

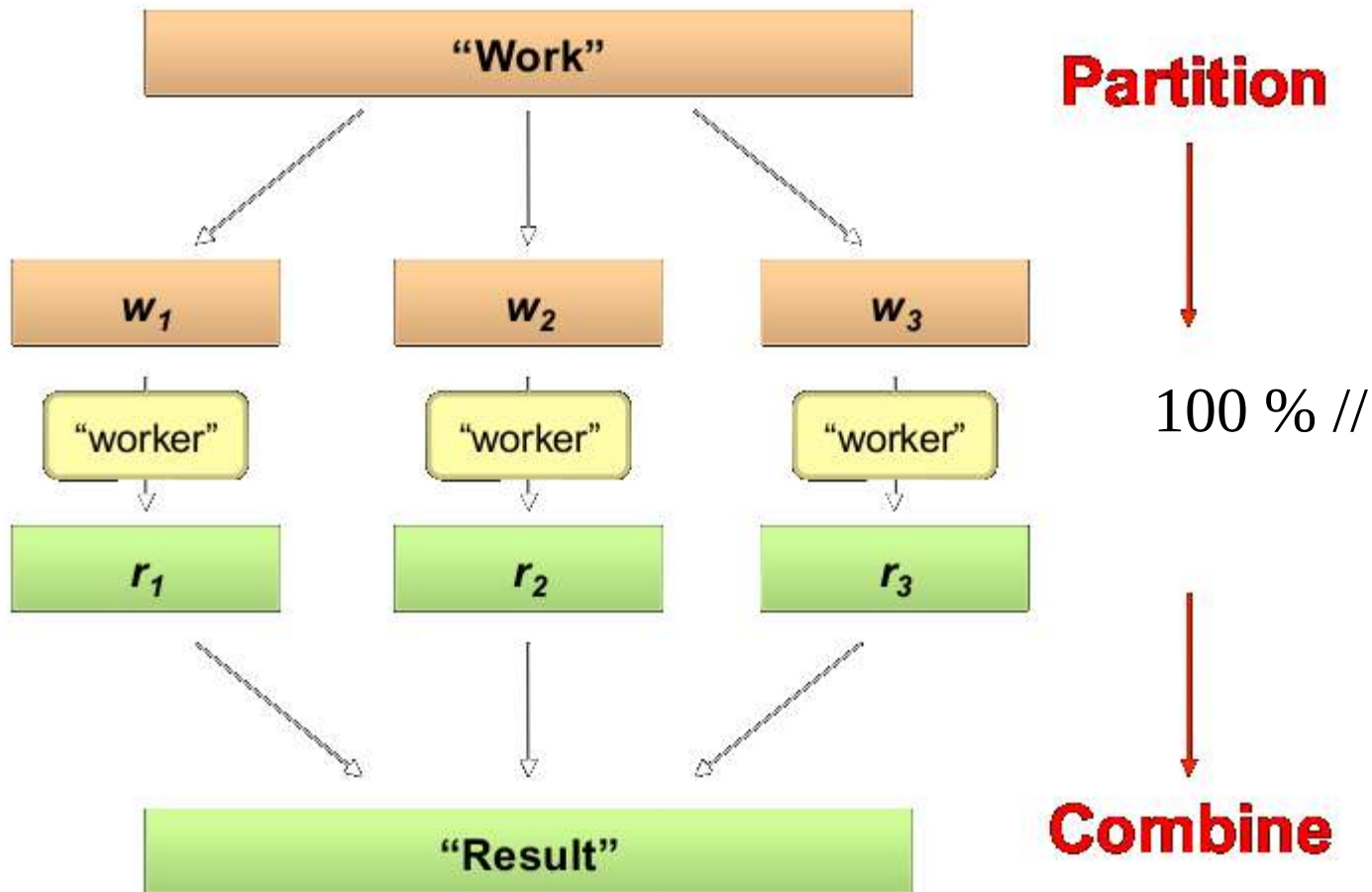


# Map Reduce

dacosta@irit.fr



## Divide and conquer at PaaS



## Typical problem

- Click to edit the outline text format
  - Iterate over a large number of records
    - Second Outline Level
  - Extract something of interest from each **MAP**
    - Third Outline Level
  - Shuffle and sort intermediate results
    - Fourth Outline Level
  - Aggregate intermediate results **Reduce**
    - Fifth Outline Level
  - Generate final output
    - Sixth Outline Level
  - Seventh Outline Level
- Key idea: functional abstraction for these two operations
- Deuxième niveau
  - Troisième niveau

