

# CLOUD COMPUTING

# MapReduce

Lorenzo Leonini  
lorenzo.leonini@unine.ch

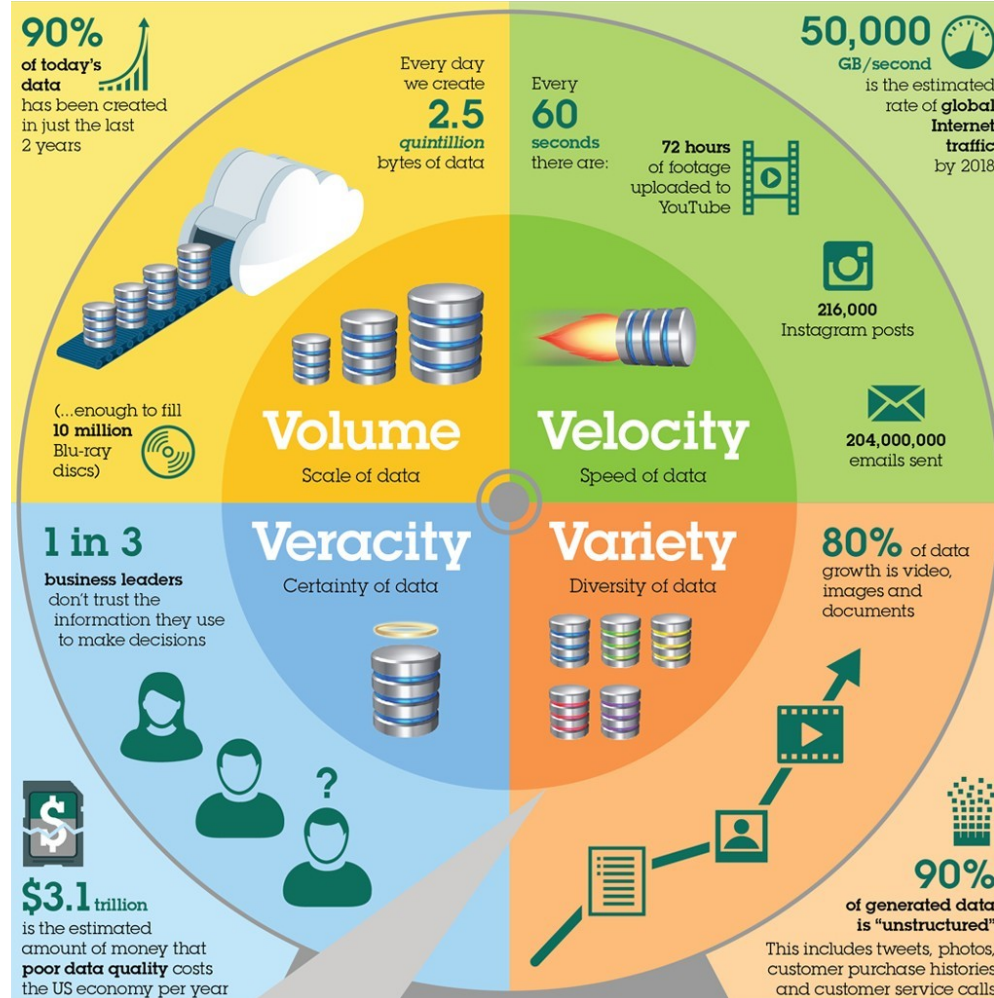
# Objectives

- Introduce the need for large-scale data processing in cloud environments
- Present the Map/Reduce principle and supporting frameworks
- Detail some representative applications

# Introduction

- Cloud computing allows larger-scale applications
- Large-scale applications attract more users
  - Users generate more data
  - Application generate more data (e.g. logs)
- How can we leverage this data to make the application better?
- "Big Data" phenomenon

# The 4 "V" of Big Data



# Data-supported services

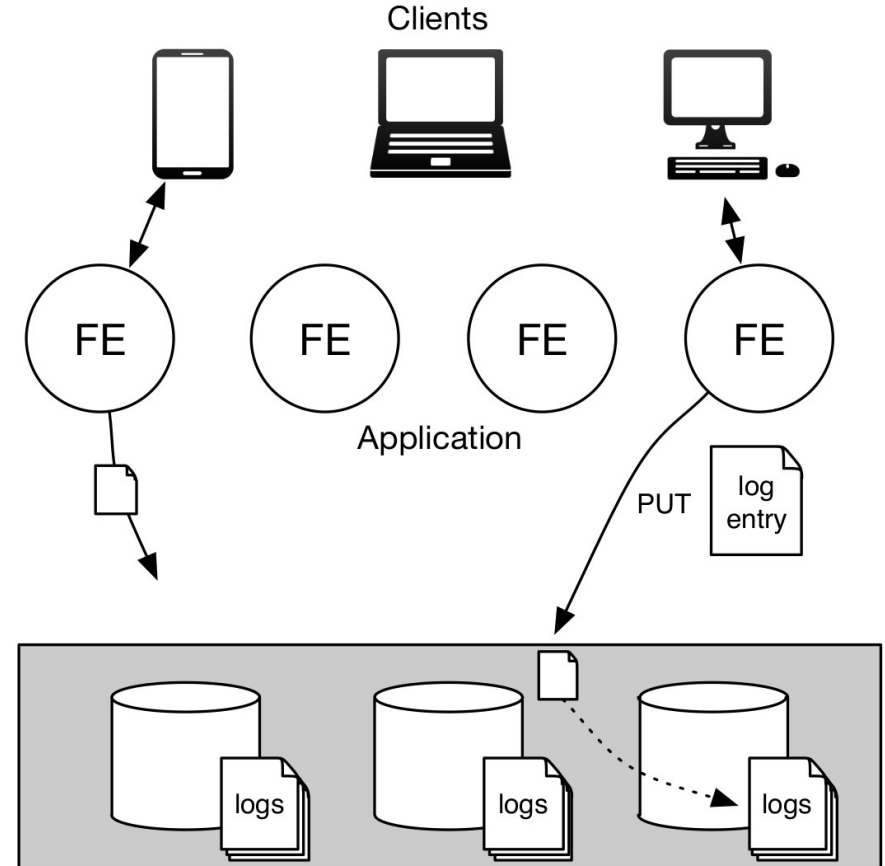
- Applications using large volumes of public data
  - Web search and indexing (initially)
- Applications using user-generated data
  - Social networks
  - Recommendation engines (Netflix, Spotify, etc.)
  - Taxi hailing predictive models (where will clients be? Where will they want to go?)
- And using both
  - Web search and indexing (now)
  - Cambridge Analytica...

# Dealing with the two first “V”s

- Need specific tools and programming models
  - (Very) large amounts of data
  - Complex environment (replication, distribution)
  - Faults and slow machines
- How to handle **V**olume & **V**elocity
  - Map/Reduce framework for large static data
  - Stream processing frameworks for dynamic data
- We will not cover Variety (data curation) and Veracity (data authenticity)

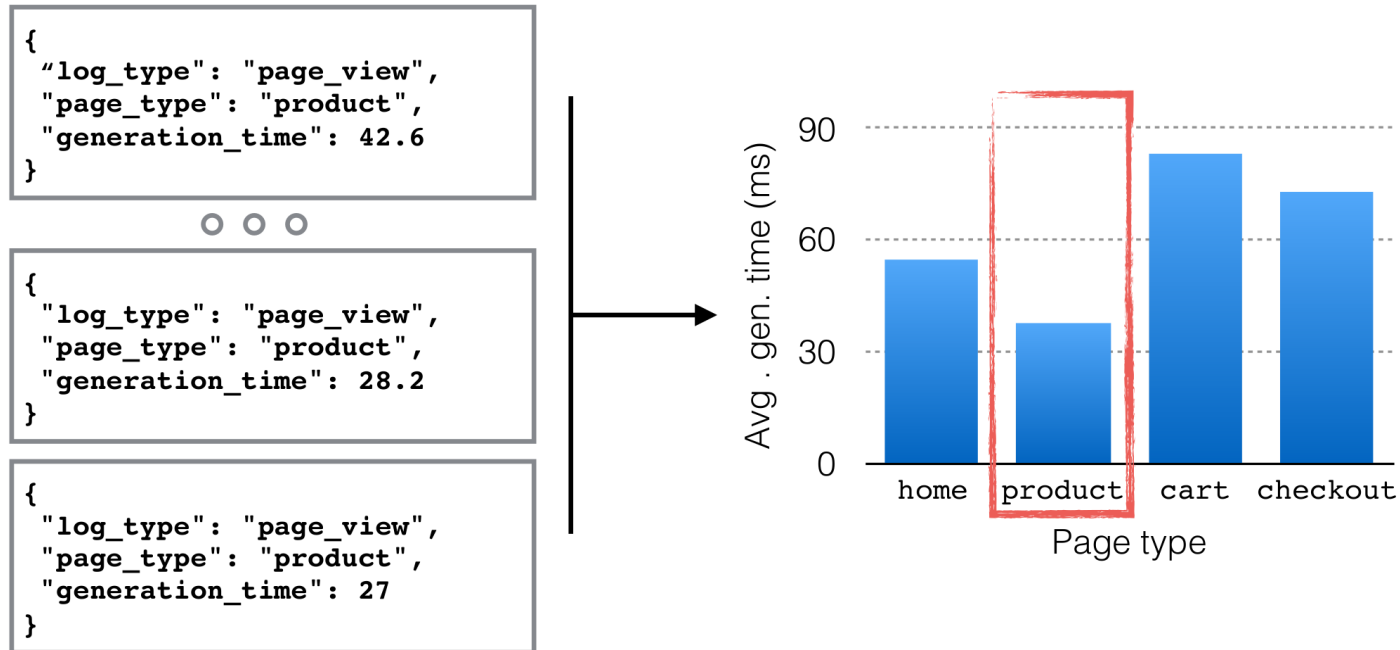
# Motivating example: Logging

- Application front-end components generate logs
  - Client information
  - Errors
  - Page generation time
  - Items accessed/purchased
- Log entries stored as JSON documents in a distributed database



