

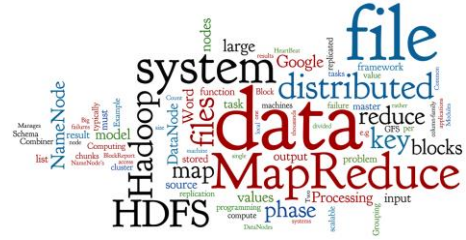


# Agenda



- Distributed Data Processing
- Google MapReduce
  - Motivation and History
  - Google File System (GFS)
  - MapReduce: Schema, Example, MapReduce Framework
- Apache Hadoop
  - Hadoop Modules and Related Projects
  - Hadoop Distributed File System (HDFS)
  - Hadoop MapReduce
- MapReduce in Other Systems

# Distributed Data Processing



What is the **best way** of doing **distributed** processing?

**Centralized** (and in memory)

**Don't do it, if don't have to**

# Big Data Processing



- Big Data **analytics** (or data mining)
    - a need to process **large** data **volumes** quickly,
    - run on a computing **cluster** instead of a super-computer
  - Communication (**sending data**) between compute nodes is **expensive**
- => model of “moving the computing to data”

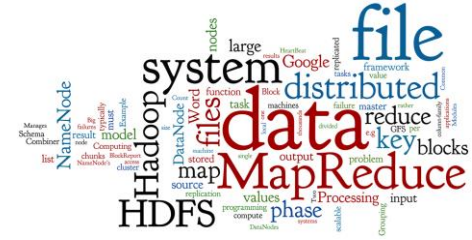


# Agenda



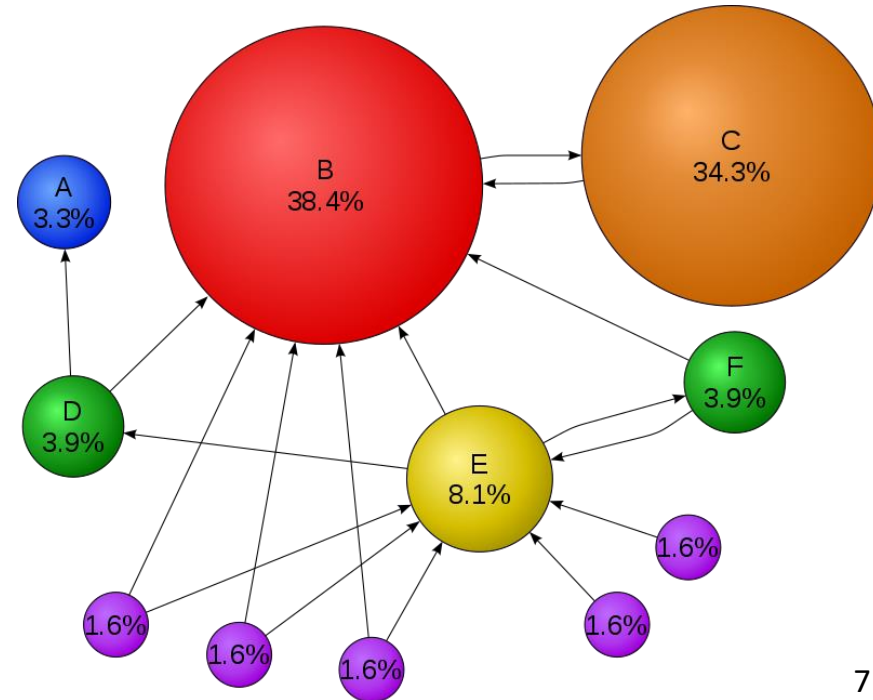
- Distributed Data Processing
- Google MapReduce
  - **Motivation** and History
  - Google File System (**GFS**)
  - **MapReduce**: Schema, Example, MapReduce Framework
- Apache Hadoop
  - Hadoop Modules and Related Projects
  - Hadoop Distributed File System (HDFS)
  - Hadoop MapReduce
- MapReduce in Other Systems

# PageRank



PageRank works by counting the **number** and **quality of links** to a page to determine a rough estimate of **how important** the website is.

The underlying assumption is that more important websites are likely to receive more links from other websites.



