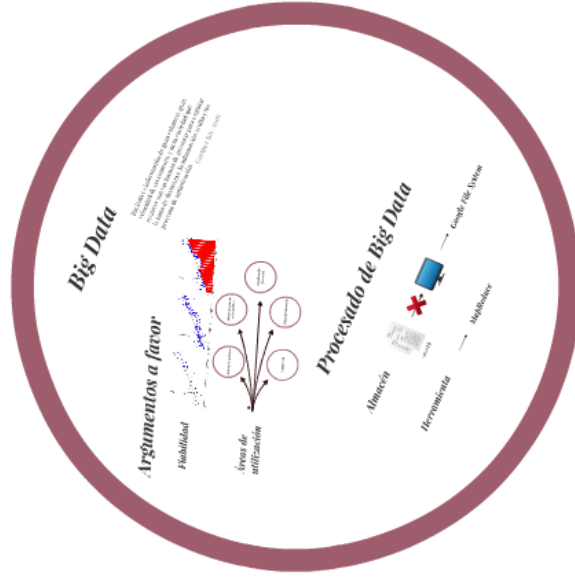


HDFS y MapReduce

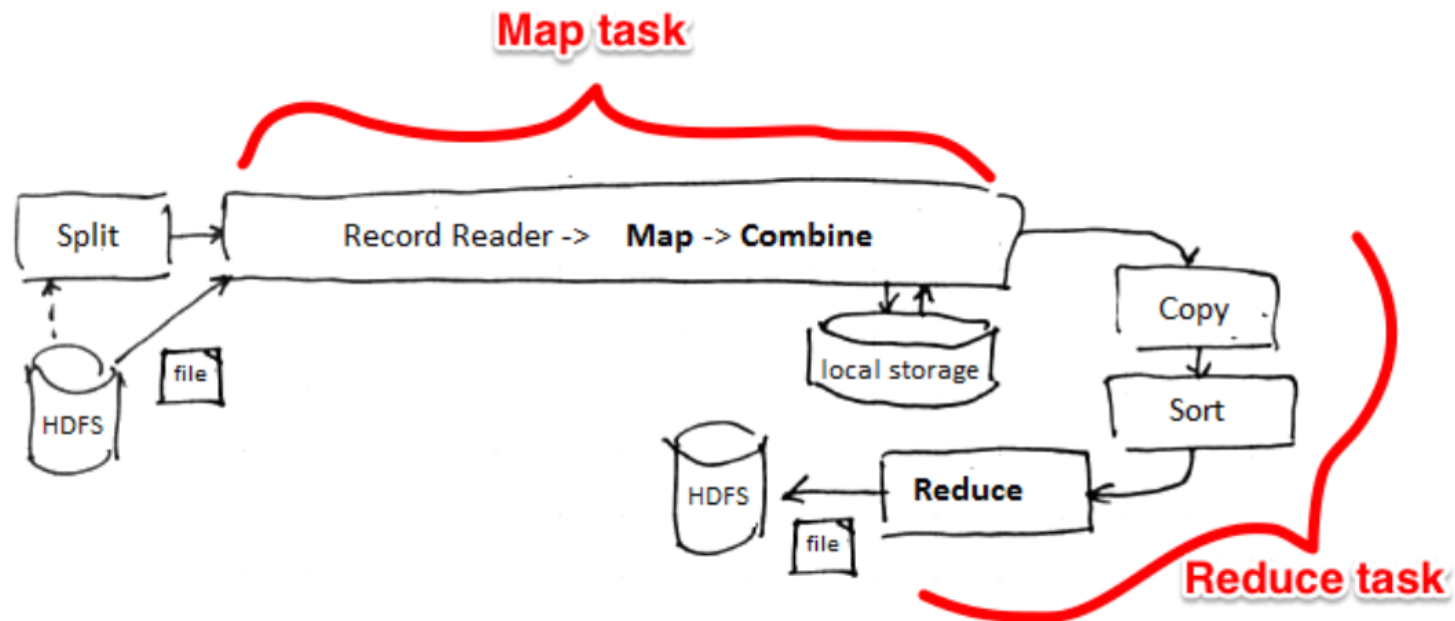


Introducción



Ejemplo práctico

MapReduce, el algoritmo de Big Data



Alejandro Tejera Pérez

Jornadas Desmitificando Big Data - 2014

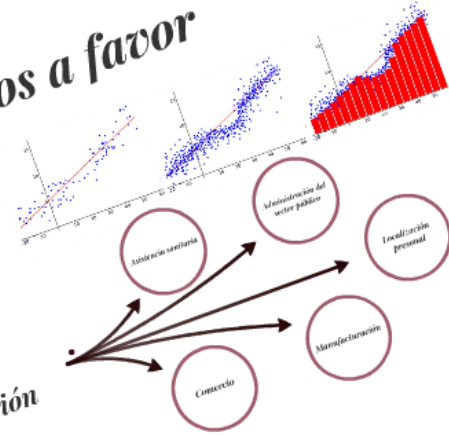
Big Data

Big Data es información de gran volumen, gran velocidad de crecimiento, y gran variedad que requiere nuevas formas de procesar para explotar la toma de decisiones, la información oculta y los procesos de optimización. Gartner Inc. 2001

