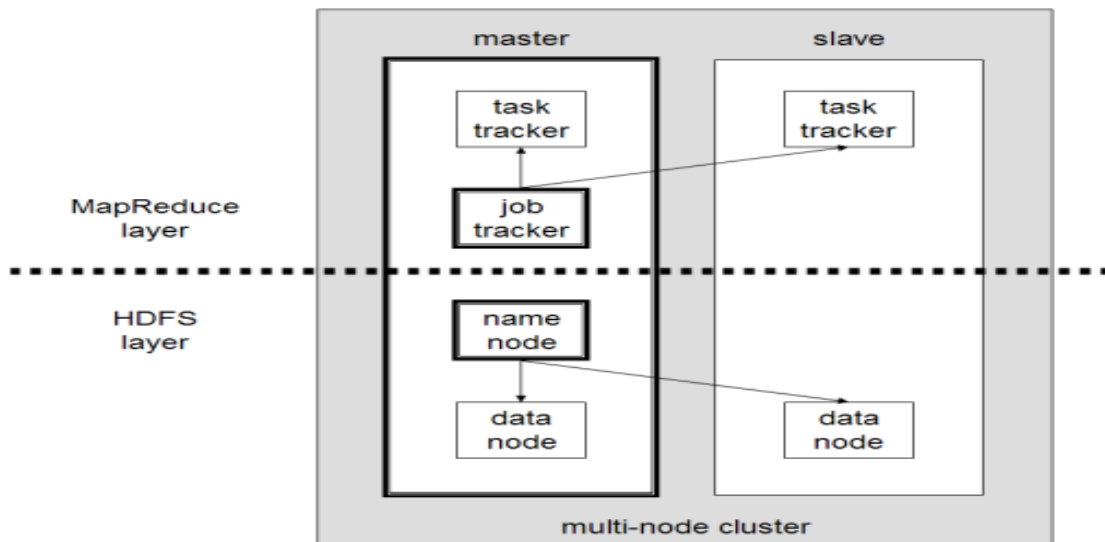


Hadoop: A Framework for Data-Intensive Distributed Computing

CS561-Spring 2012
WPI, Mohamed Y. Eltabakh

What is Hadoop?

- Hadoop is a software framework for *distributed processing* of *large datasets* across *large clusters* of computers
- Hadoop is open-source implementation for Google MapReduce
- Hadoop is based on a simple programming model called *MapReduce*
- Hadoop is based on a simple data model, *any data will fit*
- Hadoop framework consists on two main layers
 - Distributed file system (HDFS)
 - Execution engine (MapReduce)

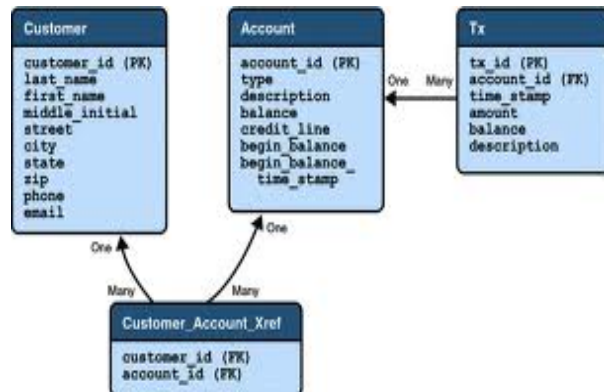


Hadoop Infrastructure

- Hadoop is a ***distributed*** system like ***distributed databases***
- **However, there are several key differences between the two infrastructures**
 - Data model
 - Computing model
 - Cost model
 - Design objectives

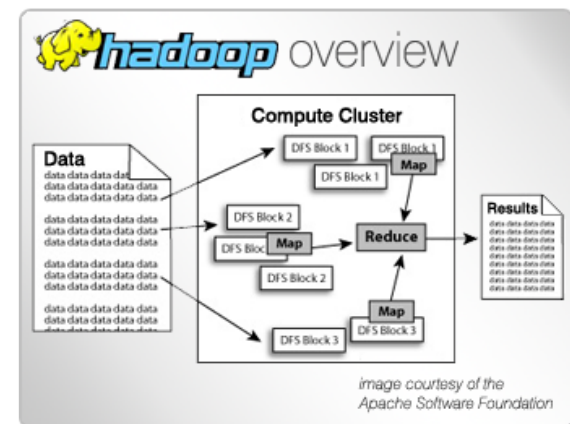
How Data Model is Different?

Distributed Databases



- Deal with tables and relations
- Must have a schema for data
- Data fragmentation & partitioning

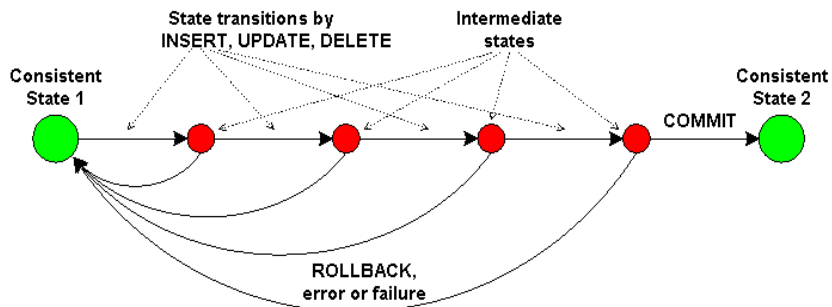
Hadoop



- Deal with flat files in any format
- No schema for data
- Files are divide automatically into blocks

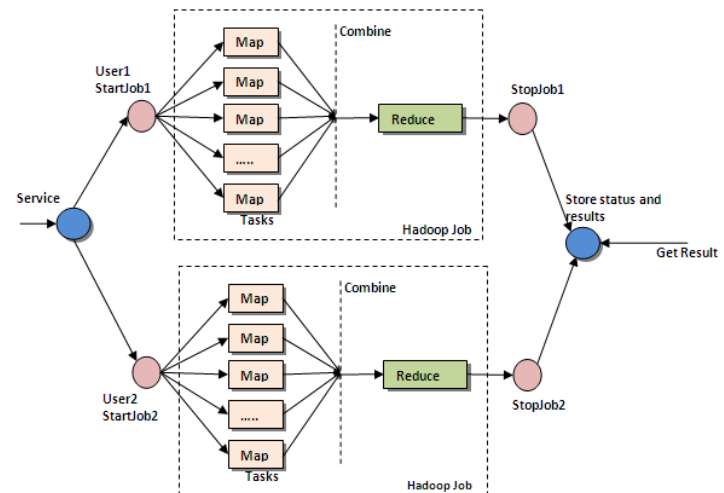
How Computing Model is Different?