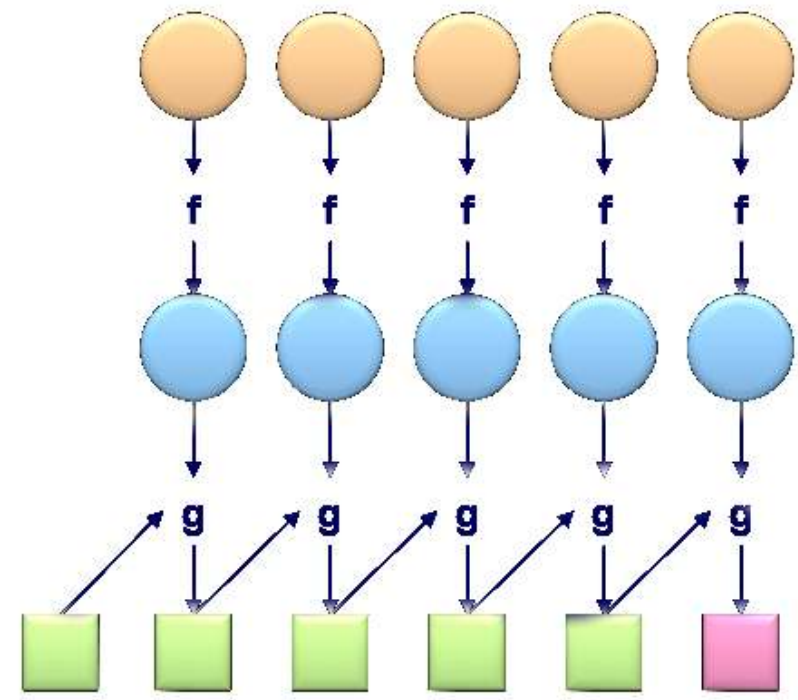
Folding

**Map**

**Map**

**Fold**

**Reduce**



– Deuxième niveau

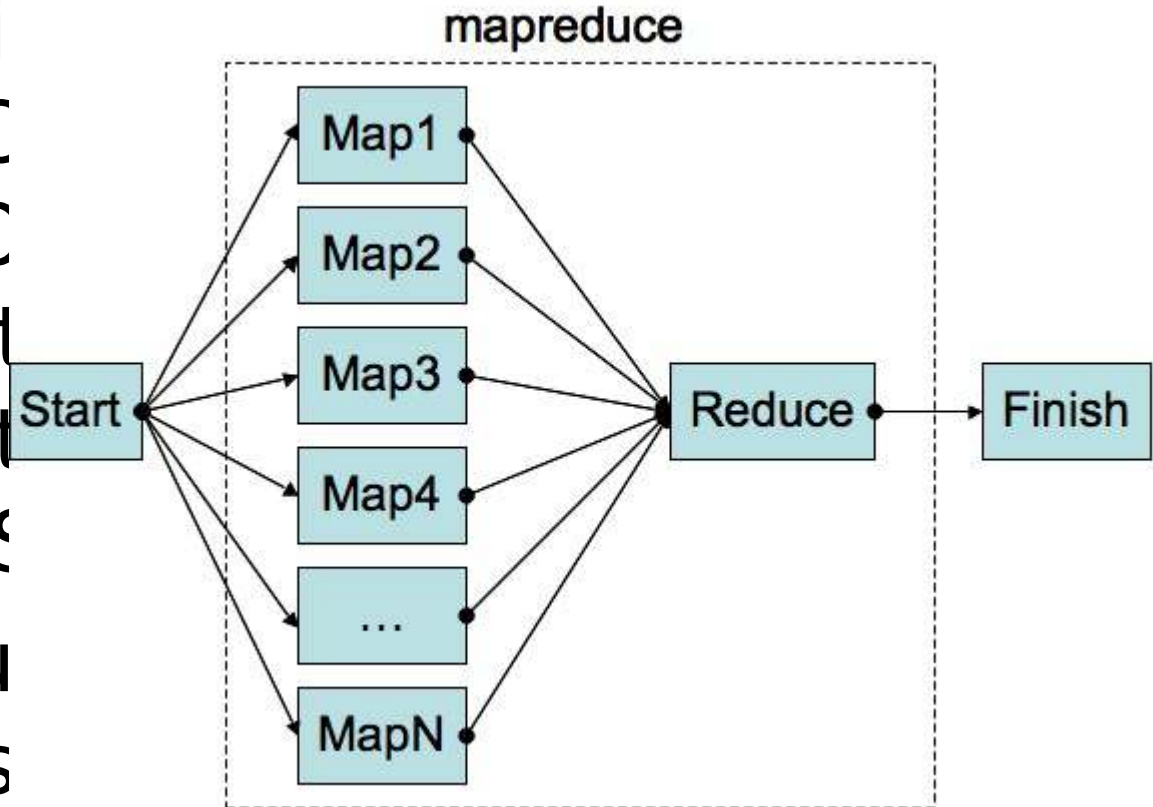
• Troisième niveau

Difficulties ?

- Huge amount of data
    - Click to edit the outline text format
    - Do not fit into memory
    - Access patterns are broad
    - Most data not accessed frequently
  - Complex data
    - links between data or the content
  - Same data can be treated in different ways
    - No pre-processing
- Exemple : crawling through internet data
- Deuxième niveau
  - Troisième niveau

Principle

- "Map" step: The master node takes the input, divides it into smaller sub-problems, and distributes them to worker nodes
- "Reduce" step: The master node then collects the answers to all the sub-problems and combines them in some way to form the output



