

CC5212-1

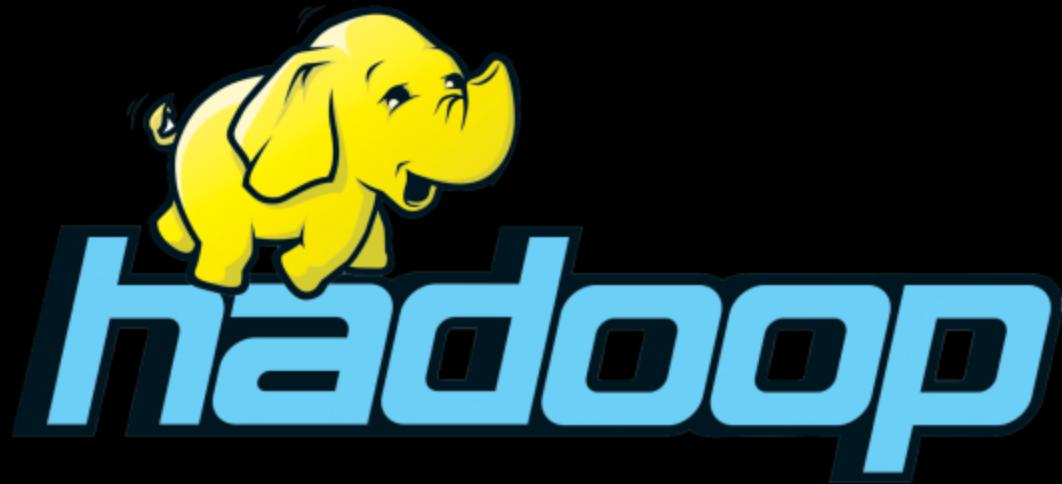
PROCESAMIENTO MASIVO DE DATOS
OTOÑO 2017

Lecture 4: MapReduce/Hadoop II

Aidan Hogan
aidhog@gmail.com

HADOOP: IN MORE DEPTH

Why “Hadoop” ...



... and why an elephant and not, e.g., a rabbit?

Why “Hadoop” ...



... and why an elephant and not, e.g., a rabbit?



... better luck next time, Tiddles.

HADOOP: UNDER THE HOOD

MapReduce Abstraction

Can we abstract any general framework?



Define **input** as a set of key value pairs $I \subseteq T_{IK} \times T_{IV}$



For example, $I = \{(1, "soy una linea"), (2, "soy otra linea")\}$

T_{IK} is the set of all **int**, T_{IV} is the set of all **string**

Define **map** as a function $I \rightarrow 2^M$ where $M \subseteq T_{MK} \times T_{MV}$

For example, $\text{map}(1, "soy una linea") := \{("soy", 1), ("una", 1), ("linea", 1)\}$

T_{MK} is the set of all **string**, T_{MV} is the set of all **int**

Define **reduce** as a function $T_{MK} \times 2^{T_{MV}} \rightarrow 2^R$ where $R \subseteq T_{RK} \times T_{RV}$

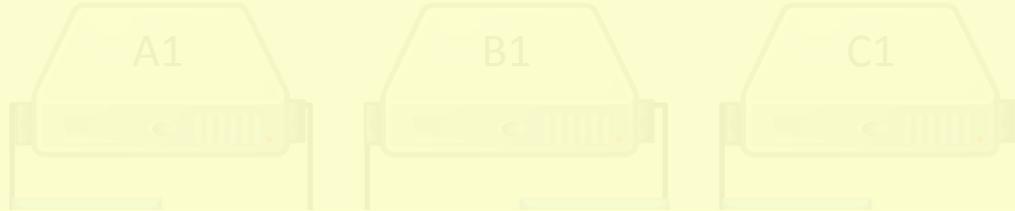
For example, $\text{reduce}("soy", \{1, 1\}) := \{("soy", 2)\}$

T_{RK} is the set of all **string**, T_{RV} is the set of all **int**

In general, we must assume bags/multisets (sets with duplicates)



MapReduce: Main Idea



Input

File on Distr. File System

Define input as a set of key value pairs $I \subseteq T_{MK} \times T_{MV}$

Given I T_{MK} is the set of all int, T_{MV} is the set of all string

... compute map over all $i \in I$

... group resulting set by map key

... apply reduce over groups

Define reduce as a function $T_{MK} \times T_{MV} \rightarrow T_{MV}$

For example, reduce["soy", {1, 1}]

