

# Technology and Application of Big Data

Qing LIAO(廖清)

School of Computer Science and Technology

HIT

# Course Details

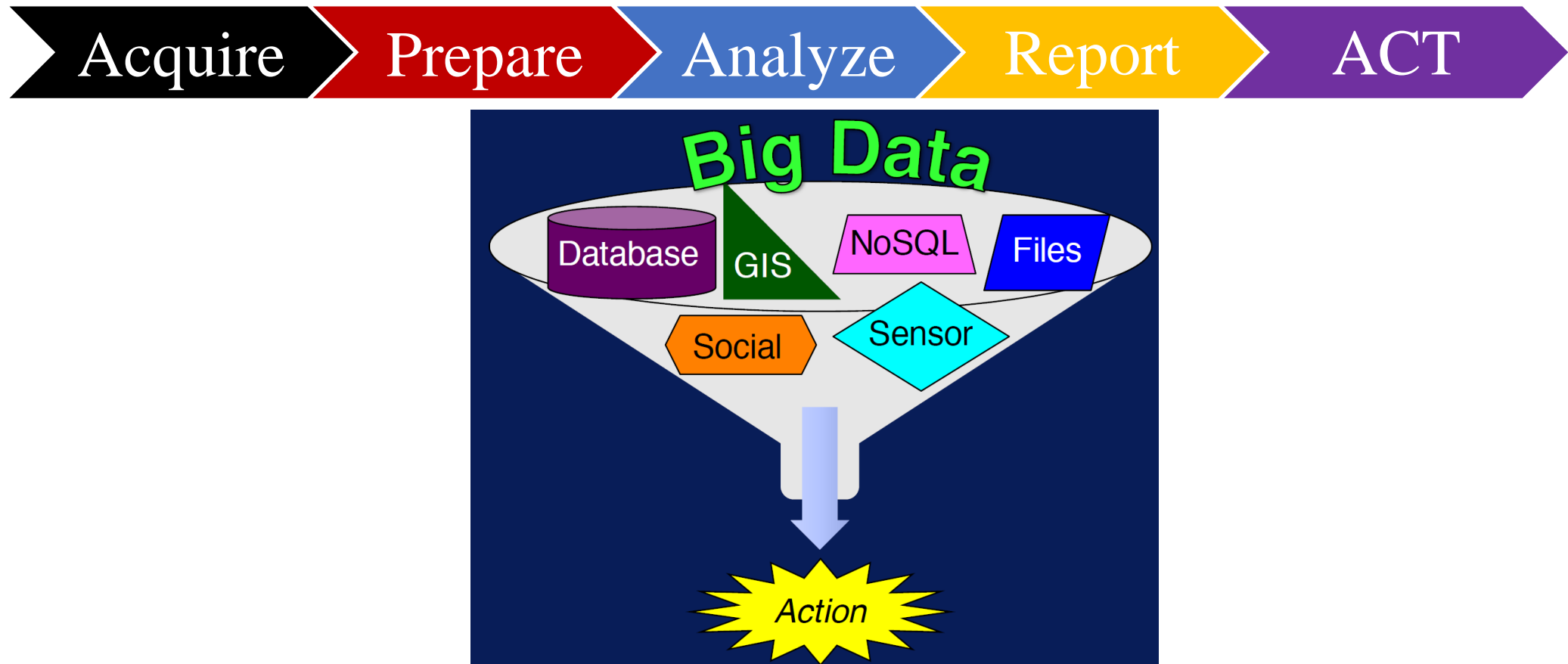
- Instructor:
  - Qing LIAO, [liaoqing@hit.edu.cn](mailto:liaoqing@hit.edu.cn)
  - Rm. 303B, Building C
  - Office hours: by appointment
- Course web site:
  - [liaoqing.me](http://liaoqing.me)
- Reference books/materials:
  - Big data courses from University of California
  - Book: BIG DATA: A Revolution That Will Transform How We Live, Work, and Think
  - Papers
- Grading Scheme:
  - Paper Report 30%
  - Final Exam 70%
- Exam:
  - 21<sup>st</sup> July(Friday), 14:00-16:00, A502

# What You Learnt: Overview

- Topics:
  - 1) Introduction of Big Data
  - 2) Characterizes of Big Data
  - 3) How to Get Value from Big Data
  - 4) Technologies of Big Data
  - 5) Applications of Big Data
- Prerequisites
  - Statistics and Probability would help
    - But not necessary
  - Machine Learning would help
    - But not necessary

# Previous Section

- How to Get Value from Big Data



# Previous Section

## ➤ Information Extraction results

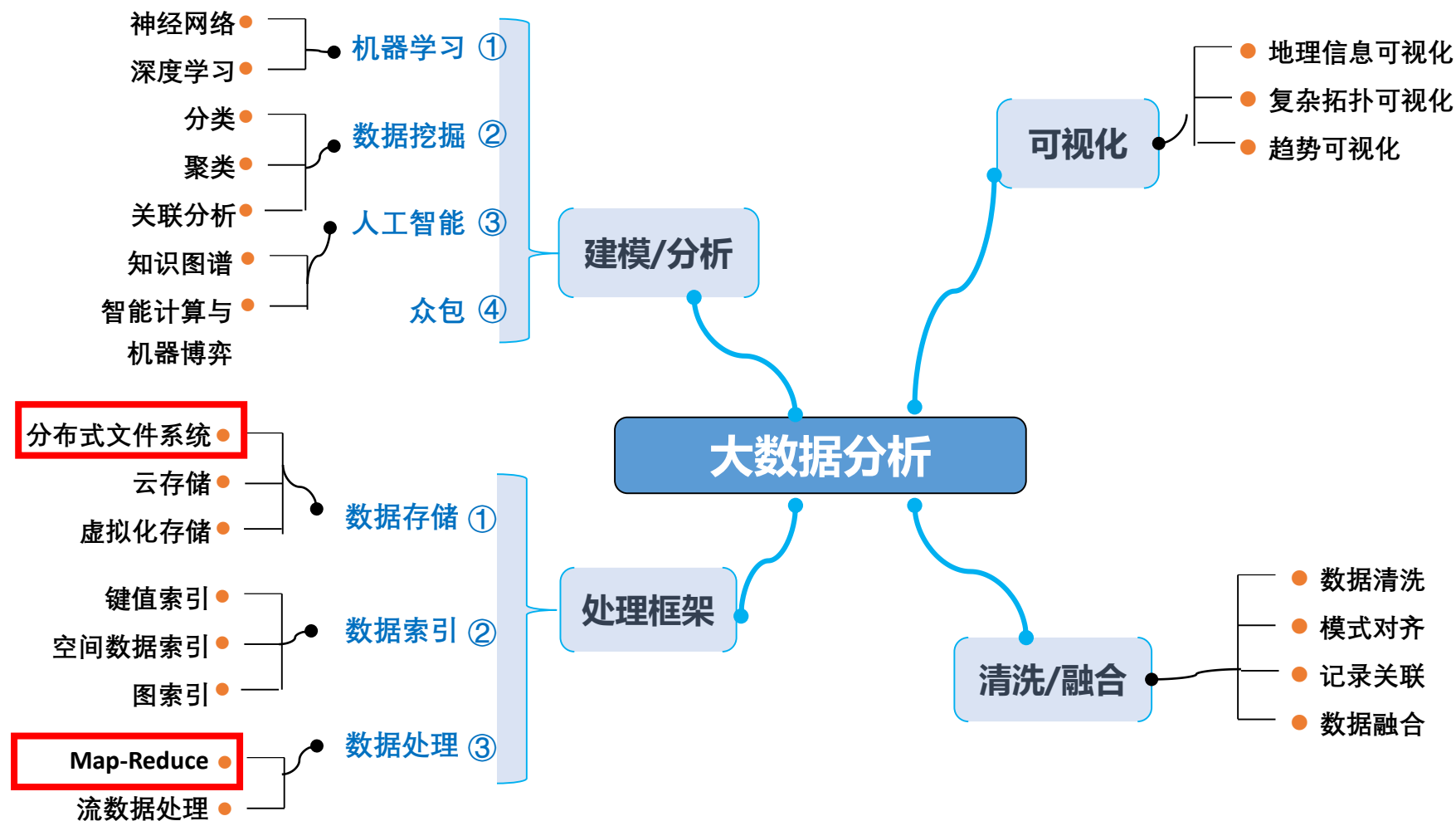


(a). An example page segment

image 1	Cabinet Organizers by Copco	9-in.	Round Turntable: White	*****	\$4.95
image 1	Cabinet Organizers by Copco	12-in.	Round Turntable: White	*****	\$7.95
image 2	Cabinet Organizers	14.75x9	Cabinet Organizer (Non-skid): White	*****	\$7.95
image 3	Cabinet Organizers	22x6	Cookware Lid Rack	*****	\$19.95

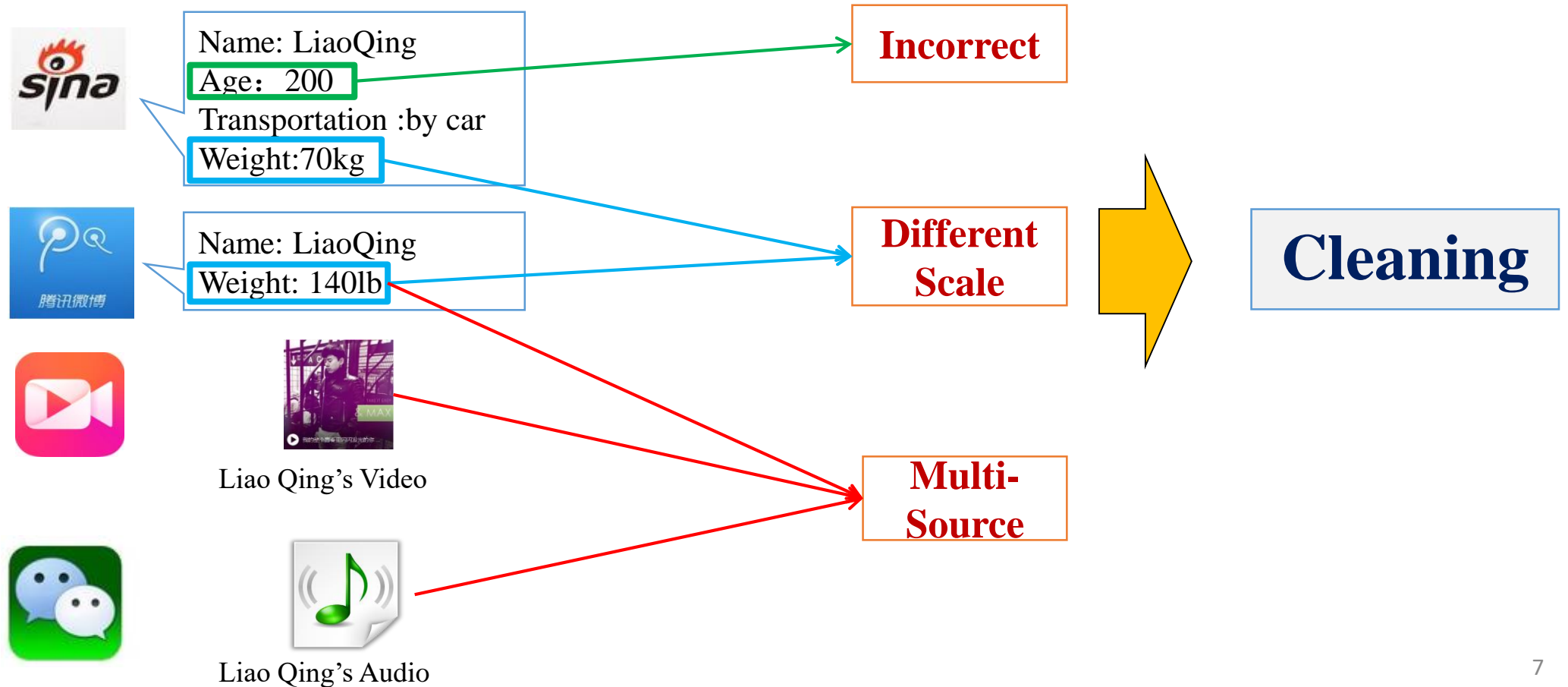
(b). Extraction results

# Technologies of Big Data



# Previous Section

- Step 2: Prepare Data - Pre-process



# Previous Section

- Step 2: Prepare Data - Pre-process

Name	Qing LIAO	✓
Gender	Female	✓
Birthplace		Empty
Birthday		Empty
ID #	220103196504032526	✓
Phone #	1502036	Anomaly
Email	xyz12596@126.com	✓
Resident	Hunan, Fuzhou	Anomaly
Height	175cm	✓
Weight	80kg	✓
Hobbies	Play Game	✓

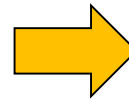
ID Card



Birthplace	Jilin, Changchun
Birthday	1965/04/03



Cleaning



Name	Qing LIAO
Gender	Female
Birthplace	Jilin, Changchun
Birthday	1965/04/03
ID #	220103196504032526
Email	xyz12596@126.com
Height	175cm
Weight	80kg
Hobbies	Play Game