Distributed Databases



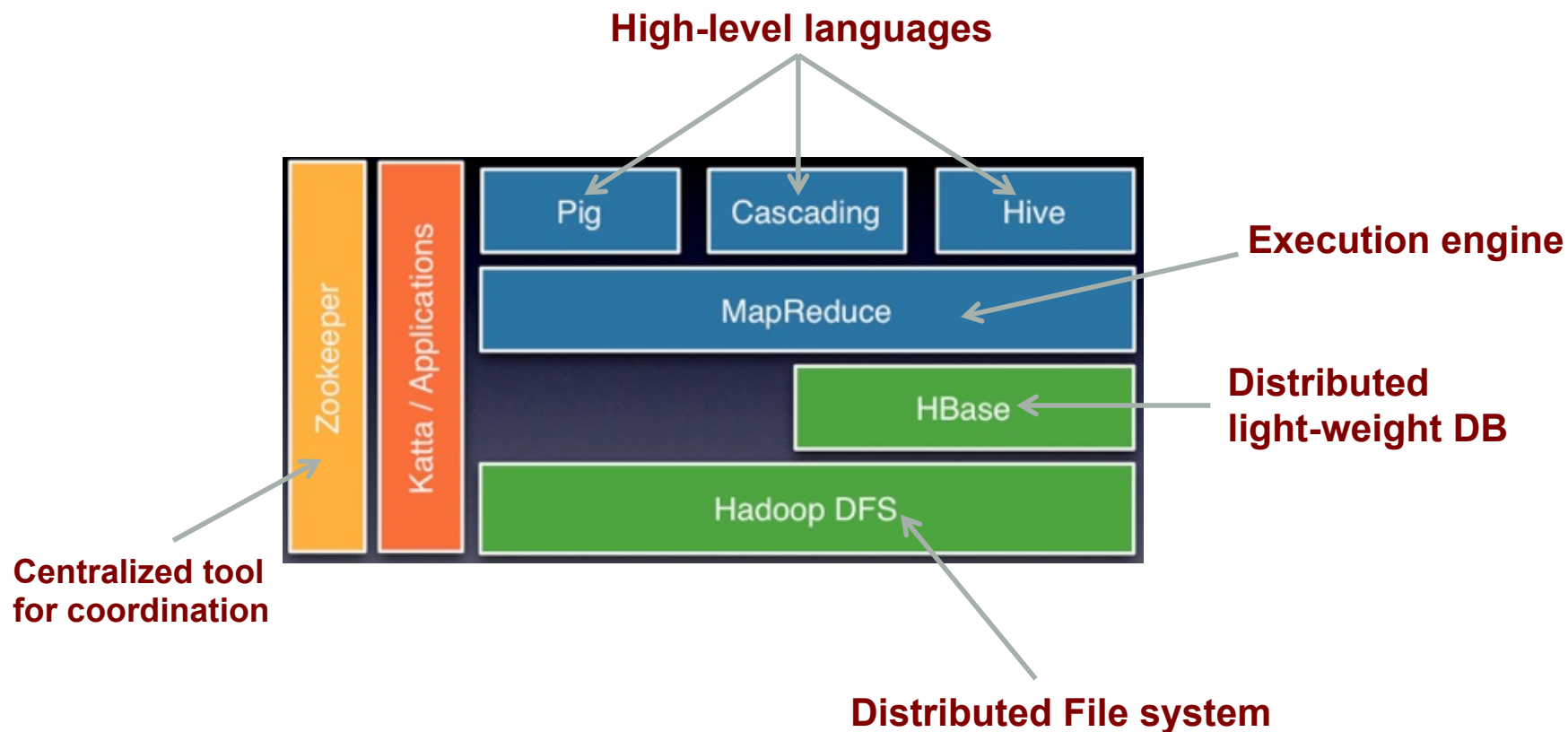
- Notion of a transaction
- Transaction properties ACID
- Distributed transaction

Hadoop



- Notion of a job divided into tasks
- Map-Reduce computing model
- Every task is either a map or reduce

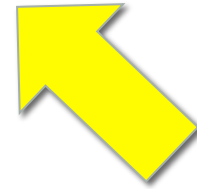
Hadoop: Big Picture



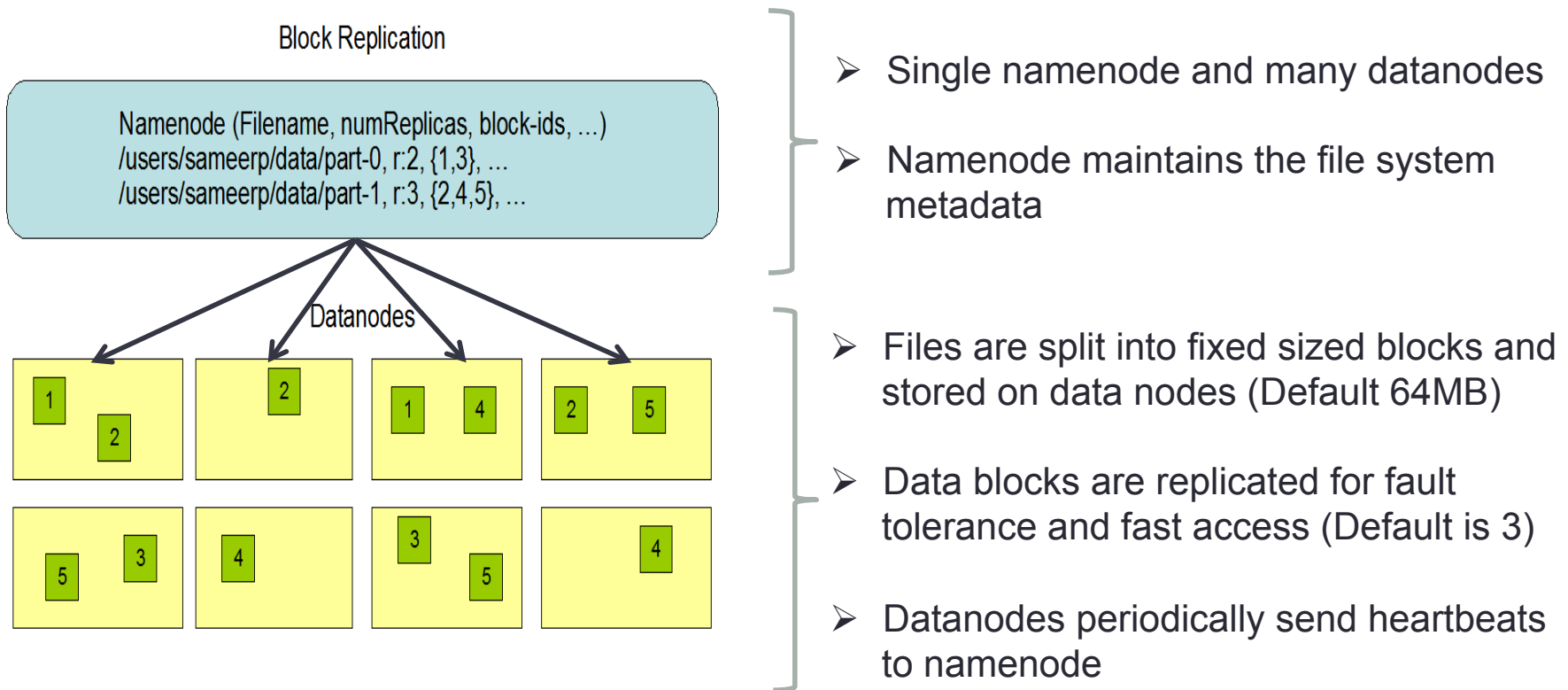
HDFS + MapReduce are enough to have things working

What is Next?

