

# Big Data Systems, Paradigms, Algorithms

Dario Colazzo

Professeur à l'Université Paris Dauphine

*Responsable du pôle Data Science*

Professeur chargé de cours à l'École Polytechnique

# Plan

- Today and tomorrow: MapReduce Hadoop and Spark
  - focus on programming
- Last class (4 hours) : Hive

# Connecting to the cluster

- Please send me an email ([dario.colazzo@dauphine.fr](mailto:dario.colazzo@dauphine.fr)) with subject 'cluster DS'
- I will replay by attaching a private key that you will store somewhere (for instance in the home)
- You will be given as login userXY (for instance user80)
- To connect

```
ssh -i <path of private key> -p 993 userXY@www.lamsade.dauphine.fr
```

# Let's run our first MR job

- WordCount : download the code, or use your own code

```
wget https://www.dropbox.com/s/3yuxnpwkv2wqpe4/mapper.py
```

```
wget https://www.dropbox.com/s/kjsugh3do9fz0jx/reducer.py
```

- Put the code in a directory of your choice, do not forget

```
chmod +x *.py
```

- Download data

```
https://www.dropbox.com/s/yq5e4gg4ztr784h/20417-8.txt
```

# Let's run our first MR job

- Prepare you HDFS directory

```
hadoop fs -mkdir /user/user81/input
```

```
hadoop fs -mkdir /user/user81/output
```

- Put the txt file into the input HDFS directory

```
hadoop fs -put local.path.to.file /user/user81/input
```

- Almost ready

# Let's run our first MR job

- Launch the job (prepare the command in a txt file before, and then copy-paste it)

```
hadoop jar /usr/local/hadoop/share/hadoop/tools/lib/hadoop-streaming-2.7.0.jar \
```

```
-input HDFS path of your input dir \
```

```
-output HDFS path of your output dir \
```

```
-file your local path to mapper.py \
```

```
-mapper your local path to mapper.py\
```

```
-file your local path to reducer.py\
```

```
-reducer your local path to reducer.py
```

# Let's run our first MR job

- When the job is completed
  - explore the output by means of the HDFS -cat command
  - also read the 'counters' at the end of the MR log.

# Add a combiner and increase the number of Reducers

```
hadoop jar /usr/local/hadoop/share/hadoop/tools/lib/hadoop-streaming-2.7.0.jar \  
  
-input HDFS path of your input dir \  
  
-output HDFS path of your output dir \  
  
-file your local path to mapper.py \  
  
-mapper your local path to mapper.py\  
  
-file your local path to reducer.py\  
  
-reducer your local path to reducer.py \  
  
-combiner your local path to reducer.py \  
  
-jobconf mapred.reduce.tasks=3
```



# Spark

## Part 1:

### *Overview and programming with Resilient Distributed Datasets*

Dario Colazzo

# Motivation

- MapReduce greatly simplified big data analysis on large, unreliable clusters.
- But as soon as it got popular, users wanted more:
  - Iterative jobs, e.g., machine learning algorithms
  - Interactive analytics

# Motivation

- Both iterative and interactive queries need one thing that MapReduce lacks

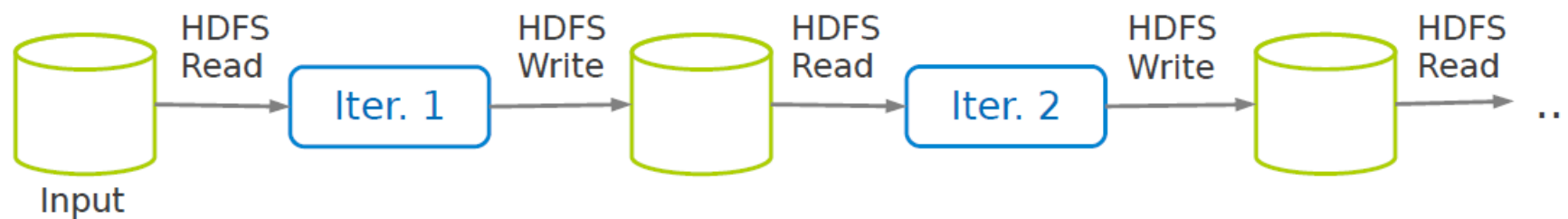
Efficient primitives for data sharing.