T_{MK} is the set of all string,

← But how to implement this part in a distributed system

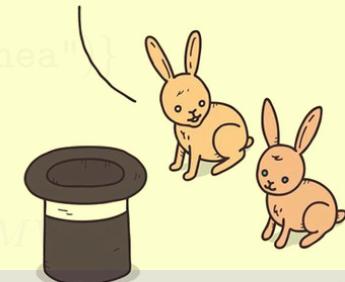
1. Partition by map key

2. Sort (in parallel) by map key

3. Apply reduce / Profit

... let's be more specific

THIS IS WHERE
THE MAGIC HAPPENS



MapReduce/Hadoop

1. Input

2. Map

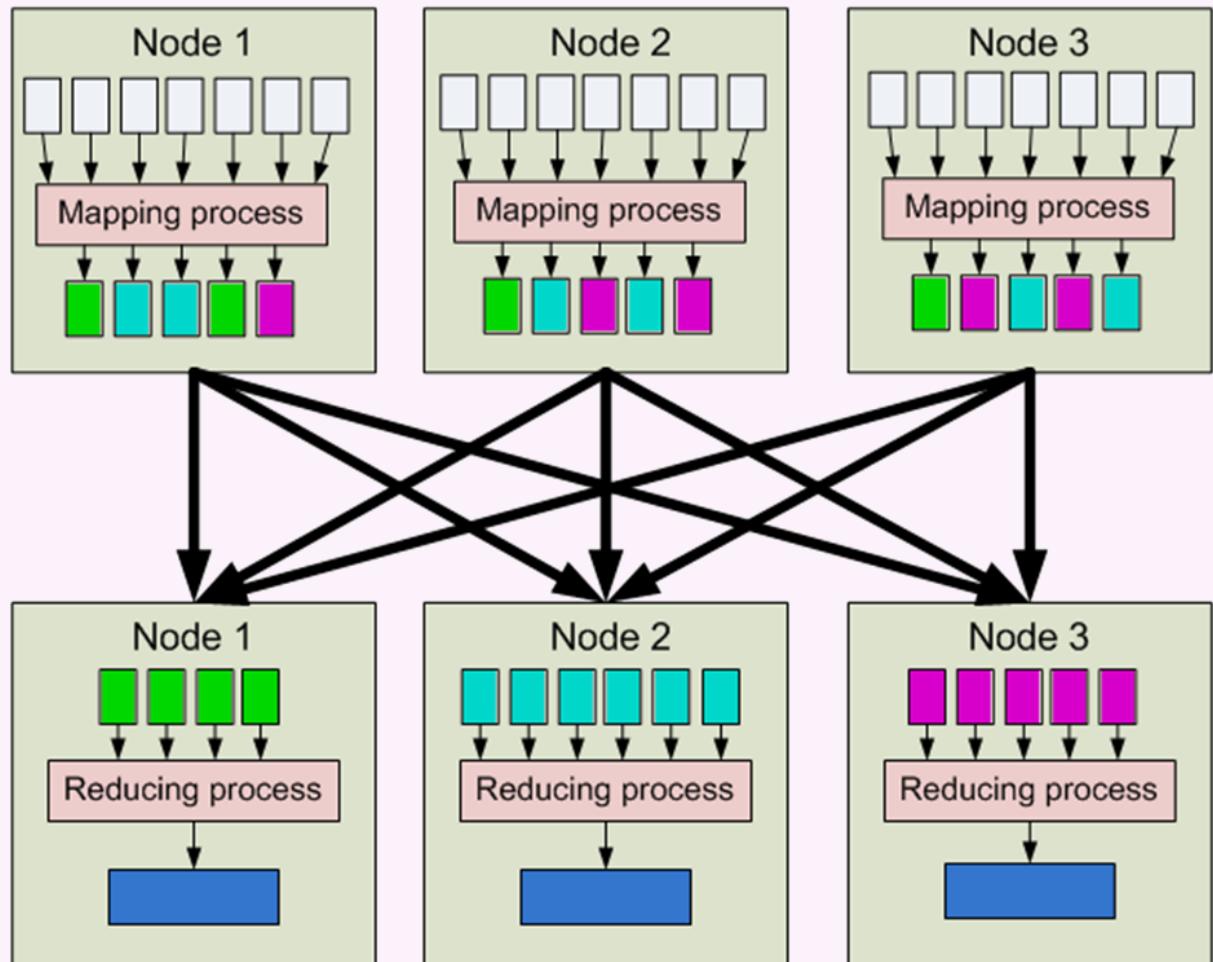
3. Partition [Sort]