# Processing logs

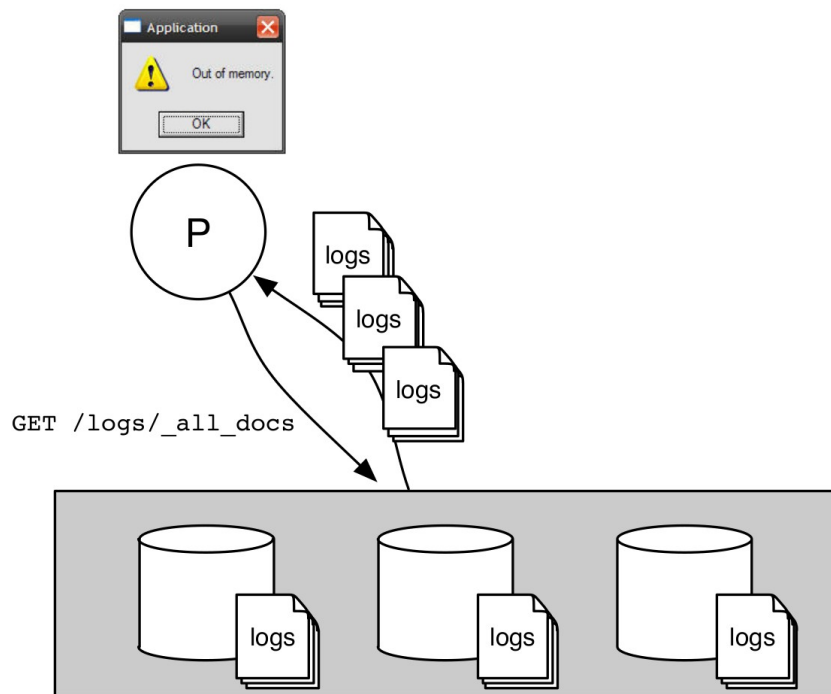
- We would like to know the *average generation time* for each type of page
  - **Types:** home, product, cart, checkout
- This information is available by *processing the logs*



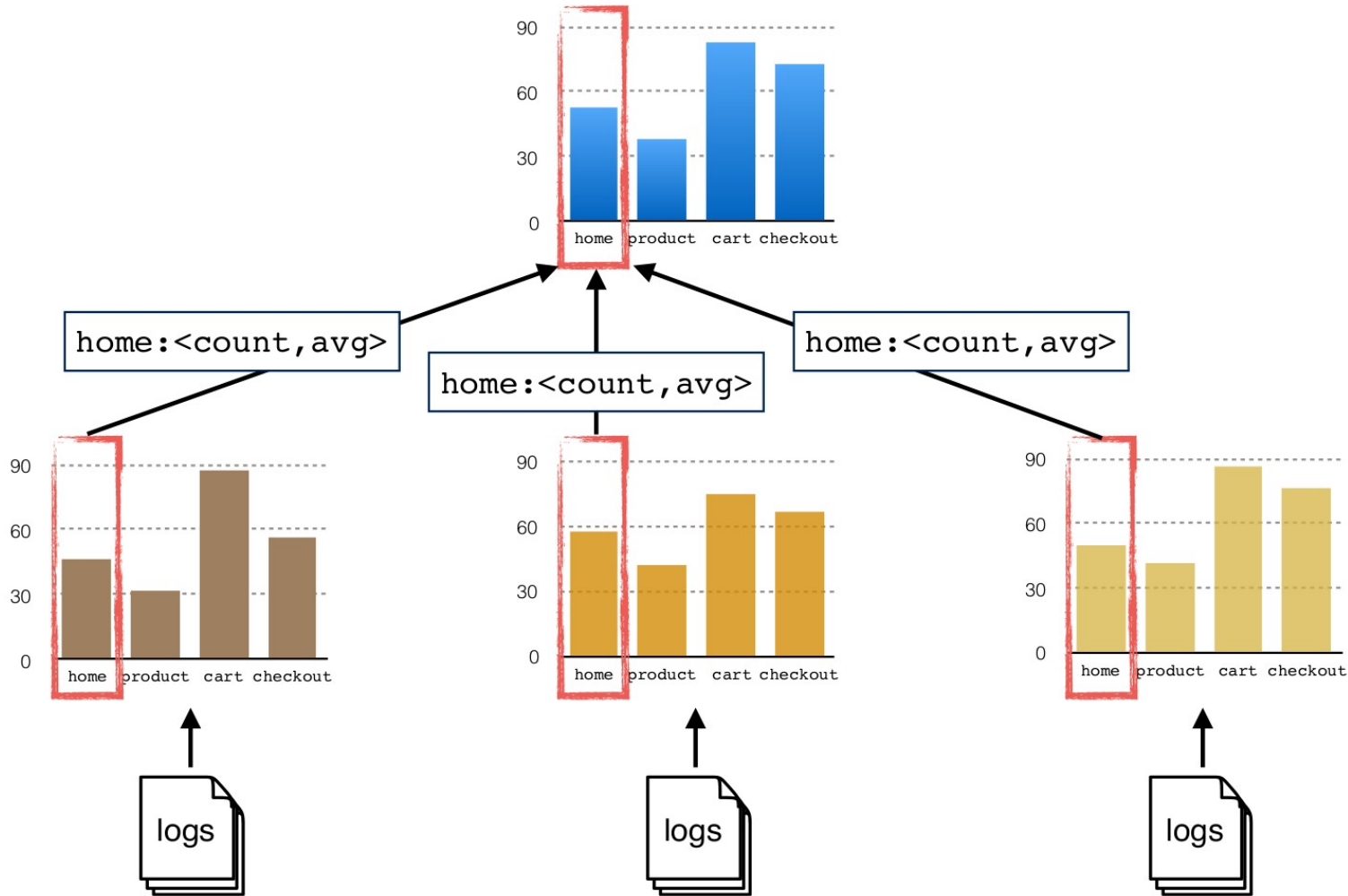


# Centralized processing?

- Should we collect all log documents and process them in a *dedicated single process*?
  - Amount of logs can be very large
    - May not fit in memory at all!
  - Not all logs entries are useful for our operation: wasted bandwidth
  - Very slow
- Need *parallel processing*



# Processing logs in parallel



# Handling Volume

- Split the dataset in multiple partitions, process each partition independently
  - Later, merge outputs for final result
- Difficult to do "manually"
  - Deploy worker processes close to the data
  - Coordinate all the workers
  - Handle faults at the workers
  - Collect all outputs and process them
  - Start again ...

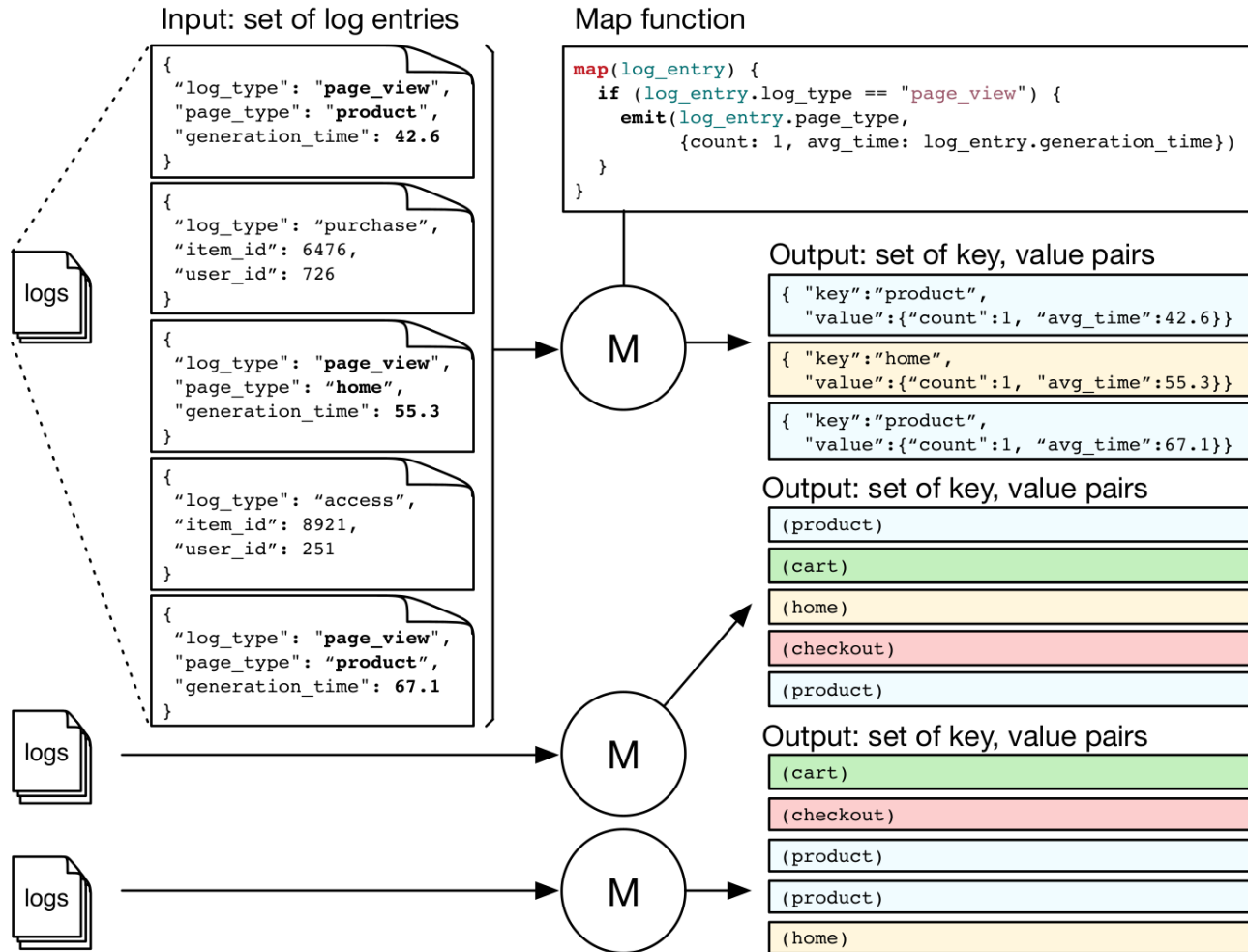
# Map/Reduce

- Big Data processing often follows a similar pattern:
  - PARTITION - Iterate in parallel over a large number of records
  - MAP - Extract information of interest from the records
  - SHUFFLE - Regroup this information in sets, one for each category
  - REDUCE - Aggregate each of the sets to obtain the final result(s)
- Map/Reduce is a programming model that allows expressing such queries and running them in parallel on many machines
- *Key idea:* provide an abstraction at the point of the two operations MAP and REDUCE, make others implicit

