



Distributed Computing with MapReduce

Lecture 2 of *NoSQL Databases* (PA195)

David Novak & Vlastislav Dohnal
Faculty of Informatics, Masaryk University, Brno

Distributed Data Processing



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

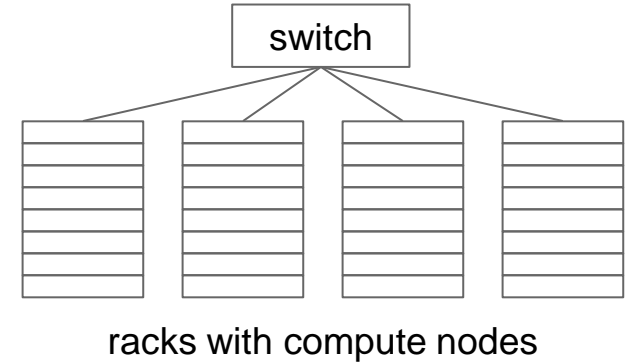
Centralized (and in memory)

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

Big Data Processing II



Computing **cluster** architecture:



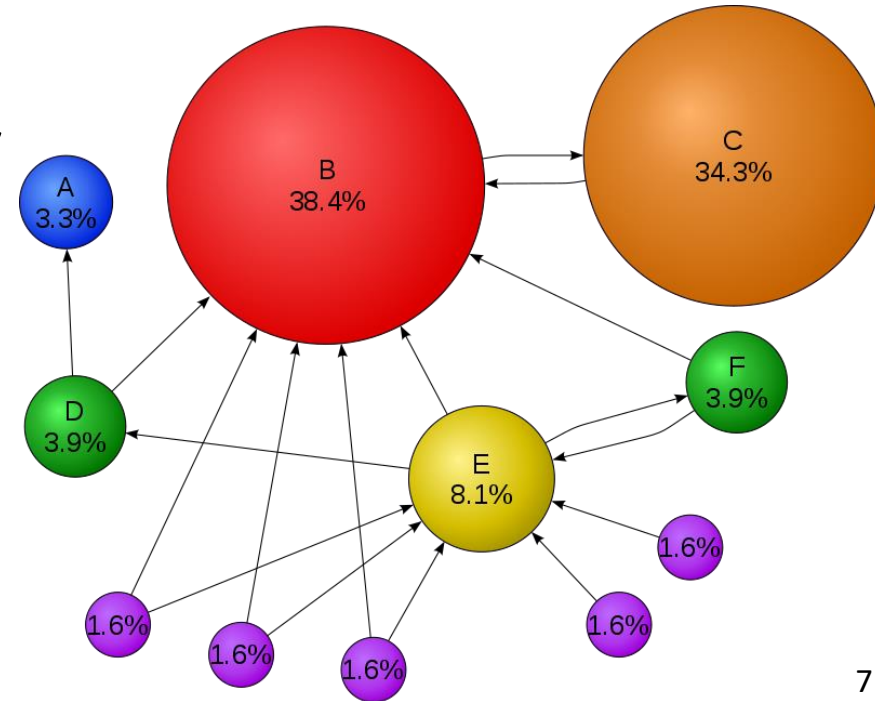
- HW **failures** are rather a rule than an exception:
 1. **Files** must be stored **redundantly**
 - over different **racks** to overcome also rack failures
 2. **Computations** must be divided into independent **tasks**
 - that can be **restarted** in case of a failure

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.