- **Hadoop Distributed File System (HDFS)**
- **MapReduce Layer**
- **Examples**
 - Word Count
 - Join
- **Fault Tolerance in Hadoop**



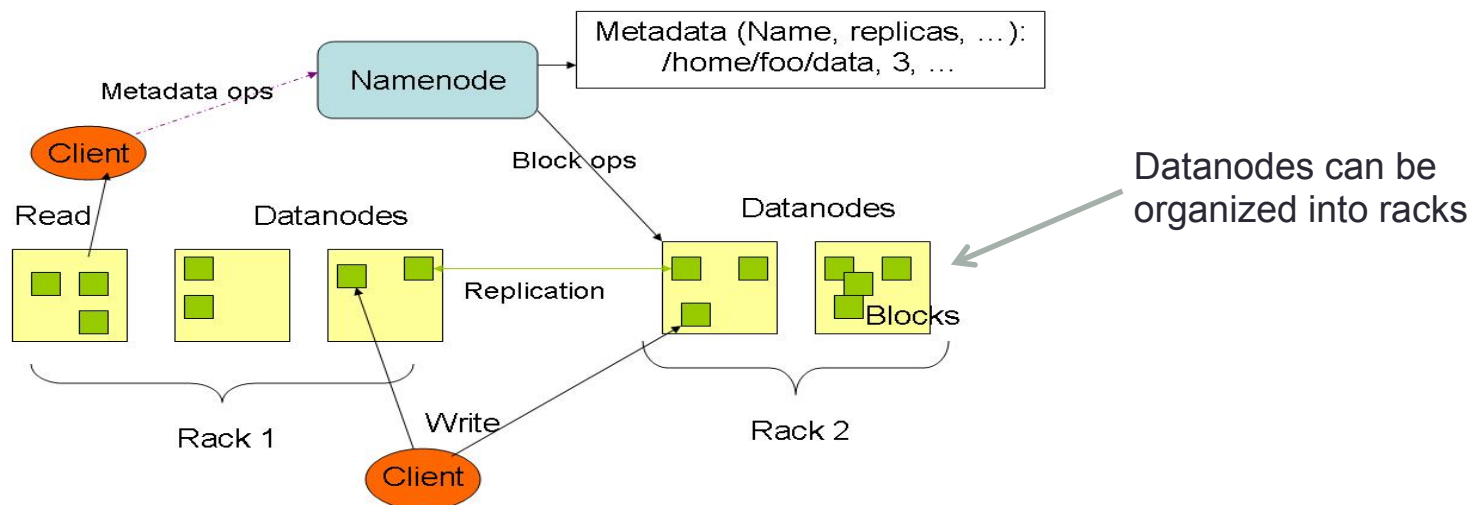
HDFS: Hadoop Distributed File System



- HDFS is a master-slave architecture
 - Master: namenode
 - Slaves: datanodes (100s or 1000s of nodes)

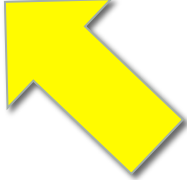
HDFS: Data Placement and Replication

HDFS Architecture

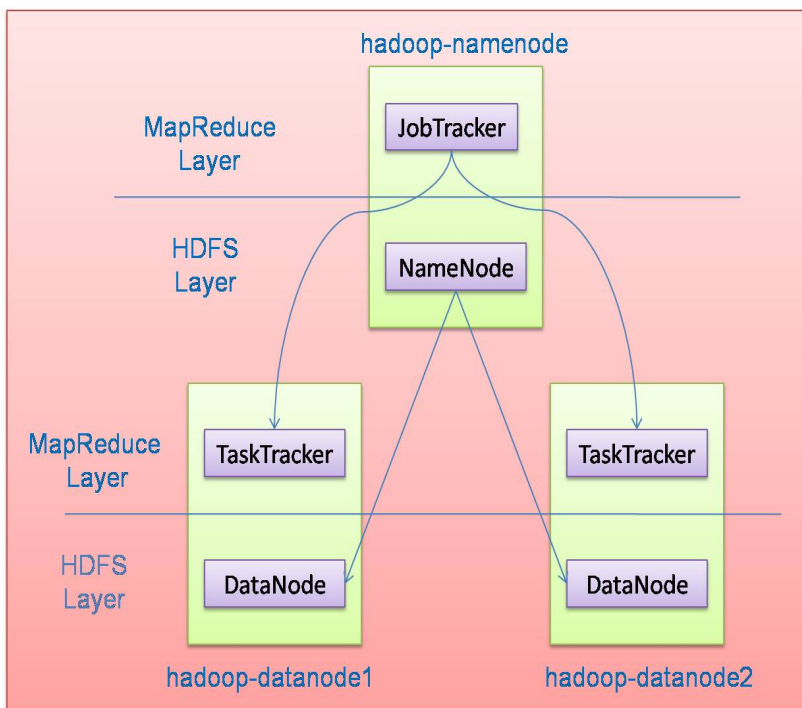


- Default placement policy: Where to put a given block?
 - **First copy** is written to the node creating the file (write affinity)
 - **Second copy** is written to a data node within the same rack
 - **Third copy** is written to a data node in a different rack
 - **Objectives:** *load balancing, fast access, fault tolerance*

What is Next?

- **Hadoop Distributed File System (HDFS)**
- **MapReduce Layer** 
- **Examples**
 - **Word Count**
 - **Join**
- **Fault Tolerance in Hadoop**