– Deuxième niveau

• Troisième niveau

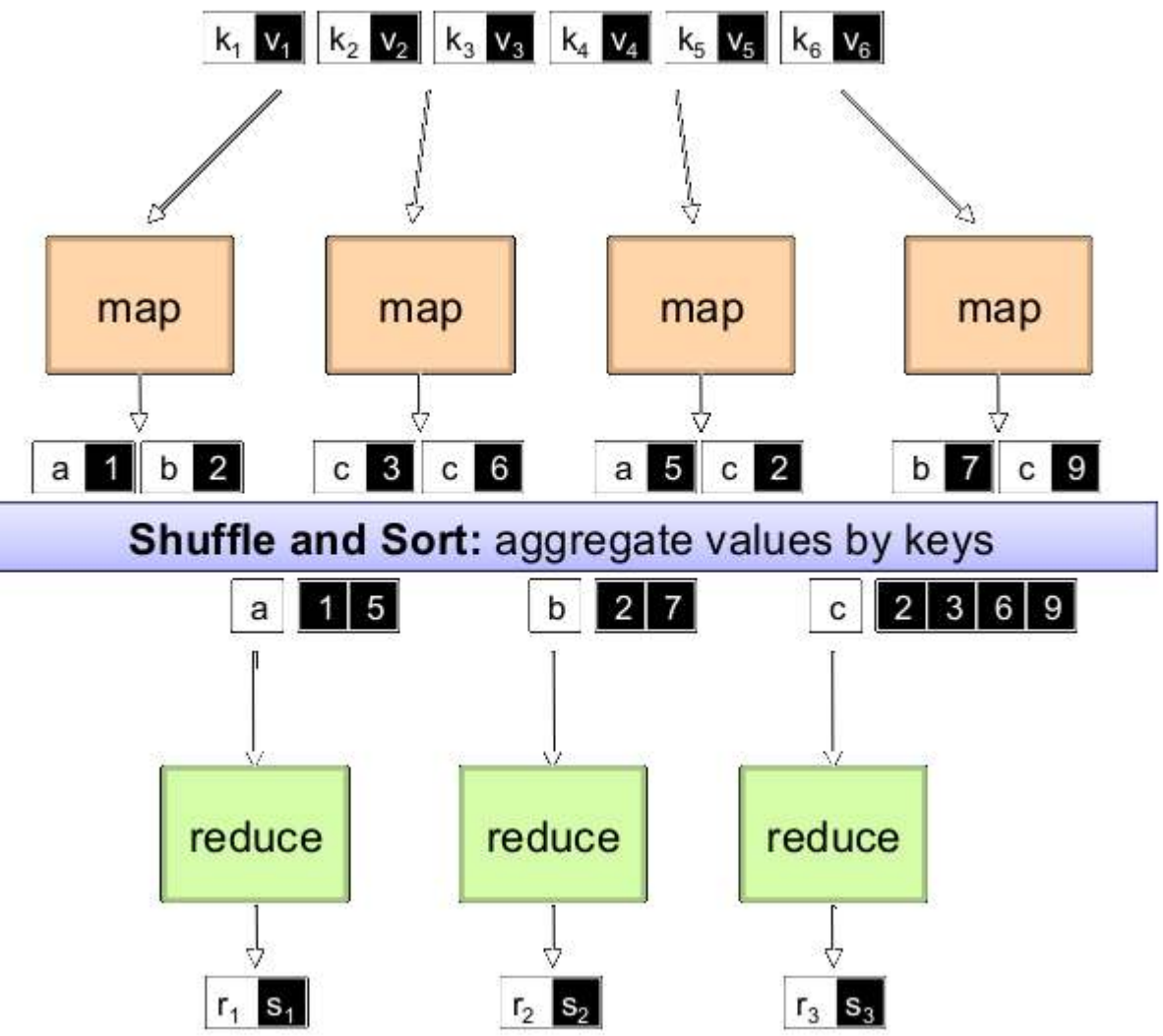
## MapReduce

- Click to edit the outline text format
- Programmers specify two functions:
  - $\text{map}(k, v) \rightarrow \langle k', v' \rangle^*$
  - $\text{reduce}(k', v') \rightarrow \langle k', v' \rangle^*$
  - All values with the same key are reduced together
- Usually, programmers also specify:
  - $\text{partition}(k', \text{number of partitions}) \rightarrow \text{partition for } k'$ 
    - Often a simple hash of the key, e.g.  $\text{hash}(k') \bmod n$
    - Allows reduce operations for different keys in parallel
  - $\text{combine}(k', v') \rightarrow \langle k', v' \rangle$
- **MapReduce** (MapReduce) that run in the map phase
  - Optimizes to reduce network traffic & disk writes
- Implementations:
  - Google has a proprietary implementation in C++
  - Hadoop is an open source implementation in Java

masque

- Deuxième niveau

• Troisième niveau



- (
- S
- r
- r

– Deuxième niveau

• Troisième niveau



## Word count

```
function map(String name, String document):  
    // name: document name  
    // document: document contents  
    for each word w in document:  
        emit (w, 1)  
  
function reduce(String word, Iterator partialCounts):  
    // word: a word  
    // partialCounts: a list of aggregated partial counts  
    sum = 0  
    for each pc in partialCounts:  
        sum += ParseInt(pc)  
    emit (word, sum)
```

modifier les styles du texte du masque

– Deuxième niveau

• Troisième niveau

## Exemple : Average number of contract by Age

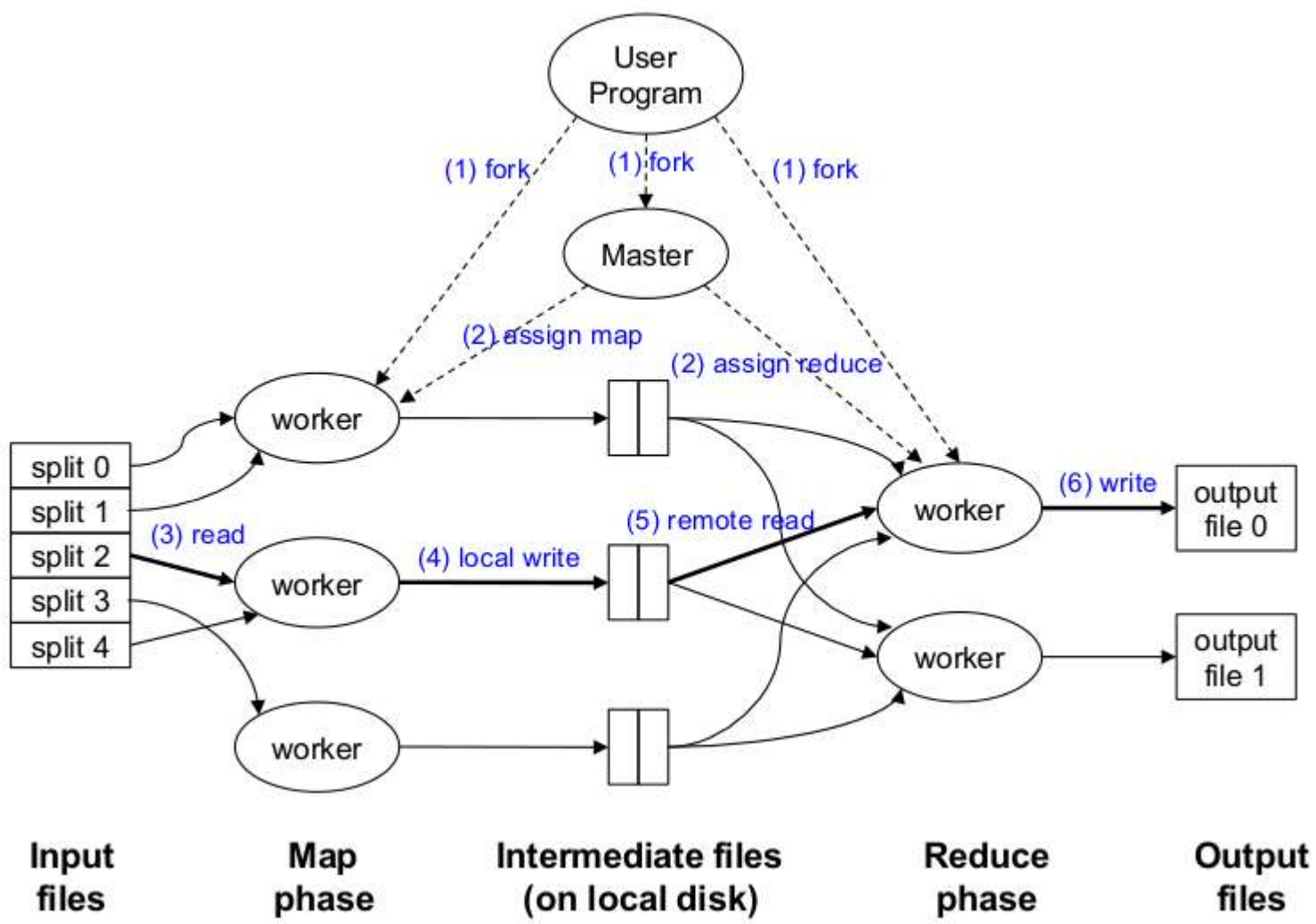
- Click to edit
- For 1 million entries
  - Batch of 1000
  - 1100 of them
- Output of Map
  - Range 8-110
- Reduce :
  - Batch of 1 Y
  - 102 of them
  - Treat 1000
- Seventh Ou
- masque
  - 102 values
  - Deuxième niveau
  - Troisième niveau

```
function Map is
  input: integer K1 between 1 and 1100
  for each social.person record in the K1 batch do
    let Y be the person's age
    let N be the number of contacts the person has
    produce one output record <Y,N>
  repeat
end function

function Reduce is
  input: age (in years) Y
  for each input record <Y,N> do
    Accumulate in S the sum of N
    Accumulate in C the count of records so far
  repeat
  let A be S/C
  produce one output record <Y,A>
end function
```