# MapReduce: Origins



- In 2003, **Google** had the following **problem**:
  1. How to **rank tens of billions** of webpages by their “importance” (PageRank) in a “reasonable” amount of time?
  2. How to **compute** these rankings **efficiently** when the data is scattered across **thousands** of **computers**?
- Additional factors:
  1. Individual data **files** can be enormous (**terabyte** or more)
  2. The files were **rarely updated**
    - the computations were **read-heavy**, but not very write-heavy
    - If **writes** occurred, they were **appended** at the end of the file





# Google File System (GFS)



- **Files** are divided **into chunks** (typically 64 MB)
  - The **chunks** are **replicated** at three different machines
    - ...in an “intelligent” fashion, e.g. never all on the same computer rack
  - The chunk size and replication factor are **tunable**
- One machine is a **master**, the other **chunkserver**s
  - The **master** keeps track of all file **metadata**
    - mappings from files to chunks and locations of the chunks
  - To find a file chunk, **client** queries the **master**, and then contacts the relevant **chunkserver**s
  - The master’s metadata files are also replicated



# MapReduce (1)



- MapReduce is a **programming model** sitting **on the top** of a Distributed File System
  - Originally: **no data model** – data stored directly in **files**
- A **distributed** computational **task** has three phases:
  1. The **map** phase: data transformation
  2. The **grouping** phase
    - done automatically by the MapReduce Framework
  3. The **reduce** phase: data aggregation
- User must define **only** map & reduce **functions**

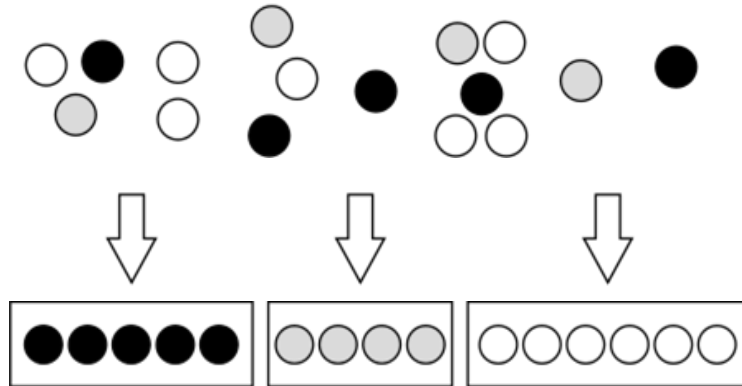




# Grouping Phase



- **Grouping** (Shuffling): The key-value **outputs** from the **map** phase are grouped by key
  - Values sharing **the same key** are sent to the same reducer.
  - These values are **consolidated** into a single list (**key, list**).
    - This is convenient for the reduce function
  - This phase is **realized by** the MapReduce **framework**.



intermediate output  
(*color indicates key*)

shuffle (grouping) phase







# Example: Word Count



Task: Calculate **word frequency** in a set of documents

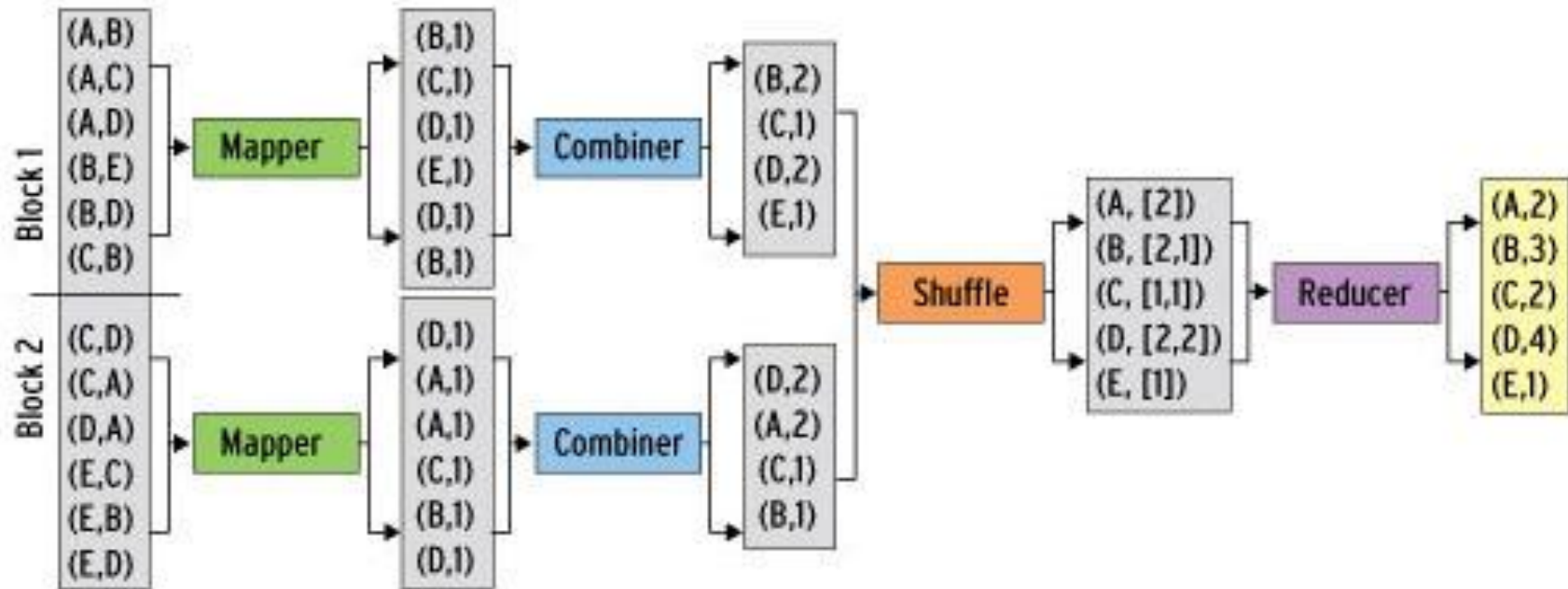
```
map(String key, Text value):  
    // key: document name (ignored)  
    // value: content of document (words)  
    foreach word w in value:  
        emitIntermediate(w, 1);  
  
    reduce(String key, Iterator values):  
        // key: a word  
        // values: a list of counts  
        int result = 0;  
        foreach v in values:  
            result += v;  
        emit(key, result);
```







