

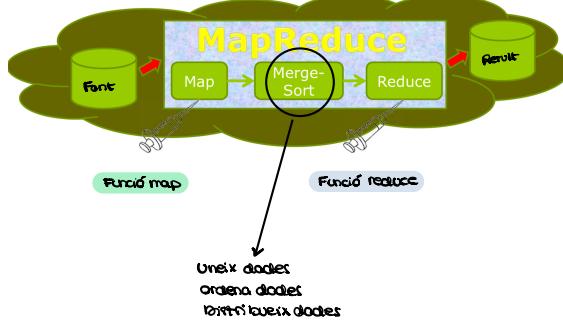
# MapReduce

- Treballa sobre Google File System (ara en més)

## Requirements

- Tu treballas, ho envies al mapReduce, i dades es distribueixen
- En comptes de per data shuffling, envies query shuffling → s'executa en local on són les dades
- Paral·lelisme:
  - Problema anterior → claus parades (per fer validacions, joins, ...) → dades estan lluny
  - Volen treballar amb petabytes.
- Escalable
- Si una màquina falla → sistema es recupera

- Basat en functional programming



- Processen parelles [key, value]
- Signature → el que hem de posar nosaltres

map (key k, value v) → [(i<sub>k1</sub>, v<sub>1</sub>), ..., (i<sub>k<sub>m</sub>(k,v)</sub>, v<sub>m(k,v)</sub>)]

reduce (key i<sub>k</sub>, v<sub>set</sub>) → [o<sub>1</sub>, ..., o<sub>v(i,k,v)</sub>]  
 ↑  
 set  
 de  
 values

## Exemple: word count

The Project Gutenberg Ebook of The Outline of Science, Vol. 1 (of 4), by J. Arthur Thomson

This work is for the use of anyone anywhere at no cost and with almost no restrictions whatsoever. You may copy it, give it away or re-use it under the terms of the Project Gutenberg license included with this eBook or online at www.gutenberg.org

Title: The Outline of Science, Vol. 1 (of 4)  
 A Plain Story Simply Told  
 Author: J. Arthur Thomson  
 Release Date: January 22, 2007 [EBook #20417]  
 Language: English  
 Character set encoding: ASCII

\*\*\* START OF THIS PROJECT GUTENBERG EBOOK OUTLINE OF SCIENCE \*\*\*

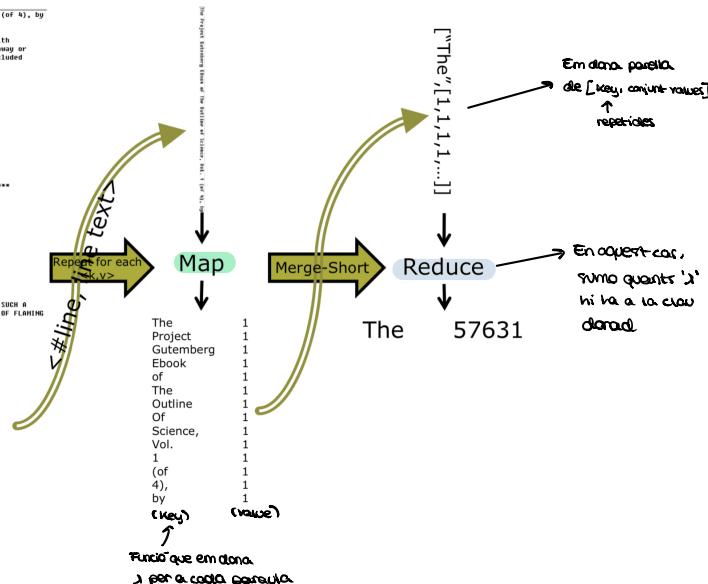
Produced by Brian James, Leonard Johnson and the Online Distributed Proofreading Team at http://www.pgdp.net

Illustration: THE GREAT SCARLET SOLAR PROMINENCES, WHICH ARE SUCH A NOTABLE FEATURE OF THE SOLAR PHENOMENA, ARE IMENSE OUTBURSTS OF FLAMING INCANDESCENT GAS RISING SOMETIMES TO A HEIGHT OF 500,000 MILES!