MapReduce: Hadoop Execution Layer



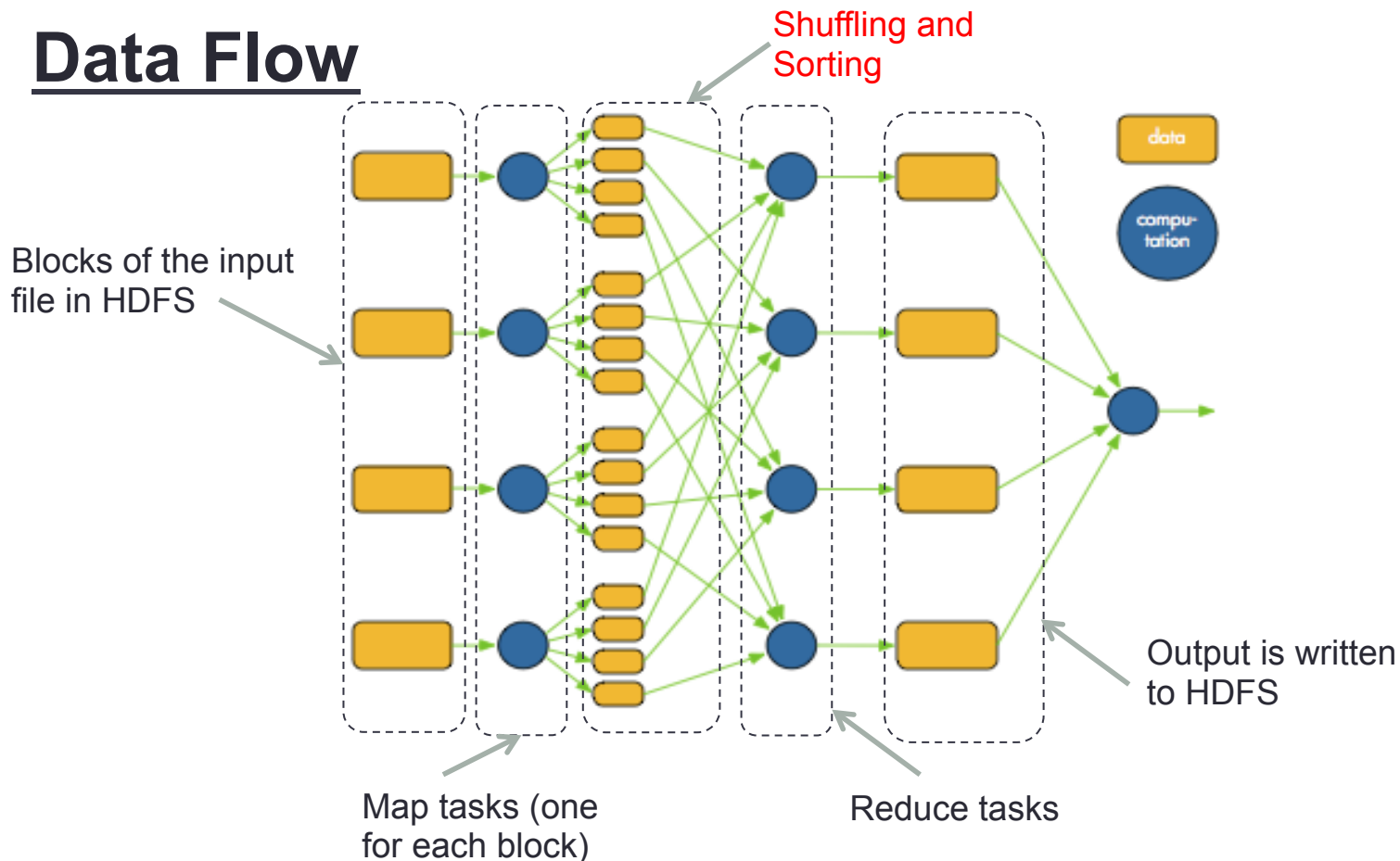
- Jobtracker knows everything about submitted jobs
- Divides jobs into tasks and decides where to run each task
- Continuously communicating with tasktrackers
- Tasktrackers execute tasks (multiple per node)
- Monitors the execution of each task
- Continuously sending feedback to Jobtracker

- MapReduce is a master-slave architecture
 - Master: JobTracker
 - Slaves: TaskTrackers (100s or 1000s of tasktrackers)
- Every datanode is running a tasktracker

Hadoop Computing Model

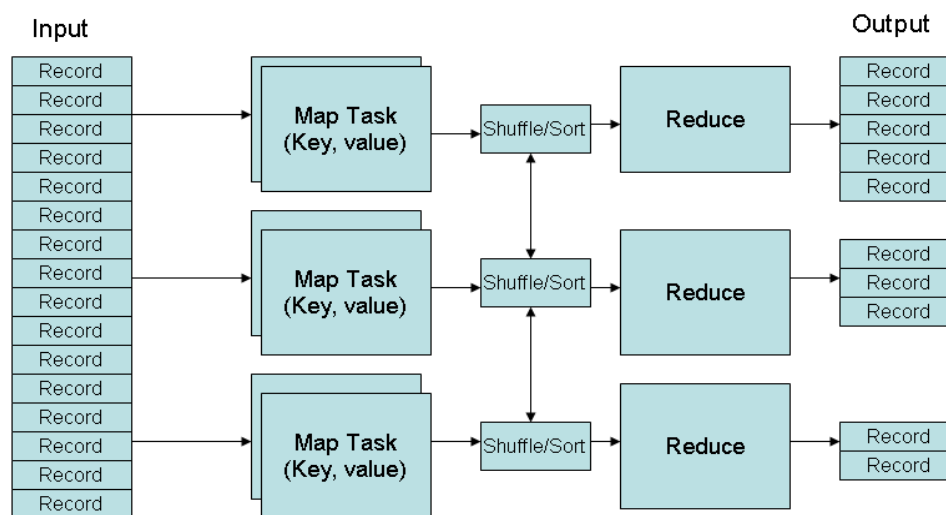
- Two main phases: Map and Reduce
- Any job is converted into *map* and *reduce* tasks
- Developers need *ONLY* to implement the *Map* and *Reduce* classes

Data Flow



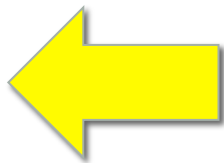
Hadoop Computing Model (Cont'd)

- Mapper and Reducers consume and produce *(Key, Value)* pairs
- Users define the data type of the *Key* and *Value*
- **Shuffling & Sorting phase:**
 - Map output is *shuffled* such that all same-key records go to the same reducer
 - Each reducer may receive multiple key sets
 - Each reducer *sorts* its records to group similar keys, then process each group