4. Shuffle

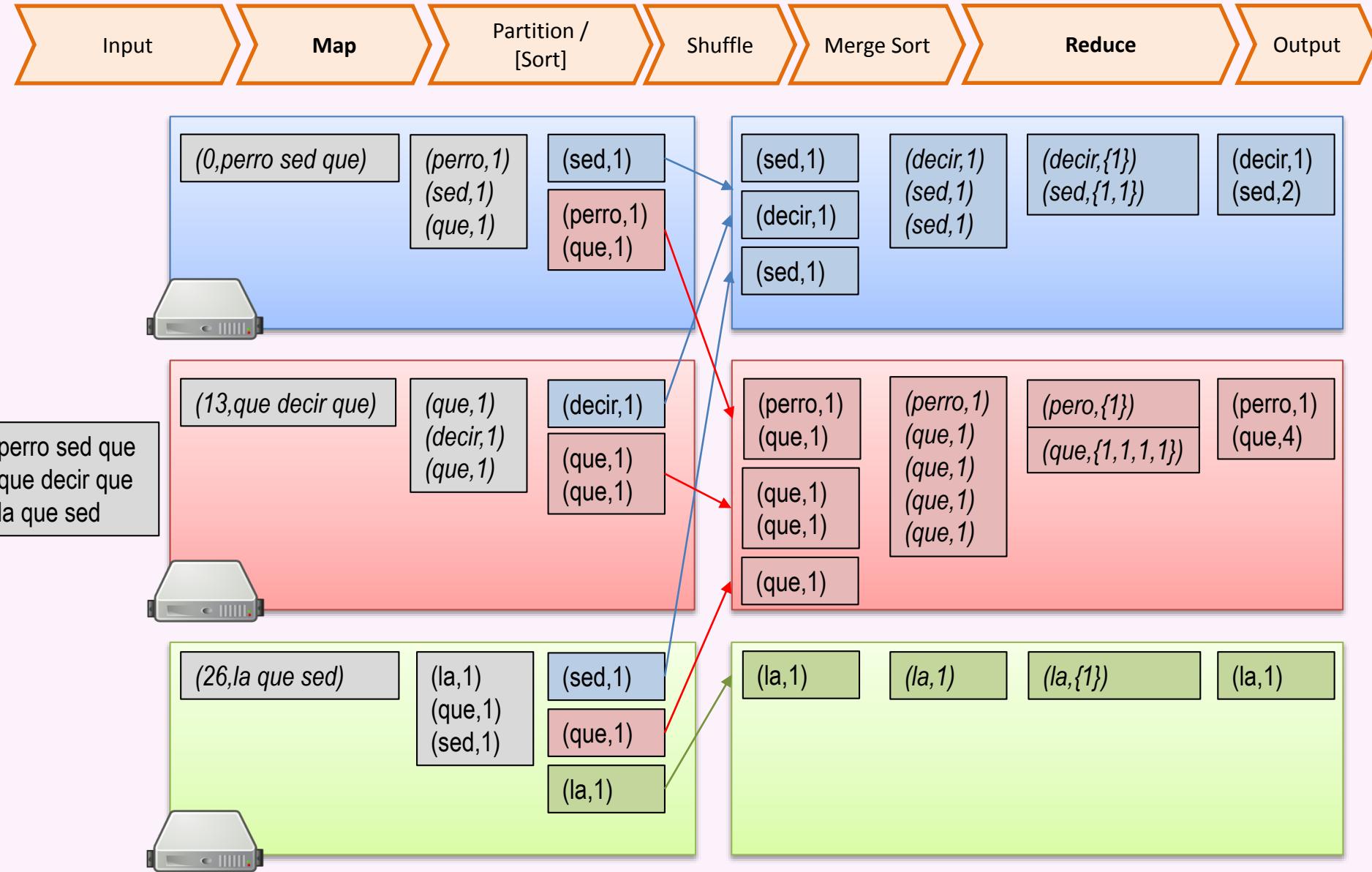
5. Merge Sort

6. Reduce

7. Output



MapReduce/Hadoop: Counting Words