Argumentos a favor

Fiabilidad

Áreas de utilización



Procesado de Big Data

Almacén



+300Tb



Google File System

Herramienta

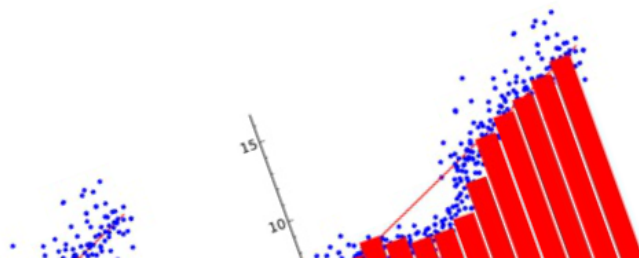


MapReduce

Big Data

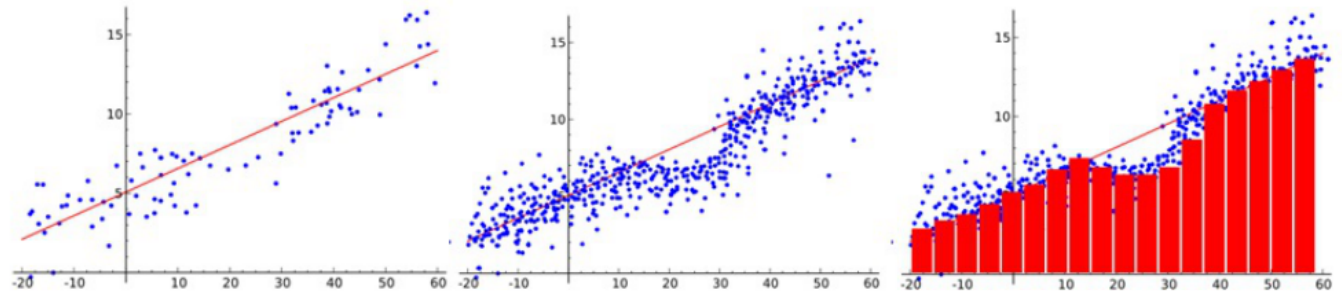
Big Data es información de gran volumen, gran velocidad de crecimiento, y gran variedad que requiere nuevas formas de procesar para explotar la toma de decisiones, la información oculta y los procesos de optimización.