Valid data



# Previous Section

- Step 2: Prepare Data - Pre-process

Data Source 1	( <b>Name</b> , <b>family phone #</b> , <b>home address</b> , <b>office phone #</b> , <b>office address</b> )
Data Source 2	( <b>Family name</b> , <b>given name</b> , <b>nick name</b> , <b>phone #</b> , <b>address</b> , <b>QQ #</b> )
Data Source 3	a: (id, <b>name</b> ); b: (id, <b>Private phone #</b> , <b>office phone #</b> ) <b>no address</b>
Data Source 4	( <b>Family name</b> , <b>given name</b> , <b>Private phone #</b> , <b>home address</b> )
Data Source 5	( <b>Nick name</b> , <b>name</b> , <b>zip code</b> , <b>phone #</b> , <b>City</b> , <b>street</b> )

*Different Format*



**Target Format**

(Nick name, **name**, **phone #**, **address**, **QQ #**)

Data Source 1	( <del>Nick name</del> , <b>name</b> , <b>phone #</b> , <b>address</b> , <b>QQ #</b> )
Data Source 2	(Nick name, <b>name</b> , <b>phone #</b> , <b>address</b> , <b>QQ #</b> )
Data Source 3	( <del>Nick name</del> , <b>name</b> , <b>phone #</b> , <b>address</b> , <b>QQ #</b> )
Data Source 4	( <del>Nick name</del> , <b>name</b> , <b>phone #</b> , <b>address</b> , <b>QQ #</b> )
Data Source 5	(Nick name, <b>name</b> , <b>phone #</b> , <b>address</b> , <b>QQ #</b> )

**Pattern  
Alignment**

# Previous Section

- Step 2: Prepare Data - Pre-process

**Web Page**



2010世界影像大赛/Supra CCD 影像传感器  
35倍长焦 25mm广角 3.0英寸 4.6万有效像素  
智能自动模式/10帧每秒/人脸检测/记忆

**Color**

**Brand**

Sony/索尼

**Model**

DSC-H300

**Type**

数码长焦照相机

35倍变焦 大广角

2010万像素

**Resolution**

13人付款

20条评论

---



限量抢购 包邮  
送 MP3  
同一购三

正品Sony/索尼DSC-H300数码小单反相机

**Feature**

35倍长焦大陆行货全国联保

48人付款

129条评论

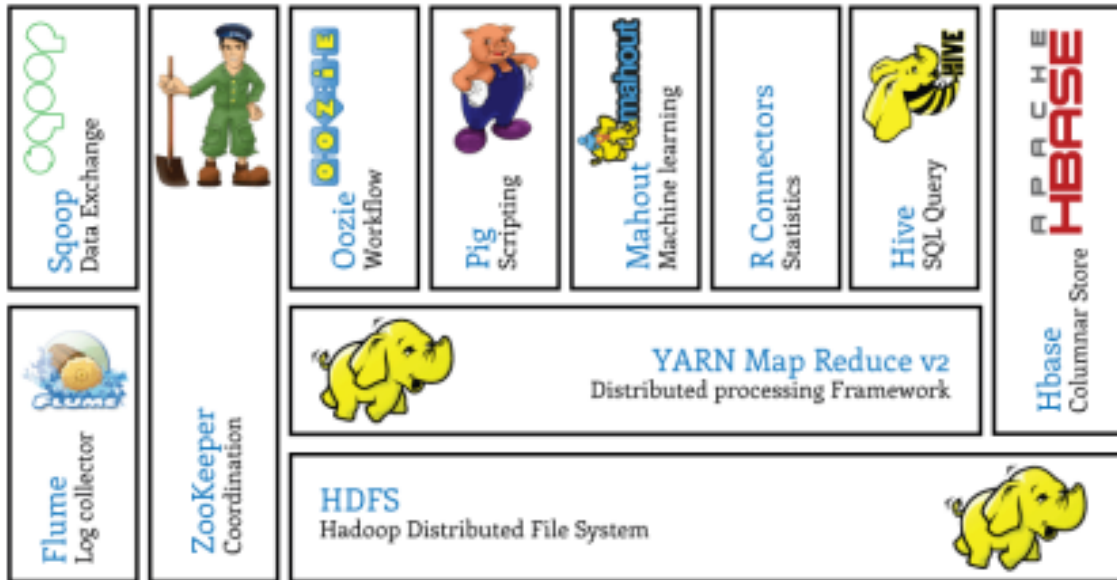
**Record Association**

**Table Format**

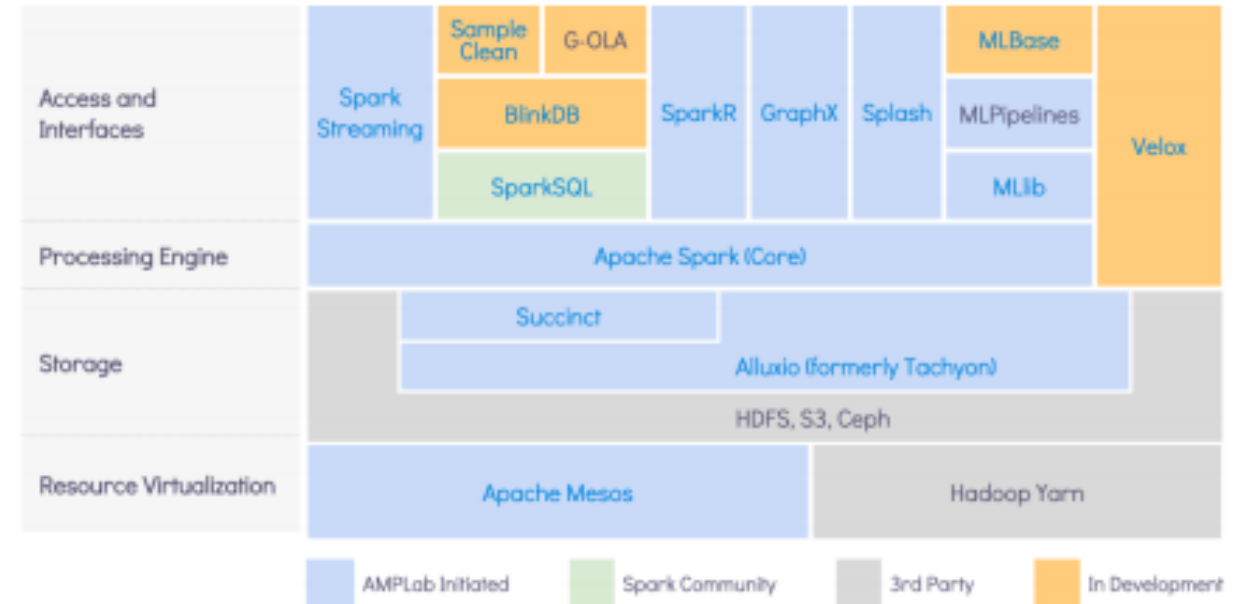
Category	Camera
Brand	Sony
Model	DSC-H300
Type	Telephoto Lens
Pixels	21 Million photo pixel
Color	Black
Feature	35X Telephoto Lens

# Hadoop Eco-System

## Hadoop Eco-System



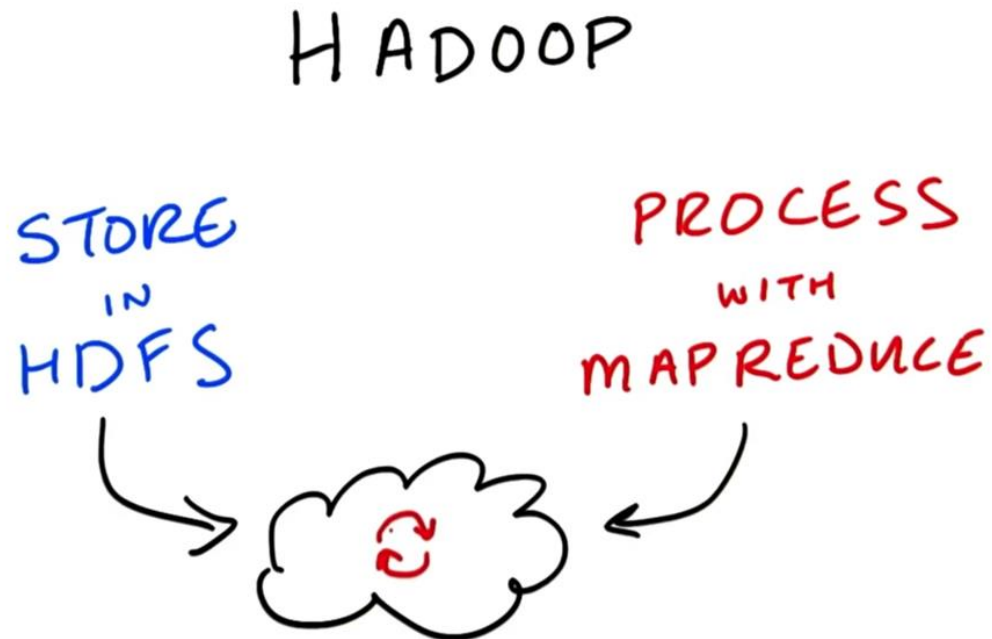
## Spark Eco-System



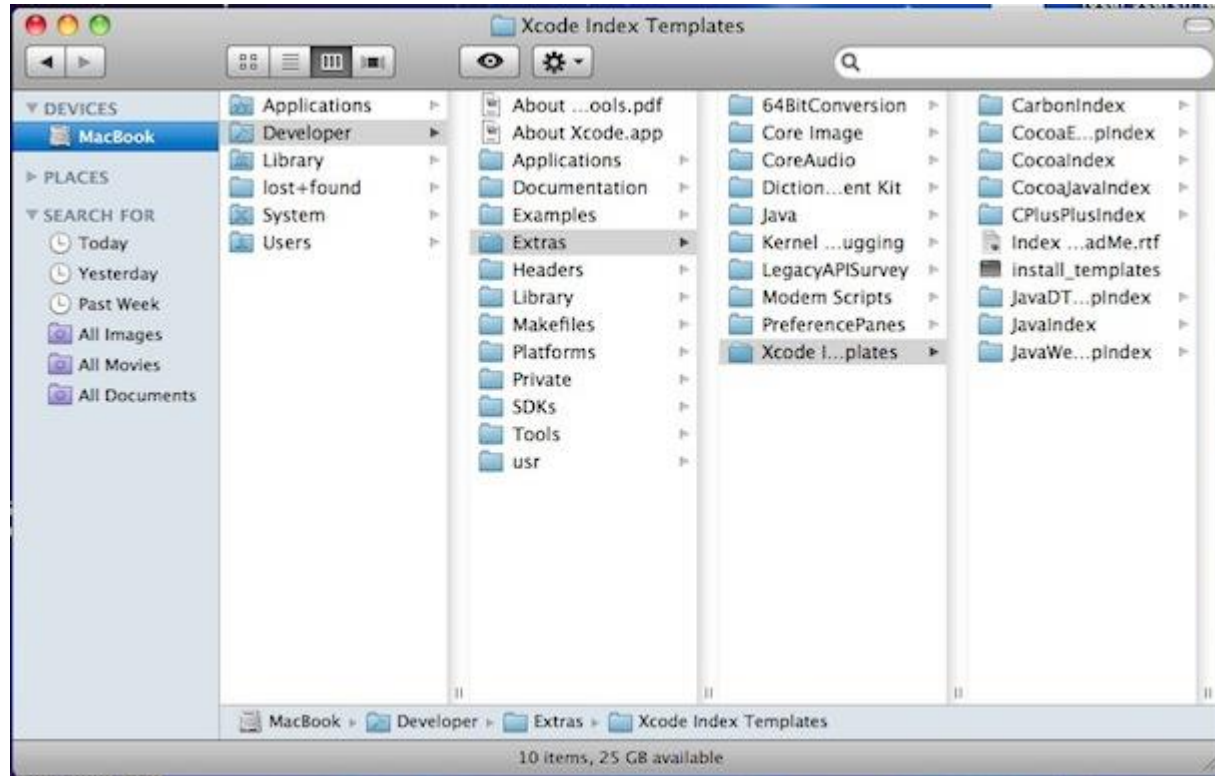
# Hadoop HDFS



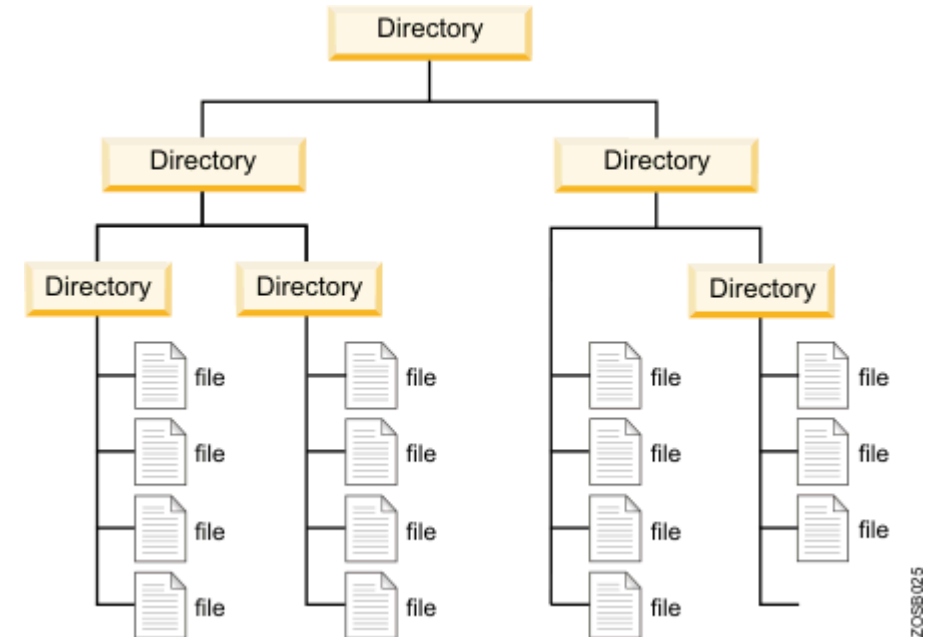
- Currently one of the best distributed file system
  - Based on GFS by Google in 2003
  - Doug Cutting first proposed by Yahoo in 2006