# Example: Word Count with Combiner









# MapReduce Framework: Details



## 1. Input reader (function)

- defines how to read data from underlying storage

## 2. Map (phase)

- master node prepares  $M$  data splits and  $M$  idle Map tasks
- pass individual splits to the Map tasks that run on workers
- these map tasks are then running
- when a task is finished, its intermediate results are stored

## 3. Combiner (function, optional)

- combine local intermediate output from the Map phase

# MapReduce Framework: Details (2)



## 4. Partition (function)

- to **partition** intermediate results for individual **Reducers**

## 5. Comparator (function)

- sort and **group** the input for each Reducer

## 6. Reduce (phase)

- **master** node creates  $R$  **idle** Reduce tasks on **workers**
- **Partition** function **defines** a data **batch** for each reducer
- each Reduce task uses **Comparator** to create **key-values pairs**
- function Reduce is **applied** on each key-values pair

## 7. Output writer (function)

- defines how the **output** key-value pairs are **written out**

# MapReduce: Example II



Task: Calculate a **graph** of web links

- what pages reference (<a href="">) each page (backlinks)

```
map(String url, Text html):  
    // url: web page URL  
    // html: HTML text of the page (linearized HTML tags)  
    foreach tag t in html:  
        if t is <a> then:  
            emitIntermediate(t.href, url);
```

```
reduce(String key, Iterator values):  
    // key: target URLs  
    // values: a list of source URLs  
    emit(key, values);
```

# Example II: Result



**Input:** (page\_URL, HTML\_code)

```
("http://cnn.com", "<html>...<a href="http://cnn.com">link</a>...</html>")
("http://ihned.cz", "<html>...<a href="http://cnn.com">link</a>...</html>")
("http://idnes.cz",
 "<html>...<a href="http://cnn.com">x</a>...
  <a href="http://ihned.cz">y</a>...<a href="http://idnes.cz">z</a>
  </html>")
```

Intermediate output **after Map** phase:

```
("http://cnn.com", "http://cnn.com")
("http://cnn.com", "http://ihned.cz")
("http://cnn.com", "http://idnes.cz")
("http://ihned.cz", "http://idnes.cz")
("http://idnes.cz", "http://idnes.cz")
```

Intermediate result **after shuffle** phase (the same as output **after Reduce** phase):

```
("http://cnn.com", ["http://cnn.com", "http://ihned.cz", "http://idnes.cz"] )
("http://ihned.cz", [ "http://idnes.cz" ])
("http://idnes.cz", [ "http://idnes.cz" ])
```





# Applicability of MapReduce



- MR is applicable if the problem is **parallelizable**.
- Two problems:
  1. The programming **model is limited**  
(only two phases with a **given schema**)
  2. There is **no data model** - it works only on “data chunks”
- Google’s **answer** to the 2nd problem was **BigTable**
  - The first **column-family** system (2005)
  - Subsequent systems: HBase (over Hadoop), Cassandra,...