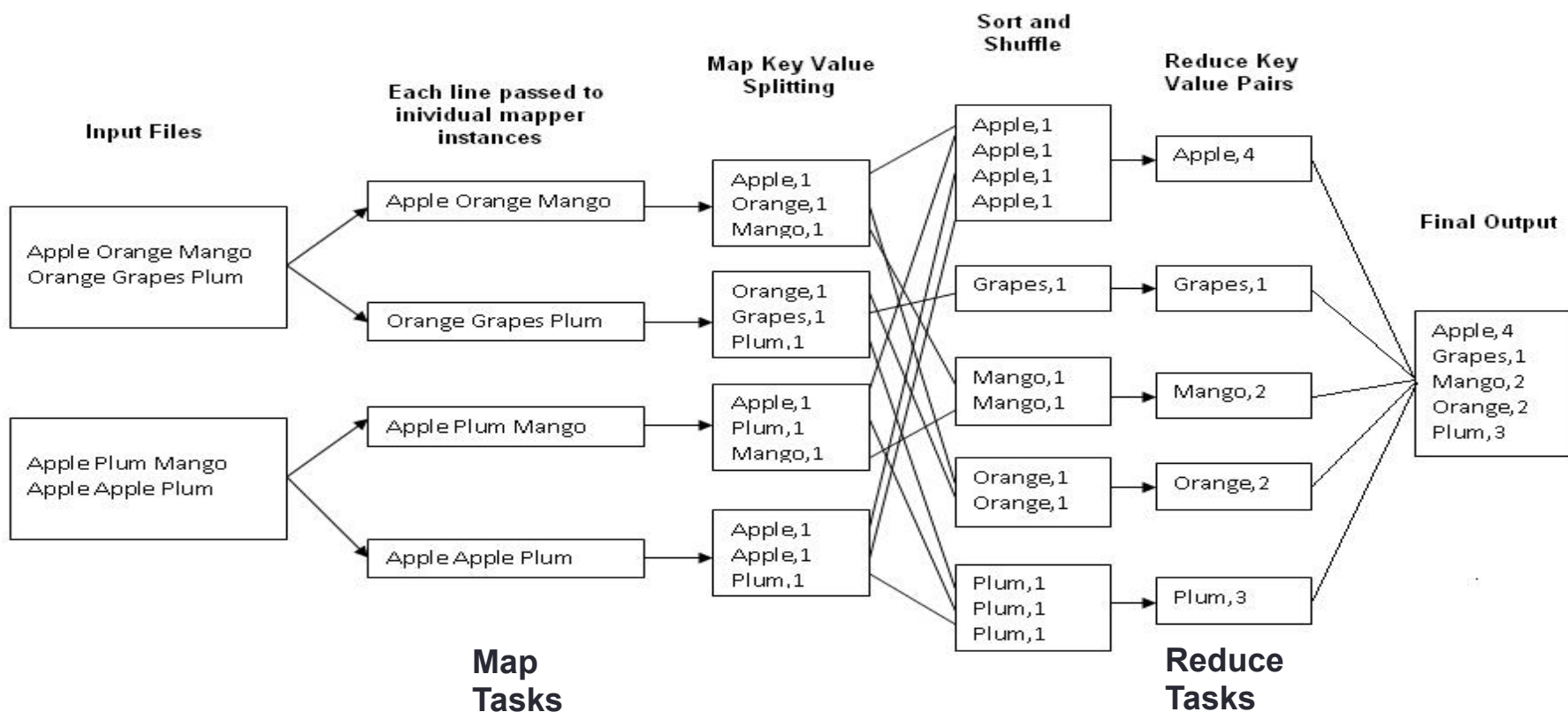
What is Next?

- **Hadoop Distributed File System (HDFS)**
- **MapReduce Layer**
- **Examples**
 - **Word Count**
 - **Join**
- **Fault Tolerance in Hadoop**