- In MapReduce, the only way to share data across processing step is stable storage (disk)
- Replication also makes the system slow, but it is necessary for fault tolerance.

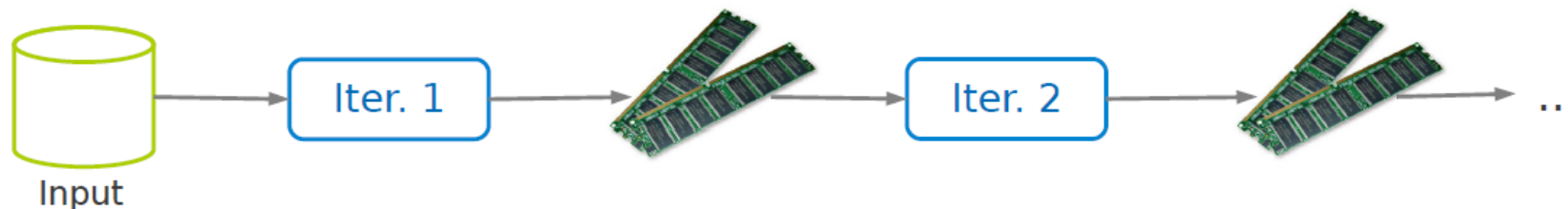
# Solution

In memory data processing and sharing

Hadoop  
MapReduce

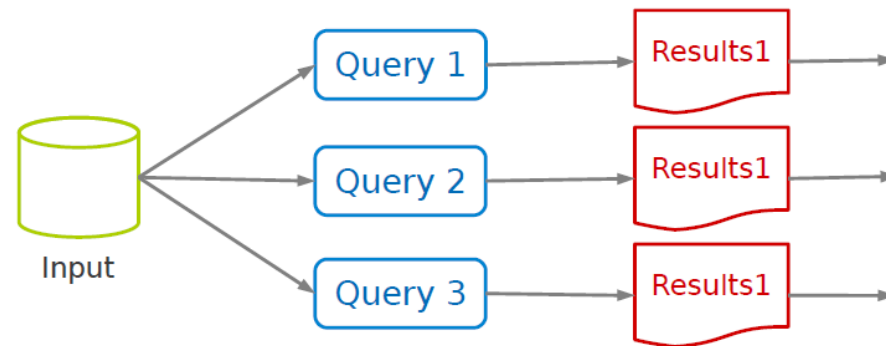


Hadoop  
Spark

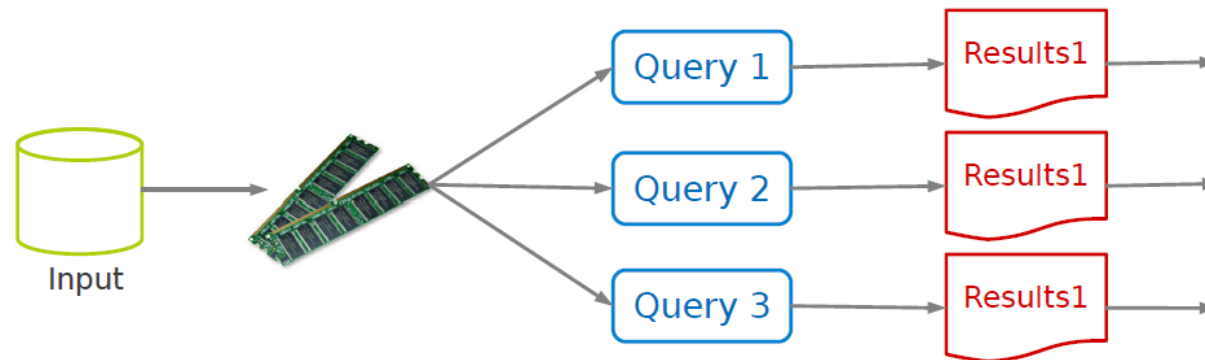


# Sharing

Hadoop  
MapReduce



Hadoop  
Spark

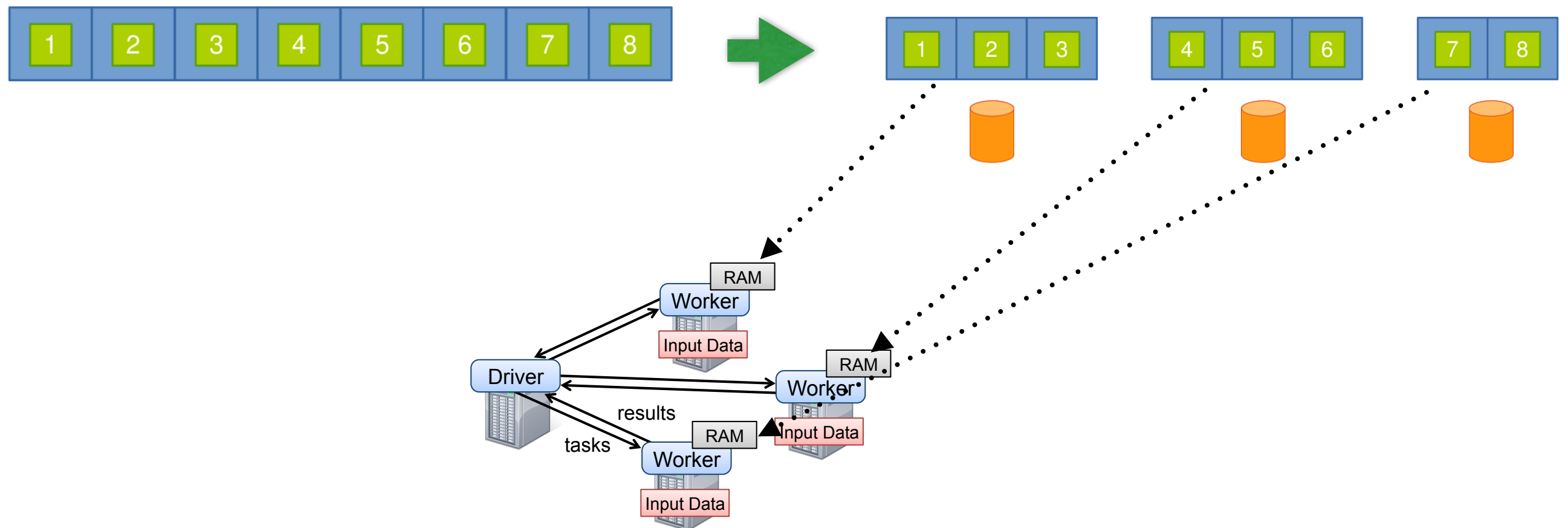


# Challenge

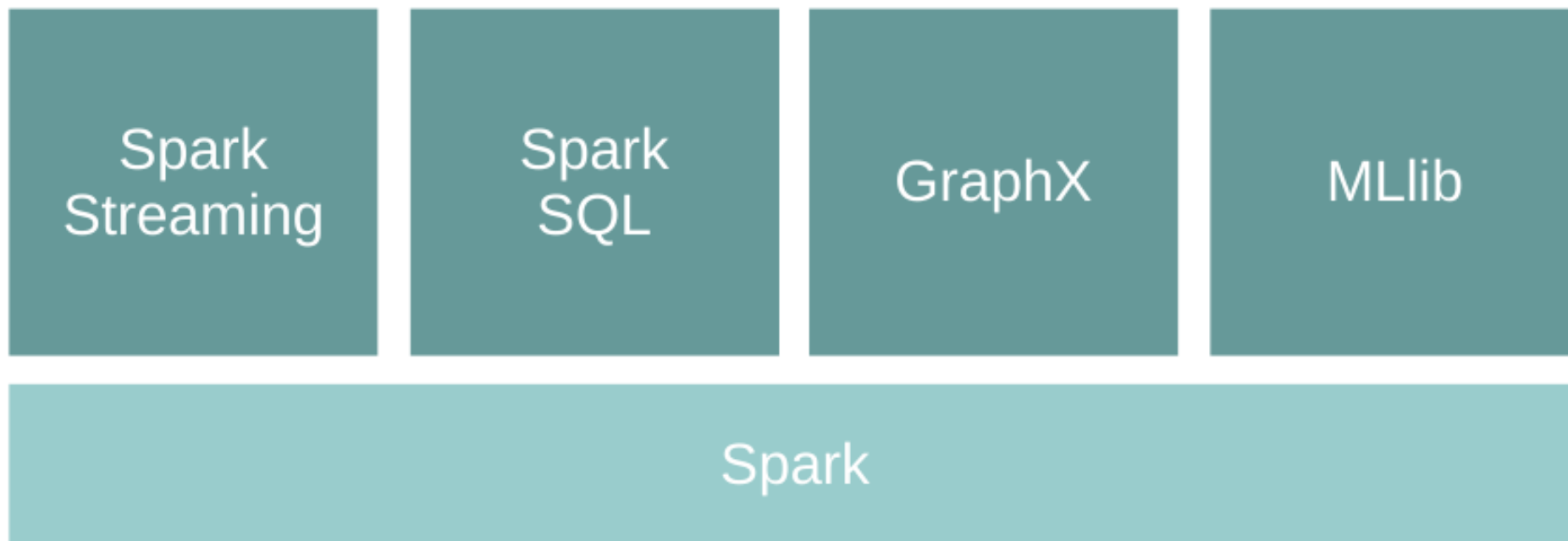
- How to design a distributed memory abstraction that is both fault **tolerant** and **efficient**?
- Solution: Resilient Distributed Datasets (RDD)
  - A **distributed** main-memory abstraction.
  - **Immutable** collections of objects spread across a cluster.
  - **Lineage** among RDDs to enable their re-evaluation in case of cluster node failures

# Resilient Distributed Datasets (RDDs)

- An RDD is a collection which divided into a number of **partitions**, which can be independently processed.