Google Solution



- Google found the following **solutions**:
 - Google File System (GFS)
 - A **distributed file system**
 - MapReduce
 - A programming **model** for **distributed data processing**

GFS: Schema

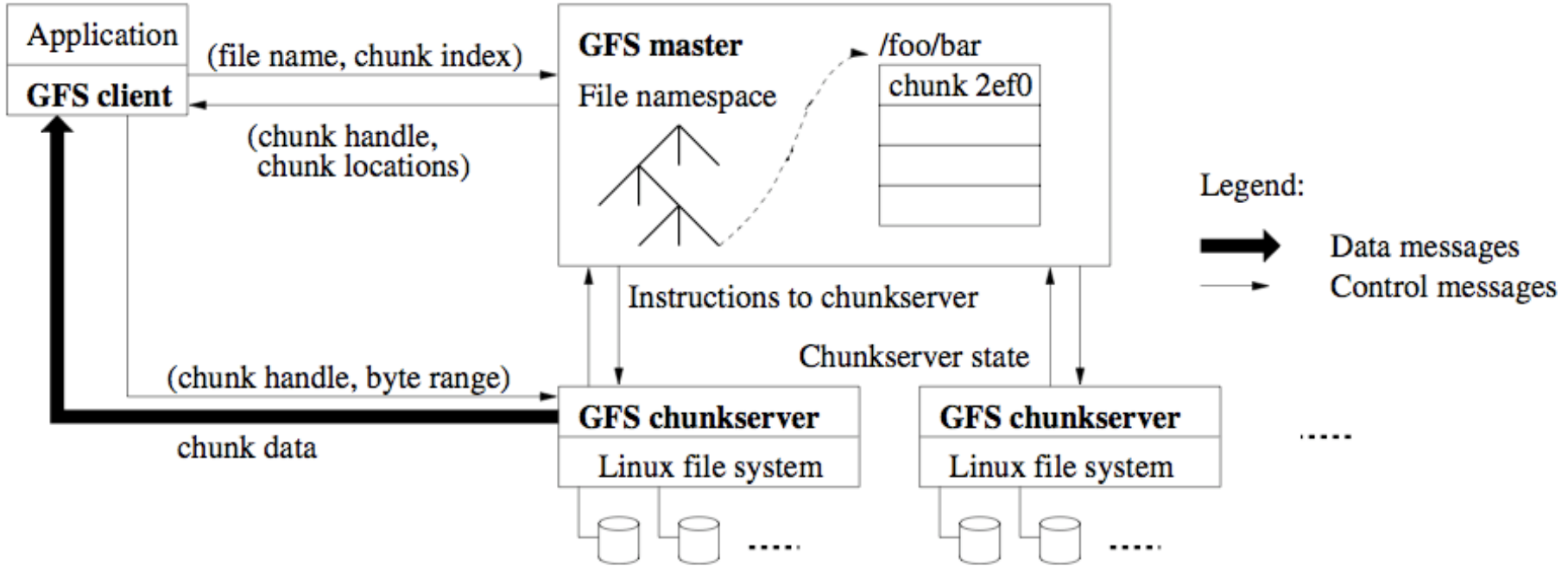
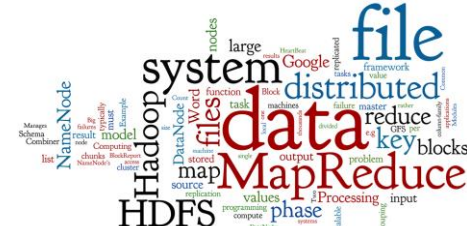


Figure 1: GFS Architecture

MapReduce (1)

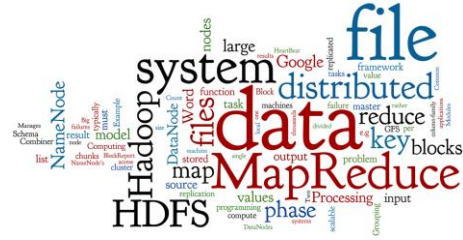


- MapReduce is a **programming model** sitting **on the top** of a Distributed File System
 - Originally: **no data model** – data stored directly in **files**
- A **distributed** computational **task** has three phases:
 1. The **map** phase: data transformation
 2. The **grouping** phase
 - done automatically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