# Programming model

- Programmer specifies **two functions**
  - **map(record)**
    - get records from a partition of the source data
    - can generate <key, value> pairs with the emit call
  - **reduce(key, [values])**
    - receives all values for the same key
    - generate aggregate results, also as <key, value> pairs
- **PARTITION** and **SHUFFLE** are handled transparently

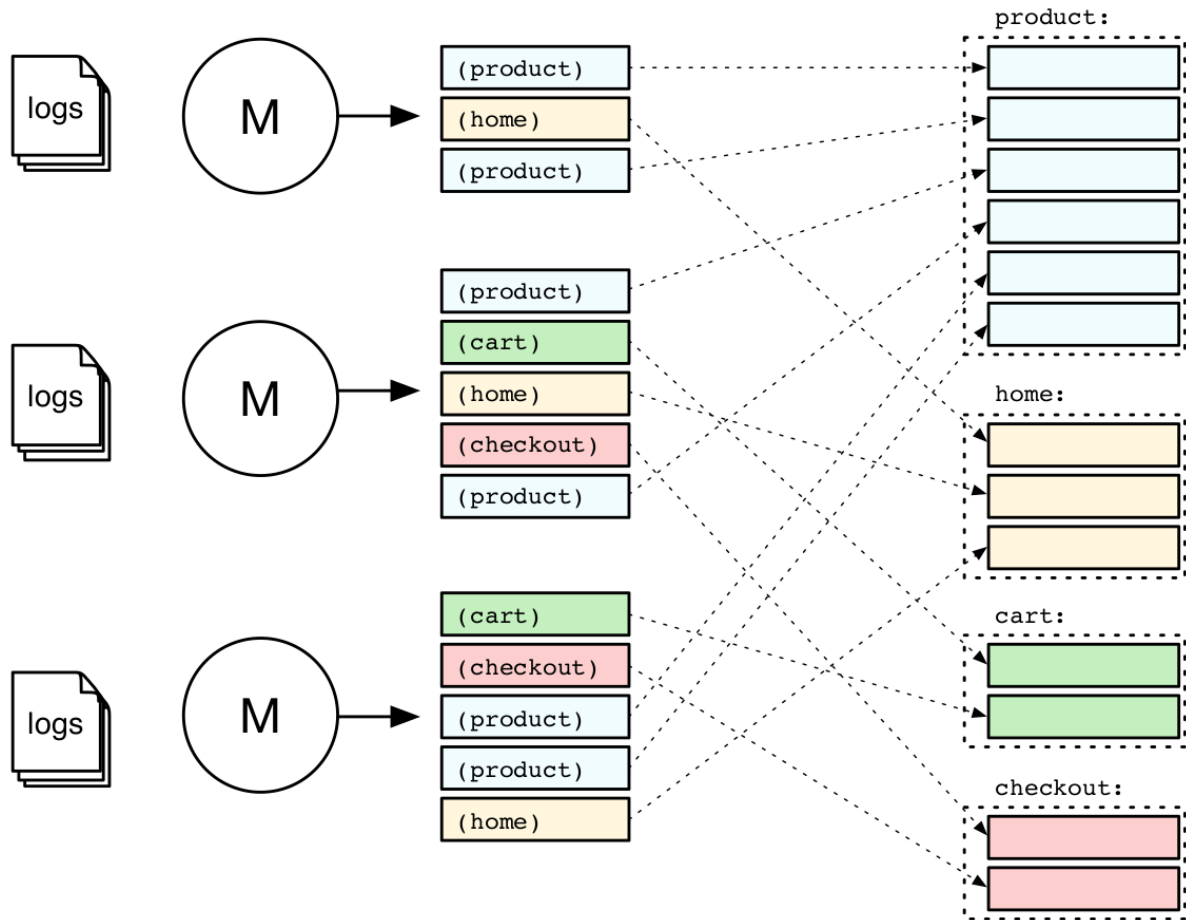
# MAP phase



# SHUFFLE phase

Map outputs:  
sets of key, value pairs

Shuffle output:  
sets of values for each key

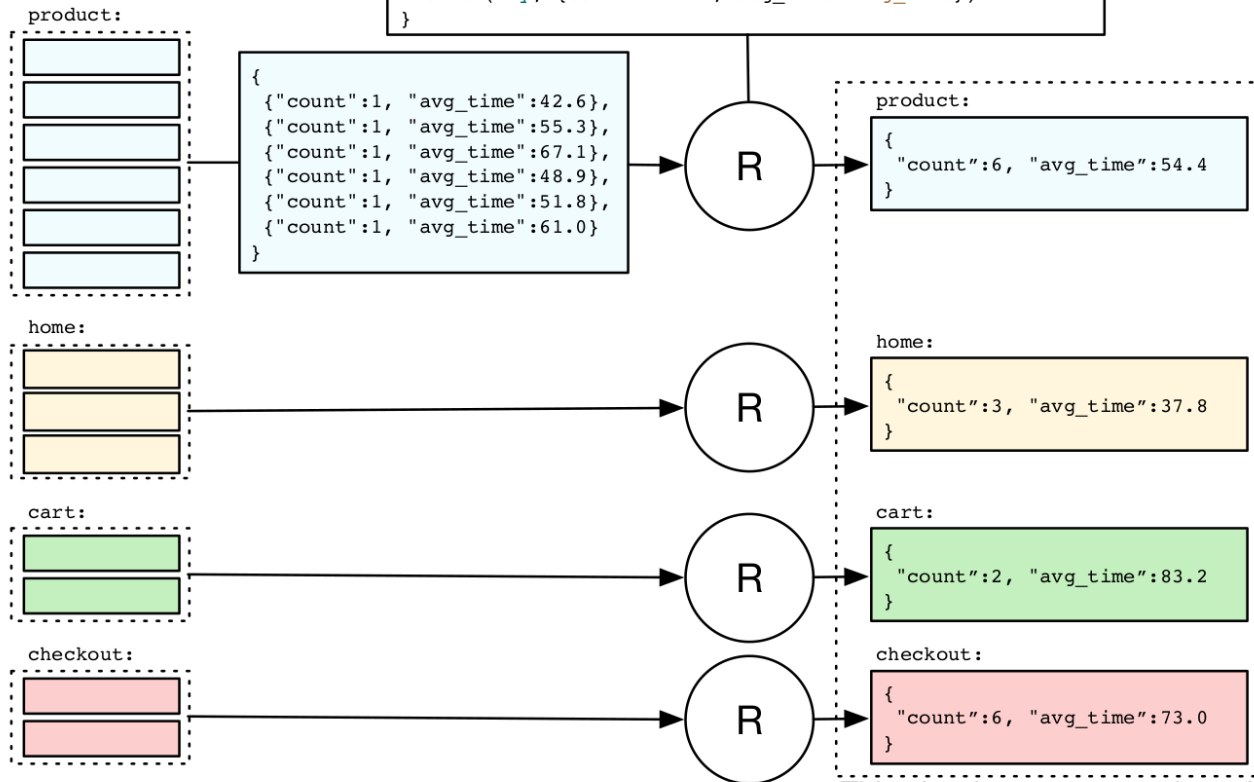


# REDUCE phase

Reduce function

```
reduce(key, values) {  
  var count = values.reduce((a,b) => a.count + b.count)  
  var avg_time = (values.reduce((a,b) =>  
    a.avg_time * a.count +  
    b.avg_time * b.count)  
    / count)  
  emit (key, {count: count, avg_time: avg_time})  
}
```

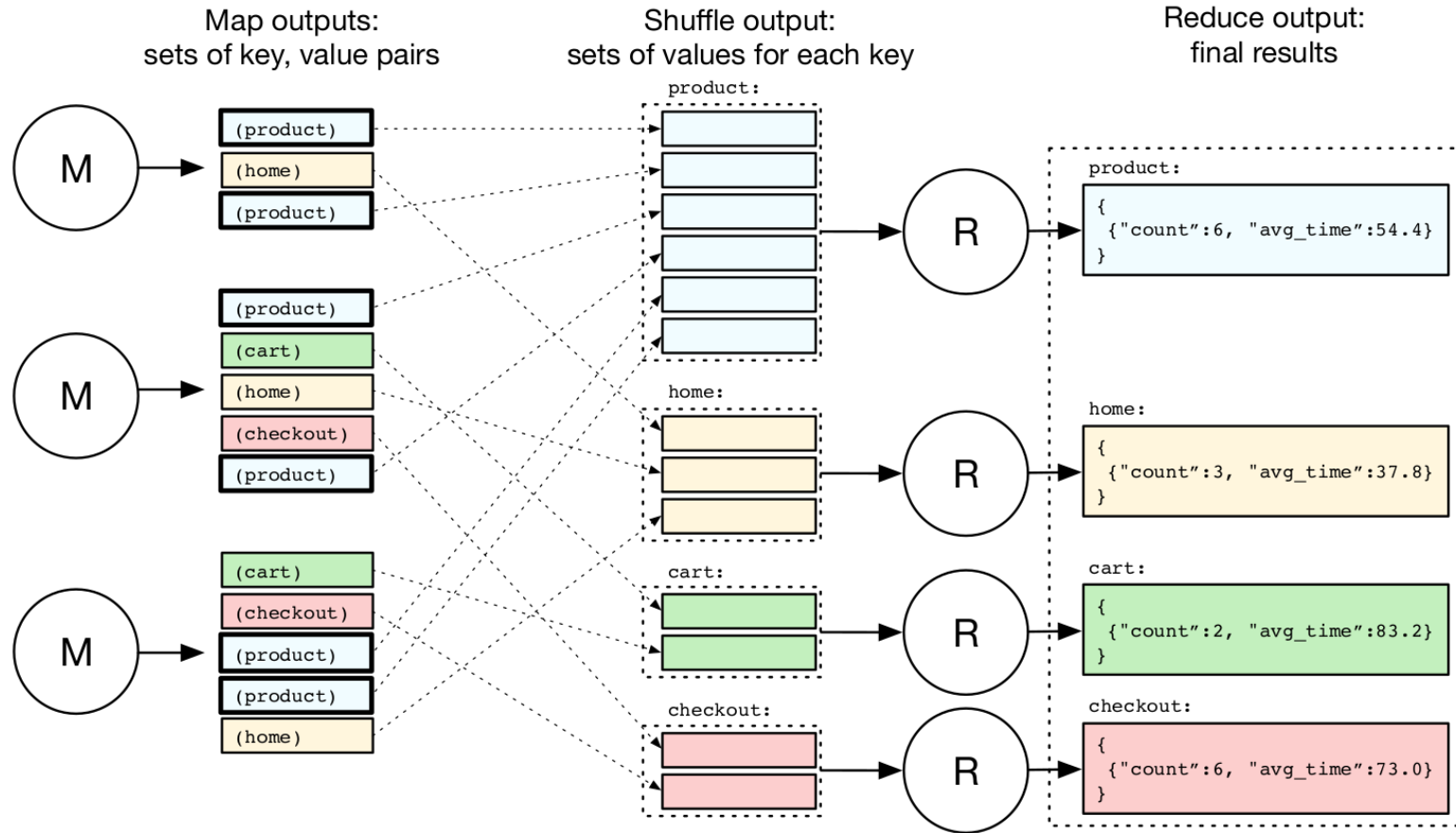
Shuffle output:  
sets of values  
for each key