MapReduce/Hadoop: Combiner

1. Input

2. Map

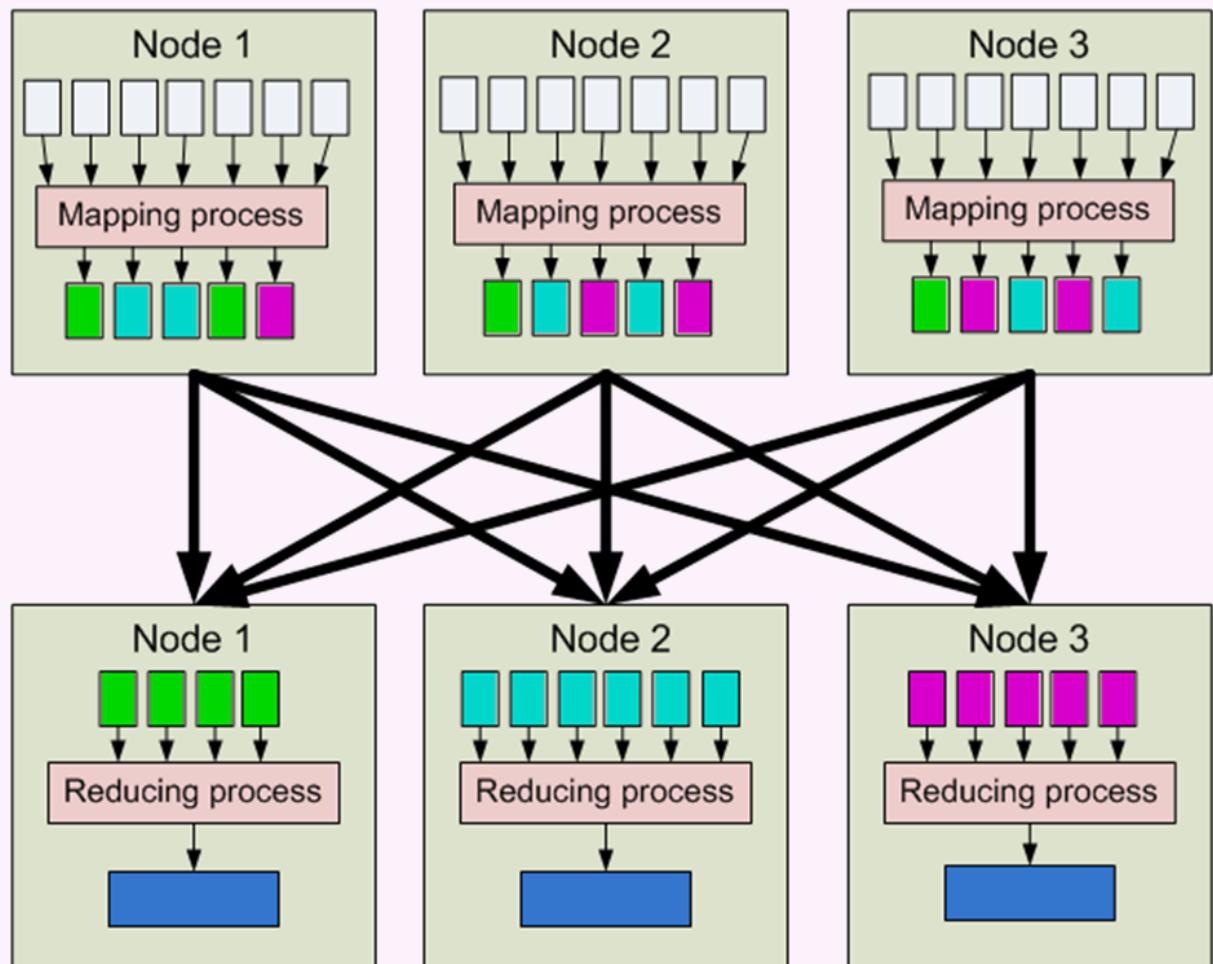
3. Partition [Sort]
("Combine")