# File System



Structured, tree-like file storage

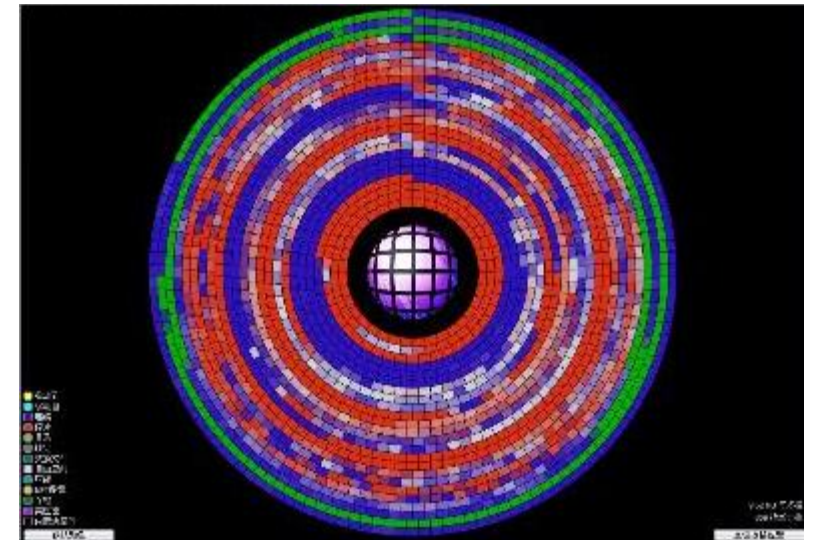
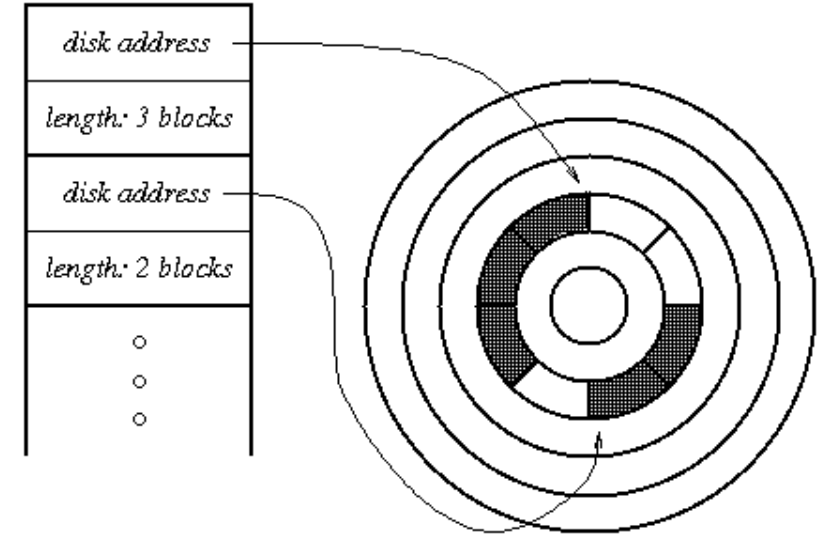


2008/02/5

# Local File System

- Files are stored **locally** (in one hard disk)
- How files are really stored:
  - File system is split into two parts:
    - 1) Header / file pointer (inode)
    - 2) File data (in blocks)
    - Formally, file system **format**
    - E.g., FAT32, NTFS, ext4, HFS, APFS (since iPhone 7)
  - Files are not stored continuously
    - Formally, **random access**

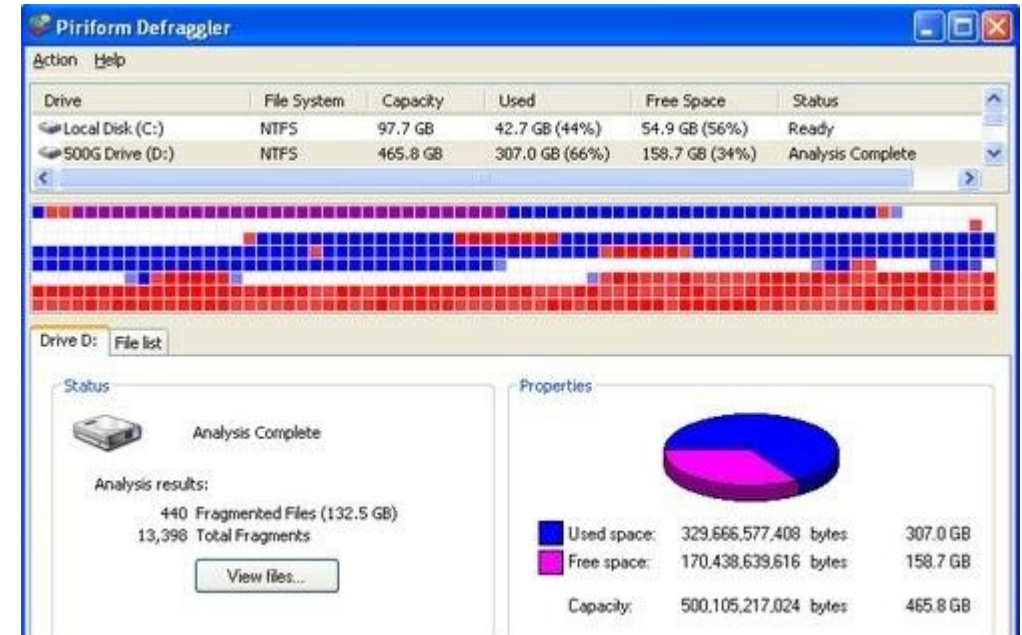
Block Pointer Allocation





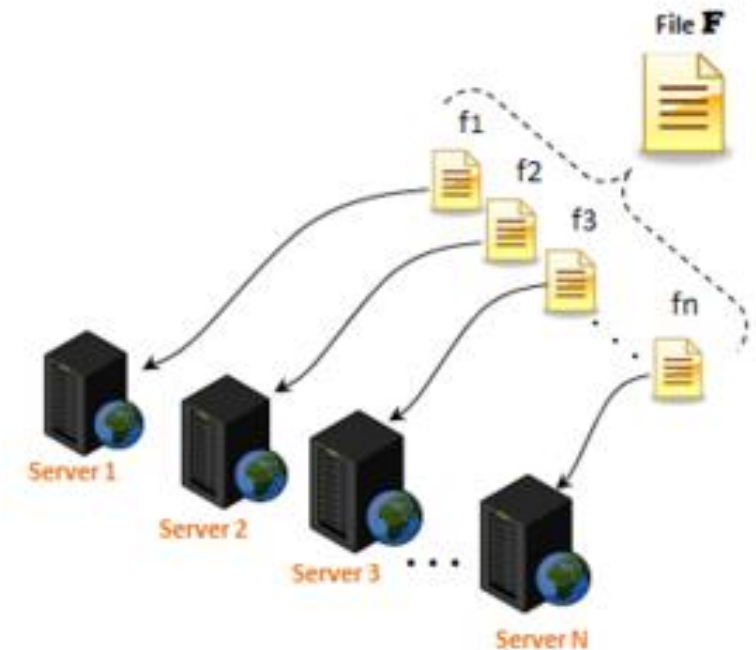
# Local File System

- In some old, bad-implemented file systems, **defragmentation** is sometimes required.
- Most of the modern file system design has eliminated such problem.



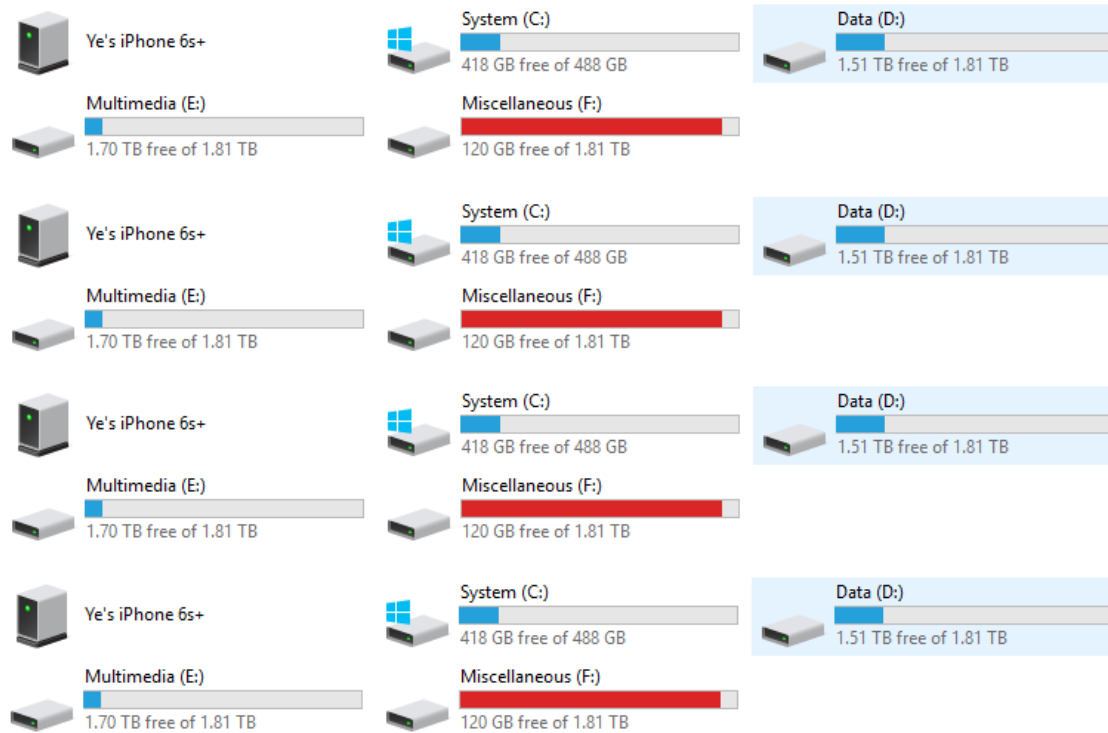
# Distributed File System

- Local file system: files are stored **locally** (in one hard disk)
- Distributed file system:
  - Files are stored in multiple hard disks
  - Sometimes even distributed across the Internet
- Pros:
  - Store up to peta level data
    - 1 SSD: up to 4 TB / 1 HDD: up to 16 TB
    - 1 PB = 1,024 TB / 64 HDDs
  - Content safe: replica / fault tolerance
  - Load balance benefits

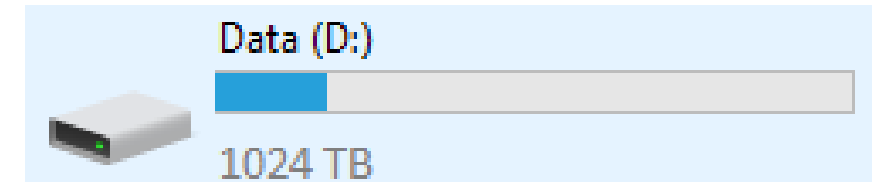




# Distributed File System



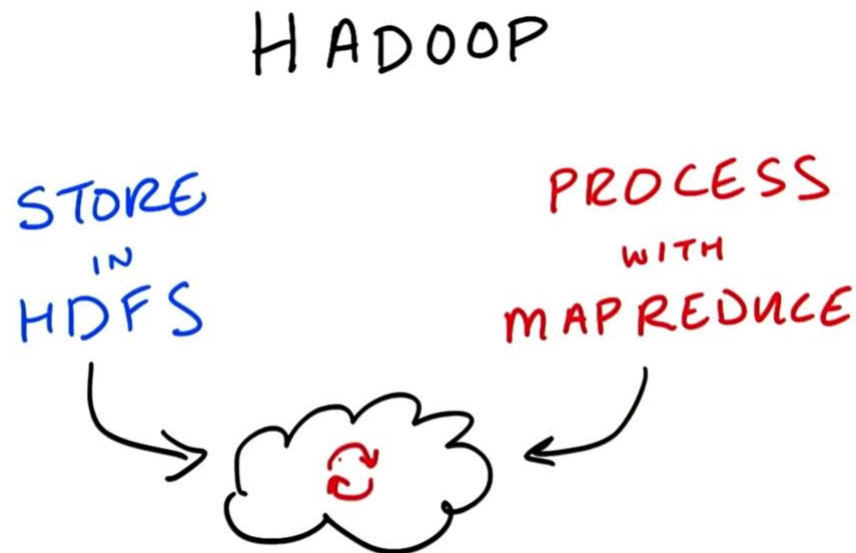
Without DFS



With DFS

# Hadoop Distributed File System

- HDFS = Hadoop Distributed File System
- One of the implementation / standard of distributed file system
  - Based on GFS by Google in 2003
  - First proposed by Yahoo in 2006



## The Hadoop Distributed File System

Konstantin Shvachko, Hairong Kuang, Sanjay Radia, Robert Chansler  
Yahoo!  
Sunnyvale, California USA  
{Shv, Hairong, SRadia, Chansler}@Yahoo-Inc.com

*Abstract*—The Hadoop Distributed File System (HDFS) is designed to store very large data sets reliably, and to stream those data sets at high bandwidth to user applications. In a large cluster, thousands of servers both host directly attached storage and execute user application tasks. By distributing storage and computation across many servers, the resource can grow with demand while remaining economical at every size. We describe the architecture of HDFS and report on experience using HDFS to manage 25 petabytes of enterprise data at Yahoo!.