This view is our expected result



# Bandwidth inefficiency



- Several <key, value> pairs for the same key are generated by each mapper and shuffled independently to the corresponding reducer!

# Local reduction

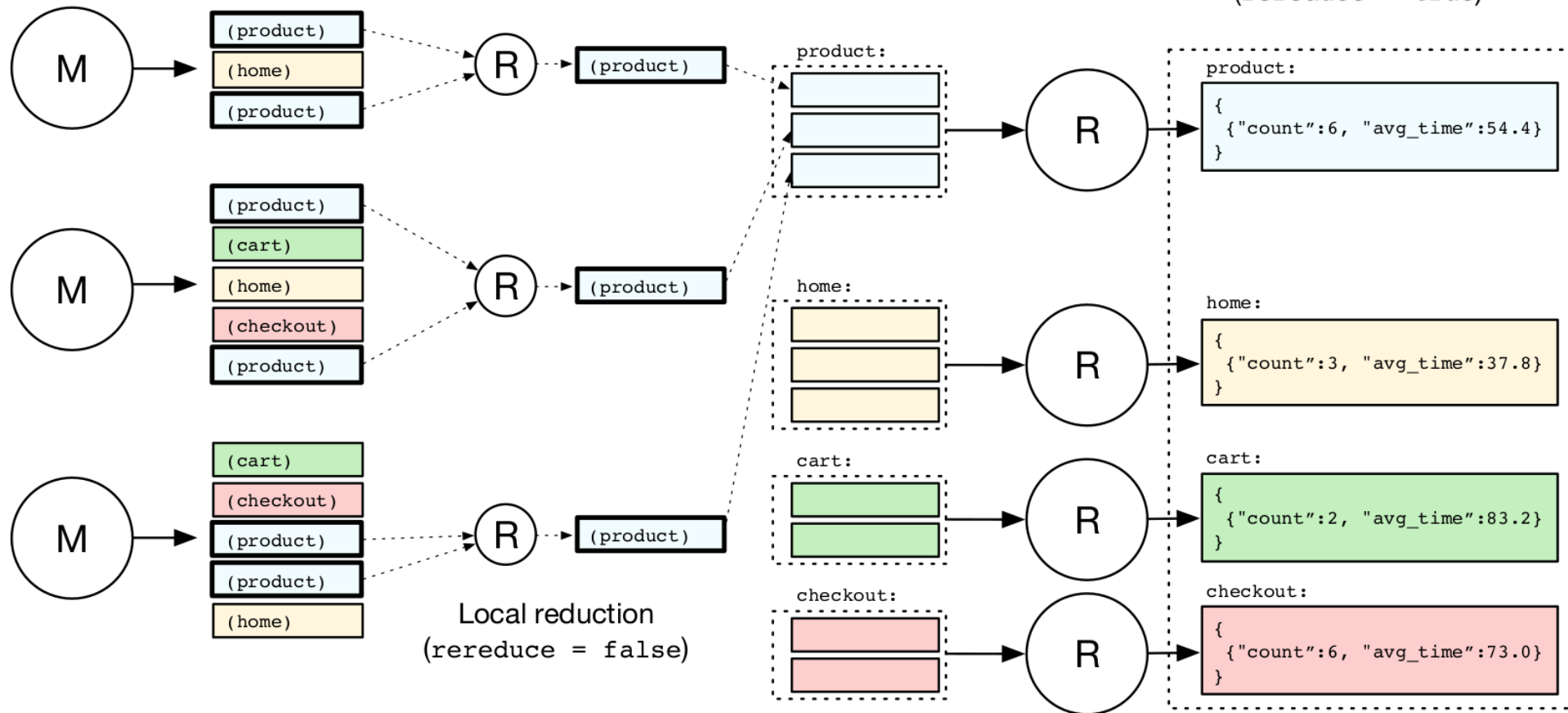
- Could we locally aggregate all <key, value> pairs for the same key at the Mapper process, before sending them to the SHUFFLE phase?
- The output of `reduce()` functions can be used in its input list
  - Local MAP output can be locally reduced before the SHUFFLE phase
  - This property holds for many aggregations: max, min, average,...
- Sometimes, we need to know if we are reducing local MAP results or the result of the SHUFFLE phase
  - Optional "rereduce" parameter to the `reduce()` call tells us if this is a local reduction or the final one:  
`reduce(key, [values], rereduce)`

# Local reduction applied

Map outputs:  
sets of key, value pairs

Shuffle output:  
sets of values for each key

Reduce output:  
final results  
(rereduce = true)



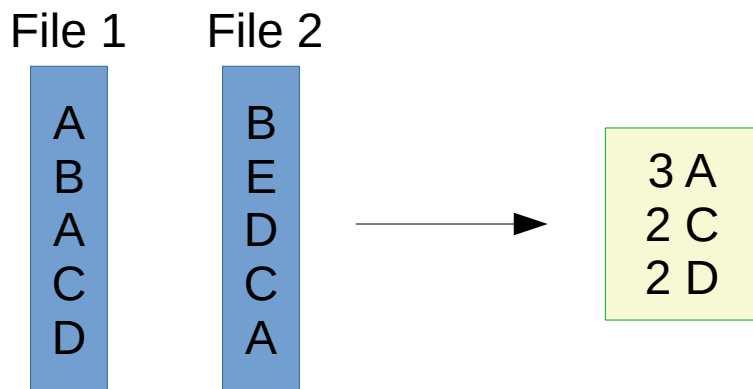
# Link with functional programming

- Higher-order functions take a function as argument
- **Map**: function applied to every element in a list
  - Result is a new list
  - Each operation is independent
- **Fold**: accumulator set to initial value, function applied to list element and the accumulator, result stored in accumulator
  - Repeated for every item in the list
  - Result is the final value in the accumulator
- Map/Reduce similar in principle but multiple output lists for Map and parallel Fold operations

# **Examples of application**

# Distributed grep

- Count lines in all files that match <regex>, show in descending order
  - `grep -Eh <regex> dir/* | sort | uniq -c | sort -nr`
- Example
  - `grep -Eh 'A|C|D' File* | sort | uniq -c | sort -nr`



# Distributed grep

- **Map** function
  - Input: files <line>
  - Emit: if line matches pattern <line, 1>
- **Reduce** function
  - Input: <line, [1,1, ...]>
  - Emit: <line, n> where n is the total of 1s

# Distributed grep

File 1

A  
B  
A  
C  
D

File 2

B  
E  
D  
C  
A

MAP

(A)  $\rightarrow$  (A, 1)

(B)  $\rightarrow$

(A)  $\rightarrow$  (A, 1)

(C)  $\rightarrow$  (C, 1)

(D)  $\rightarrow$  (D, 1)

(B)  $\rightarrow$

(E)  $\rightarrow$

(D)  $\rightarrow$  (D, 1)

(C)  $\rightarrow$  (C, 1)

(A)  $\rightarrow$  (A, 1)

REDUCE

(A, [1, 1, 1])  $\rightarrow$  (A, 3)

(C, [1, 1])  $\rightarrow$  (C, 2)

(D, [1, 1])  $\rightarrow$  (D, 2)

3 A  
2 C  
2 D



# Wordcount

- Count how many times each word appears in a text corpus

```
map(text) {  
  foreach (word: text.split(' ')) {  
    emit(word, {count: 1})  
  }  
}
```

```
reduce(key, values) {  
  var count = values.reduce((a,b) => a.count + b.count)  
  emit (key, {count: count})  
}
```

# Local reducers

- Local reducers helps save bandwidth by pre-reducing the results of a map locally and send fewer <key, value> intermediate pairs
- In the previous examples, the same code can be used as local reducer
  - Aggregate <regex, 1> counts for lines matching the pattern(s) under a single <regex, n> pair
  - Aggregate <word, 1> counts for different words under a single <word, n> pair
  - Not always the case!