Gartner Inc. 2001

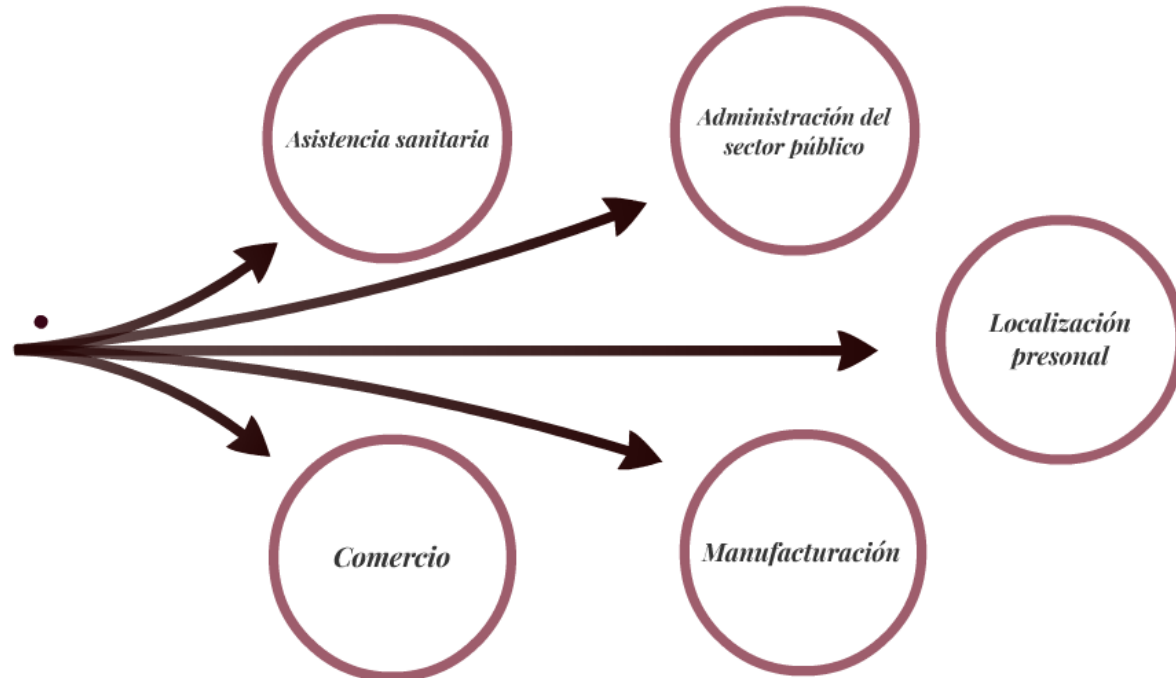


Argumentos a favor

Fiabilidad



*Áreas de
utilización*



B
ve
req
la to
proce



Procesado de Big Data

Almacén



+300Tb