## MapReduce Runtime

- Click to edit the outline text format
- Handles scheduling
  - Assigns workers to map and reduce tasks
- Handles "data distribution"
  - Moves the process to the data
- Handles synchronization
  - Gathers, sorts, and shuffles intermediate data
- Handles faults
  - Detects worker failures and restarts
- Everything happens on top of a distributed FS (later)
  - Deuxième niveau
  - Troisième niveau



– Deuxième niveau

• Troisième niveau

How do we get data to the workers

- 



NAS



Compute nodes

t  
Classical cluster  
vision

- 



SAN

What's the problem here ?  
– Deuxième niveau

• Troisième niveau

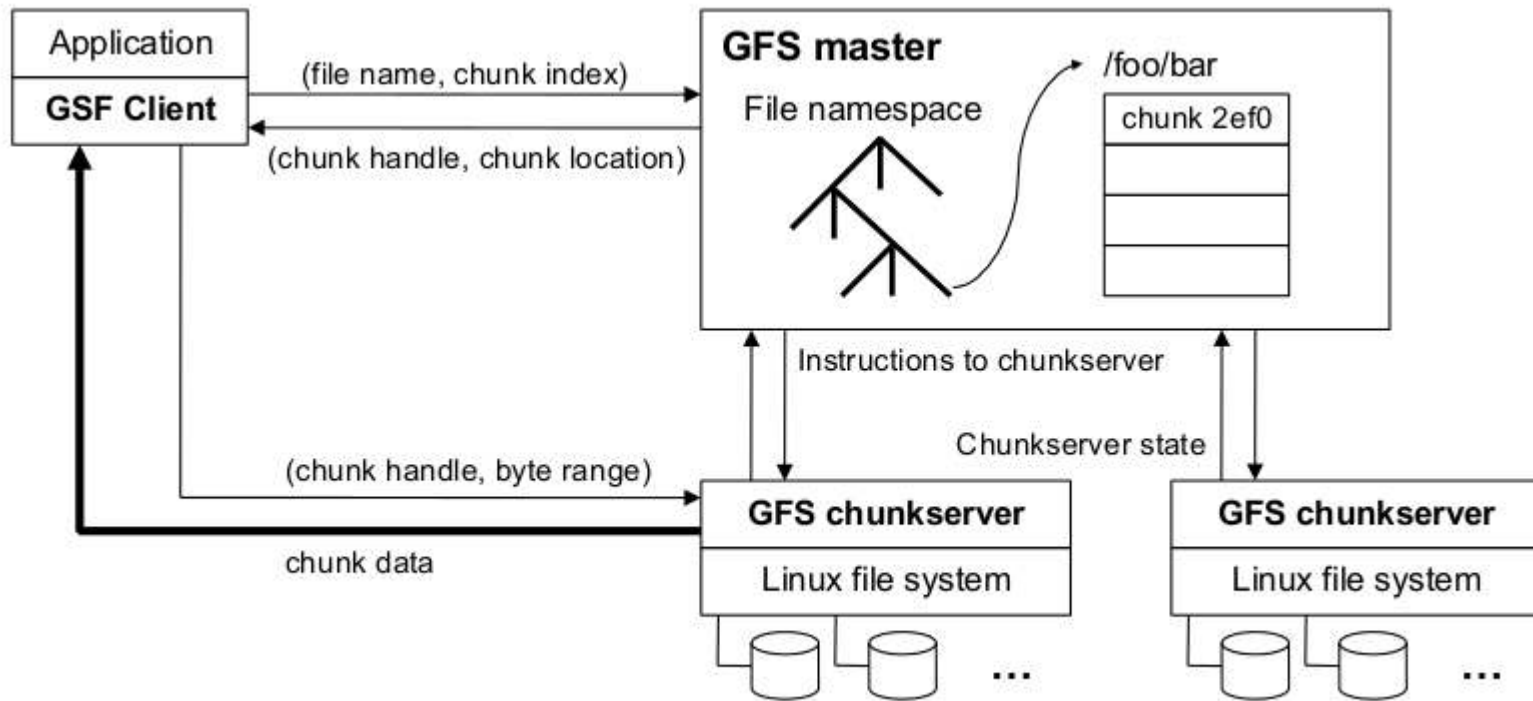
- Click to edit the outline text format
- Don't move data to workers... Move workers to the data.
  - Second Outline Level
  - Third Outline Level
    - Store data on the local disks for nodes in the cluster
    - Start up the workers on the node that has the data local
  - Fourth Outline Level
  - Fifth Outline Level
  - Sixth Outline Level
  - Seventh Outline Level
- Why ?
  - Not enough RAM to hold all the data in memory
  - Disk access is slow, disk throughput is good
- A distributed file system is the answer
  - GFS (Google File System)
  - HDFS for Hadoop (= GFS clone)
  - Deuxième niveau

- Commodity hardware over “exotic” hardware
- High Second Outline Level
  - Inexpensive commodity components fail all the time
- Third Outline Level
- “Modest Fourth Outline Level
- Files are written once, mostly appended to
  - Perhaps concurrently
- Sixth Outline Level
- Seventh Outline Level
  - Cliquez process
- High sustained throughput over low latency masque
  - Deuxième niveau
  - Troisième niveau

- Click to edit the outline text format
  - Files stored as chunks
    - Fixed size (64MB)
  - Reliability through replication
    - Each chunk replicated across 3 chunkservers
  - Single master to coordinate access, keep metadata
    - Simple centralized management
  - No data caching
  - Little benefit due to large data sets, streaming reads
  - Simplify the API
    - Push some of the issues onto the client
- masque
- Deuxième niveau
  - Troisième niveau