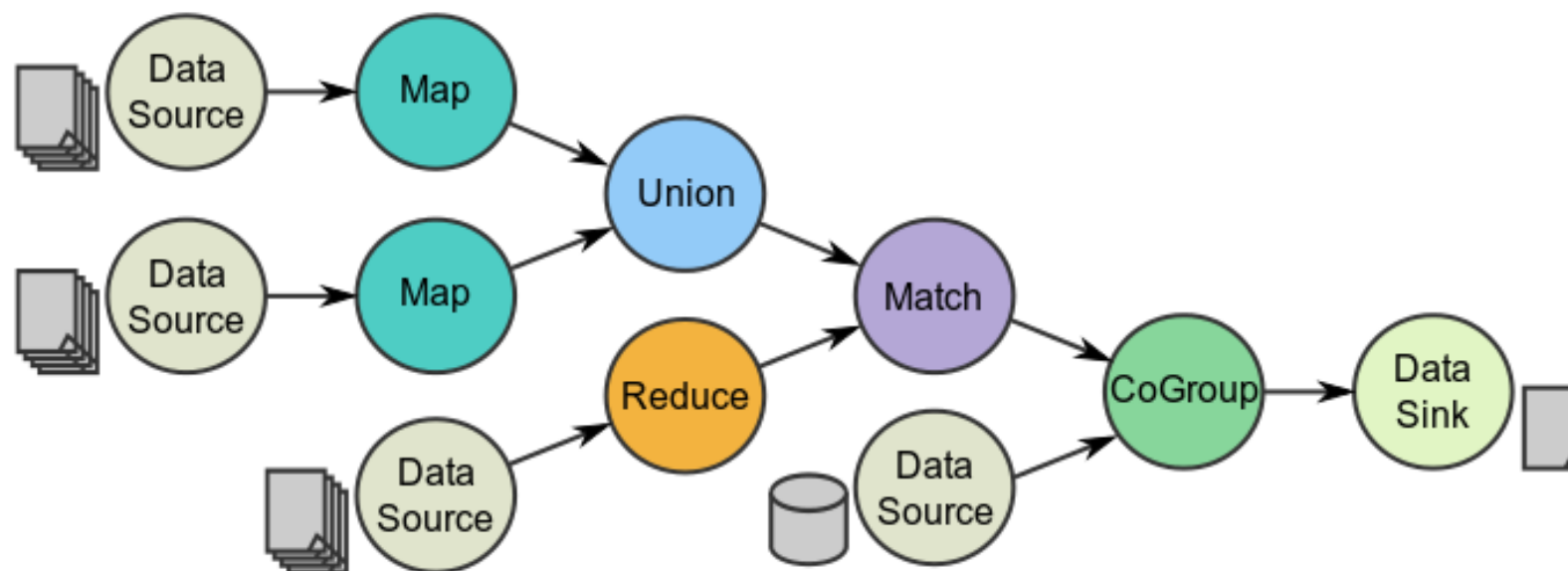
# Spark Processing engine





# Programming model

- A **data flow** is composed of any number of **data sources** and **data sinks** by connecting their inputs and outputs by means of **data operators**.

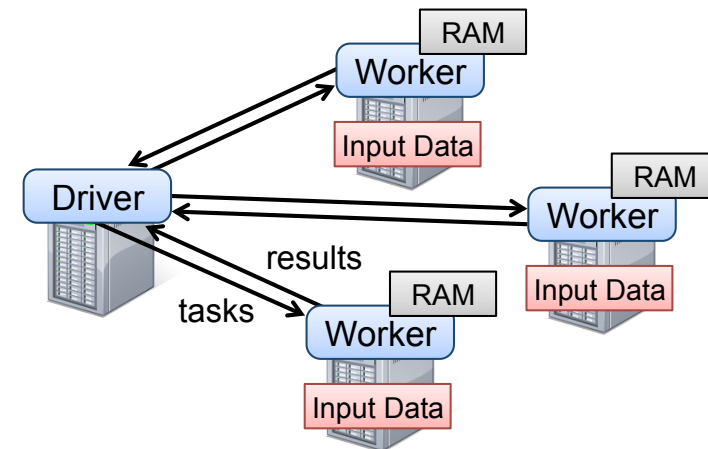


# Programming model

- Based on parallelizable operators.
- Parallelizable operators are **higher-order** functions that execute **user-defined** functions in parallel, on each partition of an RDD.
- There are two types of RDD operators : **transformations** and **actions**.

# Programming model

- **Transformations** : **lazy** operators that create **new** RDDs.
- **Actions** : launch a **computation** and return a **value** to the program driver or write data to the external storage



- Implemented in Scala:
  - a strongly and statically typed functional-OO language
  - compiled and run over the JVM
  - designed at EPFL (Switzerland).
- Java and Python can be used too for Spark programming.

# Example (1/2)

- Suppose that a web service is experiencing errors and an operator wants to search terabytes of logs in the Hadoop filesystem (HDFS) to find the cause.

- Here is Spark code in Scala (*but we will switch soon to Python*)

```
lines = spark.textFile("hdfs://...")  
errors = lines.filter(_.startsWith("ERROR"))  
errors.persist()
```

- Actions can be used to count errors:

```
errors.count()
```

- Or counting errors mentioning MySQL:

```
// Count errors mentioning MySQL:  
errors.filter(_.contains("MySQL")).count()
```

# Example (1/2)

- Suppose that a web service is experiencing errors and an operator wants to search terabytes of logs in the Hadoop filesystem (HDFS) to find the cause.

```
lines = spark.textFile("hdfs://...")
errors = lines.filter(_.startsWith("ERROR"))
errors.persist()
```

- Actions can be used to count errors:

```
errors.count()
```

- Or counting errors mentioning MySQL:

```
// Count errors mentioning MySQL:
errors.filter(_.contains("MySQL")).count()
```