Google File System

Herramienta



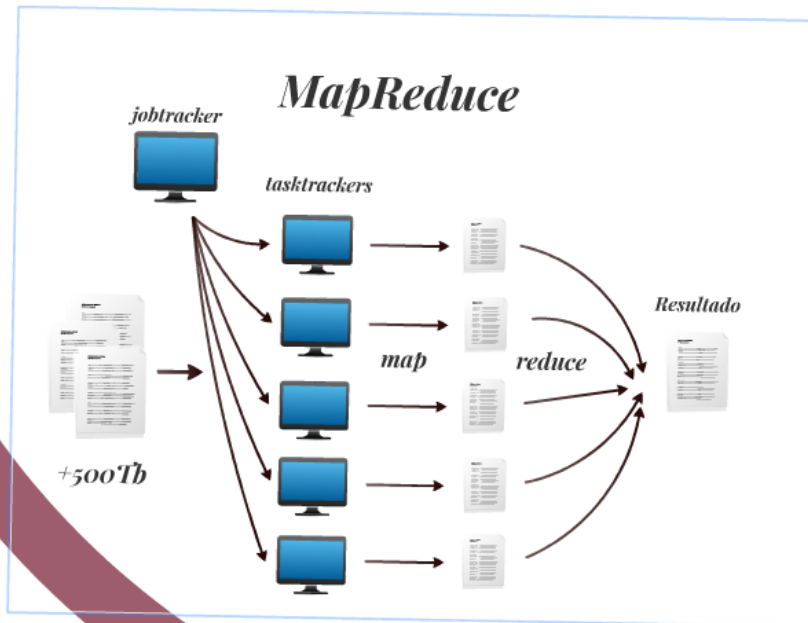
MapReduce

Google File System



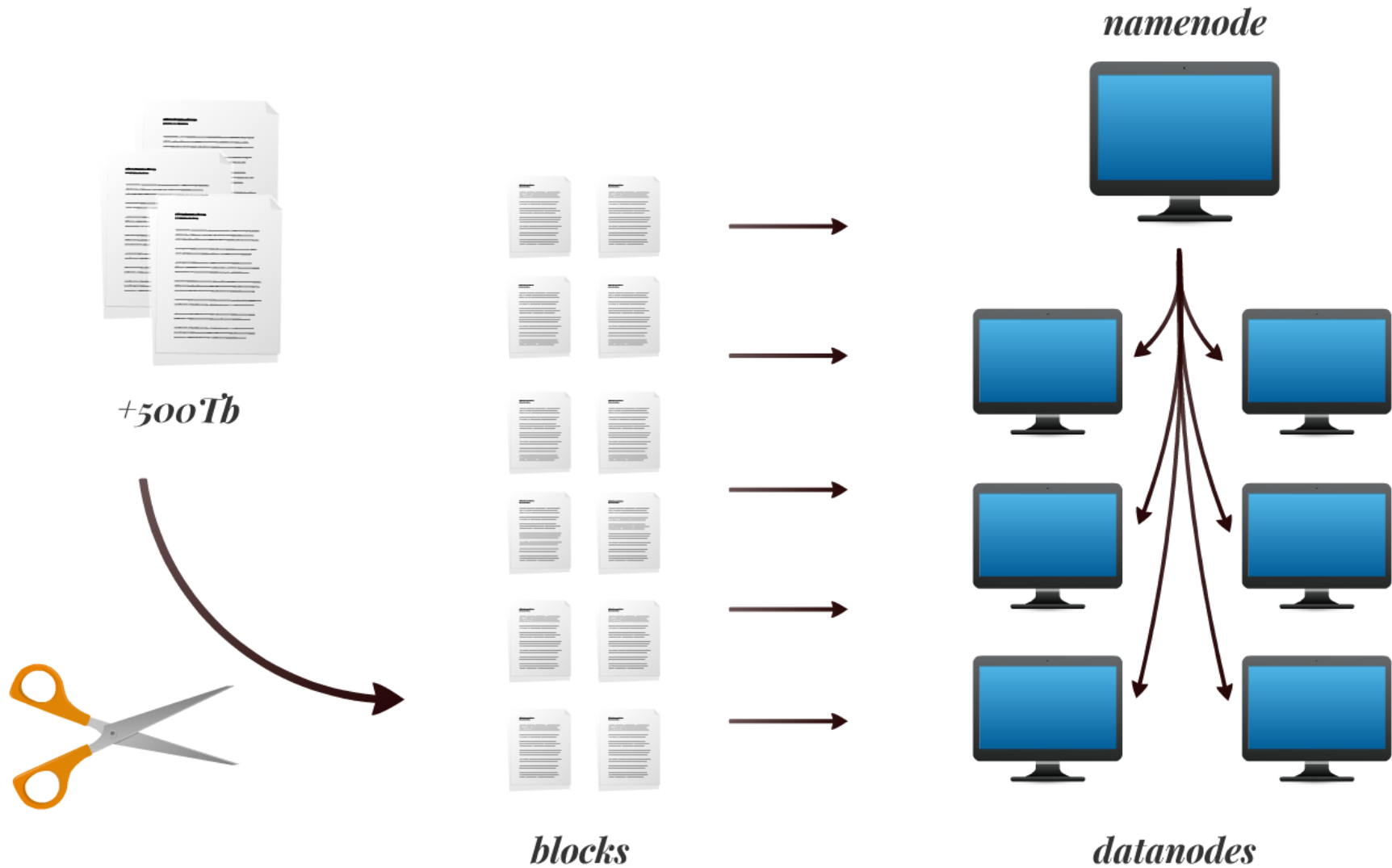
HDFS (Hadoop Distributed File System)