Map



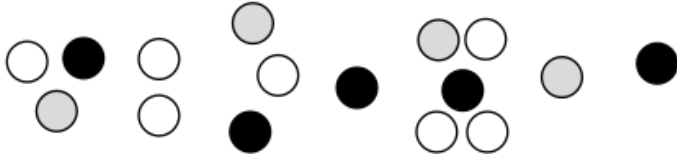
- **Map function** simplifies the problem in this way:
 - **Input**: a **single data item** (e.g., line of text) from a data file
 - **Output**: zero or more **(key, value) pairs**
- The **keys are not typical “keys”**:
 - They do **not** have to be **unique**
 - A map task can produce **several key-value pairs** with the same key (even from a single input)
- **Map phase** applies the map function to **all** items.



input data



map function



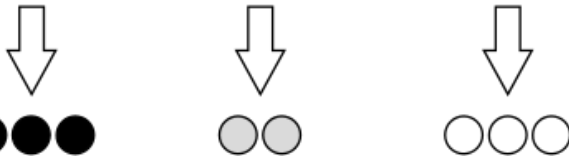
intermediate output
(color indicates key)



shuffle (grouping) phase



input data

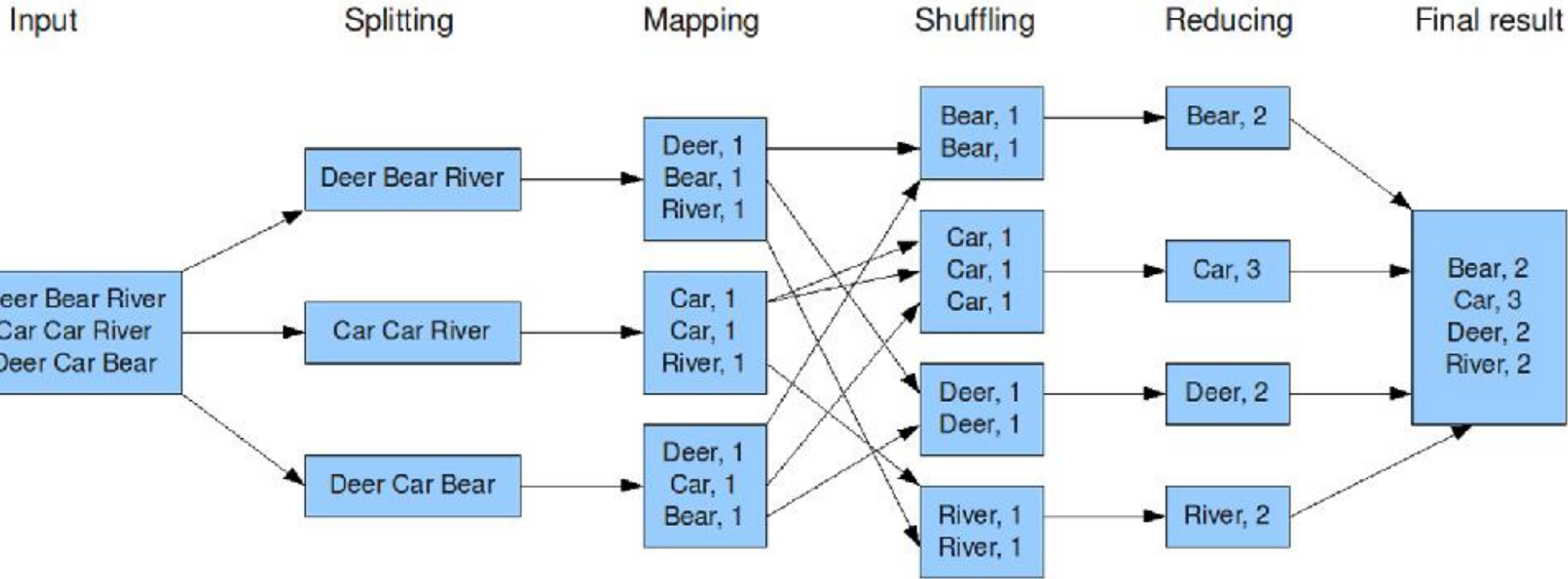


reduce function



output data

Example: Word Count (2)



MapReduce: Combiner



- If the **reduce** function is **commutative & associative**
 - The values **can be combined** in any order and **combined per partes** (grouped)
 - with the same result (e.g., Word Counts)
- ...then we can do "**partial reductions**"
 - Apply the **same reduce function** right after the map phase, **before shuffling** and redistribution to reducer nodes
- This (optional) step is known as the **combiner**
 - Note: it is still necessary to run the reduce phase.