*Keywords*: Hadoop, HDFS, distributed file system

### I. INTRODUCTION AND RELATED WORK

Hadoop [1][16][19] provides a distributed file system and a framework for the analysis and transformation of very large data sets using the MapReduce [3] paradigm. An important characteristic of Hadoop is the partitioning of data and computation across many (thousands) of hosts, and executing application computations in parallel close to their data. A Hadoop cluster scales computation capacity, storage capacity and IO bandwidth by simply adding commodity servers. Hadoop clusters at Yahoo! span 25 000 servers, and store 25 petabytes of application data, with the largest cluster being 3500 servers. One hundred other organizations worldwide report using

developed at Facebook. Pig [4], ZooKeeper [6], and Chukwa were originated and developed at Yahoo! Avro was originated at Yahoo! and is being co-developed with Cloudera.

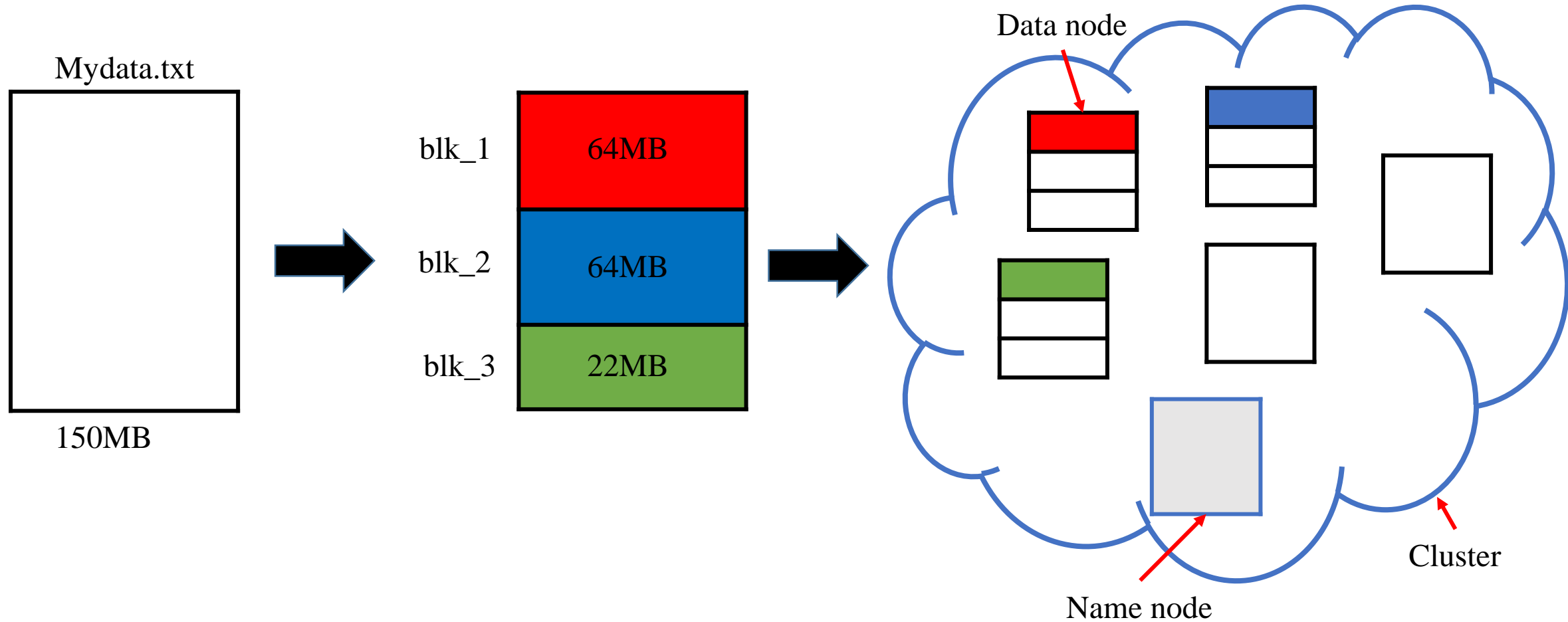
HDFS is the file system component of Hadoop. While the interface to HDFS is patterned after the UNIX file system, faithfulness to standards was sacrificed in favor of improved performance for the applications at hand.

HDFS stores file system metadata and application data separately. As in other distributed file systems, like PVFS [2][14], Lustre [7] and GFS [5][8], HDFS stores metadata on a dedicated server, called the NameNode. Application data are stored on other servers called DataNodes. All servers are fully connected and communicate with each other using TCP-based protocols.

Unlike Lustre and PVFS, the DataNodes in HDFS do not use data protection mechanisms such as RAID to make the data durable. Instead, like GFS, the file content is replicated on multiple DataNodes for reliability. While ensuring data durability, this strategy has the added advantage that data transfer bandwidth is multiplied, and there are more opportunities for locating computation near the needed data.

Several distributed file systems have or are exploring fully

# Hadoop Distributed File System



# Hadoop Distributed File System

Is there a problem?

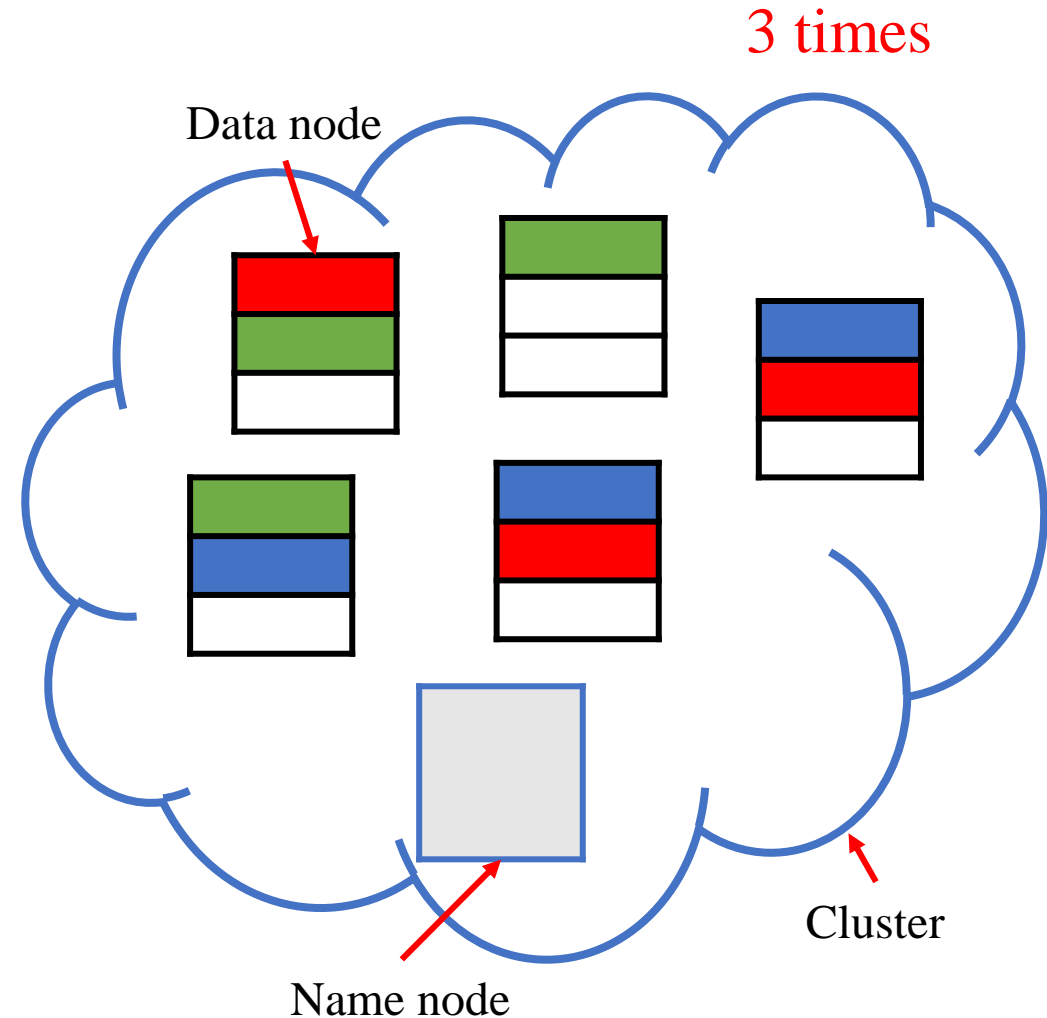
- a) Network failure
- b) Disk failure on data node
- c) Not all data node used
- d) Block size differ
- e) Disk failure on name node

# Hadoop Distributed File System

If one of our data nodes fails?  
If this data node goes away?

## Data Redundancy

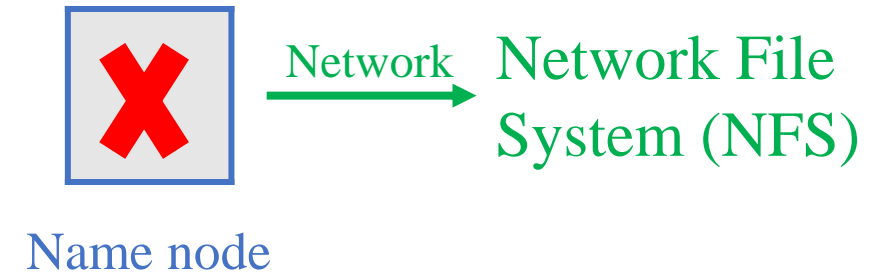
- Hadoop replicates each block three times as it's stored in HDFS
- NameNode is smart enough that if it sees that any of the blocks are under replicated, it will arrange to have those blocks re-replicated on the cluster so we're back to having three copies of them again



# Hadoop Distributed File System

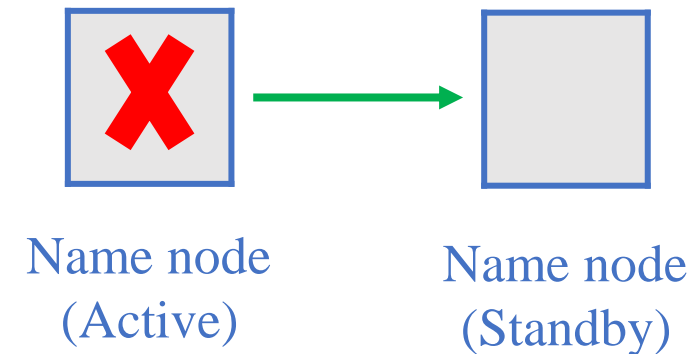
If the disk of the single NameNode would fail ?