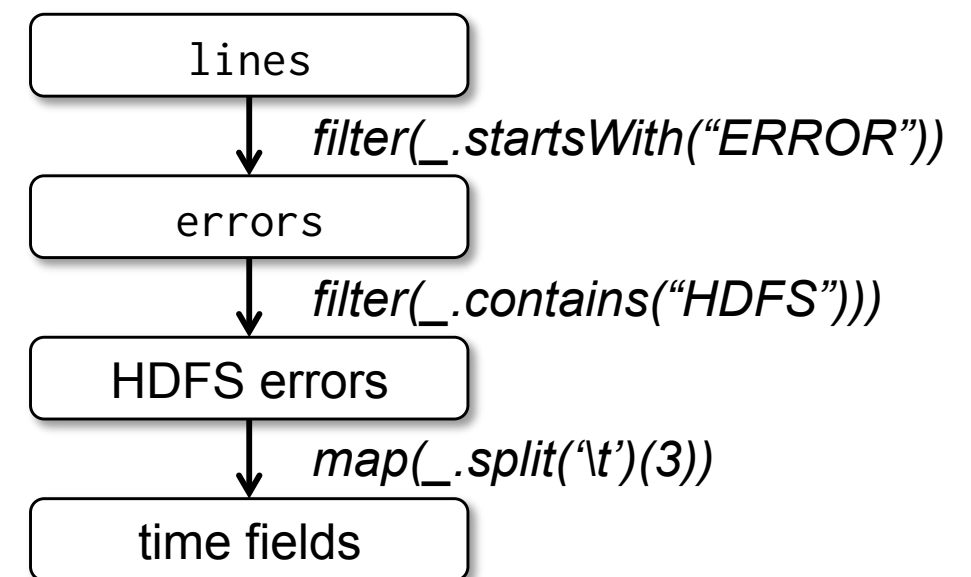
**lines** is not loaded in memory only **errors** is (simple static analysis)

*lazy evaluation:* **errors** is actually calculated and put in memory when the **count()** action is evaluated

# Fault tolerance via *lineage*

```
// Count errors mentioning MySQL:
errors.filter(_._contains("MySQL")).count()

// Return the time fields of errors mentioning
// HDFS as an array (assuming time is field
// number 3 in a tab-separated format):
errors.filter(_._contains("HDFS"))
    .map(_._split('\t')(3))
    .collect()
```



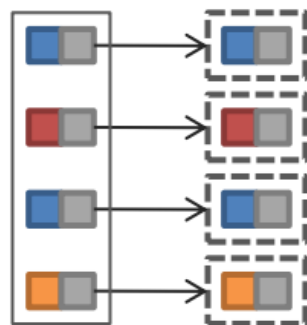
the lineage graph enables RDD re-evaluation in  
case of failure

# RDD transformations and actions

<b>Transformations</b>	$map(f : T \Rightarrow U) : RDD[T] \Rightarrow RDD[U]$ $filter(f : T \Rightarrow Bool) : RDD[T] \Rightarrow RDD[T]$ $flatMap(f : T \Rightarrow Seq[U]) : RDD[T] \Rightarrow RDD[U]$ $sample(fraction : Float) : RDD[T] \Rightarrow RDD[T]$ (Deterministic sampling) $groupByKey() : RDD[(K, V)] \Rightarrow RDD[(K, Seq[V])]$ $reduceByKey(f : (V, V) \Rightarrow V) : RDD[(K, V)] \Rightarrow RDD[(K, V)]$ $union() : (RDD[T], RDD[T]) \Rightarrow RDD[T]$ $join() : (RDD[(K, V)], RDD[(K, W)]) \Rightarrow RDD[(K, (V, W))]$ $cogroup() : (RDD[(K, V)], RDD[(K, W)]) \Rightarrow RDD[(K, (Seq[V], Seq[W]))]$ $crossProduct() : (RDD[T], RDD[U]) \Rightarrow RDD[(T, U)]$ $mapValues(f : V \Rightarrow W) : RDD[(K, V)] \Rightarrow RDD[(K, W)]$ (Preserves partitioning) $sort(c : Comparator[K]) : RDD[(K, V)] \Rightarrow RDD[(K, V)]$ $partitionBy(p : Partitioner[K]) : RDD[(K, V)] \Rightarrow RDD[(K, V)]$
<b>Actions</b>	$count() : RDD[T] \Rightarrow Long$ $collect() : RDD[T] \Rightarrow Seq[T]$ $reduce(f : (T, T) \Rightarrow T) : RDD[T] \Rightarrow T$ $lookup(k : K) : RDD[(K, V)] \Rightarrow Seq[V]$ (On hash/range partitioned RDDs) $save(path : String) : \text{Outputs RDD to a storage system, e.g., HDFS}$