## Grid Computing by the fathers of the Grid



• modifier les styles du texte au masque

– Deuxième niveau

• Troisième niveau

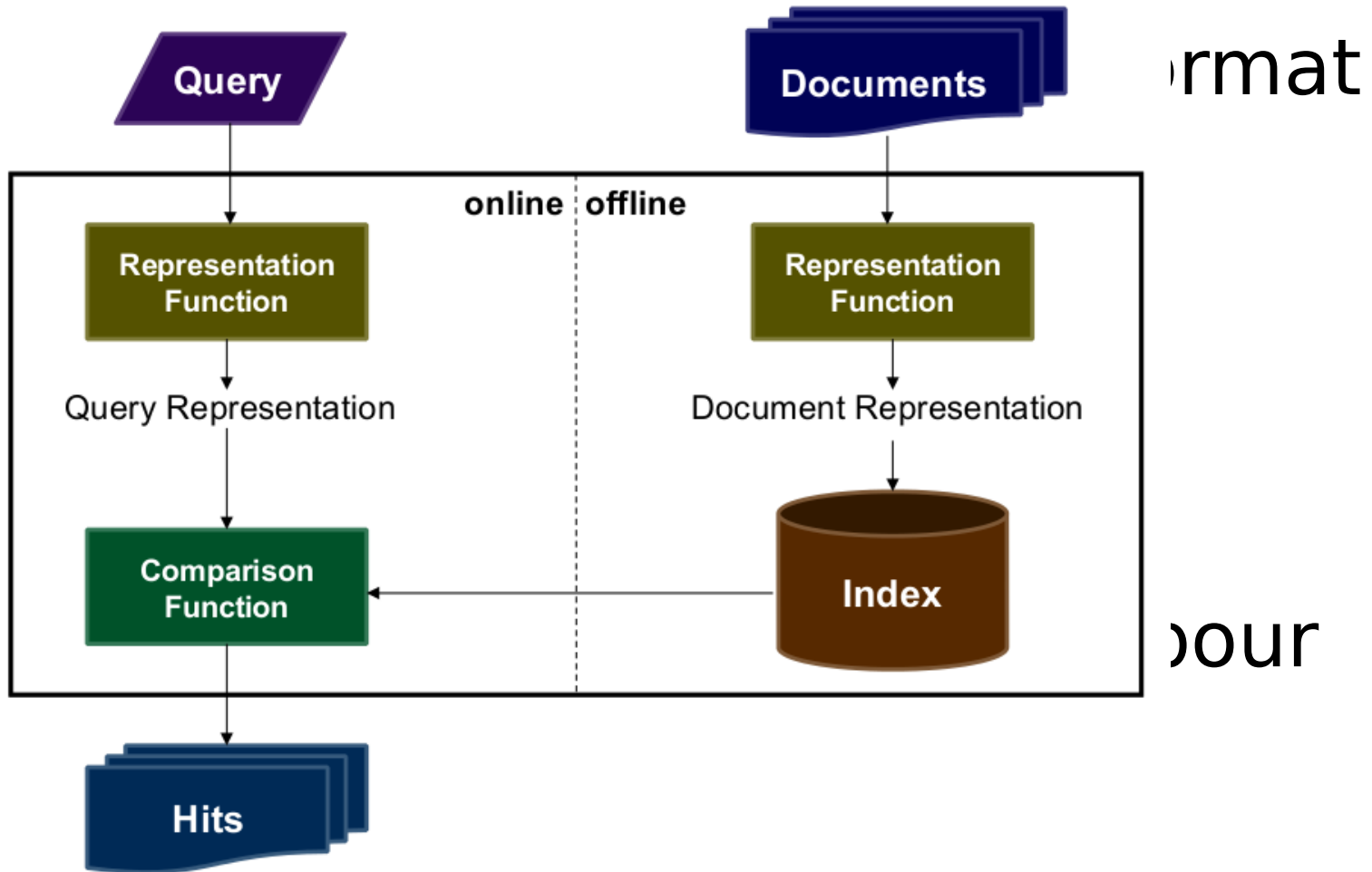


## Master's Responsibilities

- Metadata storage
  - First Outline Level
    - Click to edit the outline text format
  - Second Outline Level
- Namespace management/locking
  - Third Outline Level
- Periodic communication with chunk servers
  - Fourth Outline Level
  - Fifth Outline Level
- Chunk creation, replication
  - Sixth Outline Level
- Search
  - Seventh Outline Level
- Garbage collection
  - Cliquez pour modifier les styles du texte du masque
  - Deuxième niveau
    - Troisième niveau

- Click to edit the outline text format
  - Second Outline Level
    - Third Outline Level
      - Fourth Outline Level
        - Exemple : Inverted Indexing
          - Fifth Outline Level
            - Sixth Outline Level
- Seventh Outline Level Cliquez pour modifier les styles du texte du masque
  - Deuxième niveau
    - Troisième niveau

# Architecture of IR Systems



– Deuxieme niveau

• Troisième niveau

How do we represent text?

- “Bag of words” the outline text format
  - Treat all the words in a document as index terms for that document
    - Second Outline Level
    - Assign a weight to each term based on “importance”
      - Third Outline Level
      - Disregard order, structure, meaning, etc. of the words
        - Fourth Outline Level
        - Simple, yet effective!
          - Fifth Outline Level
          - Sixth Outline Level
- Assumptions
  - Seventh Outline Level
    - Term occurrence is independent
    - Document relevance is independent
    - Models are well-defined
  - Deuxième niveau
    - Troisième niveau

## Sample Document

- McDonald's slims down spuds
- Fast-food chain to reduce certain types of fat in its french fries with new cooking oil. "Bag of Words"
- NEW YORK (CNN/Money) - McDonald's Corp. is cutting the amount of "bad" fat in its french fries nearly in half, the fast-food chain said Tuesday as it moves to make all its fried menu items healthier.
- But does that mean the popular shoestring fries won't taste the same? The company says no. "It's a win-win for our customers because they are getting the same great french-fry taste along with an even healthier nutrition profile," said Mike Roberts, president of McDonald's USA.
- But others are not so sure. McDonald's will not specifically discuss the kind of oil it plans to use, but at least one nutrition expert says playing with the formula could mean a different taste.
- ...

Click to edit the outline text format

Second Outline Level

Third Outline Level

Fourth Outline Level

Fifth Outline Level

Sixth Outline Level

Seventh Outline Level

Cliquez pour modifier les styles du texte du masque

– Deuxième niveau

• Troisième niveau

# Representing Documents

## Document 1

- The quick brown fox jumped over the lazy dog's back.