MapReduce



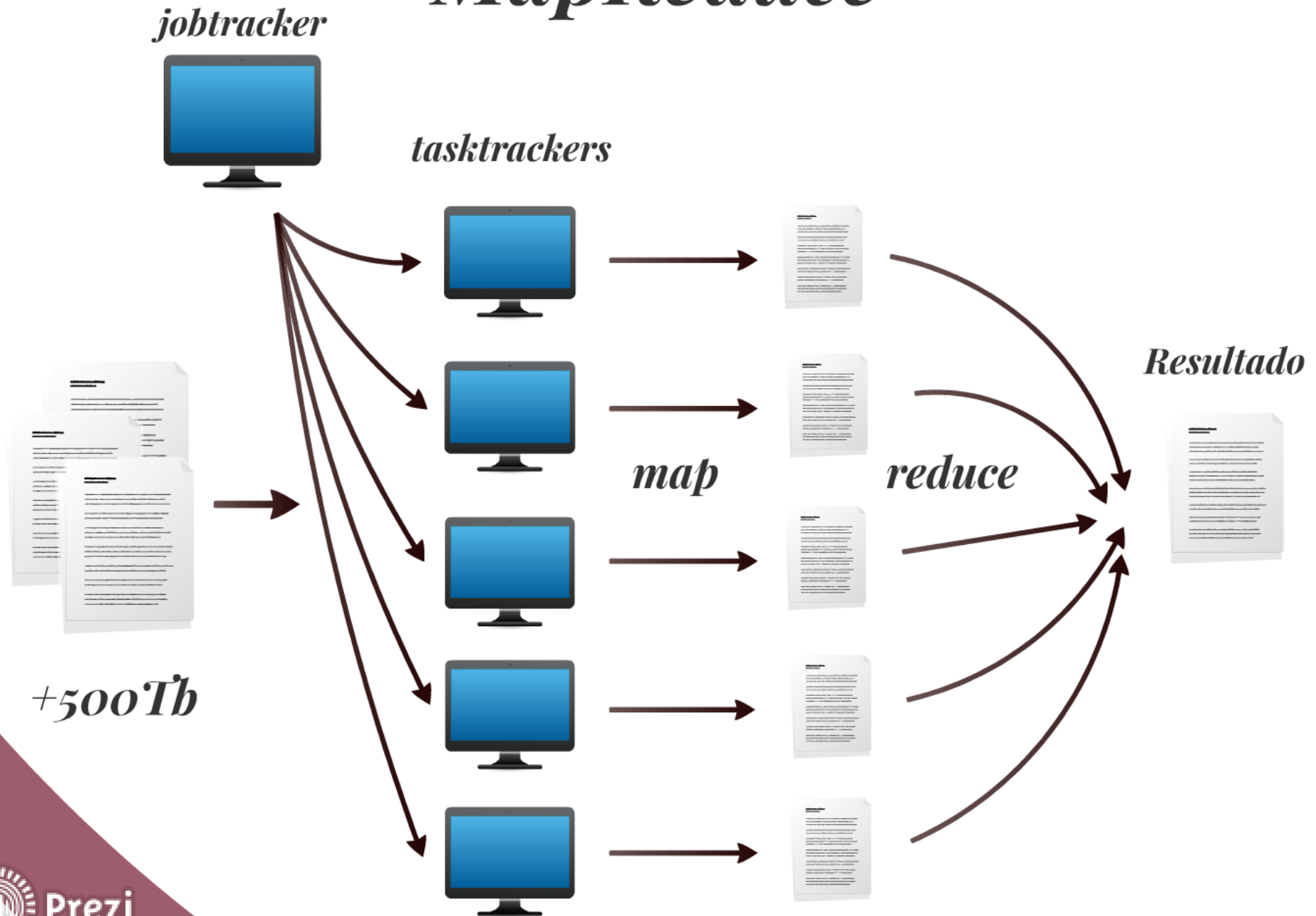
ES Y MapReduce

Google File System



HDFS (Hadoop Distributed File System)

MapReduce



Deduplicación del Registro de padrón

500000 registros

$$(500000 \times 499999) / 2$$

≈

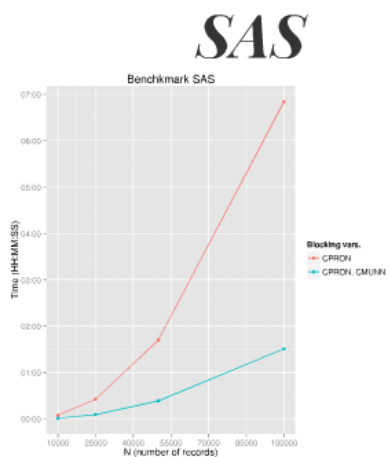
125 000 000 000

comparaciones

$$\times 2 \times 0,00014$$

≈

17 Tb



70 Mb ?

Hadoop

Clúster MapReduce de 100 nodos en el Centro de Cálculo de la ETSII

Hadoop job_201404252110_0010 on bigdata

User: hduser
 Job Name: EstadísticaRegistro
 Job File: http://bigdata:54310/ooziboo/oozie/hadoop/job_201404252110_0010/job.xml
 Submit Host: bigdata
 Submit Host Address: 192.168.101.228
 Job-ACLs: All users are allowed
 Job Setup: Successful
 Status: Running
 Started at: Sat Apr 26 17:24:55 WEST 2014
 Running for: 15m1s, 24sec
 Job Cleanup: Pending

Kind	% Complete	Num Tasks	Pending	Running	Complete	Killed	Failed/Killed Task Attempts
map	51.34%	198	0	178	22	0	0 / 2
reduce	8.53%	1	0	1	0	0	0 / 0