THE  
 OUTLINE OF SCIENCE  
 A PLAIN STORY SIMPLY TOLD

EDITED BY  
 J. ARTHUR THOMSON  
 REGIUS PROFESSOR OF NATURAL HISTORY IN THE  
 UNIVERSITY OF ABERDEEN

WITH OVER 800 ILLUSTRATIONS,  
 OF WHICH ABOUT AN AVE IN COLOUR.



```
public void map(LongWritable key, Text value) {
  String line = value.toString();
  StringTokenizer tokenizer = new StringTokenizer(line);
  while (tokenizer.hasMoreTokens()) {
    write(new Text(tokenizer.nextToken()), new IntWritable(1));
  }
}

public void reduce(Text key, Iterable<IntWritable> values) {
  int sum = 0;
  for (IntWritable val : values) {
    sum += val.get();
  }
  write(key, new IntWritable(sum));
}
```

## PROBLEMA.

Com que no sap què hi ha al "whatever", no pot optimitzar res.

```
public void map (key, value) {
  whatever
  key, value
}
```

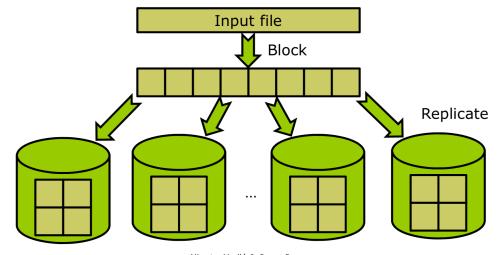
```
public void reduce (key, value) {
  whatever
  key, value
}
```

Verràs 2nd spark - intenta mirar el "whatever" per optimitzar una mica.

## Algoritme

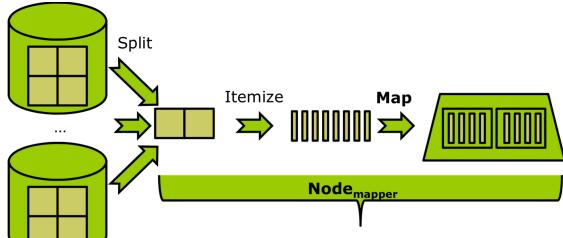
### ① Fixar fitxer a HDFS (cloud)

- HDFS ho divideix en chunks
- Per defecte, ho replica a 3 més.



Alberto Abelló & Oscar Romero

1



### ② A cada un dels workers, se li assigna un split (subset de blocks)

### ③ Worker el divideix en records

### ④ Executa la funció del map i els posa a memòria (spill)

↳ conte resultats  
d'haver executat  
lo més cap "map"

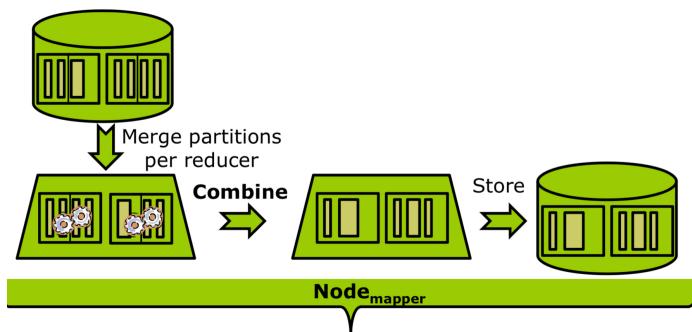
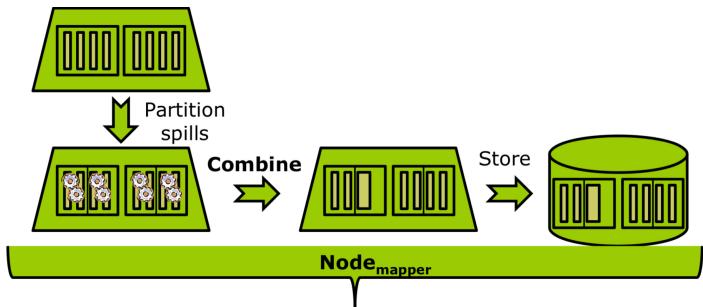
### ⑤ En cada spill, el partitiona amb tants trastos com reducers tenim

### ⑥ Cada partició de "spill" s'ordena de forma independent.

### ⑦ Cada cop que ordenem, s'executa "combinar"

↑  
reduce parcial → pel merge.

### ⑧ Guardem les particions de spill a disc (massive writing)

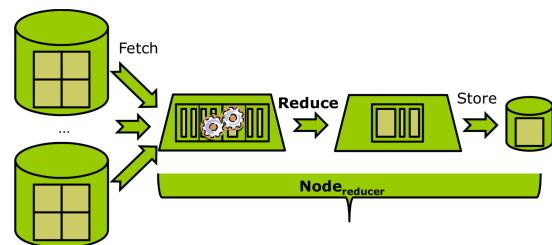
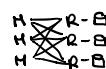