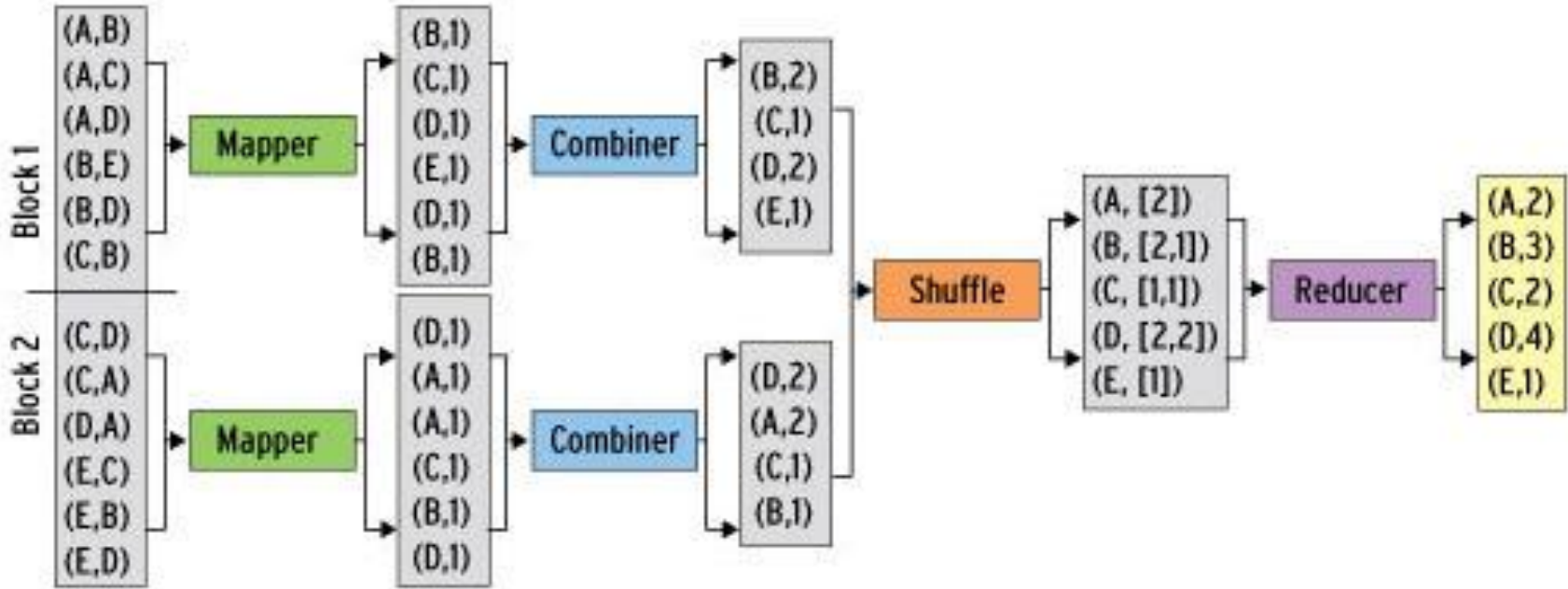
Example: Word Count, Combiner



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

```
combine(String key, Iterator values):  
    // key: a word  
    // values: a list of local counts  
int result = 0;  
foreach v in values:  
    result += v;  
emit(key, result);
```