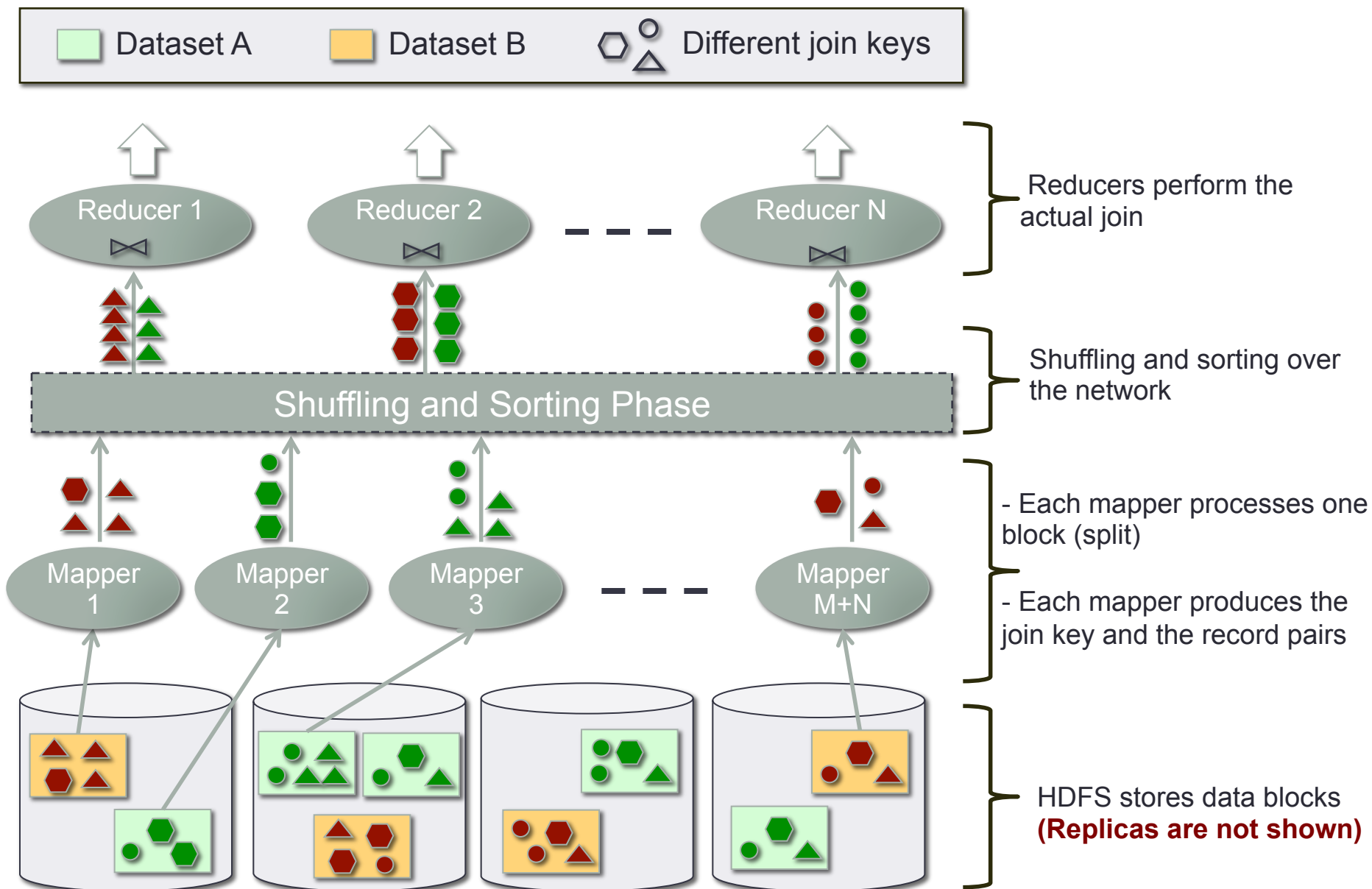
Word Count

- Job: Count the occurrences of each word in a data set

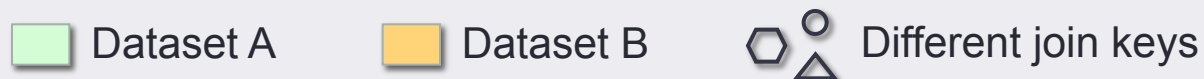


Reduce phase is optional: Jobs can be MapOnly

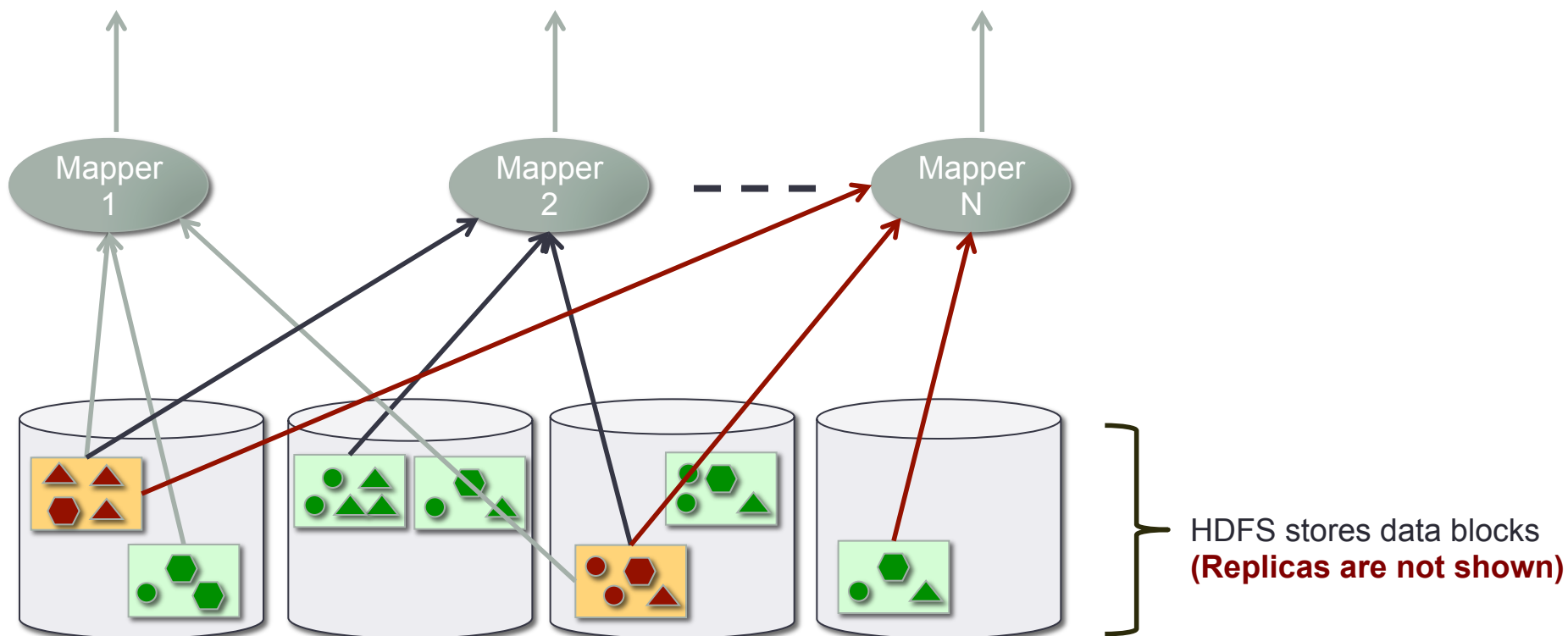
Joining Two Large Datasets



Joining Large Dataset (A) with Small Dataset (B)

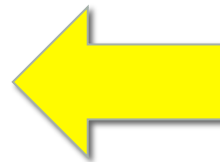


- Every map task processes one block from A and the entire B
- Every map task performs the join (**MapOnly job**)
- Avoid the shuffling and reduce *expensive* phases



What is Next?

- **Hadoop Distributed File System (HDFS)**
- **MapReduce Layer**
- **Examples**
 - **Word Count**
 - **Join**
- **Fault Tolerance in Hadoop**



Hadoop Fault Tolerance

- Intermediate data between mappers and reducers are *materialized* to simple & straightforward fault tolerance

- **What if a task fails (map or reduce)?**

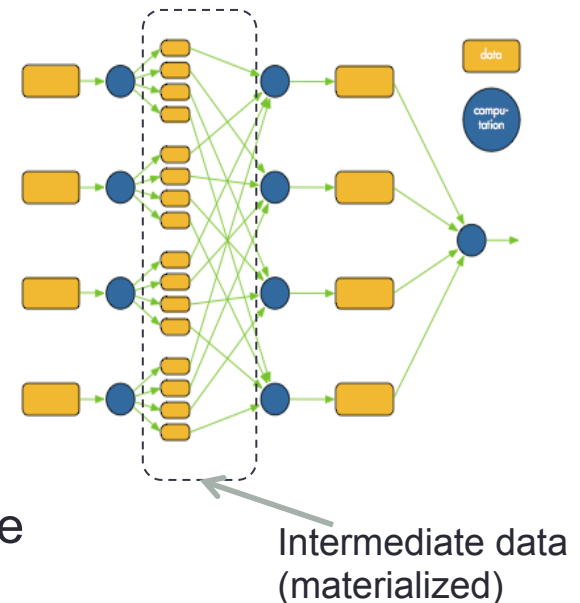
- Tasktracker detects the failure
- Sends message to the jobtracker
- Jobtracker re-schedules the task

- **What if a datanode fails?**

- Both namenode and jobtracker detect the failure
- All tasks on the failed node are re-scheduled
- Namenode replicates the users' data to another node

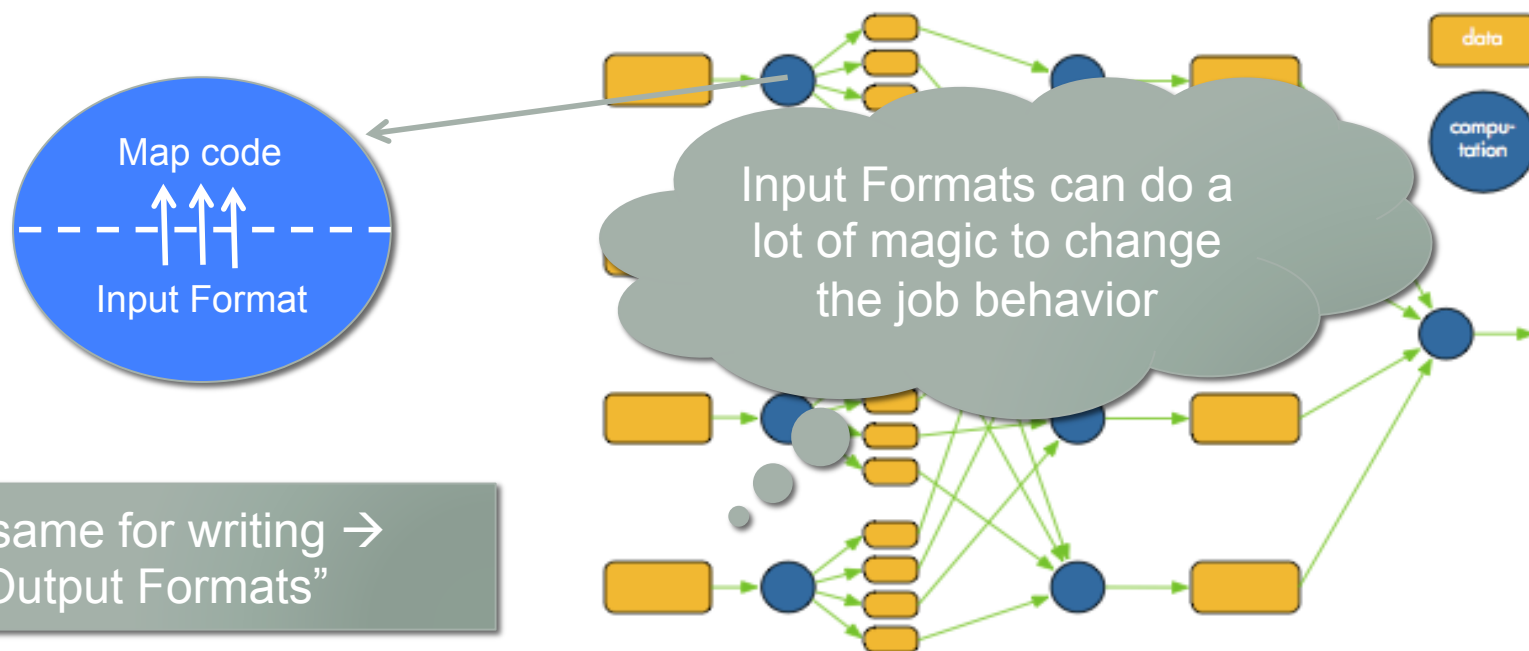
- **What if a namenode or jobtracker fails?**