## Document 2

- Now is the time for all good men to come to the aid of their party.



Term	Document 1	Document 2
aid	0	1
all	0	1
back	1	0
brown	1	0
come	0	1
dog	1	0
fox	1	0
good	0	1
jump	1	0
lazy	1	0
men	0	1
now	0	1
over	1	0
party	0	1
quick	1	0
their	0	1
time	0	1

format

Stopword List

for
is
of
the
to

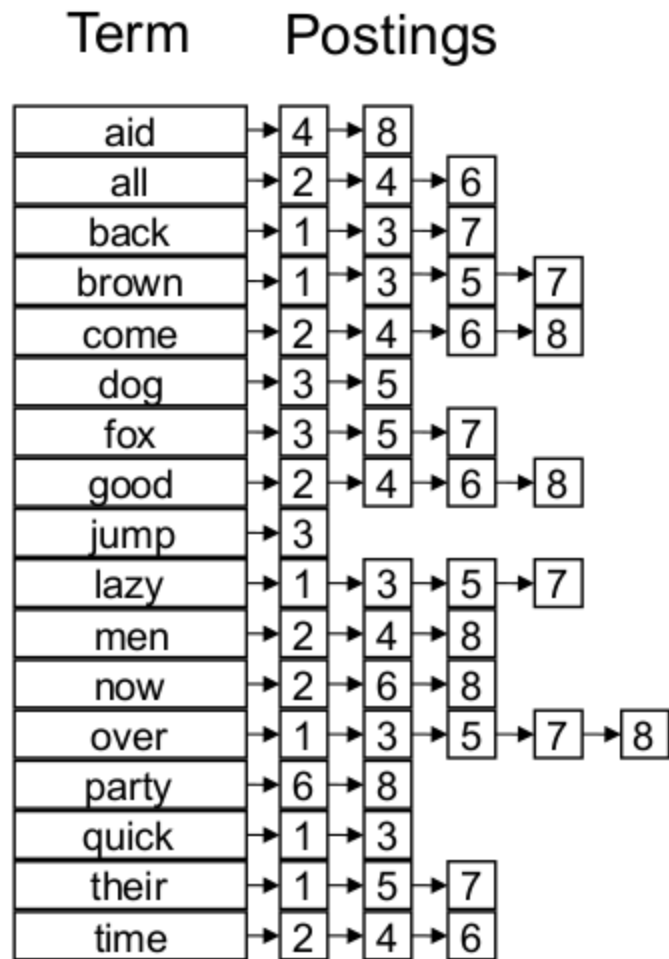
jour

– Deuxieme niveau

• Troisième niveau

## Inverted Index

Term	Doc 1	Doc 2	Doc 3	Doc 4	Doc 5	Doc 6	Doc 7	Doc 8
aid	0	0	0	1	0	0	0	1
all	0	1	0	1	0	1	0	0
back	1	0	1	0	0	0	1	0
brown	1	0	1	0	1	0	1	0
come	0	1	0	1	0	1	0	1
dog	0	0	1	0	1	0	0	0
fox	0	0	1	0	1	0	1	0
good	0	1	0	1	0	1	0	1
jump	0	0	1	0	0	0	0	0
lazy	1	0	1	0	1	0	1	0
men	0	1	0	1	0	0	0	1
now	0	1	0	0	0	1	0	1
over	1	0	1	0	1	0	1	1
party	0	0	0	0	0	1	0	1
quick	1	0	1	0	0	0	0	0
their	1	0	0	0	1	0	1	0
time	0	1	0	1	0	1	0	0



– Deuxième niveau

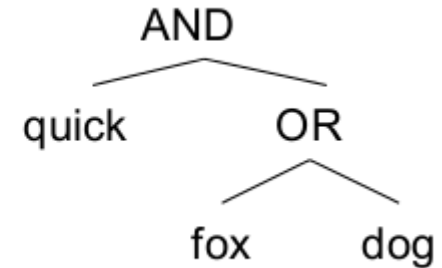
• Troisième niveau



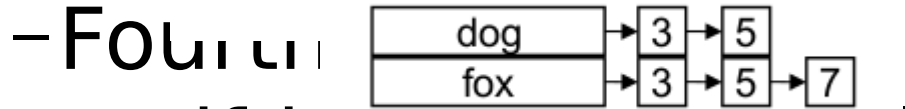
- To execute a Boolean query:

- **Click to e**

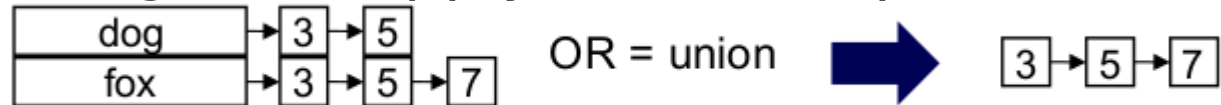
- Build query syntax tree
- **Second** ( fox or dog ) and quick



- For each phrase, look up postings



- Traverse postings and apply Boolean operator



- Efficiency analysis

- Postings traversal is linear (assuming sorted postings)
- Start with shortest posting first

masque

- Deuxième niveau

- Troisième niveau

- Click to edit the outline text format

$$w_{i,j} = \text{tf}_{i,j} \cdot \log \frac{N}{n_i}$$

$w_{i,j}$  weight assigned to term  $i$  in document  $j$

$\text{tf}_{i,j}$  number of occurrence of term  $i$  in document  $j$

$N$  number of documents in entire collection

$n_i$  number of documents with term  $i$

- Se  
m

our

masque

- Deuxième niveau

- Troisième niveau

	<i>tf</i>				
	1	2	3	4	<i>idf</i>
complicated			5	2	0.301
contaminated	4	1	3		0.125
fallout	5		4	3	0.125
information	6	3	3	2	0.000
interesting		1			0.602
nuclear	3		7		0.301
retrieval		6	1	4	0.125
siberia	2				0.602



complicated	→ 0.301	→ 3,5	4,2		
contaminated	→ 0.125	→ 1,4	2,1	3,3	
fallout	→ 0.125	→ 1,5	3,4	4,3	
information	→ 0.000	→ 1,6	2,3	3,3	4,2
interesting	→ 0.602	→ 2,1			
nuclear	→ 0.301	→ 1,3	3,7		
retrieval	→ 0.125	→ 2,6	3,1	4,4	
siberia	→ 0.602	→ 1,2			