4. Shuffle

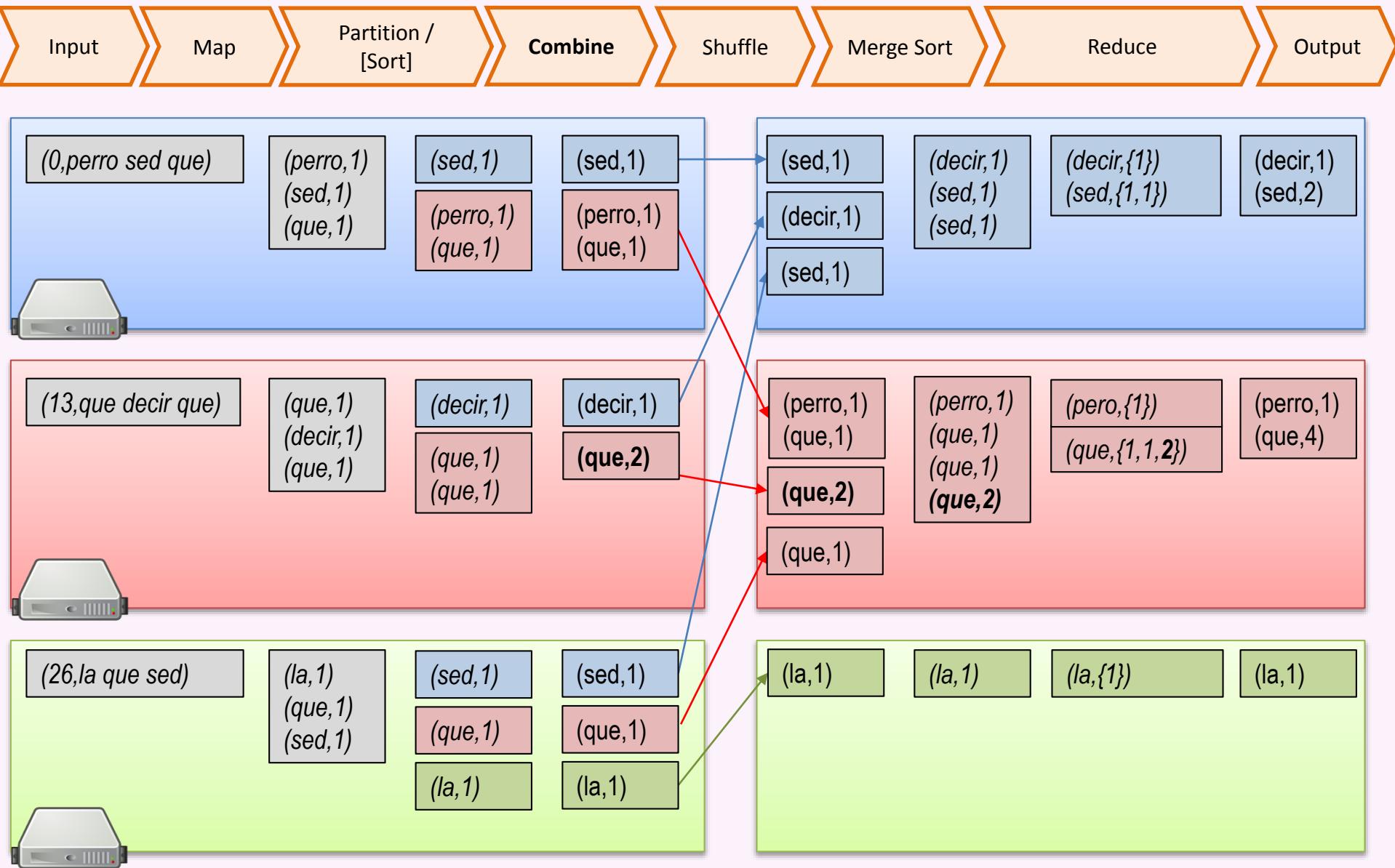
5. Merge Sort

6. Reduce

7. Output



MapReduce/Hadoop: Combiner