- The entire cluster is down

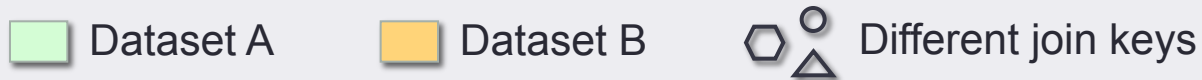


Reading/Writing Files

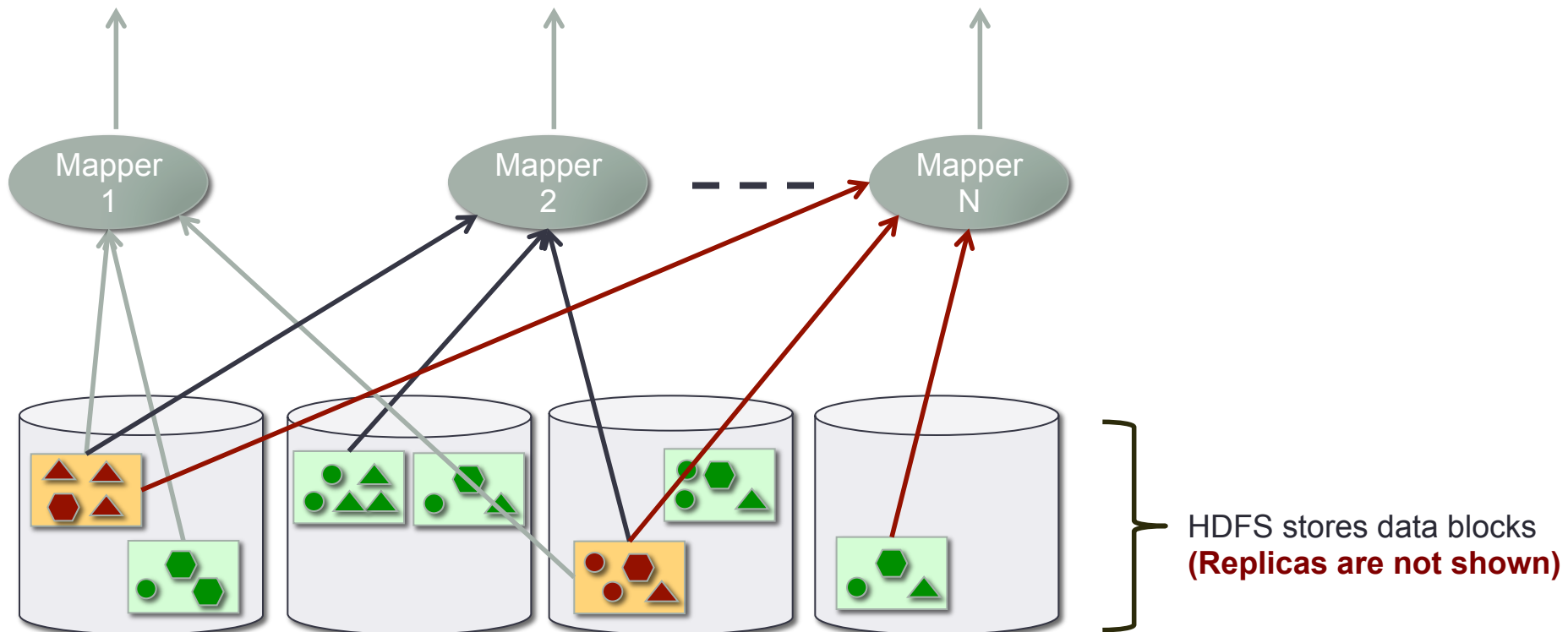
- **Recall: Any data will fit in Hadoop, so how does Hadoop understand/read the data?**
- User-pluggable class *“Input Format”*
 - Input formats know how to parse and read the data (convert byte stream to records)
 - Each record is then passed to the mapper for processing
- Hadoop provides *built-in* Input Formats for reading *text* & *sequence* files



Back to Joining Large & Small Datasets



- Every map task processes one block from A and the entire B
- How does a single mapper reads multiple splits (from different datasets)?
 - Customized input formats



Using Hadoop

- **Java language**
- **High-Level Languages**
 - Hive (Facebook)
 - Pig (Yahoo)
 - Jaql (IBM)

Java Code Example

Source Code

```

WordCount.java
1. package org.myorg;
2.
3. import java.io.IOException;
4. import java.util.*;
5.
6. import org.apache.hadoop.fs.Path;
7. import org.apache.hadoop.conf.*;
8. import org.apache.hadoop.io.*;
9. import org.apache.hadoop.mapred.*;
10. import org.apache.hadoop.util.*;
11.
12. public class WordCount {
13.
14.     public static class Map extends MapReduceBase implements Mapper<LongWritable, Text, Text, IntWritable> {
15.         private final static IntWritable one = new IntWritable(1);
16.         private Text word = new Text();
17.
18.         public void map(LongWritable key, Text value, OutputCollector<Text, IntWritable> output, Reporter reporter) throws IOException {
19.             String line = value.toString();
20.             StringTokenizer tokenizer = new StringTokenizer(line);
21.             while (tokenizer.hasMoreTokens()) {
22.                 word.set(tokenizer.nextToken());
23.                 output.collect(word, one);
24.             }
25.         }
26.     }
27.
28.     public static class Reduce extends MapReduceBase implements Reducer<Text, IntWritable, Text, IntWritable> {
29.         public void reduce(Text key, Iterator<IntWritable> values, OutputCollector<Text, IntWritable> output, Reporter reporter) throws IOException {
30.             int sum = 0;
31.             while (values.hasNext()) {
32.                 sum += values.next().get();
33.             }
34.             output.collect(key, new IntWritable(sum));
35.         }
36.     }
37.
38.     public static void main(String[] args) throws Exception {
39.         JobConf conf = new JobConf(WordCount.class);
40.         conf.setJobName("wordcount");
41.
42.         conf.setOutputKeyClass(Text.class);
43.         conf.setOutputValueClass(IntWritable.class);
44.
45.         conf.setMapperClass(Map.class);
46.         conf.setCombinerClass(Reduce.class);
47.         conf.setReducerClass(Reduce.class);
48.
49.         conf.setInputFormat(TextInputFormat.class);
50.         conf.setOutputFormat(TextOutputFormat.class);
51.
52.         FileInputFormat.setInputPaths(conf, new Path(args[0]));
53.         FileOutputFormat.setOutputPath(conf, new Path(args[1]));
54.
55.         JobClient.runJob(conf);
56.     }
57. }
58. }