- Data inaccessible
- Data on HDFS would be lost permanently



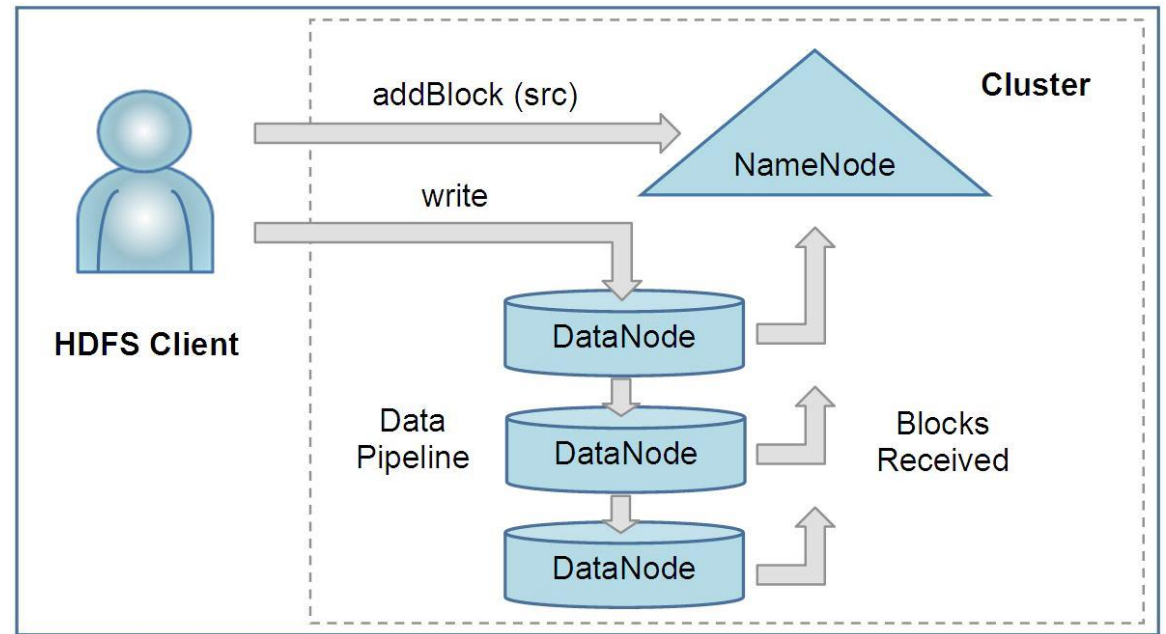
Two plans of Name Node

- NFS
- Standby Name node



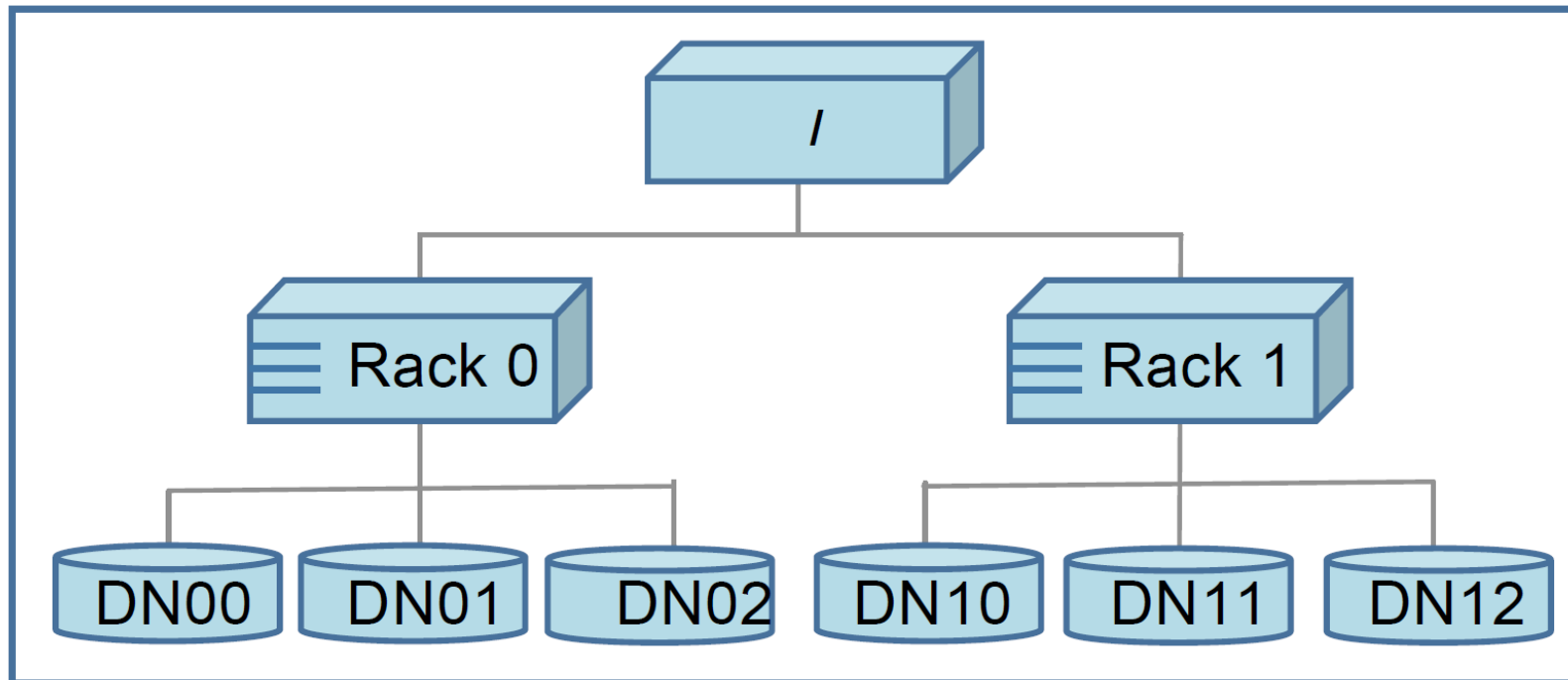
# Hadoop Distributed File System

- Composed by **NameNodes** and **DataNodes**
- NameNode = header / file pointer in FS
- DataNode = file data in FS
- Write a file:
  - In FS: 1) header -> data blocks
  - In HDFS: 1) record in NameNode -> data blocks in DataNode



# Hadoop Distributed File System

- File System: data blocks in one single disk
- HDFS: data blocks in multiple DataNodes contained in server racks across the network





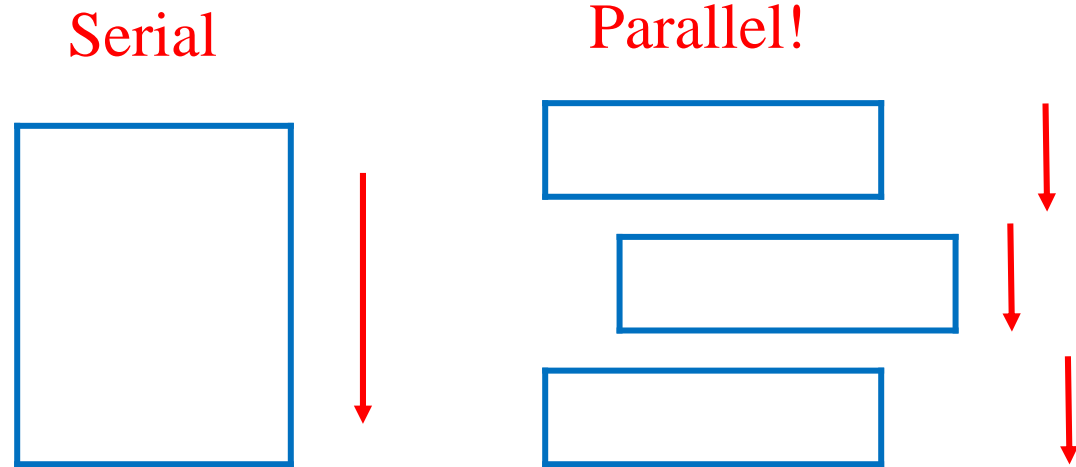
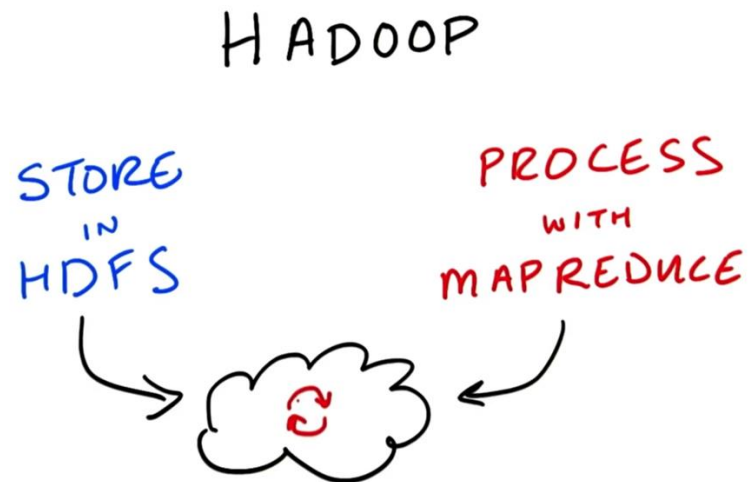
# Hadoop Distributed File System

## Summary

- One form of distributed file system
- Belongs to Hadoop eco system
- Competitors:
  - GFS
  - Ceph
  - GlusterFS
- Pros and cons: difficult to compare through standards
  - Is Chinese law good or Western law good?

# Mapreduce

Mapreduce is designed to by a very parallelized way of managing data, meaning that you input data is split into many pieces, and each piece is processed simultaneously.



# Mapreduce

Ledger:

2012-01-01 London Clothes 25.99  
2012-01-01 Miami Music 12.15  
2012-01-02 NYC Toys 3.10  
2012-01-02 Miami Clothes 50.00  
.....



London 25.99  
Miami 62.15  
NYC 3.10  
.....

Mappers



Reducers



Ledger 1:  
2012-01-01 London Clothes 25.99  
2012-01-01 Miami Music 12.15  
...



Ledger 2:  
2012-01-02 NYC Toys 3.10  
2012-01-02 Miami Clothes 50.00  
...



Ledger 3:  
2012-01-03 NYC Toys 6.10  
2012-01-05 Miami Clothes 30.00  
...

Key: City Name  
Value: Money



# Mapreduce

Ledger:

2012-01-01 London Clothes 25.99

2012-01-01 Miami Music 12.15

2012-01-02 NYC Toys 3.10

2012-01-02 Miami Clothes 50.00

.....

Key: City Name  
Value: Money



London 25.99

Miami 62.15

NYC 3.10

.....

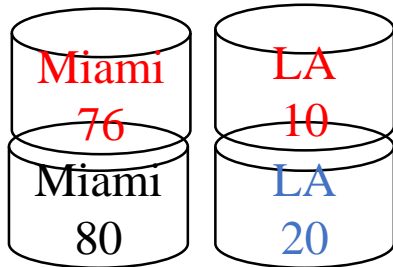
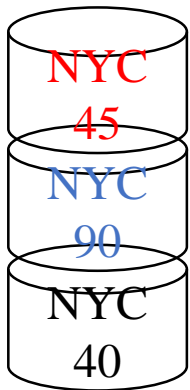
Mappers