# Agenda



- Distributed Data Processing
- Google MapReduce
  - Motivation and History
  - Google File System (GFS)
  - MapReduce: Schema, Example, MapReduce Framework
- Apache **Hadoop**
  - Hadoop **Modules** and Related Projects
  - Hadoop Distributed File System (**HDFS**)
  - Hadoop **MapReduce**
- MapReduce in Other Systems















# HDFS: NameNode



- **NameNode** has a structure called **FsImage**
  - Entire **file system** namespace + mapping of **blocks** to files + file system properties
  - Stored in a file in NameNode's local file system
  - Designed to be **compact**
    - Loaded in NameNode's memory (4 GB of RAM is sufficient)
- **NameNode** uses a **transaction log** called **EditLog**
  - to **record** every **change** to the file system's meta data
    - E.g., creating a new file, change in replication factor of a file, ..
  - EditLog is stored in the NameNode's local file system

# HDFS: DataNode



- Stores **data in files** on its local file system
  - Each HDFS **block** in a **separate file**
  - Has **no knowledge** about **HDFS** file system
- When the DataNode **starts up**:
  - It **generates** a list of all HDFS blocks = **BlockReport**
  - It sends the report to NameNode





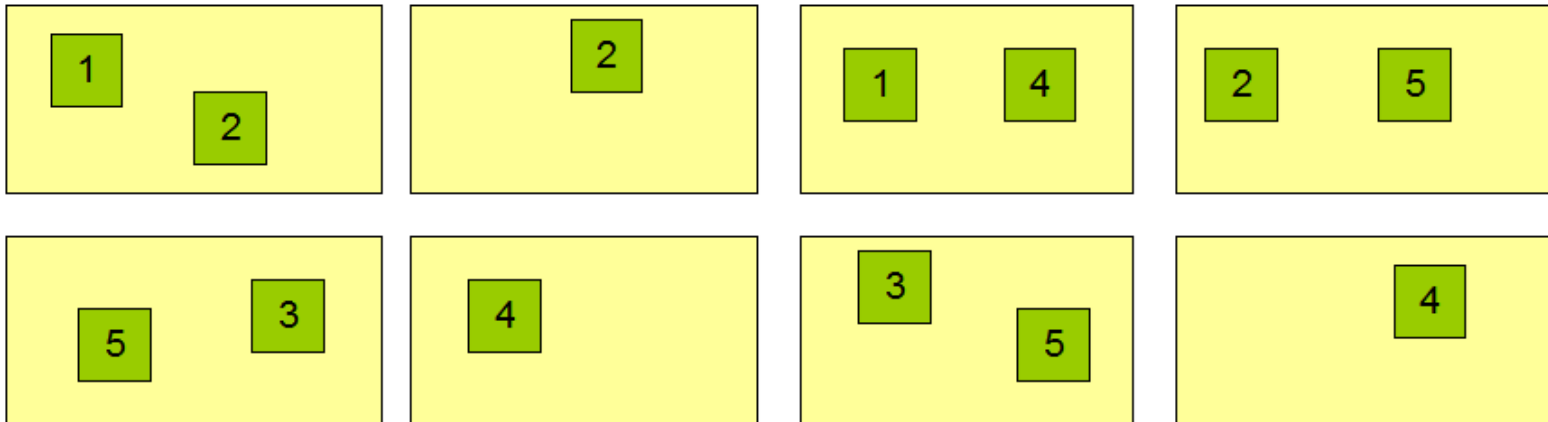
# HDFS: Block Replication



## Block Replication

Namenode (Filename, numReplicas, block-ids, ...)  
/users/sameerp/data/part-0, r:2, {1,3}, ...  
/users/sameerp/data/part-1, r:3, {2,4,5}, ...