### ⑨ Reducer fa fetch de les particions que hi ha al mapper

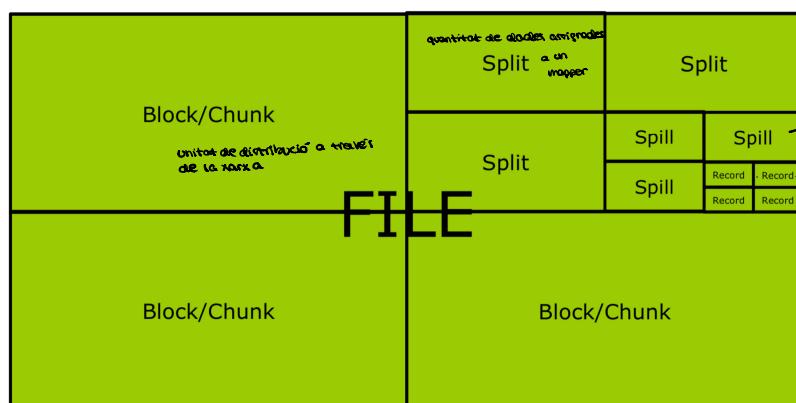
### ⑩ L'output del mapper es torna a fer merge

### ⑪ S'executa el reduce

### ⑫ El guarda en disc un altre cop.



## Objetos



memòria assignada  
a un worker → equivalent a  
un buffer

register : key, value → cicle map()

Record=Key-Value pair

## Combiner

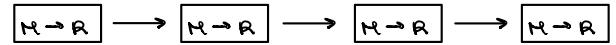
- Permet optimitzar ( $\downarrow$  lectures de disc)
- Redueix trànsit de xarxa
- Important  $\rightarrow$  commutativa  $\rightarrow$  igual que el reducer  
 $\hookrightarrow$  associativa
- $\frac{1 \text{ input}}{1 \text{ output}} \gg \# \text{ CPU}$

## Primeres mapReduce

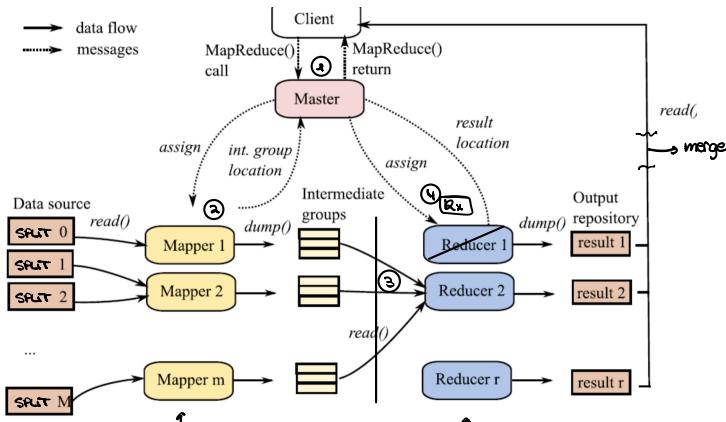
- Massa baix nivell (java)  $\rightarrow$  tant eficient com tu el fagis

## Avantatges mapReduce

- Pots fer el que vulguis
- Tancat: pots fer cicles entre maps-reduce



## Tasques i dataflow



**PROBLEMES**  $\rightarrow$  Reassignments needed (4)  
 $\hookleftarrow$  Single-point failure (1 master)

Sempre que s'escriu, es fa en local

És el "reader" qui va a buscar les dades

} no perdem fitxer mai

## Costos

- (1) Start-up time
- (2) Writing intermediate results  
(necessari per no perdre dades)
- (3) Transferència de dades (xarxa)
- (4) Fer pings a master per dir que estan vius  
(necessari per tenir availability)

## Què fa Spark?

- menys escriptures a disc: romés quan ho veure necessari
- permet més operacions
- et facilita controlar les particions

## Exercici mapreduce

Arremem que tenim el següent codi:

```
public void map(LongWritable key, Text value) {
    String line = value.toString();
    StringTokenizer tokenizer = new StringTokenizer(line);
    while (tokenizer.hasMoreTokens()) {
        write(new Text(tokenizer.nextToken()), new IntWritable(1));
    }
}

public void reduce(Text key, Iterable<IntWritable> values) {
    int sum = 0;
    for (IntWritable val : values) {
        sum += val.get();
    }
    write(key, new IntWritable(sum));
}
```

#split + #màniques virtuales

Fill the gaps:

- 1) Machine0 contains 3 blocks. (block0)
- Machine1 contains 3 blocks. (block1)
- 2) We keep 2 replicas (including the master copy) per block.
- 3) We have 2 splits per machine. (2 split = 1 block)
- 4) Mapper0 reads 2 records.
- Mapper1 reads 2 records.