Reducers

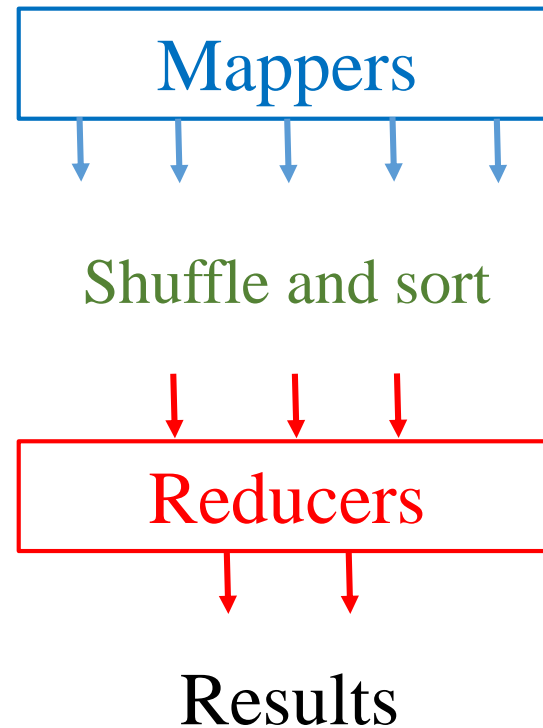


Reducers



# Mapreduce

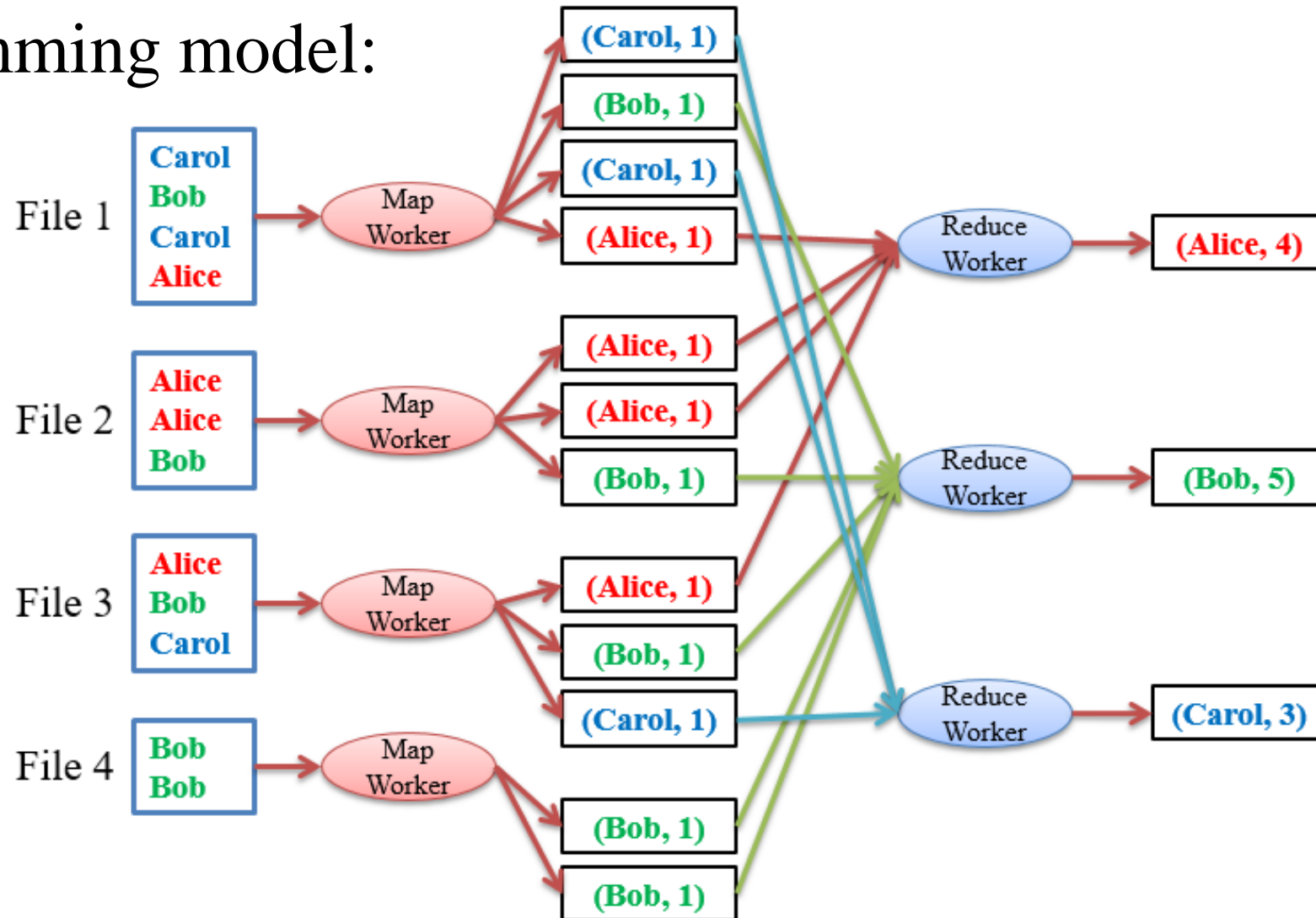
- Shuffle is the movement of the intermediate data from the **Mappers** to the **Reducers** and the combination of all the small sets of records together.
- Sort is the fact that the Reducers will organize the sets of records
  - For example, the piles of city cards into order
  - Reduce works on one set of records



Intermediate records  
(Key, Value)

# Mapreduce

The programming model:



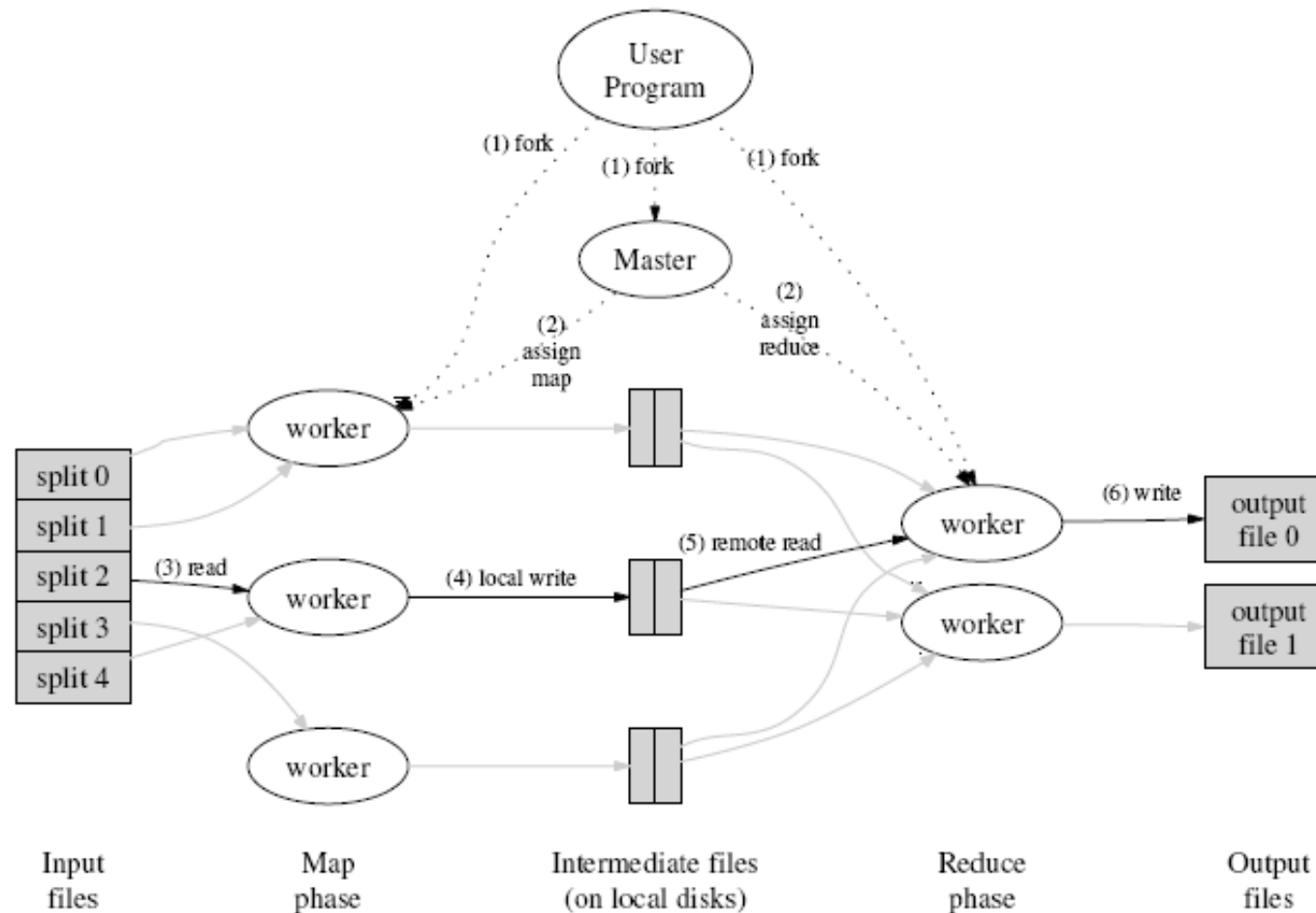
Word count: read text files and count how often words occur.

Key?

Value?

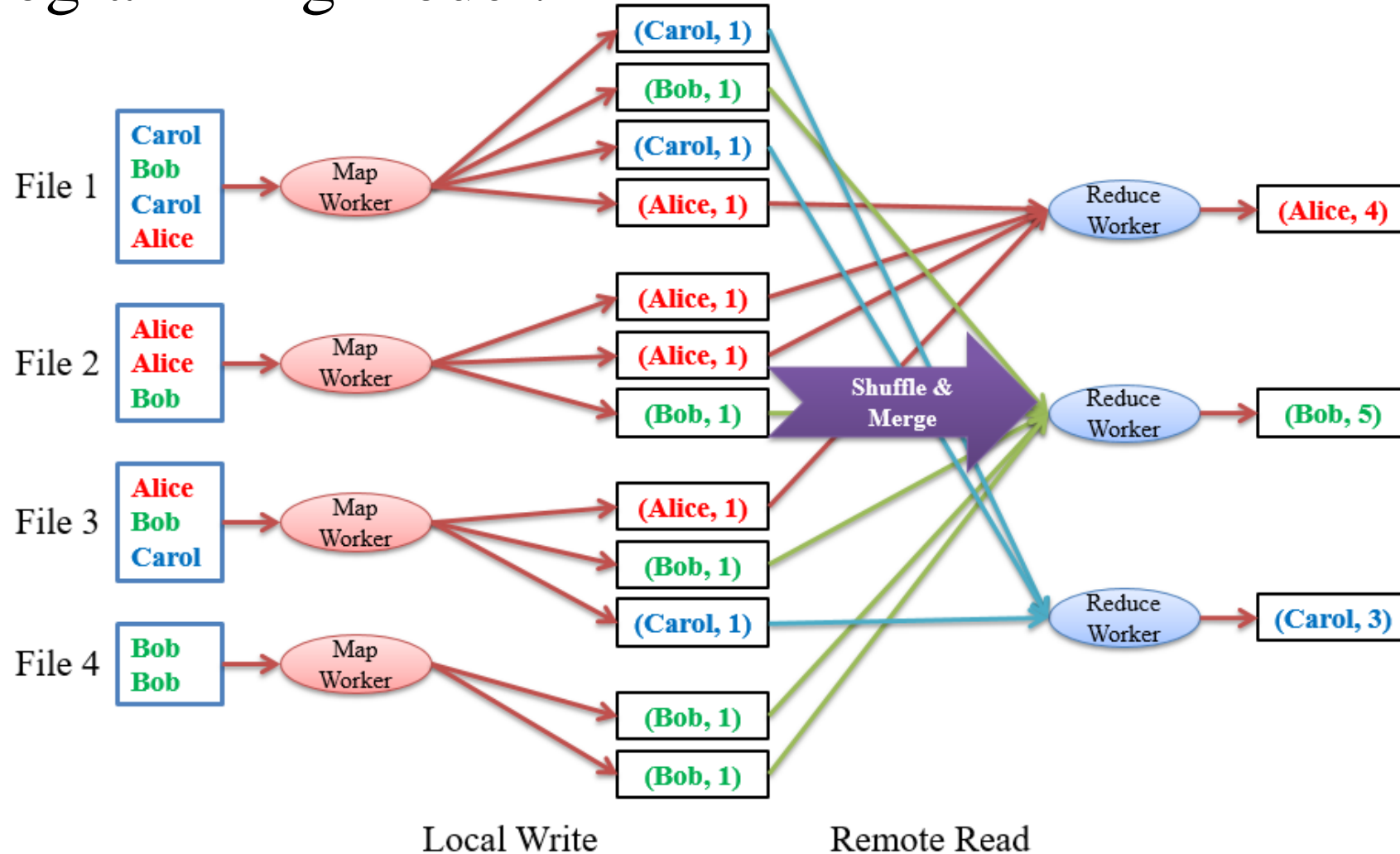
# Mapreduce

The programming model:



# Mapreduce

The programming model:



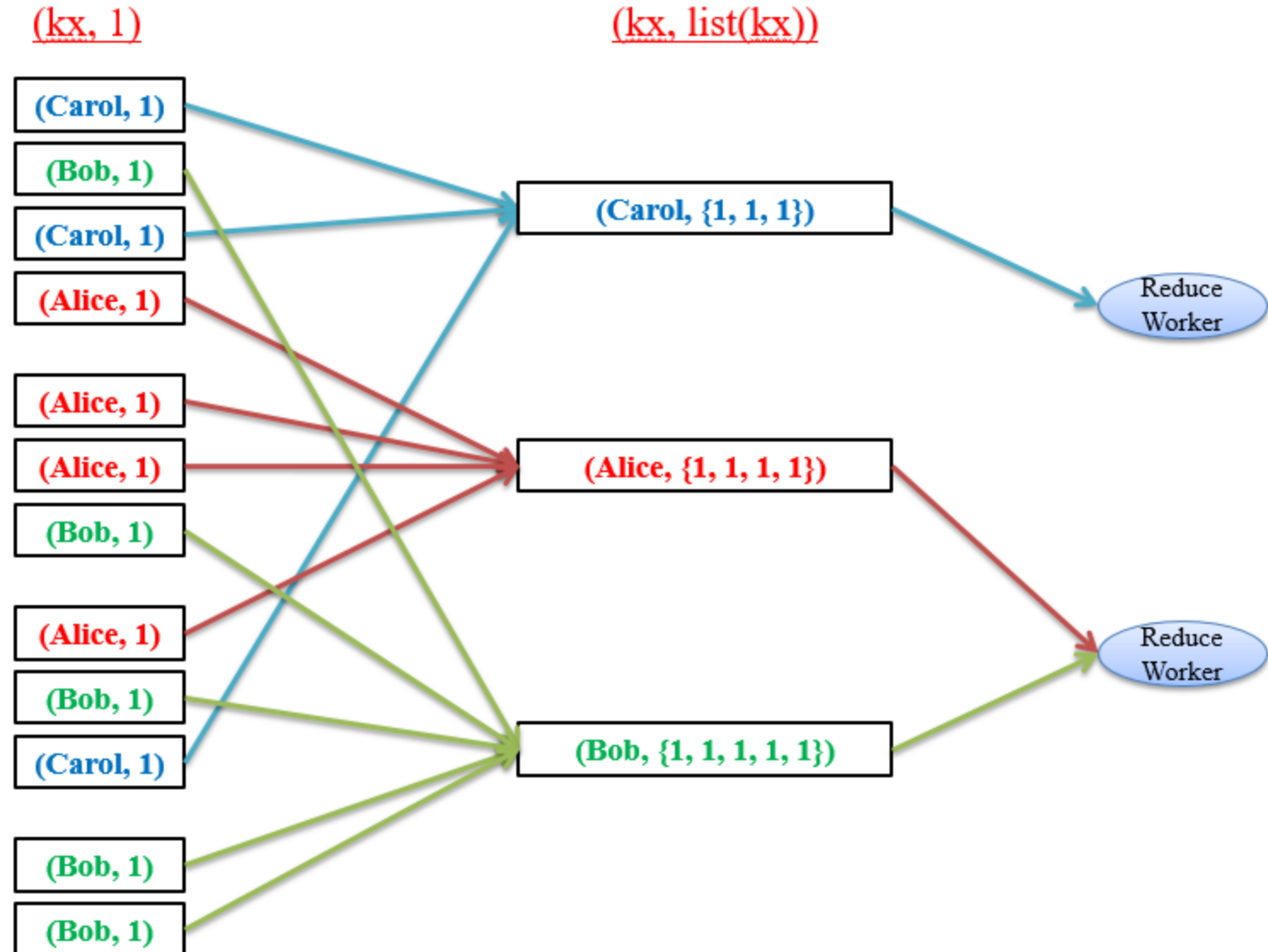


# Mapreduce

- Map:
  - $(k1, 1) \rightarrow \text{list}(k1) \rightarrow \{1, 1, \dots\}$
- Reduce:
  - $(k1, \text{list}(k1)) \rightarrow (k1, n)$
- For word count example:
  - $(k1, 1): (\text{word}, 1)$
  - $\text{list}(k1): \{(the\ 1^{st}\ occurrence\ of\ \text{word}\ in\ file, 1), (the\ 2^{nd}\ occurrence\ of\ \text{word}\ in\ file, 1), (the\ 3^{rd}\ occurrence\ of\ \text{word}\ in\ file, 1), \dots\}$
  - $(k1, \text{list}(k1)): (word, \{1, 1, 1, 1, \dots\})$
  - $(k1, n): (word, \text{sum}\{1, 1, 1, 1, \dots\})$