## masque

– Deuxième niveau

• Troisième niveau

- The indexing problem
  - Second Outline Level
    - Must have sub-second response
  - Third Outline Level
    - For Web, incremental updates are important
  - Fourth Outline Level
    - Crawling is a phase in itself
  - Fifth Outline Level
    - Sixth Outline Level
- The retrieval problem
  - Seventh Outline Level
    - Cliquez pour modifier les styles du texte du masque
    - Deuxième niveau
      - Troisième niveau

- Click to edit the outline text format
- Fundamentally, a large sorting problem
  - Second Outline Level
    - Third Outline Level
      - Terms usually fit in memory
      - Fourth Outline Level
        - Postings usually don't
        - Fifth Outline Level
- How is it done on a single machine?
- Seventh Outline Level
  - How large is the inverted index?
    - Size of vocabulary
    - Size of postings
    - Deuxième niveau
      - Troisième niveau

- Click to edit the outline text format
- Map over all documents
  - Second Outline Level
    - Third Outline Level
  - Emit *term* as key, (*docid*, *tf*) as value
    - Fourth Outline Level
      - Fifth Outline Level
    - Sixth Outline Level
  - Trivial: each value represents a posting!
- Seventh Outline Level
  - Might want to sort the postings (e.g., by *docid* or *tf*)
- Reduce
  - Cliquez pour modifier les styles du texte du masque
- MapReduce does all the heavy lifting!
  - Deuxième niveau
    - Troisième niveau

- **MapReduce** is meant for text & data at batch processing Level
  - Not suitable for interactive operations requiring low latency
- **Third Outline Level**
- **Fourth Outline Level**
- **Fifth: “The secret sauce”**
- **Sixth Outline Level**
  - Most likely involves document partitioning
- **Seventh Outline Level** Cliquez pour modifier les styles du texte du masque
  - Lots of system engineering: e.g., caching, load balancing, etc.
  - Deuxième niveau
  - Troisième niveau

- Click to edit the outline text format
  - Second Outline Level
    - Third Outline Level
      - Fourth Outline Level
        - Fifth Outline Level
  - Sixth Outline Level
- Seventh Outline LevelCliquez pour modifier les styles du texte du masque
  - Deuxième niveau
  - Troisième niveau

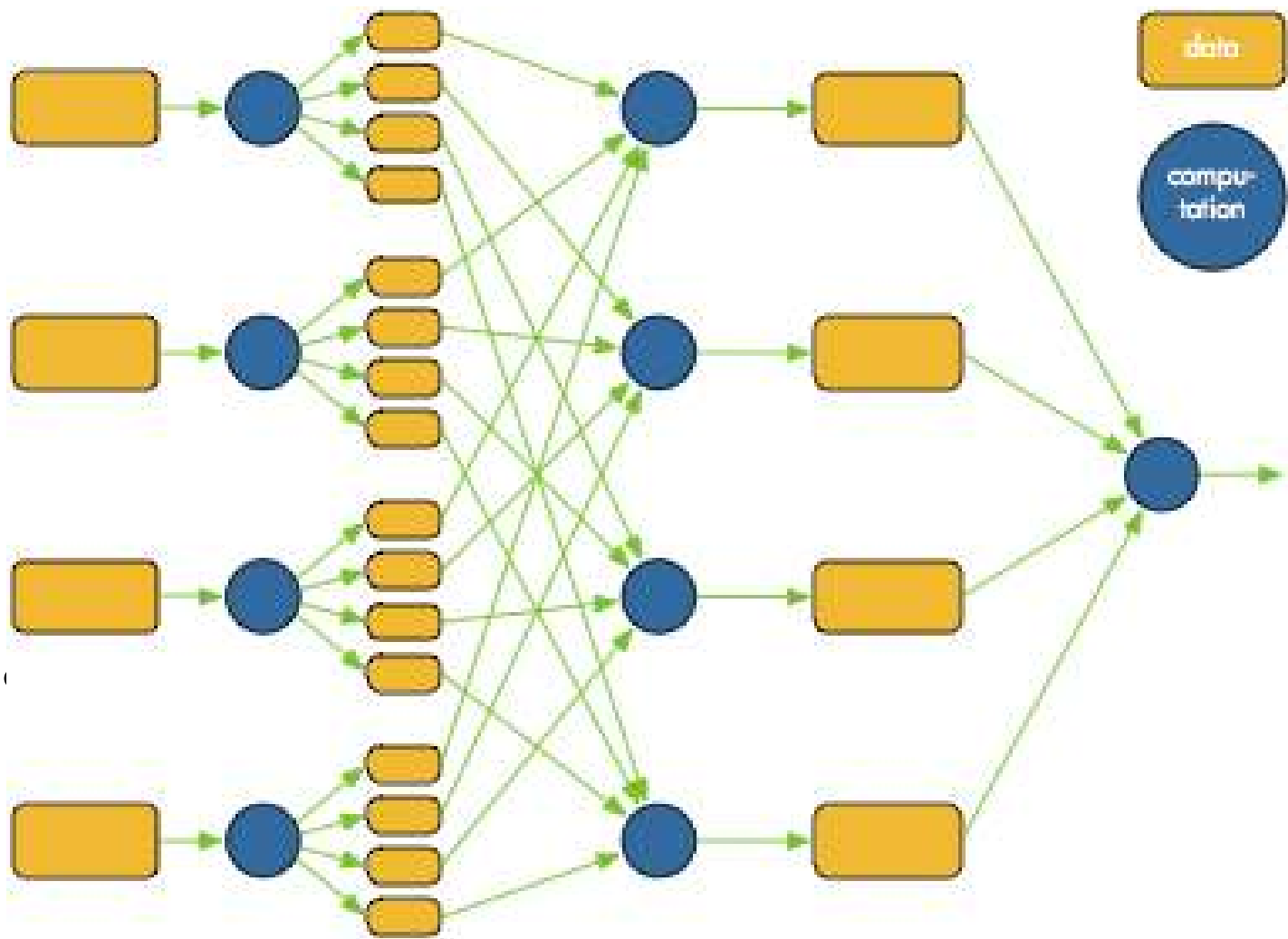


- Click to edit the outline text format
- Remember: Mappers run in isolation
  - Second Outline Level
    - Third Outline Level
    - You have no idea on what node the mappers run
  - Fourth Outline Level
    - Fifth Outline Level
    - You have no idea when each mapper finishes
  - Sixth Outline Level
- Tools for synchronization:
  - Ability to hold state in reducer across multiple key-value pairs
- Seventh Outline Level
  - Cliquez pour modifier les styles du texte du masque
  - Partitioner
  - Cleverly-constructed data structures
  - Deuxième niveau
    - Troisième niveau