## Datanodes











# Hadoop MR: WordCount Example



```
public class Map
    extends Mapper<LongWritable, Text, Text, IntWritable> {

    private final static IntWritable one = new IntWritable(1);
    private final Text word = new Text();

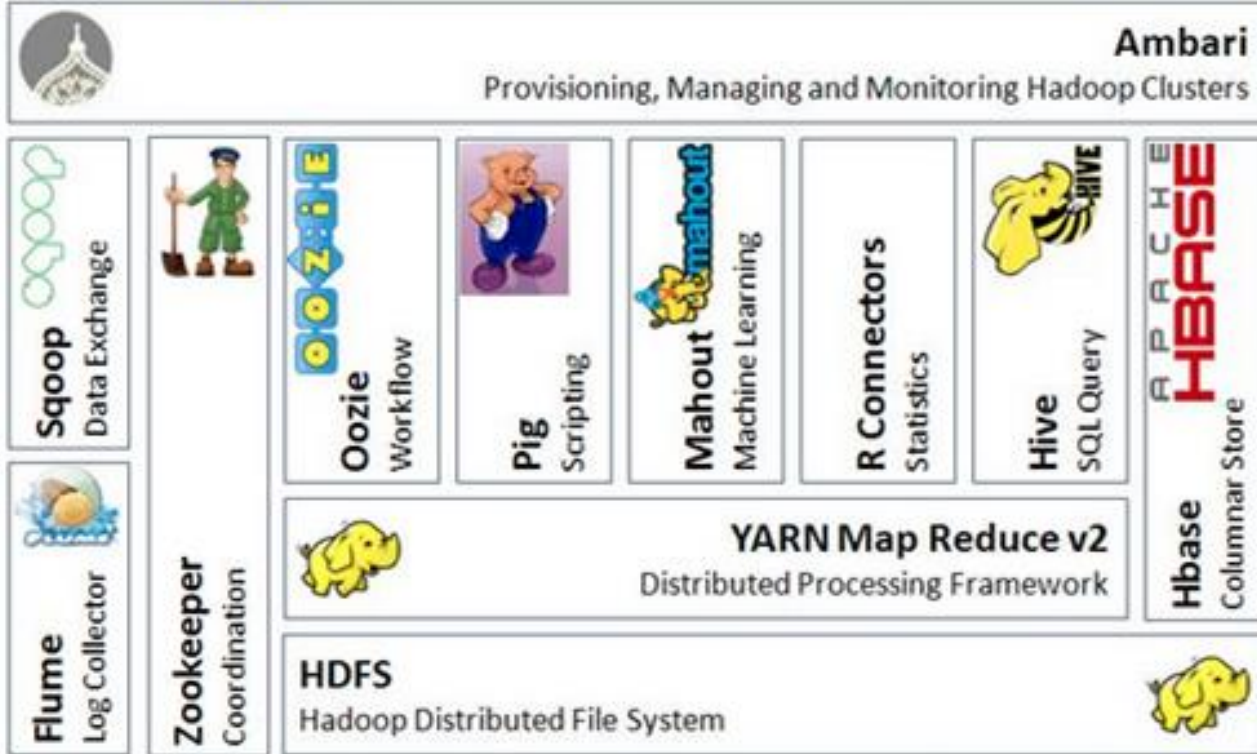
    @Override protected void map(LongWritable key, Text value,
        Context context) throws ... {
        String string = value.toString()
        StringTokenizer tokenizer = new StringTokenizer(string);
        while (tokenizer.hasMoreTokens()) {
            word.set(tokenizer.nextToken());
            context.write(word, one);
        }
    }
}
```







# Apache Hadoop Ecosystem







# MapReduce: Implementation



Amazon Elastic  
MapReduce



redis



Cassandra

# Apache Spark



- **Engine** for **distributed** data processing
  - **Runs** over Hadoop Yarn, Apache Mesos, standalone, ...
  - Can **access data** from HDFS, Cassandra, HBase, AWS S3
- Can do **MapReduce**
  - Is much **faster** than pure Hadoop
    - They say 10x on the disk, 100x in memory
  - The main reason: **intermediate** data in **memory**
- Different **languages** to write MapReduce tasks
  - Java, Scala, Python, R





# References



- RNDr. Irena Holubova, Ph.D. MMF UK course NDBI040: Big Data Management and NoSQL Databases
- Dean, J. & Ghemawat, S. MapReduce: Simplified Data Processing on Large Clusters. In OSDI 2004 (pp 137-149)
- Firas Abuzaid, Perth Charernwattanakul (2014). Lecture 8 “NoSQL” of Stanford course CS145. [link](#)
- J. Leskovec, A. Rajaraman, and J. D. Ullman, Mining of Massive Datasets. 2014.
- I. Holubová, J. Kosek, K. Minařík, D. Novák. Big Data a NoSQL databáze. Praha: Grada Publishing, 2015. 288 p.