- 5) Spills in Machine0:

Spill1	Spill2	Spill3	Spill4
[a,1][b,1]	[c,1][d,1]	[a,1][e,1]	[ , ][ , ]
[b,1][a,1]	[d,1][c,1]	[ , ][ , ]	[ , ][ , ]

Spills in Machine1:

Spill1	Spill2	Spill3	Spill4
[a,1][b,1]	[a,1][b,1]	[c,1][d,1]	[ , ][ , ]
[d,1][a,1]	[b,1][c,1]	[ , ][ , ]	[ , ][ , ]

amico  
work

- 6) Partitions in machine0:

Spill 1		Spill 2		Spill 3	
Partition 0	Partition 1	Partition 0	Partition 1	Partition 0	Partition 1
[b,1][b,1]	[a,1][a,1]	[d,1][ , ]	[c,1][c,1]	[ , ][ , ]	[a,1][e,1]
[ , ][ , ]	[ , ][ , ]	[ , ][ , ]	[c,1][ , ]	[ , ][ , ]	[ , ][ , ]

Partitions in Machine1:

Spill 1		Spill 2		Spill 3	
Partition 0	Partition 1	Partition 0	Partition 1	Partition 0	Partition 1
[b,1][d,1]	[a,1][ , ]	[b,1][b,1]	[a,1][c,1]	[f,1][ , ]	[c,1][ , ]
[a,1][ , ]	[ , ][ , ]	[ , ][ , ]	[ , ][ , ]	[ , ][ , ]	[ , ][ , ]

ta Set:

Da. Block0: "a b b a c | c d c a e"  
 Da. Block1: "a b d d a | b b c c p"

configuration

- combine = reduce
  - 1 split = 1 block
  - '1' divideix els records del block més 2 records/block
  - dfs.replication = 1 (hadoop)
  - 1 spill = 4 [key,value]
  - 2 mappers, 2 reducers
- |                     |                     |
|---------------------|---------------------|
| ↳ machine0 → block0 | , machine1 → block1 |
| ↳ mapper0           | ↳ mapper1           |
| ↳ reducer0          | ↳ reducer1          |
- hash for reducers → {b,d,f} → 0  
 ↳ {a,c,e} → 1

7) Partitions in machine0:

Spill 1		Spill 2		Spill 3	
Partition 0	Partition 1	Partition 0	Partition 1	Partition 0	Partition 1
[b,2][ , ]	[a,2][ , ]	[d,1][ , ]	[e,2][ , ]	[ , ][ , ]	[a,1][e,1][ , ]
[ , ][ , ]	[ , ][ , ]	[ , ][ , ]	[ , ][ , ]	[ , ][ , ]	[ , ][ , ]

← tenir  
combinar  
que suma)

Partitions in Machine1:

Spill 1		Spill 2		Spill 3	
Partition 0	Partition 1	Partition 0	Partition 1	Partition 0	Partition 1
[b,1][d,2]	[a,1][ , ]	[b,2][ , ]	[a,1][e,1]	[g,1][ , ]	[c,1][ , ]
[ , ][ , ]	[ , ][ , ]	[ , ][ , ]	[ , ][ , ]	[ , ][ , ]	[ , ][ , ]

+ ordenar

8) Files in machine0:

File0	File1	File2
[b,2][ , ][ , ][ , ]	[a,2][ , ][ , ][ , ]	[d,1][ , ][ , ][ , ]
File3	File4	File5
[c,3][ , ][ , ][ , ]	[ , ][ , ][ , ][ , ]	[a,1][e,1][ , ][ , ]

Files in Machine1:

File0	File1	File2
[b,2][d,2][ , ][ , ]	[a,1][ , ][ , ][ , ]	[b,2][ , ][ , ][ , ]
File3	File4	File 5
[a,1][c,1][ , ][ , ]	[g,2][ , ][ , ][ , ]	[c,1][ , ][ , ][ , ]

9) Merges in machine0:

Merge0	Merge1
[b,2][d,1][ , ][ , ][ , ]	[a,2][c,3][e,1][ , ][ , ][ , ]

Merges in Machine1:

Merge0	Merge1
[b,3][d,2][g,2][ , ][ , ][ , ]	[a,2][c,2][ , ][ , ][ , ][ , ]

10) Files in machine0:

File0	File1
[b,2][d,1][ , ][ , ][ , ]	[a,3][c,3][e,1][ , ][ , ][ , ]

Files in Machine1:

File0	File1
[b,3][d,2][g,1][ , ][ , ][ , ]	[a,2][c,2][ , ][ , ][ , ][ , ]

11) Reducer0 reads ..... File 0..... Files from machine0

and ..... File 0..... Files from machine1. (answer which Files)

Reducer1 reads ..... File 1..... Files from machine0  
and ..... File 1..... Files from machine1. (answer which Files)

12) Merge in machine0:

File0	File1
[b,12,34][d,11,24][g,3,24]	[ , ][ , ][ , ][ , ][ , ]

Merge in Machine1:

File0	File1
[a,13,45][e,43,27][c,11,7]	[ , ][ , ][ , ][ , ][ , ]

13) Reduce function is executed 3.... times in machine0.