# Top-k page frequency

- **Input**: log of web page requests
- **Output**: the top- $k$  accessed webpages
- **Map** function
  - Parse log, output  $\langle \text{URL}, 1 \rangle$  pairs
  - Similar to word count but goes to a **single key**
- **Reduce** function (for single key)
  - Aggregate  $\langle \text{URL}, 1 \rangle$  pairs for individual URLs
  - Sort by decreasing frequency
  - Output top- $k$   $\langle \text{URL}, n \rangle$  pairs

# Top-k page frequency

- Going to a single key is vastly inefficient
  - A single worker will receive all individual  $\langle \text{URL}, 1 \rangle$  pairs
  - Only to aggregate them and get the top-k URLs
- Use local reducer?
  - Here we cannot use the same code as in previous examples
    - It includes the selection of the top-k URLs
    - The union of the local top-k is not the global top-k
    - Sometimes OK to do local trimming for approximate solution

# Reverse Web-link graph

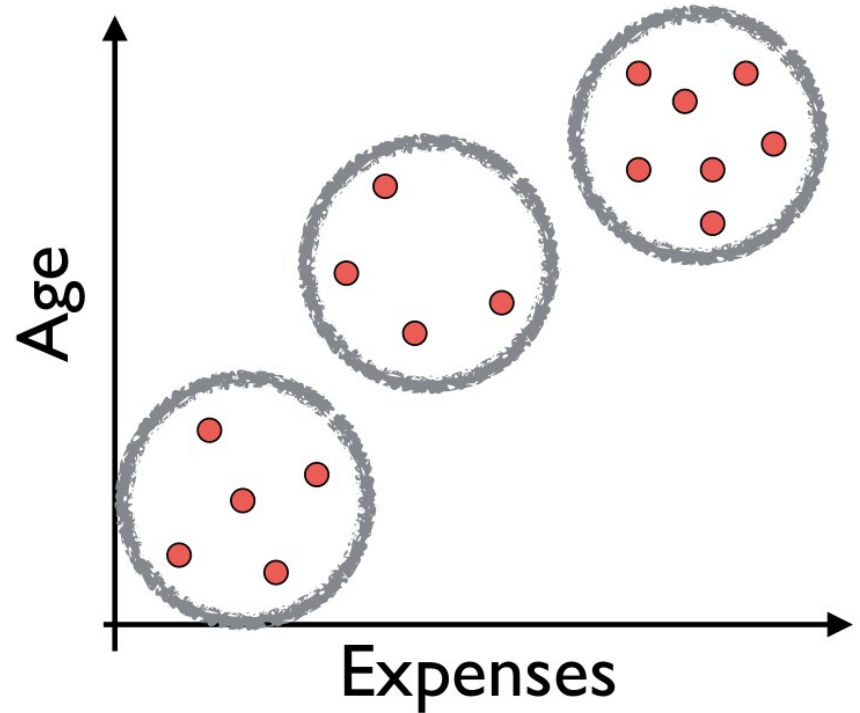
- Get all the links pointing to some page
  - This is the basis for the [PageRank algorithm](#)!
- **Map** function
  - Input: web pages (fetched by a *web crawler*)
  - Output a <target, source> pair for each link to *target* URL in a page named source
- **Reduce** function
  - Concatenate the list of all source URLs associated with a given target URL and emits the pair:  
<target, list(sources)>

# Inverted index

- Get all documents containing some particular keyword
  - Used by the search mechanisms of Google, Yahoo!, etc.
  - Second input for PageRank
- **Map** function
  - Input: (e.g.) web pages (fetched by a *web crawler*)
  - Parse each document and emit a set of pairs <word, documentID>
- **Reduce** function
  - Take all pairs for a given word
  - Sort the documents IDs
  - Emit a final <word, list(document IDs)> pair

# K-Means clustering

- Typical data mining / database problem
- Group items into  $k$  clusters
  - Clusters must minimize distance of contained points (relative to the center)
- NP hard problem
  - Heuristic algorithm



# K-Means algorithm

- Goal: obtain  $m_1, m_2, \dots, m_k$  as representative points of the  $k$  clusters
  - These will be the *centroids* of the clusters
- Initialize  $m_1, m_2, \dots, m_k$  to random values
- Iterate ( $t$  = iteration)

- Assign each point to the closest centroid

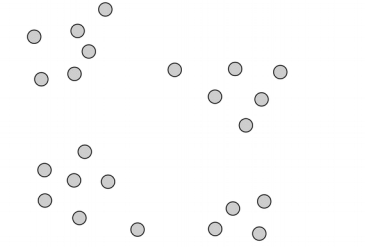
$$S_i^{(t)} = \left\{ \mathbf{x}_j : \|\mathbf{x}_j - \mathbf{m}_i^{(t)}\| \leq \|\mathbf{x}_j - \mathbf{m}_{i^*}^{(t)}\| \forall i^* = 1, \dots, k \right\}$$

- $m_i$  becomes the new centroid for its points

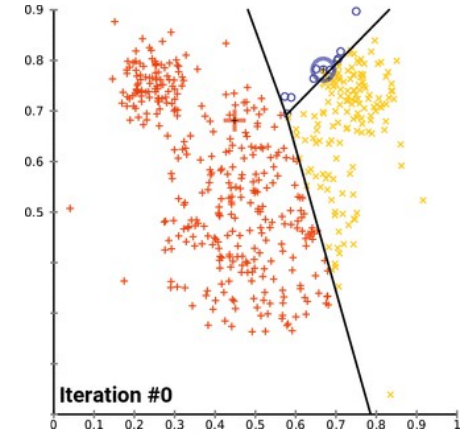
$$m_i^{(t+1)} = \frac{1}{|S_i^{(t)}|} \sum_{x_j \in S_i^{(t)}} x_j$$



# K-Means iterations



0a. Données d'entrée



# K-Means in Map/Reduce

- Initialize
  - Choose centroids randomly

- Classify

- Map each point to the closest centroid

$$S_i^{(t)} = \left\{ \mathbf{x}_j : \|\mathbf{x}_j - \mathbf{m}_i^{(t)}\| \leq \|\mathbf{x}_j - \mathbf{m}_{i^*}^{(t)}\| \forall i^* = 1, \dots, k \right\}$$

- Recenter

- $m_i$  becomes new centroid for its points

$$m_i^{(t+1)} = \frac{1}{|S_i^{(t)}|} \sum_{x_j \in S_i^{(t)}} x_j$$

- Repeat until no change

- Centroids have converged

# Classification step as Map

- Initialize
  - Set global var centroids from file
    - Initially  $k$  random points
- `map(point)`
  - Access to global var centroids
    - Can also be provided as data directly inside the map function definition to keep the functional paradigm
  - Compute nearest centroid based on centroids:
    - `emit(nearest centroid, point)`

# Recenter step as Reduce

- Initialize
  - Set global var centroids = []
- `reduce(centroidt, [points])`
  - Recompute centroid from points in it
    - `emit(centroidt+1)`
- Cleanup (after all calls to reduce are made):
  - Save new centroids to file
  - Check if change since last iteration (if yes → re-iterate)
- Does not completely "fit" Map/Reduce model that assume no global shared state
  - But widely used in practice

# Origin and implementation

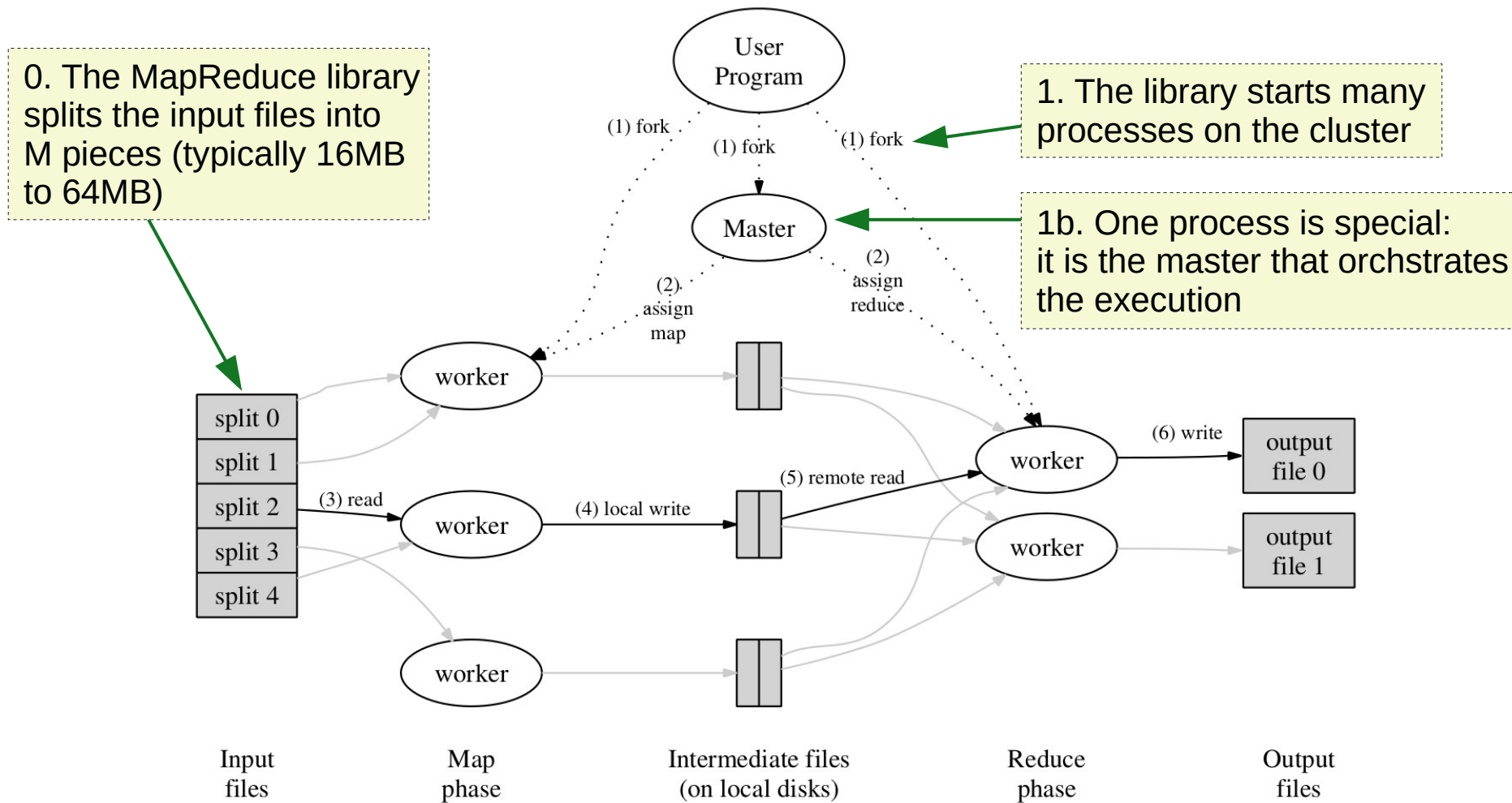
# Origin

- Map/Reduce originally proposed by Google in 2004
  - Necessary due to their unprecedented scale
  - One (if not the) most cited paper in computer science
  - Computation of the PageRank for webpages
    - Also, extraction of keywords, sanitation, etc.
    - Many other uses: log parsing, network monitoring, etc.
- Proprietary implementation
- Huge and immediate success
  - Simple programming model
  - Distributed computation complexity left to framework
  - Usable by non-CS majors