MapReduce/Hadoop: Combiner

1. Input

2. Map

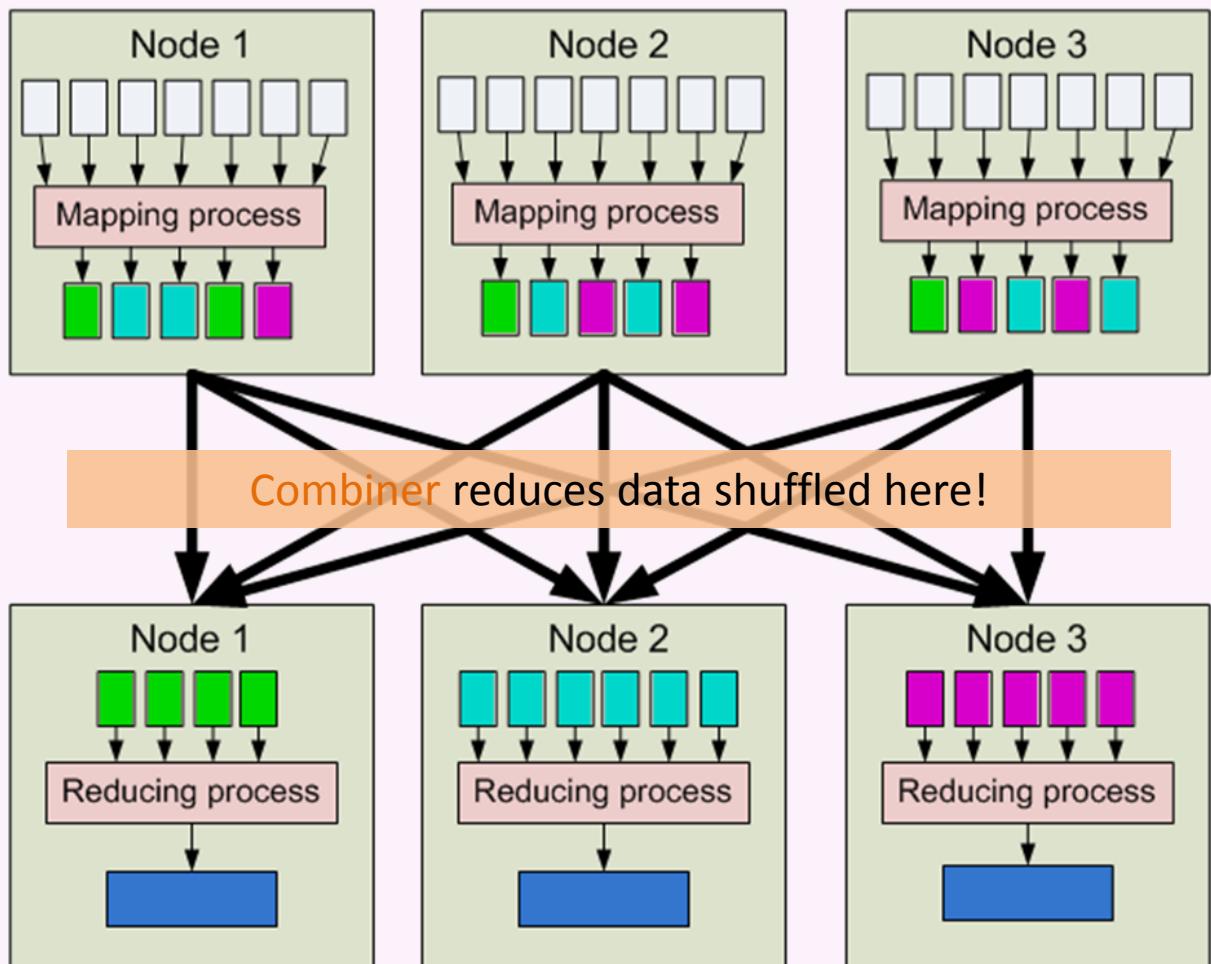
3. Partition [Sort]
("Combine")