For the programmer

- Input reader
  - Click to edit the outline text format
    - The input reader reads data from stable storage and generates key/value pairs.
- Map function
  - Second Outline Level
    - Third Outline Level
      - Takes a series of key/value pairs, processes each, and generates zero or more output key/value pairs
    - Fourth Outline Level
- Partition function
  - Fifth Outline Level
    - Each Map function output is allocated to a particular reducer by the application's partition function
  - Sixth Outline Level
- Compare function
- Reduce function
  - Seventh Outline Level
    - Cliquez pour modifier les styles du texte du masque
    - The framework calls the application's Reduce function once for each unique key in the sorted order
- Output writer
  - Deuxième niveau
    - writes the output of the Reduce to the stable storage
  - Troisième niveau

Input -> Map -> Copy/Sort -> Reduce -> Output



nat

ur

– Deuxième niveau

• Troisième niveau

## Use cases

- Word count in the text
- Before text format
- Second Outline Level message per word
- Third Outline Level in the text

```
class Mapper
  method Map(docid id, doc d)
    H = new AssociativeArray
    for all term t in doc d do
      H{t} = H{t} + 1
    for all term t in H do
      Emit(term t, count H{t})
```

Fourth Outline Level  
 Here  
 line Level  
 I message per  
 Outline Level word in the  
 text

- Seventh Outline Level Cliquez pour modifier les styles du texte du masque

– Deuxième niveau

• Troisième niveau

- Count the number of co-occurrence of  $n$  elements in sets
  - Second Outline Level
- Exemple
  - Third Outline Level
    - Words appears in same sentence
    - Customer who buy this also buy that
  - Fourth Outline Level
- If there are  $n$  elements
  - Fifth Outline Level
    - Report occurrence of  $N \times N$  couples
  - Sixth Outline Level
- On a single node, quite simple
  - Foreach set
    - Foreach i in set
      - Foreach j in set
        - `res[i][j]++`
  - Deuxième niveau
    - Troisième niveau

Cliquez pour modifier les styles du texte du masque

Map Reduce version ?

```

• class Mapper
  method Map(null, items [i1, i2, ...] )
    for all item i in [i1, i2, ...]
      for all item j in [i1, i2, ...]
        Emit(pair [i j], count 1)
  
```

```

class Reducer
  method Reduce(pair [i j], counts [c1, c2, ...])
    s = sum([c1, c2, ...])
    Emit(pair[i j], count s)
  
```

• Fifth Outline Level

• Sixth Outline Level

- Too many intermediary keys
- Easy and straightforward implementation
- Optimize using local accumulation of counts of  $[i,j]$ 
  - Easy optimization
  - Only few improvement (large space)

– Deuxième niveau

• Troisième niveau

format

Cliquez pour modifier les styles du texte du masque



## Stripes Approach

- ```

class Mapper
  method Map(null, items [i1, i2,...] )
    for all item i in [i1, i2,...]
      H = new AssociativeArray : item -> counter
      for all item j in [i1, i2,...]
        H{j} = H{j} + 1
      Emit(item i, stripe H)

```

```

class Reducer
  method Reduce(item i, stripes [H1, H2,...])
    H = new AssociativeArray : item -> counter
    H = merge-sum( [H1, H2,...] )
    for all item j in H.keys()
      Emit(pair [i j], H{j})

```

format

- Faster, lower number of intermediate keys
- Can lead to memory problems
- More complex implementation
  - Deuxième niveau

- Troisième niveau

Other exemples

- Grep
  - 10<sup>10</sup> 100-byte records
  - Second Outline Level
  - Seek a rare 3 letters word
    - Third Outline Level
  - 1800 machines
    - Fourth Outline Level
  - Peak performance 30 GB/s with 1764 workers
    - Fifth Outline Level
  - 150s
    - Sixth Outline Level
    - 1 minute startup
- Sort
- Cliquez pour modifier les styles du texte du masque
  - Same environment and dataset
  - 50 lines of code
  - 80 seconds
  - Deuxième niveau
    - Troisième niveau





## Characteristics

- Manage well failure
  - Click to edit the outline text format
  - Just send the keys again
- Heavy on the file system
  - Second Outline Level
  - Third Outline Level
  - Need dedicated and adapted filesystem
  - Fourth Outline Level
- Scale well
  - Fifth Outline Level
  - In term of data, workflow
  - Sixth Outline Level
- Easy to use
- Some translation tools from SQL are available
- Cliquez pour modifier les styles du texte du
- Middleware manages data- and computing-  
Infrastructure
  - Deuxième niveau
  - Troisième niveau

Some users

- Google
  - Click to edit the outline text format
    - They normalized it
    - Second Outline Level
    - They use it internally
      - Third Outline Level
        - large-scale machine learning problems,
        - clustering problems for the Google News and Froogle products,
        - extracting data to produce reports of popular queries (e.g. Google Zeitgeist and Google Trends).
      - Fourth Outline Level
      - Fifth Outline Level
      - Sixth Outline Level
  - Seventh Outline Level
    - Cliquez pour modifier les styles du texte du masque
      - language models processing for statistical machine translation, and large-scale graph computations.
    - Deuxième niveau
    - Troisième niveau

Other users

- Facebook
  - Click to edit the outline text format
    - Hadoop
    - Now use Corona (own implementation)
- Yahoo
  - Third Outline Level
    - Fourth Outline Level
    - More than 100,000 CPUs in more than 40,000 computers
    - Hadoop
  - Fifth Outline Level
- Linkedin
  - Sixth Outline Level
    - 5000 servers on hadoop
- Ebay
  - Seventh Outline Level
    - 532 nodes cluster (8 \* 532 cores, 5.3PB)

masque

- Deuxième niveau

• Troisième niveau



Some links

- Google
  - **Click to edit the outline text format**
  - **Second Outline Level**
    - **MapReduce: Simplified Data Processing on Large Clusters** by Jeffrey Dean and Sanjay Ghemawat
  - **Third Outline Level**
    - Technical report
- Apache
  - **Fourth Outline Level**
    - **Hadoop: The definitive guide**
    - Book
  - **Sixth Outline Level**
- Microsoft
  - **Seventh Outline Level**
    - **Google's MapReduce Programming Model – Revisited**
    - **masque**
    - **Technical report**
    - **Deuxième niveau**
    - **Troisième niveau**