# RDD transformations : Map

- All pairs are independently processed



```
# passing each RDD element through a function
nums = sc.parallelize([1,2,3])
squares = nums.map(lambda x: x * x)
```

```
# selecting elements making a boolean function returning true
even = squares.filter(lambda x : x % 2 == 0)
```

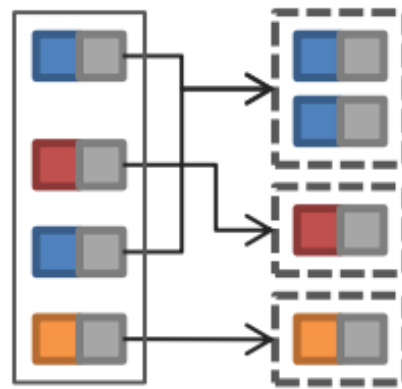
```
# map + flattening
m = nums.map(lambda x: range(x))
# [[0], [0, 1], [0, 1, 2]]
fm = nums.flatMap(lambda x: range(x))
# [0, 0, 1, 0, 1, 2]
```



# RDD transformations :

## Reduce

- Pairs with identical key are grouped
- Each group is independently processed



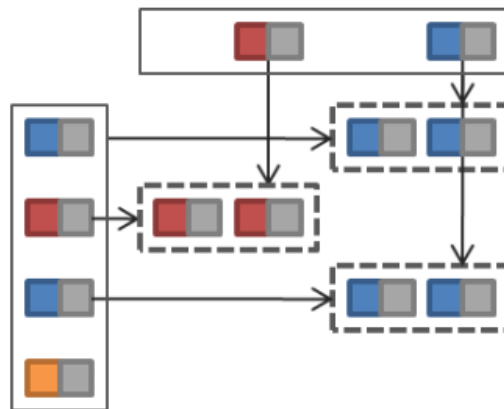
```
pets = sc.parallelize([("cat", 1), ("dog", 1), ("cat", 2), ("dog", 3) ])
```

```
pets.reduceByKey(lambda x, y : x +y)  
# [('dog', 4), ('cat', 3)]
```

```
pets.groupByKey()  
pets.groupByKey().map(lambda x : (x[0], list(x[1])))  
# [('dog', [1,3]), ('cat', [1, 2])]
```

# RDD transformations : Join

- Equi-join on the key



```
visits = sc.parallelize( [("h", "1.2.3.4"), ("a", "3.4.5.6"), ("h", "1.3.3.1")] )
```

```
pageNames = sc.parallelize( [("h", "Home"), ("a", "About")] )
```

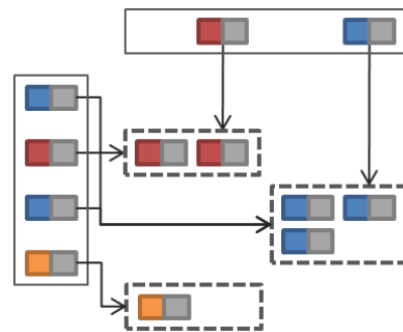
```
visits.join(pageNames)
```

```
# [('a', ('3.4.5.6', 'About')), ('h', ('1.2.3.4', 'Home')),  
   ('h', ('1.3.3.1', 'Home'))]
```