4. Shuffle

5. Merge Sort

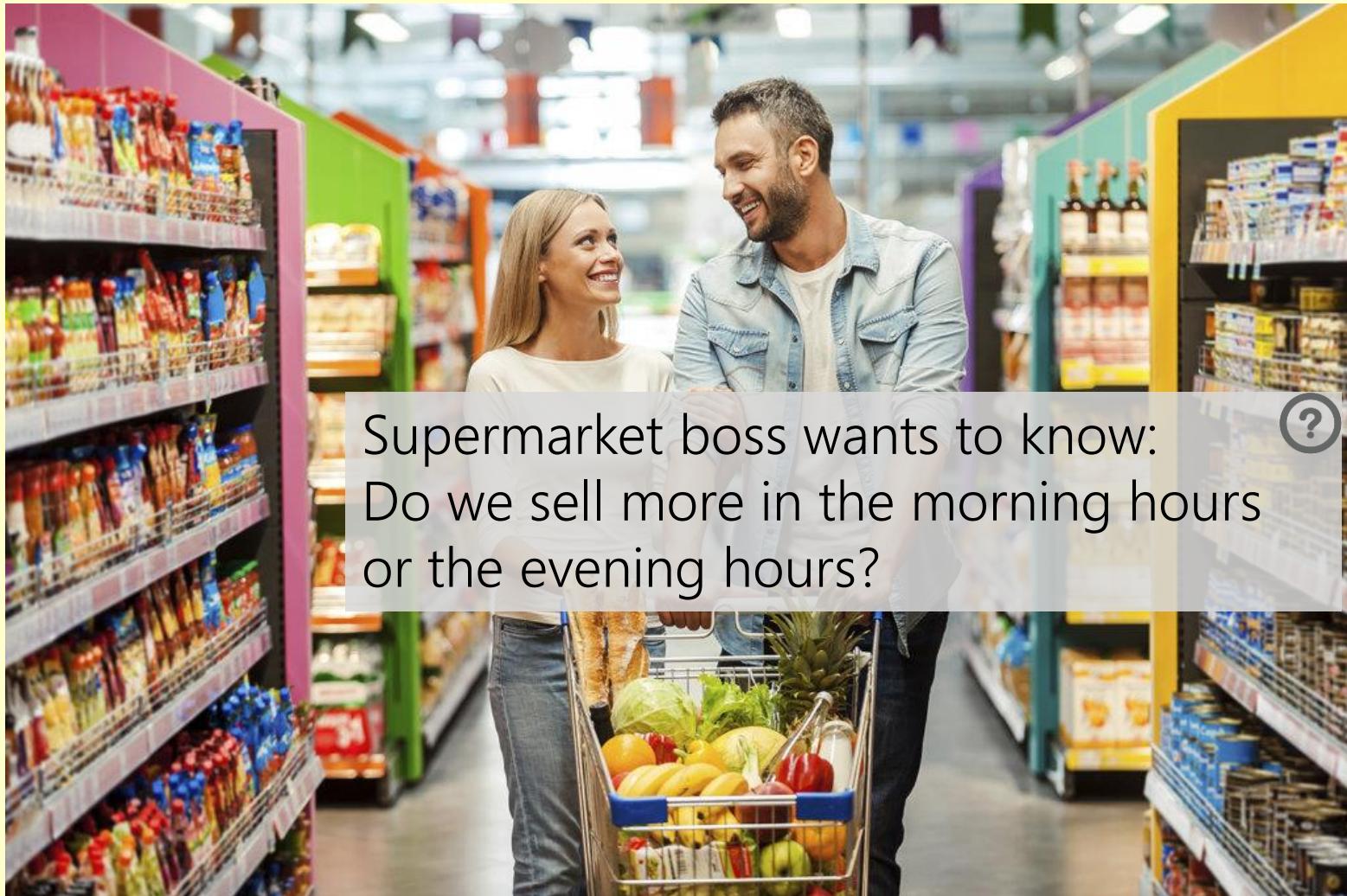
6. Reduce

7. Output



IN-CLASS EXERCISES

Supermarket Example



Supermarket boss wants to know:
Do we sell more in the morning hours
or the evening hours?

MapReduce: Supermarket Example

ReceiptItems	
RECEIPT ID	ITEM ID
R1401	I306
R1401	I306
R1401	I504
R1402	I007
R1402	I306
R1403	I306
R1403	I504
...	...

ReceiptTimes	
RECEIPT ID	TIME
R1403	19:00
R1401	18:59
R1402	19:01
...	...

ItemDetails		
ITEM ID	NAME	PRICE (\$)
I306	Zanahoria 500g	500
I504	CocaCola 3L	1400
I007	Comfort	1200
...

Compute total sales per hour of the day?



Output	
HOUR	TOTAL
...	...
18:00–18:59	\$2400
19:00–19:59	\$3600
...	...