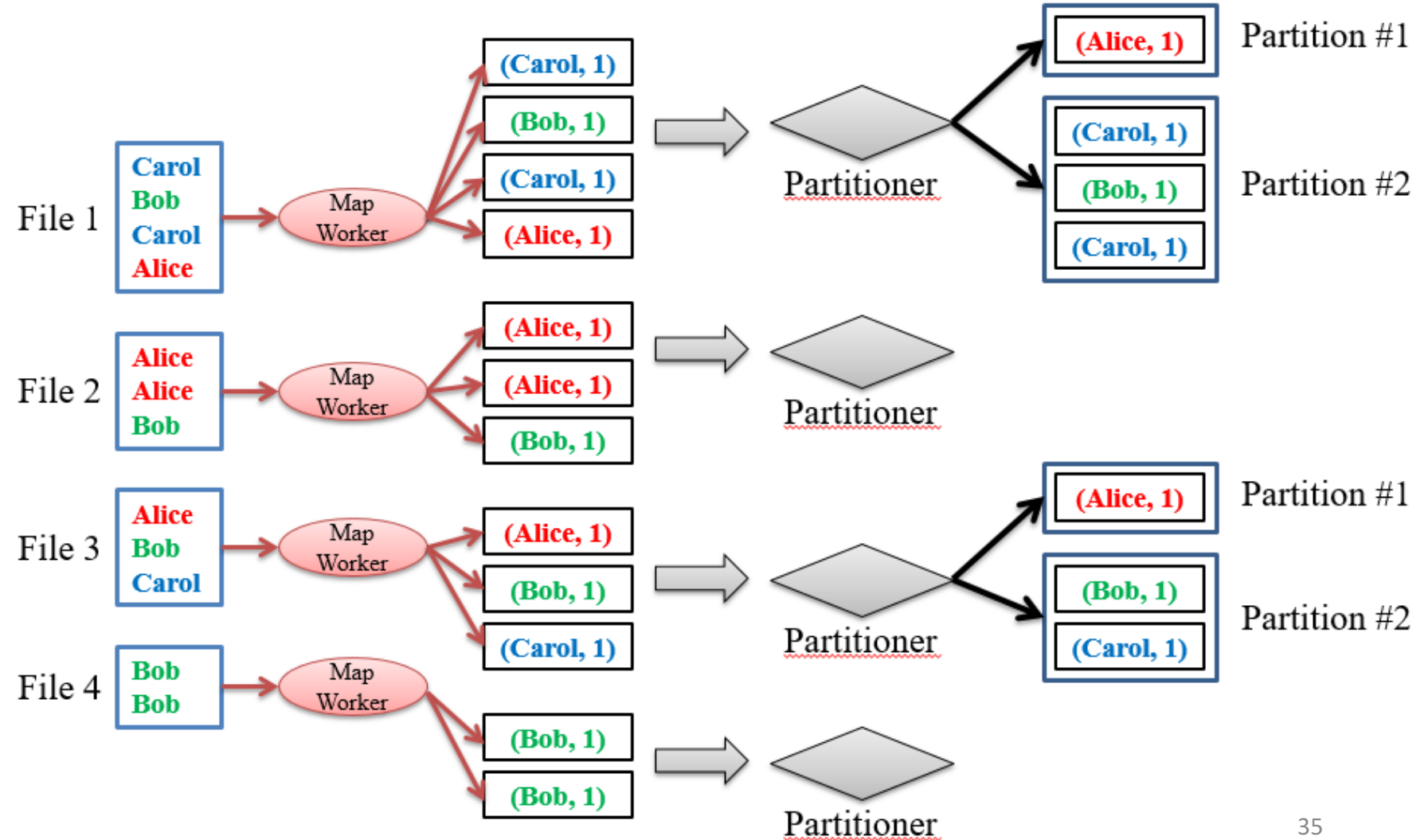
# Mapreduce

- Rule No. 1: The same key goes to the same reduce worker



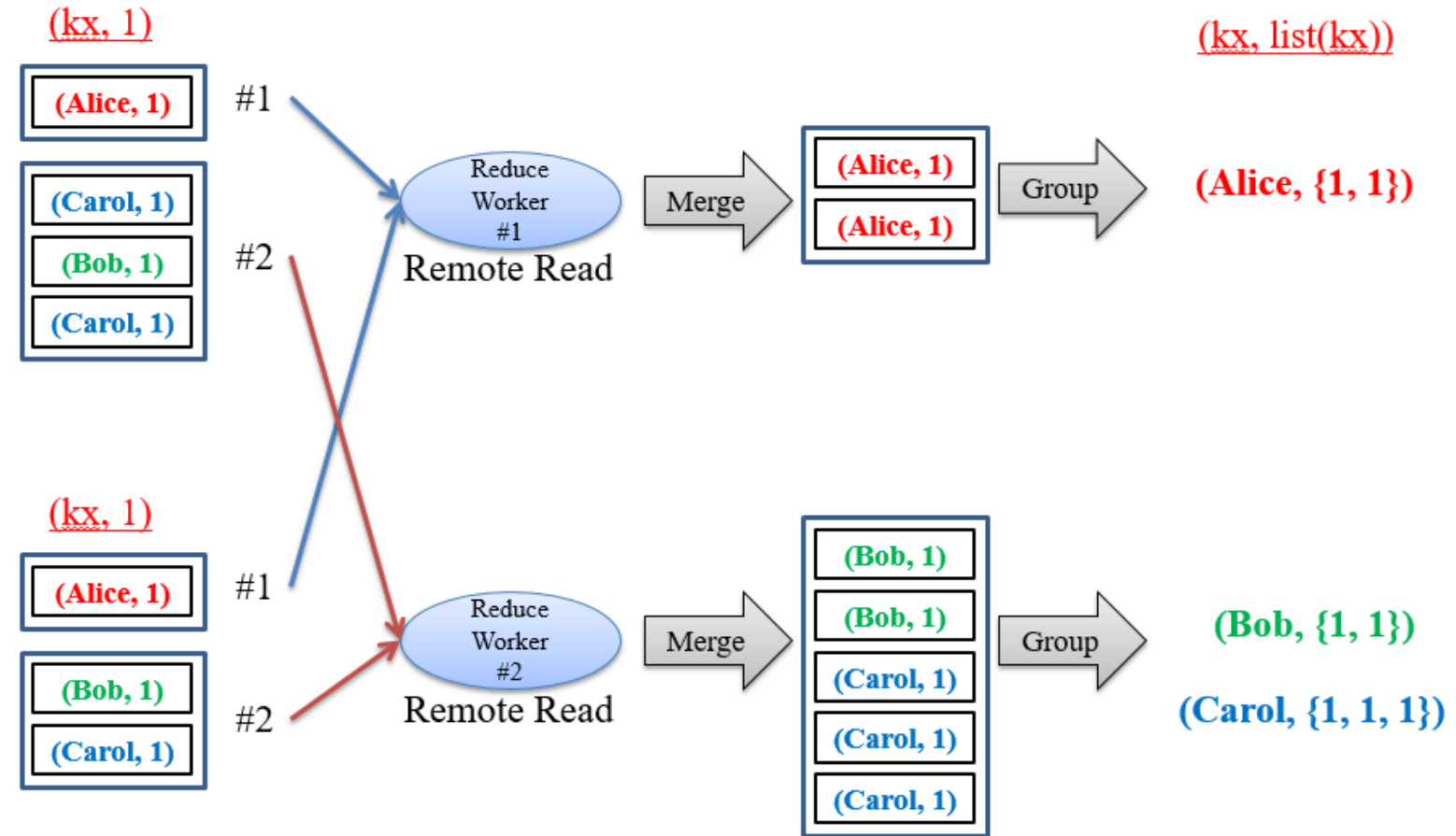
# Mapreduce

- How Rule No. 1 is done (1)



# Mapreduce

- How Rule No. 1 is done (2)



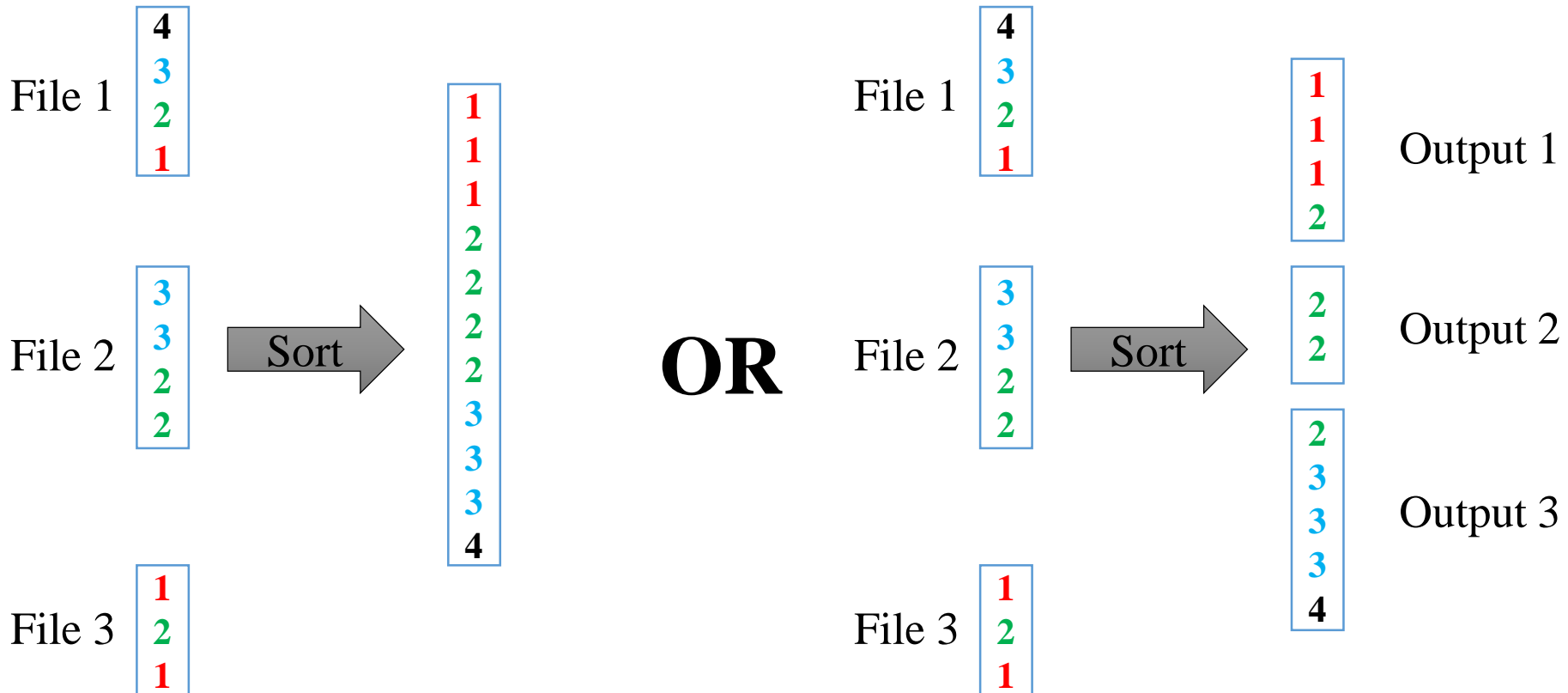
# of Partitions == # of Reduce Workers

# Mapreduce

- **Rule No. 2:** Keys are sorted on each reduce worker
  - A side effect of the **Shuffle and Merge** phase
  - Important to advanced tasks such as **Sort** and **Join**