Reduce function is executed ...3.. times in machine1.

14) Files in machine0:

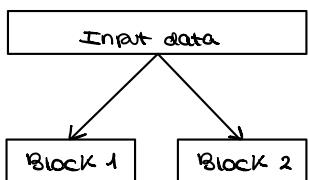
File0	File1
[b,5][d,3][g,2][ , ][ , ][ , ]	[ , ][ , ][ , ][ , ][ , ][ , ]

Files in Machine1:

File0	File1
[a,4][c,5][e,4][ , ][ , ][ , ]	[ , ][ , ][ , ][ , ][ , ][ , ]

Algorithm per se.

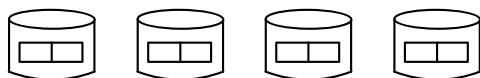
- ① Input data partitioned into blocks



Partition les données en cloud par blocks → blocks  
Block 1, Block 2

- ② Replicate them into different nodes

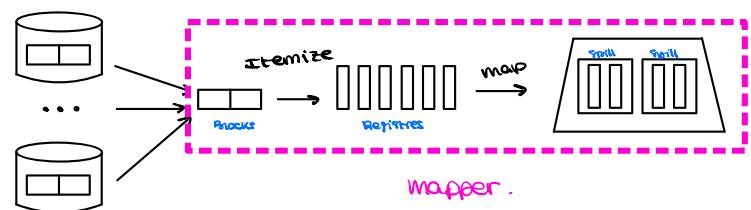
& replica / block



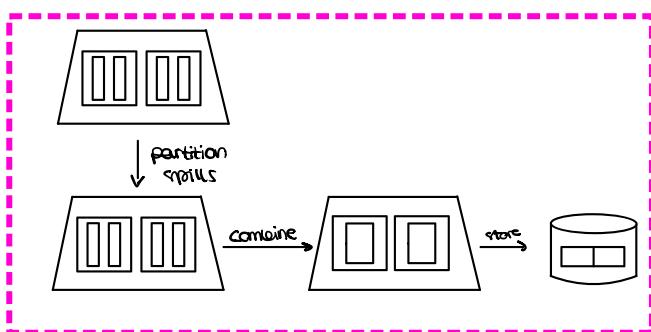
- ③ Each map function reads a subset of blocks (split)

& mapper and block = 1 split

- ④ Divides into records. and 2 records / block (split)



- ⑤ Executes the map for each record and 4 records and 4 maps → 4 map segments map



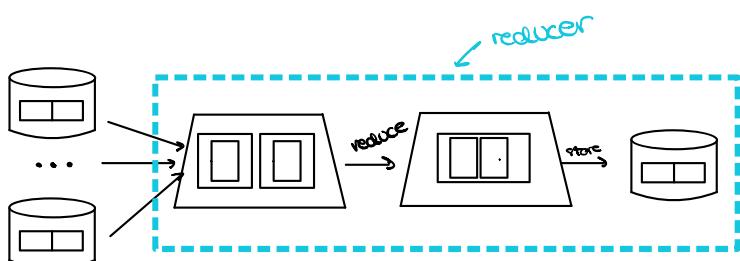
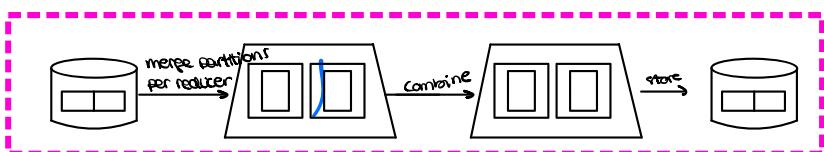
- ⑥ Each spill is partitioned per reducer and 2 reducer and 2 spills

- ⑦ Each spill is stored independently (+combine)

- ⑧ Store spill partitions into disk

- ⑨ Spill partitions are merged and stored independently (+combine)

- ⑩ store into disk.



- ⑪ Reducers fetch data from mappers

- ⑫ Mappers output is merged and stored

- ⑬ Reduce function is executed per key

- ⑭ Store the result into disk.