# RDD transformations :

## CoGroup

- Groups each input on key
- Groups with identical keys are processed together

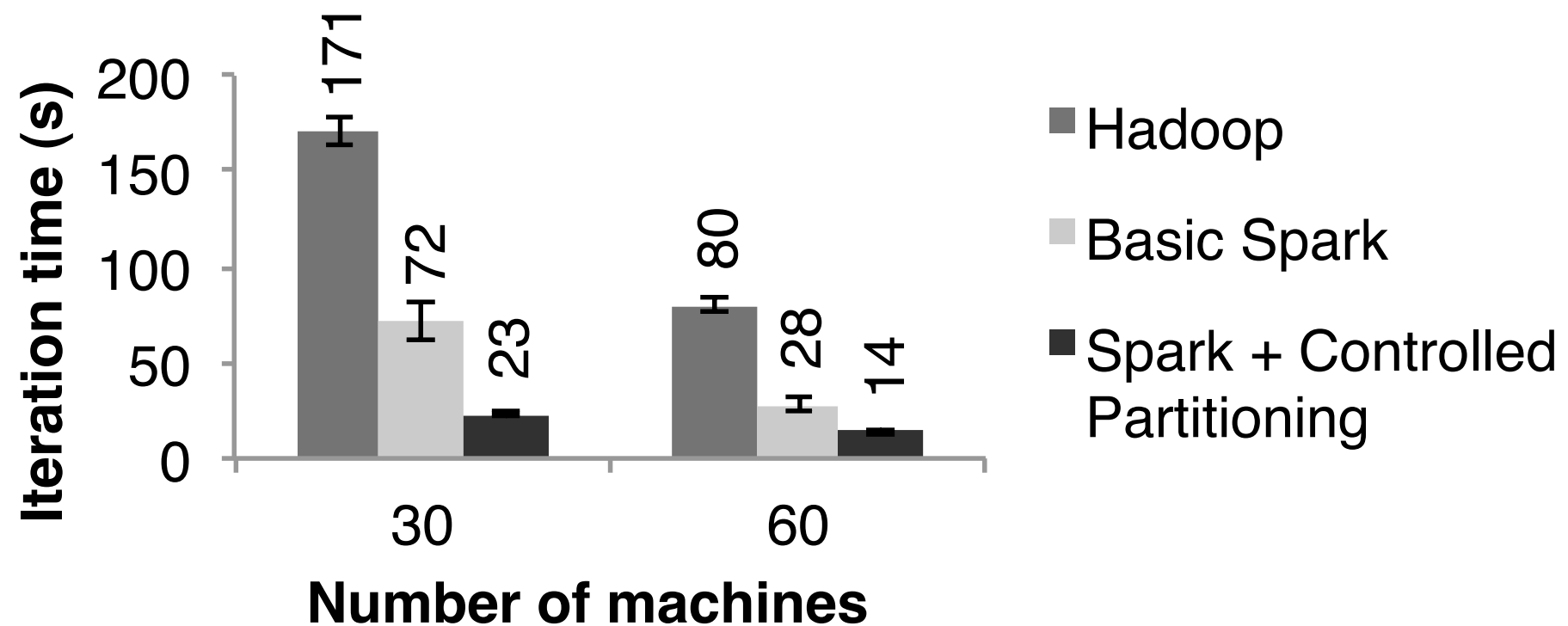


```
visits = sc.parallelize([("h", "1.2.3.4"), ("a", "3.4.5.6"), ("h", "1.3.3.1")])
pageNames = sc.parallelize([("h", "Home"), ("a", "About"), ("o", "Other")])
visits.cogroup(pageNames)

visits.cogroup(pageNames).map(lambda x : (x[0], (list(x[1][0]), list(x[1][1]))))

# [('a', ([ '3.4.5.6' ], [ 'About' ])), ('h', ([ '1.2.3.4', '1.3.3.1' ], [ 'Home' ])),
#  ('o', ([ ], [ 'Other' ]))]
```

# Some experiments on PageRank



Borrowed from *Resilient Distributed Datasets: A Fault-Tolerant Abstraction for In-Memory Cluster Computing*.  
Matei Zaharia et al, NSDI 2012.

# Remarks

- MapReduce makes important abstraction step that greatly helps rapid development of efficient and robust Big Data data flows.
- But:
  - we still need some ‘acking’ to ensure good performances
  - problems with iterative analyses
  - MapReduce programming is not easy
- Spark overcomes these limitations in a large extent, at the cost of more RAM needed.
- Makes a one more step towards ‘The data center is the computer’ scenario.

# Lab Session

- Connect to the cluster and launch the Python Spark shell, by means of `>pyspark`
- There is a predefined `sc` object that allows you to create RDDs
- Previous examples can be copy-pasted and eventually changed.
- Exercice 1 : identify what is the path leading to the directory containing the Hadoop installation, by means of `>echo $HADOOP_HOME`
  - the directory contains a LICENCE textual file
  - create an RDD from this file by means of

```
>>>t= sc.textFile("file: ...path to ...LICENSE.txt")
```
  - write and run Spark code to perform word counting on the LICENSE.txt file
- Exercices 2 : write and run Spark code to calculate the average of integers values associated to pets (use the previously seen RDD for pets)