Example: Word Count with Combiner

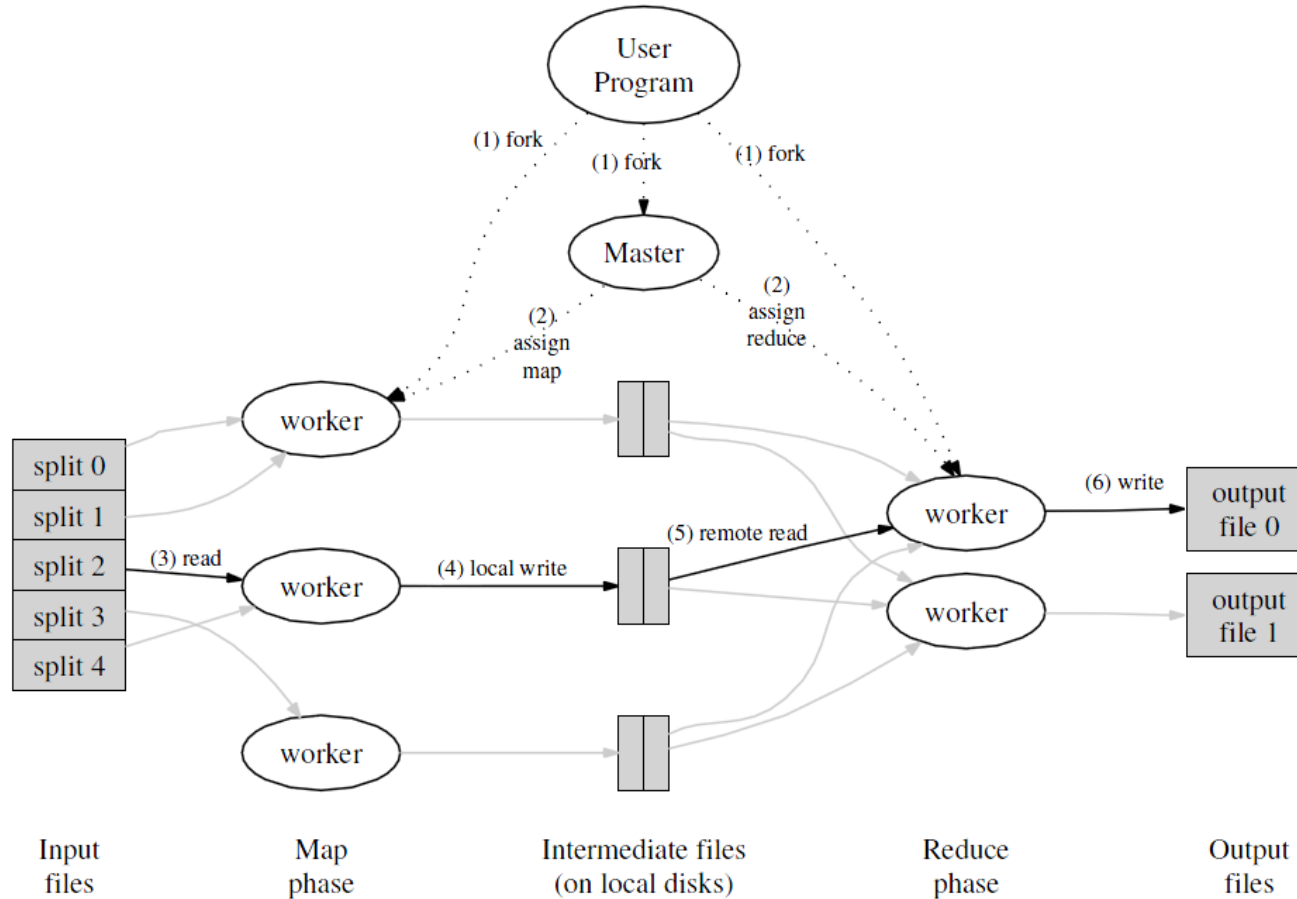


MapReduce Framework



- MapReduce **framework** takes care about
 - **Distribution** and parallelizing of the computation
 - **Monitoring** of the whole distributed task
 - The **grouping** (shuffling) phase
 - putting together intermediate results
 - **Recovering** from any failures
- User must define **only** map & reduce **functions**
 - but can define also other additional functions (see below)