# Original Map/Reduce

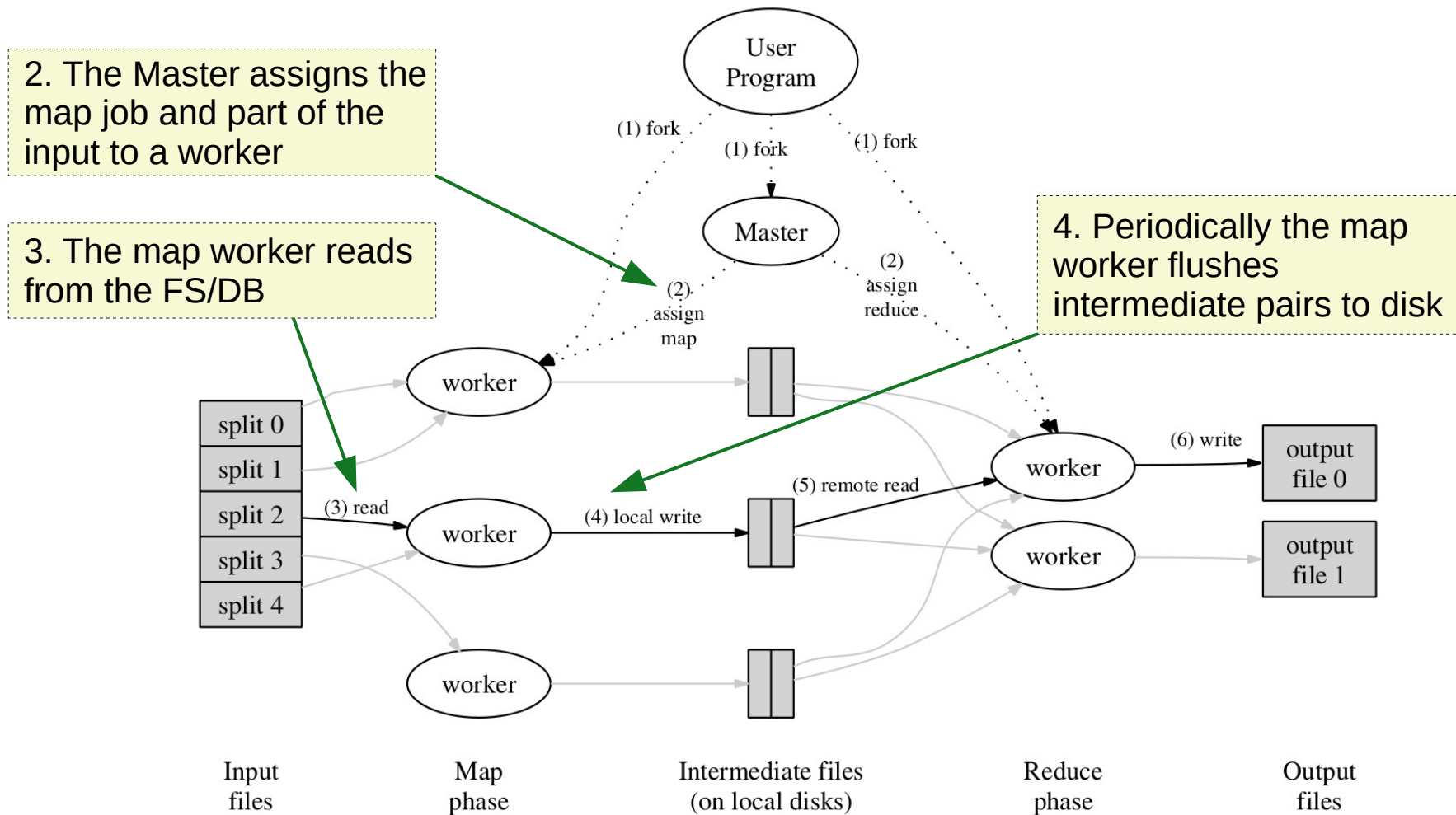
- The job is submitted to a **master** process
- The master orchestrates its execution
- Each node supports one or more workers
- Each worker can handle a map or reduce job when instructed by the master
- Map jobs get the data from file system
  - Google File System: optimized for storing large append-only files
  - Also possible from NoSQL database: Google BigTable
- Communication is based on key/values pairs
- Output written to file system