Actualización del Registro de padrón

500000 registros

$$(500000 \times 499999) / 2$$

~ =

125 000 000 000

comparaciones

$$\times 2 \times 0,00014$$

~ =

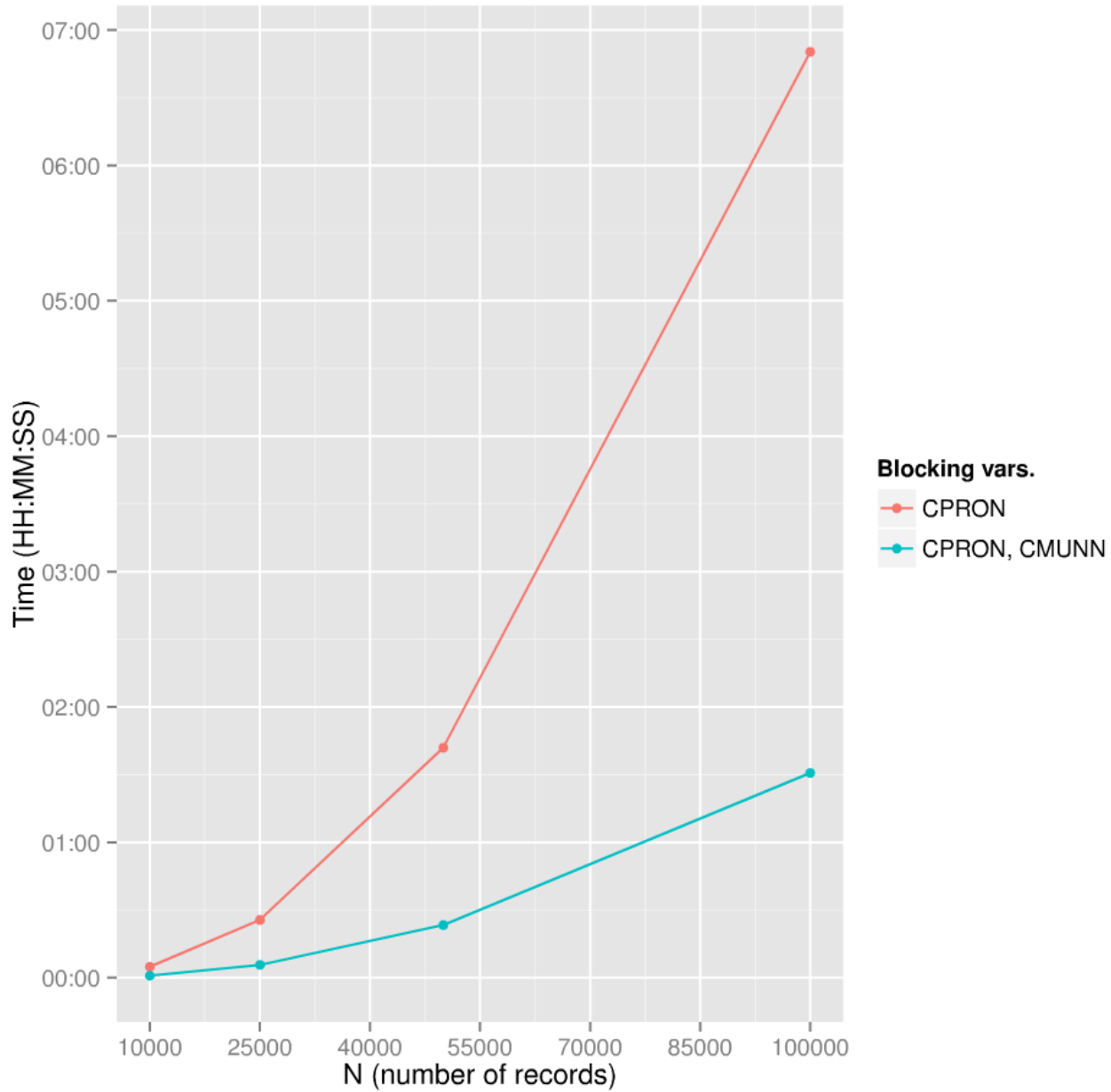
17 Tb



70 Mb ?

SAS

Benchmark SAS



Hadoop

Clúster MapReduce de 100 nodos en el Céntrro de Cálculo de la ETSII

Hadoop job_201404252110_0010 on [bigdata](#)

User: hduser

Job Name: EstadisticasRegistro

Job File: hdfs://bigdata:54310/app/hadoop/tmp/mapred/staging/hduser/.staging/job_201404252110_0010/job.xml

Submit Host: bigdata

Submit Host Address: 192.168.101.228

Job-ACLs: All users are allowed

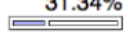
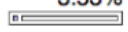
Job Setup: [Successful](#)

Status: Running

Started at: Sat Apr 26 17:24:55 WEST 2014

Running for: 15mins, 24sec

Job Cleanup: Pending

Kind	% Complete	Num Tasks	Pending	Running	Complete	Killed	Failed/Killed Task Attempts
map	31.34% 	198	0	176	22	0	0 / 2
reduce	3.53% 	1	0	1	0	0	0 / 0