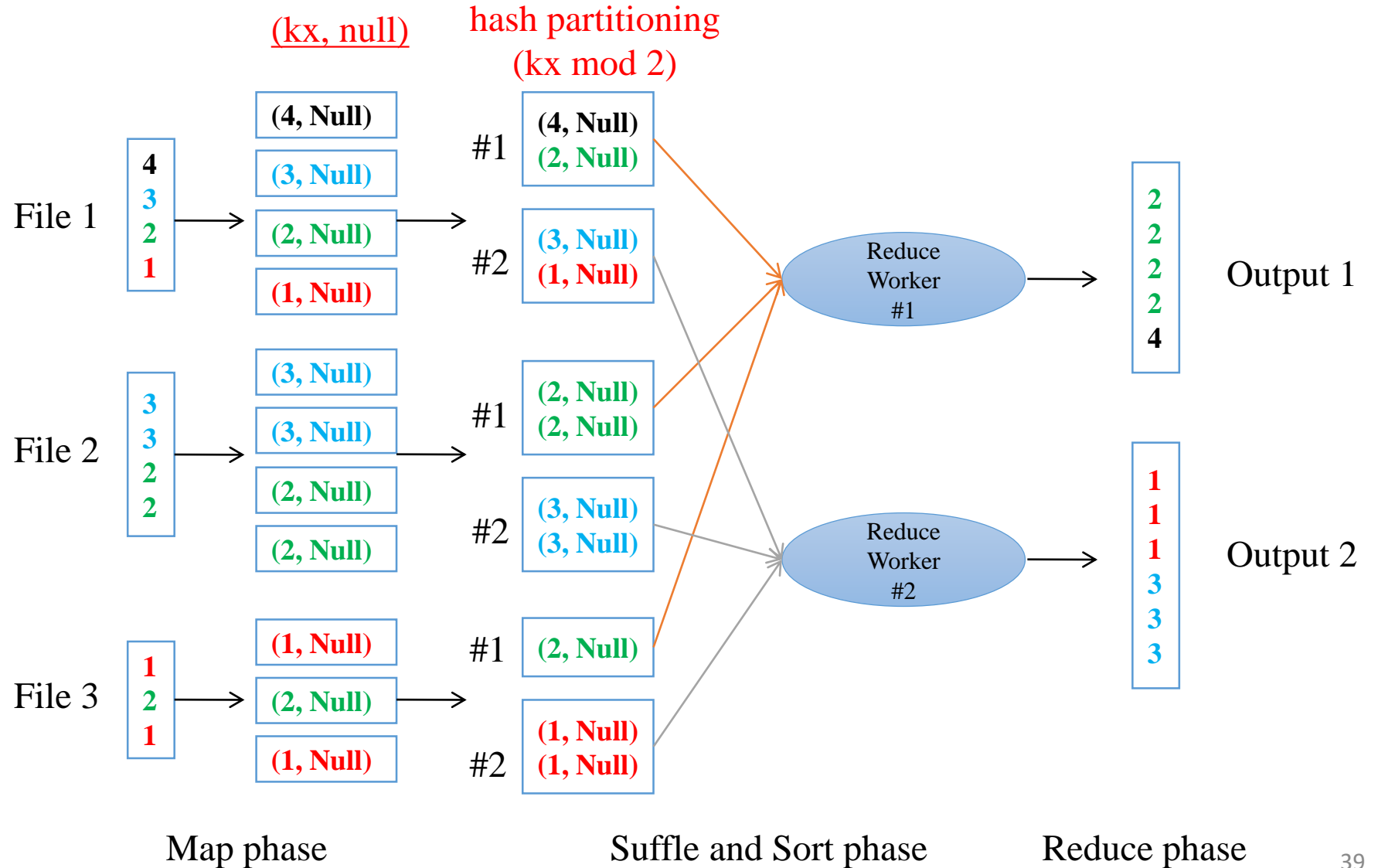
# Mapreduce

- **Rule No. 2:** Keys are sorted on each reduce worker



# Mapreduce

- Keys are already sort of sorted (partially sorted)



# Mapreduce

- However, to obtain a good balance partitioner is not trivial
- **Load balance is crucial** since the job execution time is determined by the **slowest** reduce worker.
- Range split points depend on **data distribution**.
  - {1, 2, 3, 4, 5, 5, 5, 5, 7, 7, 7, 7}: [1, 5), [5, 7), [7, +Infinity)
  - {1, 1, 1, 2, 3, 4, 5, 7, 8, 8, 9, 9}: [1, 3), [3, 8), [8, +Infinity)
  - .....
- It is feasible to use **sampling** to estimate the data distribution.



# Mapreduce

- Join

Name	Color
Alice	Red
Alice	Green
Bob	Blue

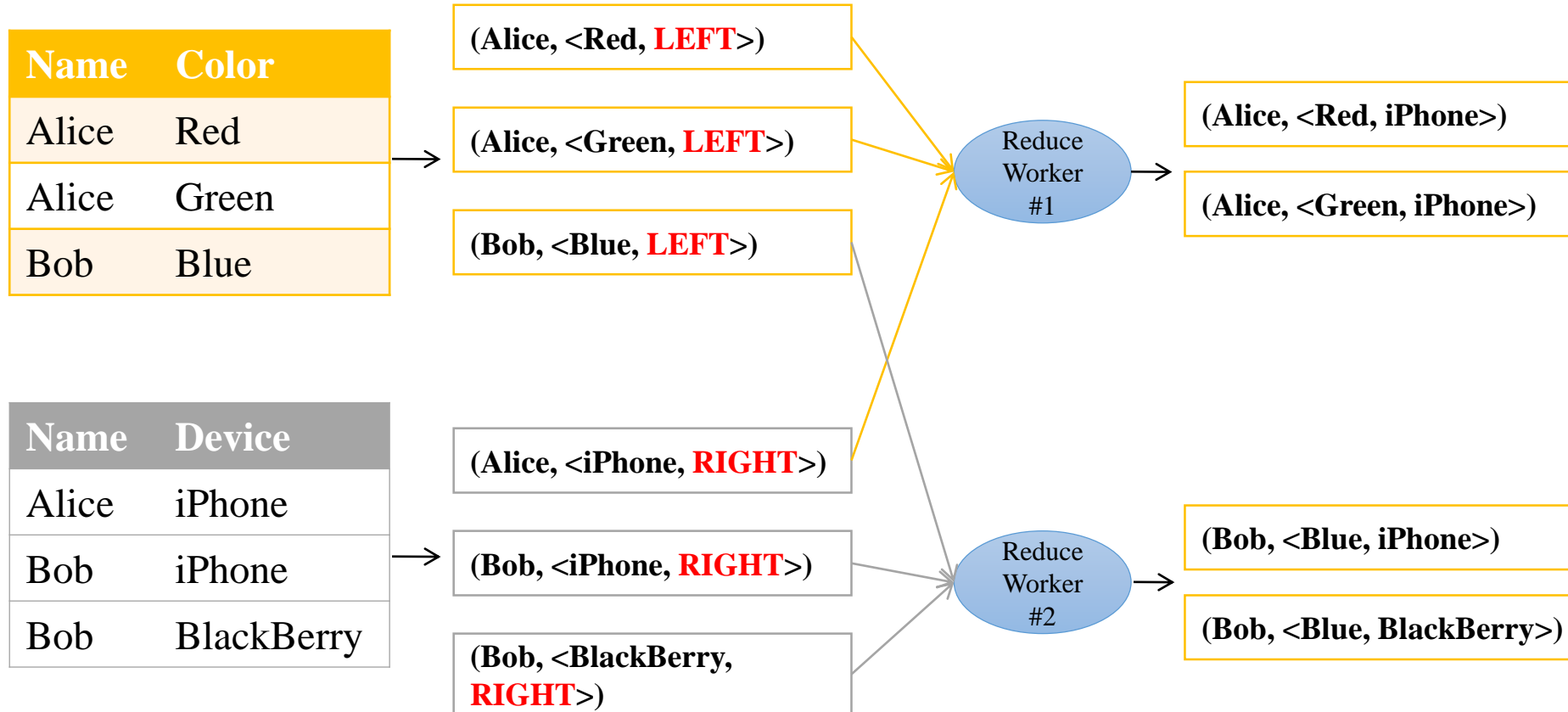
Name	Device
Alice	iPhone
Bob	iPhone
Bob	BlackBerry



Name	Color	Device
Alice	Red	iPhone
Alice	Green	iPhone
Bob	Blue	iPhone
Bob	Blue	BlackBerry

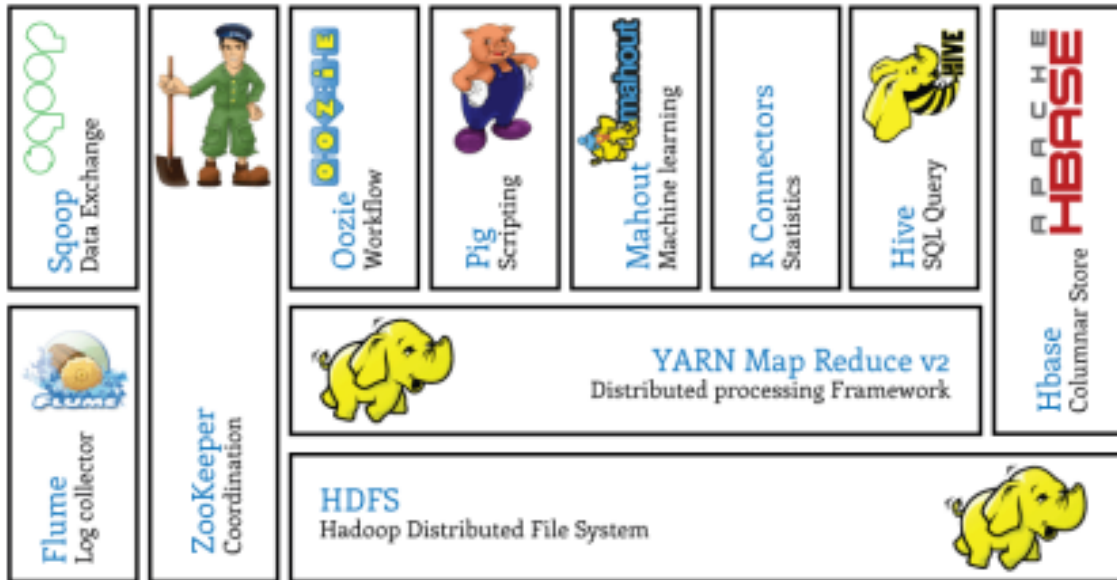
# Mapreduce

- Join

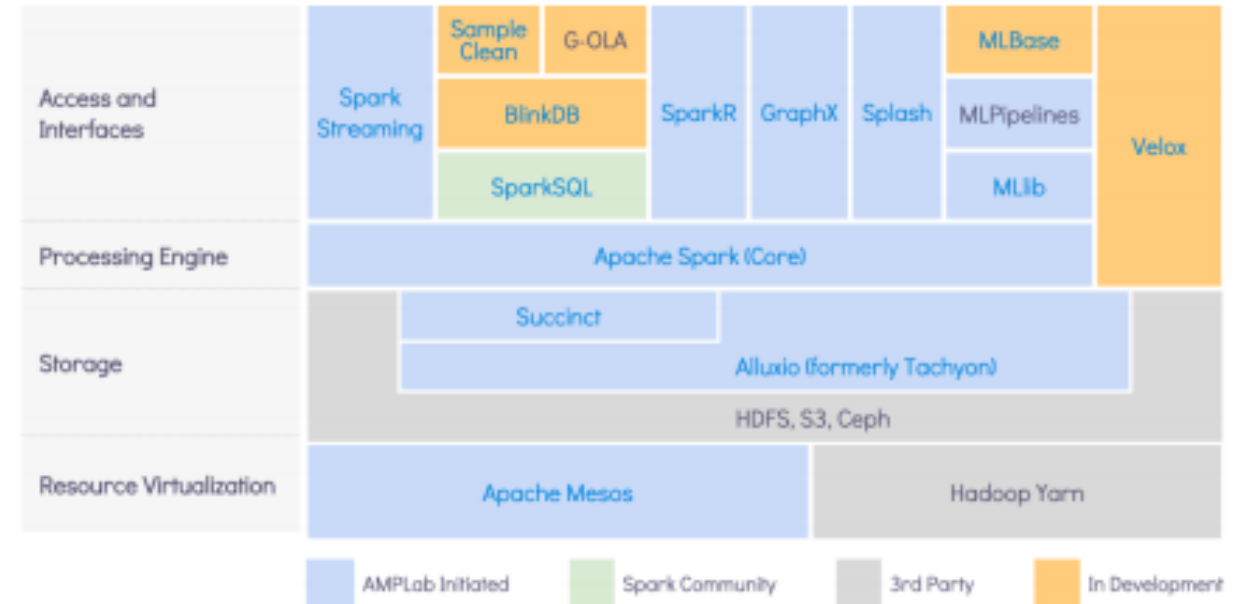


# Hadoop Eco-System

## Hadoop Eco-System



## Spark Eco-System



# Hadoop Eco-System

- Hbase:
  - Column Store
  - Column vs row store
- Hive
  - SQL Query
- R Connectors
  - R Programming Language
- Mahout
  - Machine Learning Library
  - Mahout is great / is bad?
- Pig
  - Pig scripting