# Execution overview





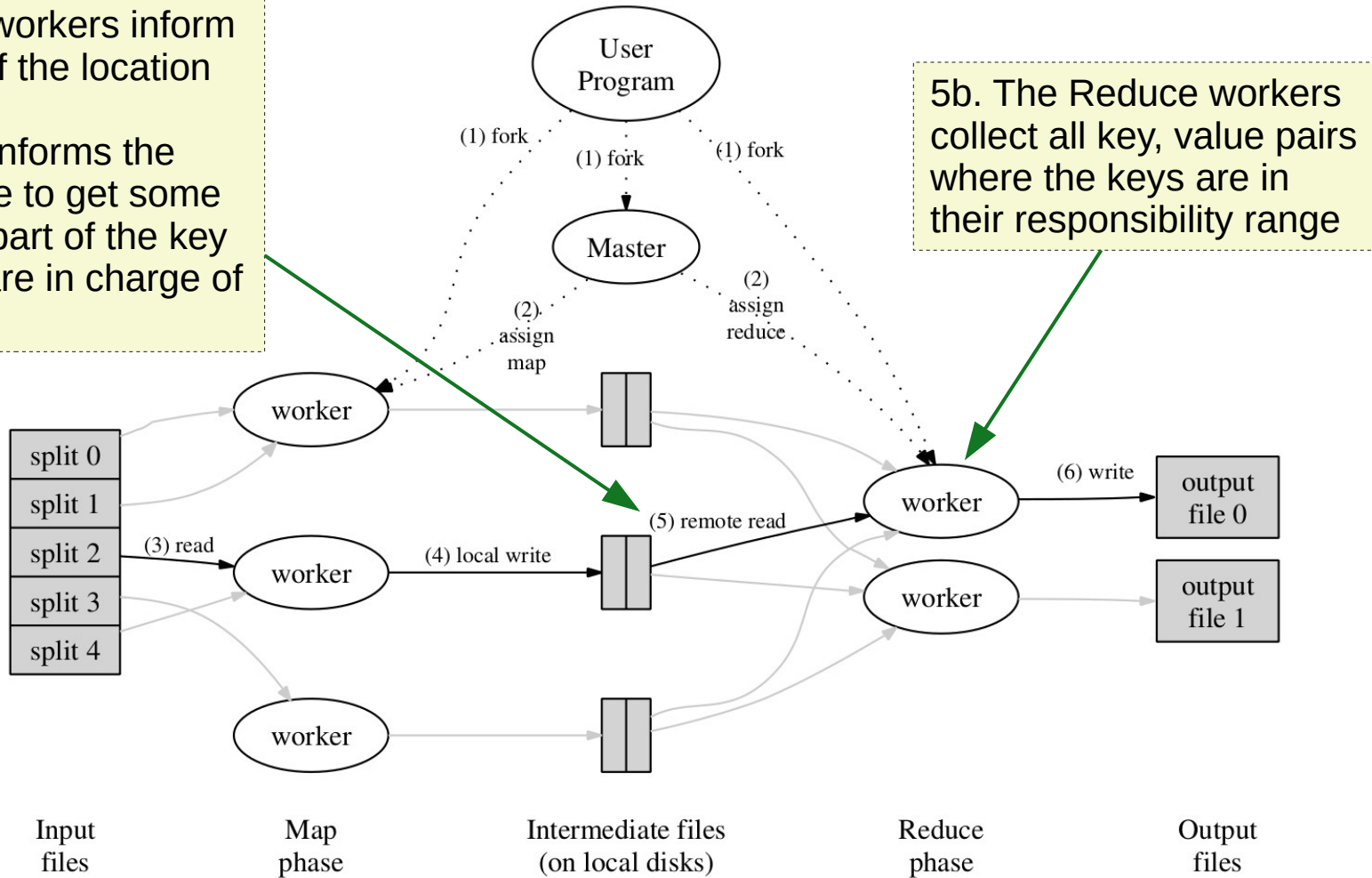
# Execution overview



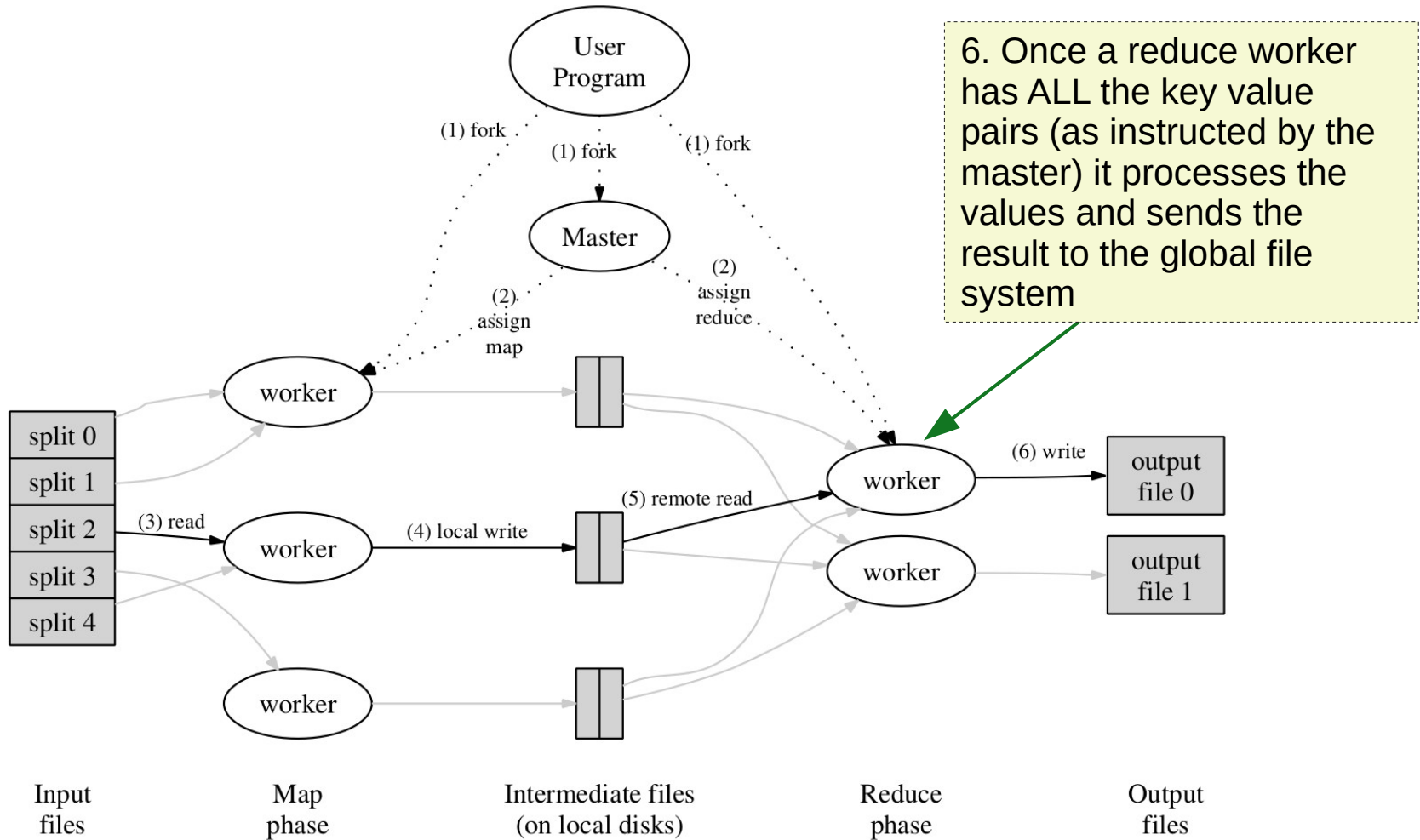
# Execution overview

5. The Map workers inform the Master of the location of fresh data

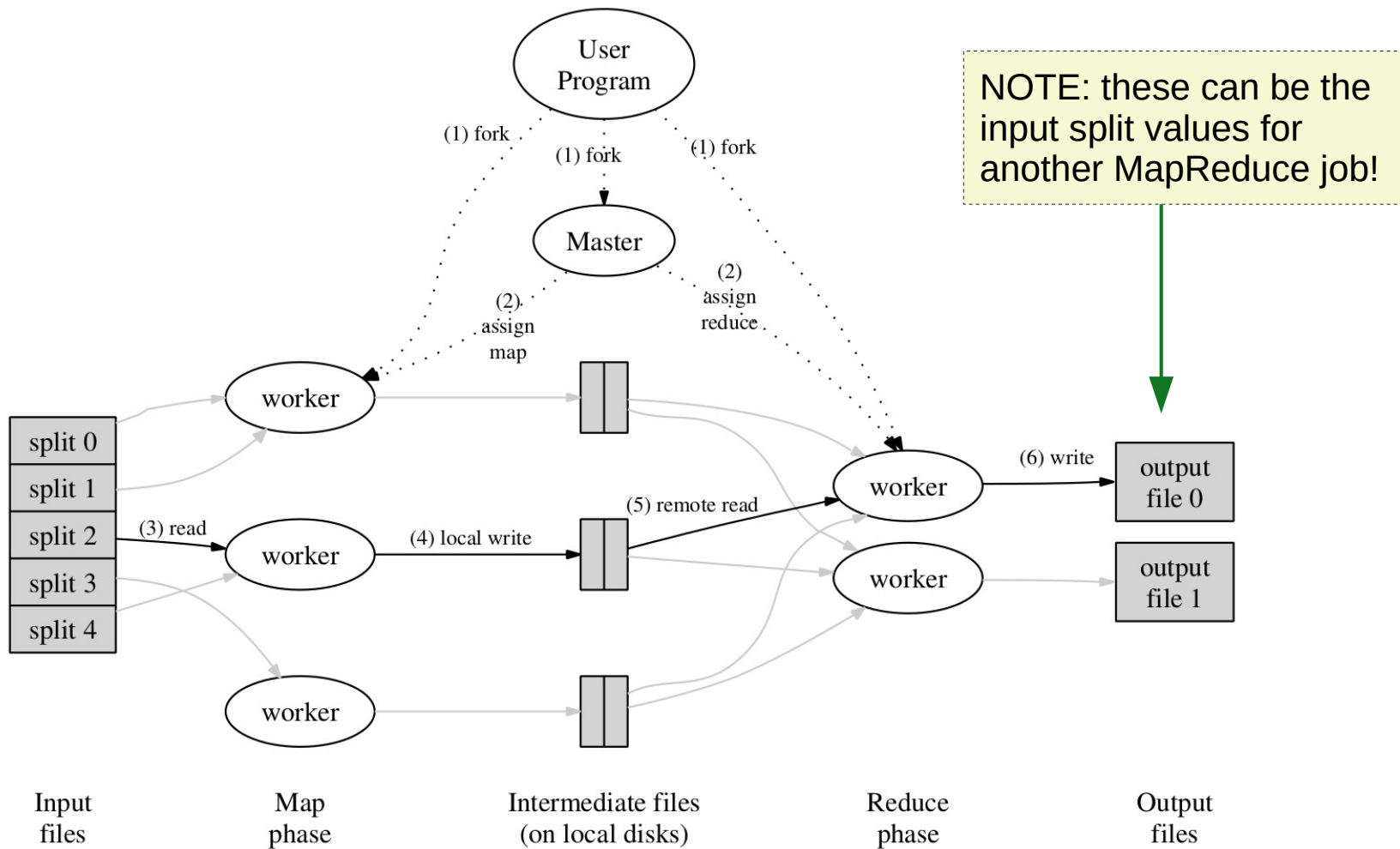
The Master informs the worker where to get some data for the part of the key space they are in charge of (locally)



# Execution overview



# Execution overview



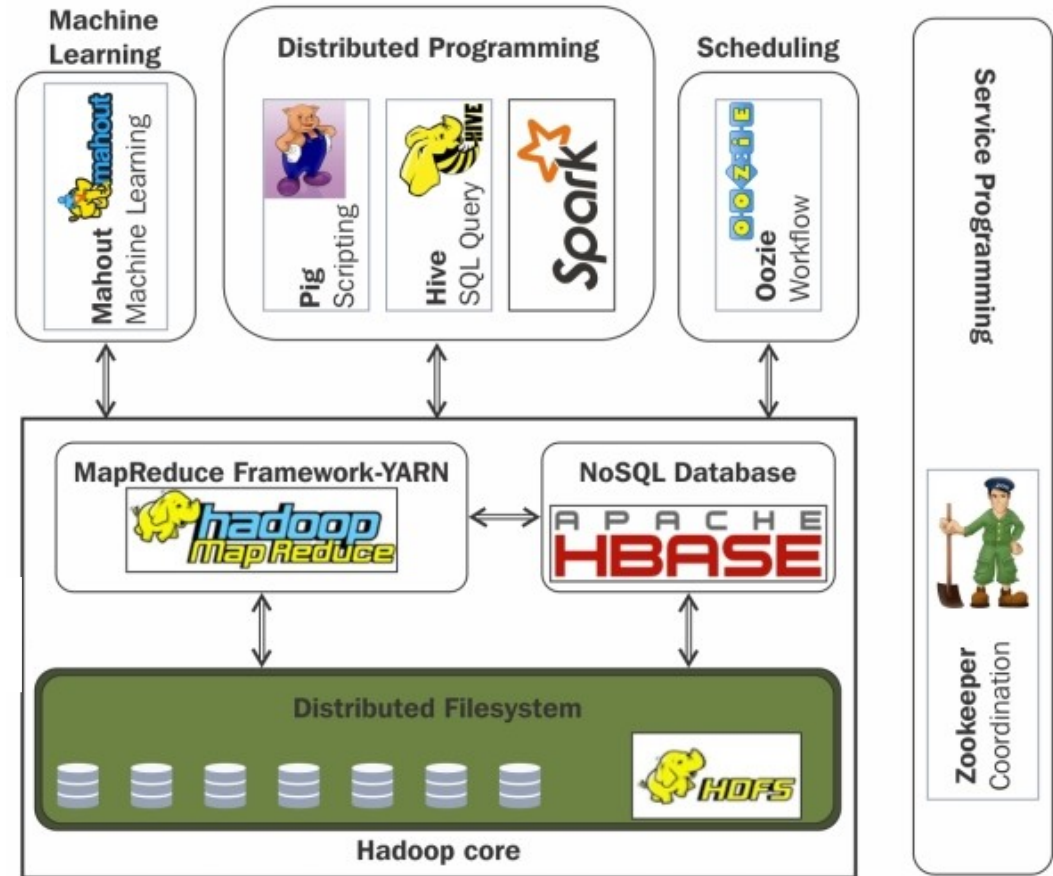
# Implementation challenges

- Failing workers
  - Must detect and re-assign to another worker
  - Partitioning helps: partially processed data simply discarded
- Slow workers
  - Alive but slow worker can slow down entire job execution
  - Monitor work, assign fast workers work done by slow ones
    - Redundant computations, keep only first finisher
- Failing master
  - Original Map/Reduce: snapshot and retry
  - Today: use a coordination kernel!