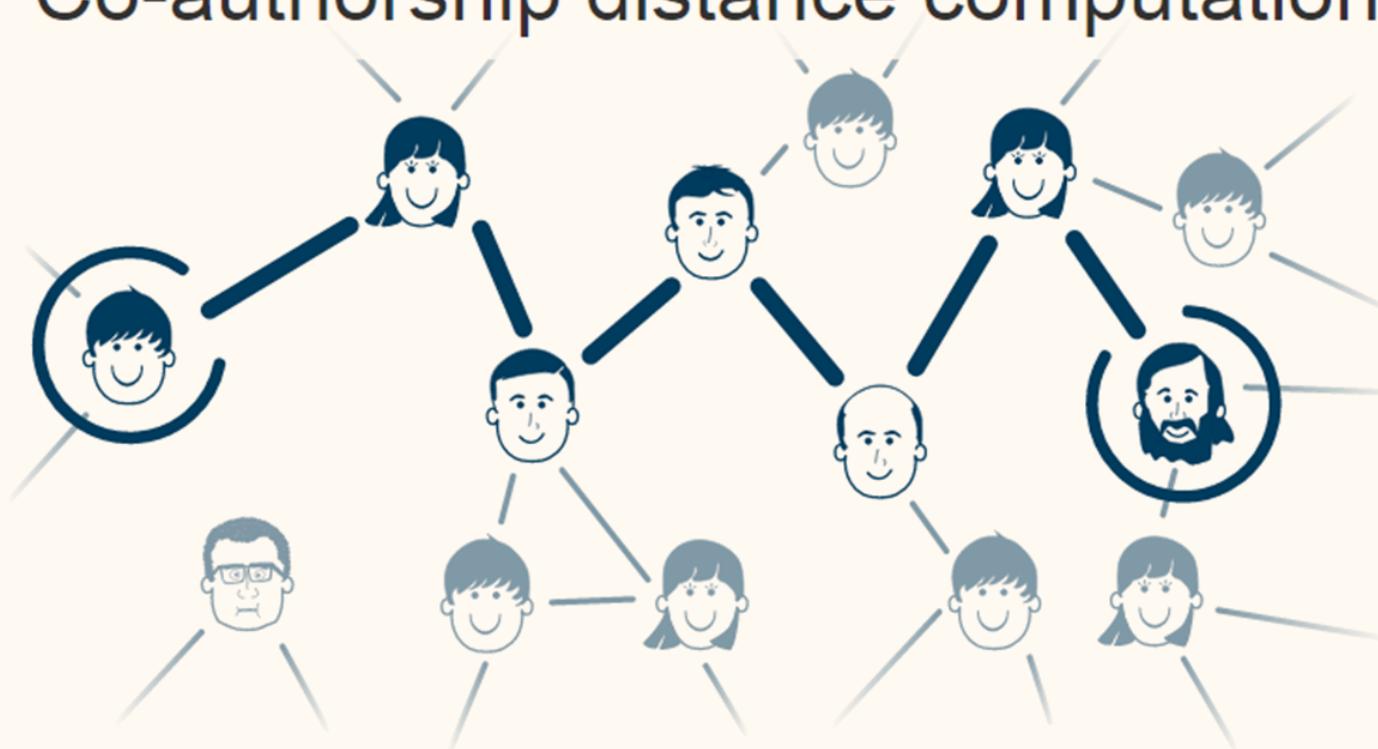
Erdős-number Example



Everybody wants to know:
What's my Erdős-number?



Co-authorship distance computation



Find the path between two authors:

Paul Erdős

Erdős-number Example



Input: Authors.tsv

AUTHOR	PAPER
Aidan Hogan	IMGpedia: Enriching the Web of Data with Image Content Analysis
Benjamin Bustos	IMGpedia: Enriching the Web of Data with Image Content Analysis
Benjamin Bustos	Scalability of Non-Rigid 3D Shape Retrieval
H. Li	A new class of Ramsey-Turán problems
H. Li	Scalability of Non-Rigid 3D Shape Retrieval
Paul Erdős	Random induced graphs
Richard H. Schelp	A new class of Ramsey-Turán problems
Richard H. Schelp	Random induced graphs
Sebastian Ferrada	IMGpedia: Enriching the Web of Data with Image Content Analysis
...	...

Who has an Erdős-number of 3?



Output

AUTHOR

Sebastián Ferrada
Aidan Hogan

...



Questions?