MapReduce Framework (2)



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: Example II



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

- what pages reference () 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" ])
```

MapReduce: Example III



Task: What are the **lengths** of words in the input text

- output = **how many** words are in the text for **each length**

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

MapReduce: Features



- MapReduce uses a “**shared nothing**” architecture
 - Nodes operate **independently**, sharing no memory/disk
 - Common feature of many NoSQL systems

- Data **partitioned** and **replicated** over many nodes
 - Pro: **Large** number of **read/write** operations per second
 - Con: **Coordination** problem – which nodes have my data, and when?

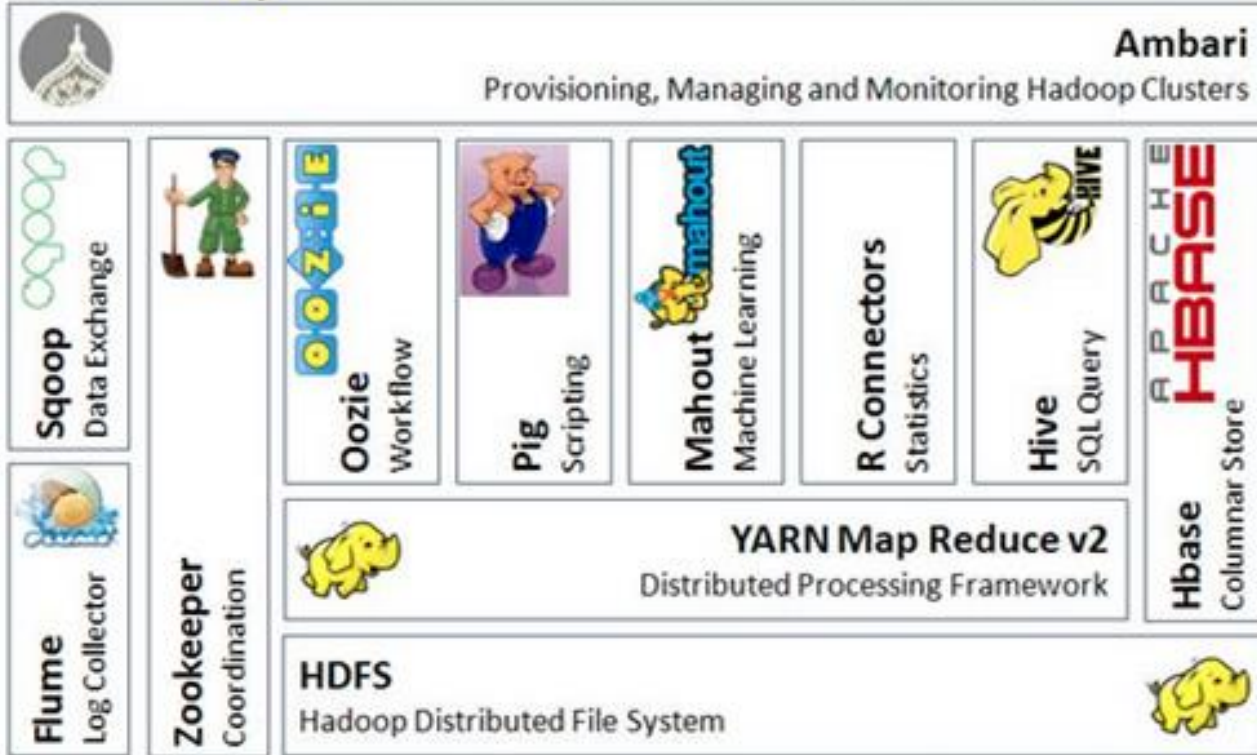
Applicability of MapReduce



- 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,...



Apache Hadoop Ecosystem



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



MapReduce: Implementation



Amazon Elastic
MapReduce



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.