# **Map/Reduce frameworks**

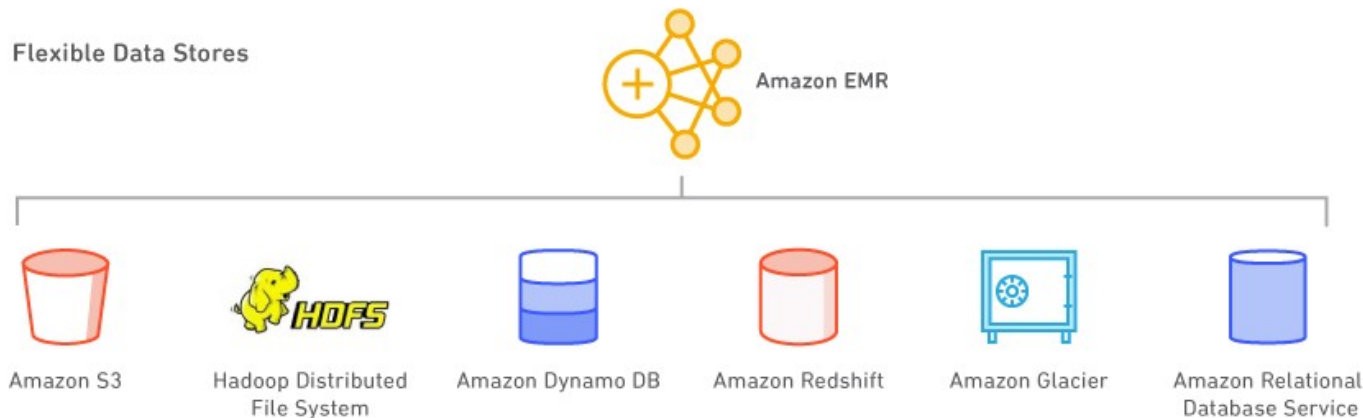
# Map/Reduce implementations

- Original Google Map/Reduce proprietary
- **Hadoop**: open source implementation of MapReduce and associated tools
  - File systems, workflow management



# Map/Reduce as a service

- Amazon Elastic Map Reduce: pre-configured versions of BigData frameworks including Hadoop, etc.
  - Can feed in data from S3, Dynamo DB, etc.



- Similar offers at other providers



# Map/Reduce in NoSQL databases

- Hadoop is a heavy machinery
  - Useful for processing vast amounts of unstructured data
  - But Cloud application data is stored in NoSQL databases
- Map/Reduce calls supported by most NoSQL databases
  - Computations performed directly on top of the data

# mySQL

# MongoDB

SELECT

```
Dim1, Dim2,
SUM(Measure1) AS MSum,
COUNT(*) AS RecordCount,
AVG(Measure2) AS MAvg,
MIN(Measure1) AS MMin
MAX(CASE
  WHEN Measure2 < 100
  THEN Measure2
END) AS MMax
FROM DenormAggTable
WHERE (Filter1 IN ('A','B'))
AND (Filter2 = 'C')
AND (Filter3 > 123)
GROUP BY Dim1, Dim2
HAVING (MMin > 0)
ORDER BY RecordCount DESC
LIMIT 4, 8
```

- ① Grouped dimension columns are pulled out as keys in the map function, reducing the size of the working set.
- ② Measures must be manually aggregated.
- ③ Aggregates depending on record counts must wait until finalization.
- ④ Measures can use procedural logic.
- ⑤ Filters have an ORM/ActiveRecord-looking style.
- ⑥ Aggregate filtering must be applied to the result set, not in the map/reduce.
- ⑦ Ascending: 1; Descending: -1

```
db.runCommand({
  mapreduce: "DenormAggCollection",
  query: {
    filter1: { '$in': [ 'A', 'B' ] },
    filter2: 'C',
    filter3: { '$gt': 123 }
  },
  map: function() { emit(
    { d1: this.Dim1, d2: this.Dim2 },
    { msum: this.measure1, recs: 1, mmin: this.measure1,
      mmax: this.measure2 < 100 ? this.measure2 : 0 }
  );},
  reduce: function(key, vals) {
    var ret = { msum: 0, recs: 0, mmin: 0, mmax: 0 };
    for(var i = 0; i < vals.length; i++) {
      ret.msum += vals[i].msum;
      ret.recs += vals[i].recs;
      if(vals[i].mmin < ret.mmin) ret.mmin = vals[i].mmin;
      if((vals[i].mmax < 100) && (vals[i].mmax > ret.mmax))
        ret.mmax = vals[i].mmax;
    }
    return ret;
  },
  finalize: function(key, val) {
    val.mavg = val.msum / val.recs;
    return val;
  },
  out: 'result1',
  verbose: true
});
db.result1.
  find({ mmin: { '$gt': 0 } }).
  sort({ recs: -1 }).
  skip(4).
  limit(8);
```

Revision 4, Created 2010-03-06  
Rick Osborne, rickosborne.org

# MongoDB - Reduce function

- The reduce function should not affect the outside system
- MongoDB will not call the reduce function for a key that has only a single value
- MongoDB can invoke the reduce function more than once for the same key
  - In this case, the previous output from the reduce function for that key will become one of the input values to the next reduce function invocation for that key
  - The type of the return object must be identical to the type of the value emitted by the map function
  - The reduce function must be **associative**
    - `reduce(key, [ C, reduce(key, [ A, B ]) ]) == reduce( key, [ C, A, B ])`
  - The reduce function must be **idempotent**
    - `reduce( key, [ reduce(key, valuesArray) ]) == reduce( key, valuesArray )`
  - The reduce function must be **commutative**
    - `reduce( key, [ A, B ]) == reduce( key, [ B, A ])`

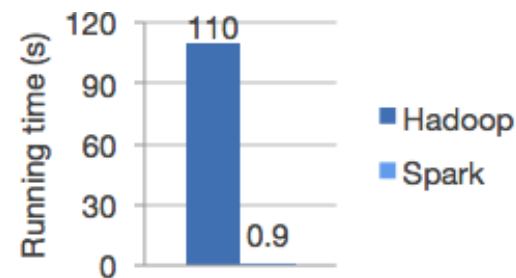
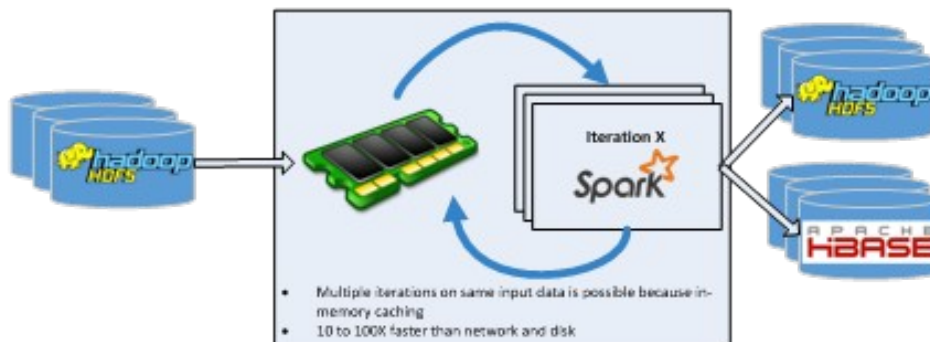
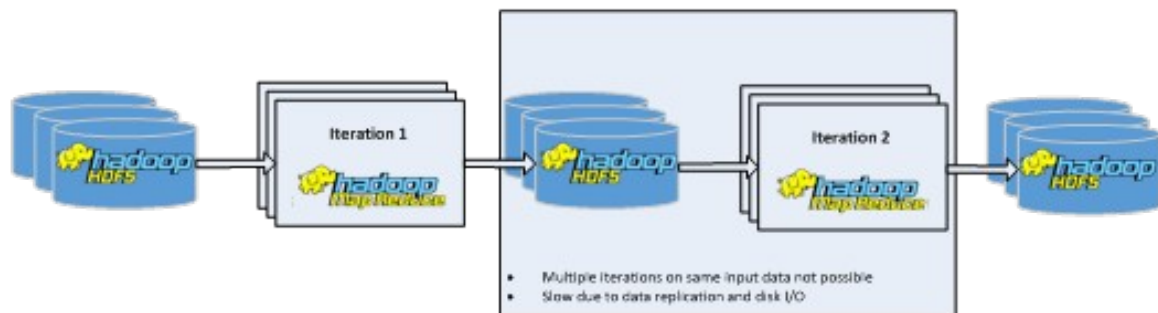
# Beyond Map/Reduce

- Many algorithms require to iterate until stopping condition
  - K-Means: no point changes to new cluster
  - PageRank: ranking of webpage based on incoming links
  - Many *machine learning* algorithms
- And new usages
  - Algorithms applied on data streams
  - Interactive queries
- Map/Reduce intermediate storage (on distributed FS) is not well suited for these use cases
  - New abstractions have emerged to keep intermediary results in memory



- Writing to disk is slow and costly
  - Especially if we need to re-read immediately!
- Can we keep data in-memory instead?
- Apache Spark introduces Resilient Distributed Datasets
  - Shared memory abstraction
  - Programming model remains the same as Map/Reduce
  - Avoids re-reading from disk systematically
  - Builds different querying models on top of core engine
  - Can use pipeline of queries under different models

# Spark VS Map/Reduce



Logistic regression in Hadoop and Spark

# Conclusion

- Cloud-based web applications collect vast amounts of information about their clients, servers, infrastructure
  - Big Data processing
    - Necessary for Cloud scale
    - Enables new usages and applications
  - Process it in the background to derive information useful for business
- Map/Reduce: game changer for large-scale Big Data processing
  - Pioneered the use of a functional approach
  - Handles large-scale orchestration and fault tolerance
  - Specific (and simple) programming model
- Apache Spark improves/generalizes Map/Reduce principle
  - Support for iterative and interaction computations

# References

- MapReduce: Simplified Data Processing on Large Clusters
  - Jeffrey Dean and Sanjay Ghemawat
    - OSDI'04: Sixth Symposium on Operating System Design and Implementation
- The Anatomy of a Large-Scale Hypertextual Web Search Engine
  - Sergey Brin, Lawrence Page
    - Computer Networks 30 - 1998
- Resilient Distributed Datasets: A Fault-Tolerant Abstraction for In-Memory Cluster Computing
  - Matei Zaharia, Mosharaf Chowdhury, Tathagata Das, Ankur Dave, Justin Ma, Murphy McCauley, Michael J. Franklin, Scott Shenker, Ion Stoica
    - NSDI'12
- Scaling Spark in the Real World: Performance and Usability
  - Michael Armbrust, Tathagata Das, Aaron Davidson, Ali Ghodsi, Andrew Or, Josh Rosen, Ion Stoica, Patrick Wendell, Reynold Xin, Matei Zaharia
    - Proceedings of the VLDB Endowment - 2015