```

Import Hadoop libs

Map class

Reduce class

Job configuration

Hive Language

- High-level language on top of Hadoop
 - Like SQL on top of DBMSs
- Support structured data, e.g., creating tables, as well as extensibility for un-structured data

page_view			user			pv_users	
pageid	userid	time	userid	age	gender	pageid	age
1	111	9:08:01	111	25	female	1	25
2	111	9:08:13	222	32	male	2	25
1	222	9:08:14				1	32

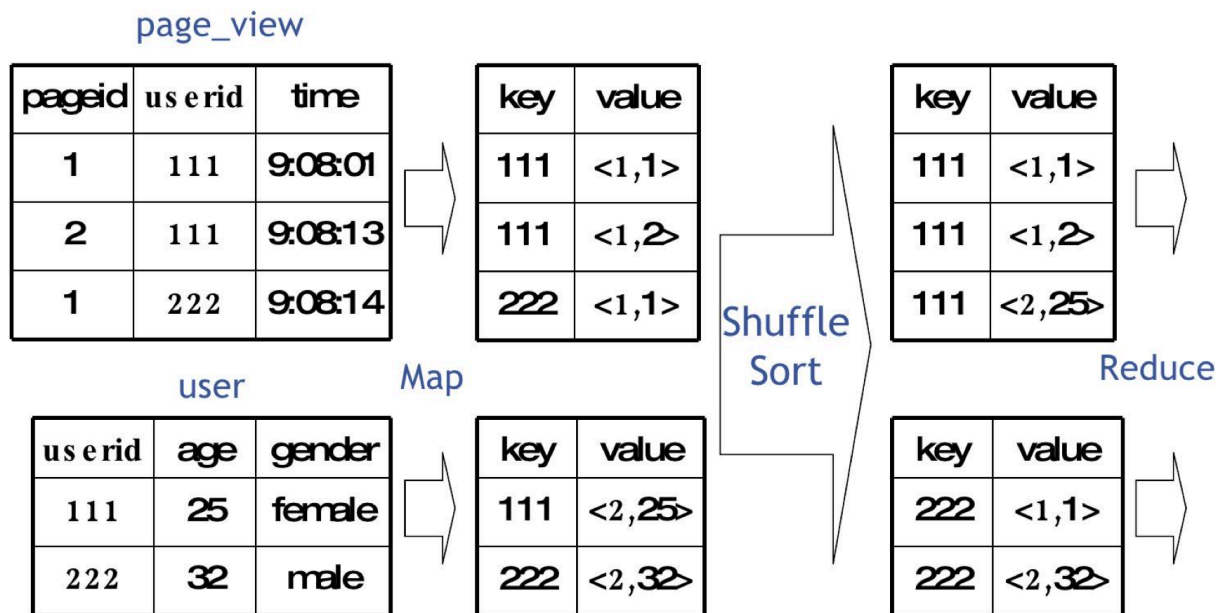
```
INSERT INTO TABLE pv_users
SELECT pv.pageid, u.age
FROM page_view pv JOIN user u ON (pv.userid = u.userid);
```

Create Table user (userID int,
age int,
gender char)
Row Format Delimited Fields;

Load Data Local Inpath '/user/
local/users.txt' into Table user;

From Hive To MapReduce

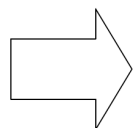
```
INSERT INTO TABLE pv_users
SELECT pv.pageid, u.age
FROM page_view pv JOIN user u ON (pv.userid = u.userid);
```



Hive: Group By

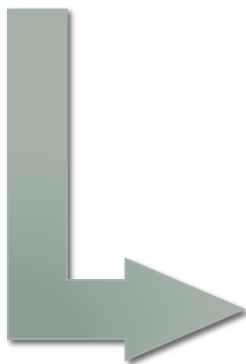
pv_users

pageid	age
1	25
2	25
1	32
2	25



pageid	age	count
1	25	1
2	25	2
1	32	1

```
SELECT pageid, age, count(1)
FROM pv_users
GROUP BY pageid, age;
```



pv_users

pageid	age
1	25
2	25



key	value
<1,25>	1
<2,25>	1

Map



pageid	age
1	32
2	25

key	value
<1,32>	1
<2,25>	1

Shuffle
Sort

key	value
<1,25>	1
<1,32>	1

Reduce



key	value
<2,25>	1
<2,25>	1

Summary

- **Hadoop is a distributed systems for processing large-scale datasets**
 - Scales to thousands of nodes and petabytes of data
- **Two main layers**
 - HDFS: Distributed file system(NameNode is centralized)
 - MapReduce: Execution engine (JobTracker is centralized)
- **Simple data model, any format will fit**
 - At query time, specify how to read (write) the data using input (output) formats
- **Simple computation model based on Map-Reduce phases**
 - Very efficient in aggregation and joins
- **Higher-level Languages on top of Hadoop**
 - Hive, Jaql, Pig

Summary: Hadoop vs. Other Systems

	Distributed Databases	Hadoop
Computing Model	<ul style="list-style-type: none"> - Notion of transactions - Transaction is the unit of work - ACID properties, Concurrency control 	<ul style="list-style-type: none"> - Notion of jobs - Job is the unit of work - No concurrency control
Data Model	<ul style="list-style-type: none"> - Structured data with known schema - Read/Write mode 	<ul style="list-style-type: none"> - Any data will fit in any format (un)(semi)structured - ReadOnly mode
Cost Model	<ul style="list-style-type: none"> - Expensive servers 	<ul style="list-style-type: none"> - Cheap commodity machines
Fault Tolerance	<ul style="list-style-type: none"> - Failures are rare - Recovery mechanisms 	<ul style="list-style-type: none"> - Failures are common over thousands of machines - Simple yet efficient fault tolerance
Key Characteristics	<ul style="list-style-type: none"> - Efficiency, optimizations, fine-tuning 	<ul style="list-style-type: none"> - Scalability, flexibility, fault tolerance

• Cloud Computing

- A computing model where any computing infrastructure can run on the cloud
- Hardware & Software are provided as remote services
- Elastic: grows and shrinks based on the user's demand
- Example: Amazon EC2

