

MapReduce

EE412: Foundation of Big Data Analytics

Announcements

- Todo reminder
 - Register to Classum
 - (Next week) Login to Haedong machine and change your password
 - Please check later an announcement at Classum
- Lecture videos
 - Will be recorded and uploaded from today
- Homeworks
 - Will post HW0 and HW1 next Thursday

↓
apadhe
spark

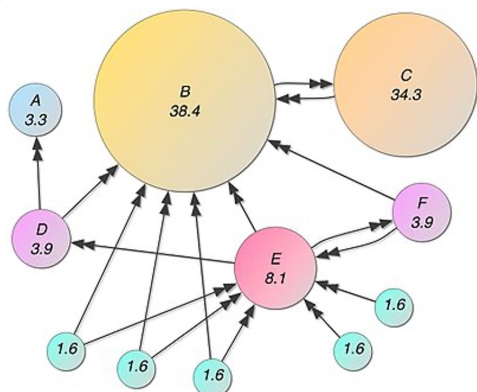
Te 10 assignment

Outline

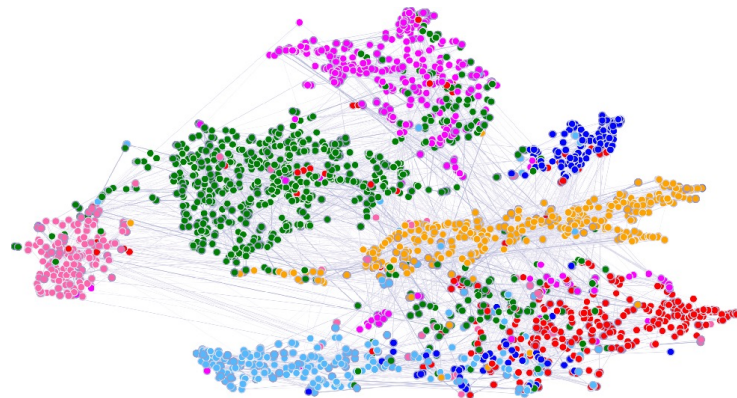
1. **Distributed Computing for Data Mining**
2. MapReduce: Basics
3. MapReduce: Details

Managing Large Amounts of Data Quickly

- Ranking Web pages by importance
 - Iterated matrix-vector multiplication where dimension is many billions
- Search friends in social networks
 - Graphs with hundreds of millions of nodes and many billions of edges



Source: Spatial



Source: Velickovic et al. (2018)

Large-Scale Computing

- Supercomputers are expensive and cannot scale infinitely
- Instead, we have a large **collection of commodity hardware**
 - Connected by Ethernet cables or inexpensive switches

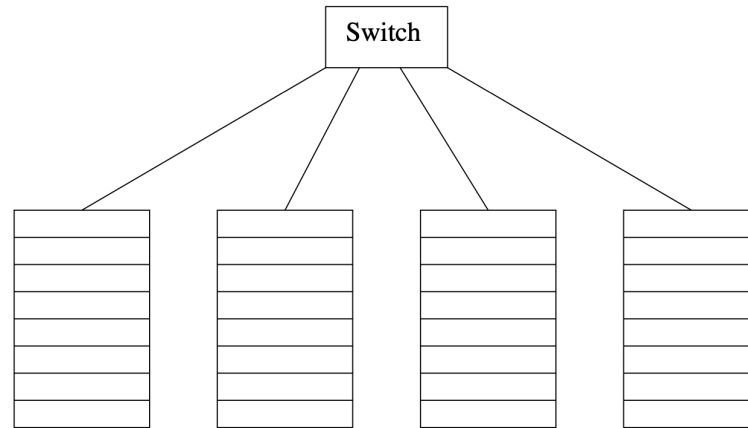


Source: UMBC

Physical Organization of Nodes

- **Parallel-computing architecture**

- Compute nodes are stored on racks (perhaps 8 – 64 on a rack)
- The nodes on a single rack are typically connected by gigabit Ethernet
- There can be many racks connected by a switch



Source: Textbook

Challenges

- How do you distribute computation?
- How can we make it easy to write distributed programs?
- Machines fail:
 - One server may stay up 3 years ($\approx 1,000$ days)
 - If you have 1,000 servers, expect to lose 1/day
 - With 1M machines, 1,000 machines fail every day!

Idea and Solution

- **Issue:** Copying data over a network takes time
 - **Idea:** *Bring computation to data* result 만 공유하고, raw data exchange X
모든 곳에 모든 data 있는 (subset이 존재)
 - **Hadoop/Spark** address these problems
 - **Storage infrastructure (file system):** GFS and HDFS
 - **Programming model:** MapReduce and Spark
- ↳ philosophy ↳ real SW
- ↑
not the
real
software

Storage Infrastructure //

- **Problem:**

- If nodes fail, how to store data persistently?
 - Loss of single node (e.g., disk crashes)
 - Loss of an entire rack (e.g., network fails)

→ 어떻게 해결하죠?

- **Answer: Distributed File System (DFS)**

- **Typical usage pattern:**

- Huge files (100s of GB to TB)
- Data is rarely updated in place → 드문드문 rare한
- Reads and appends are common → to dependent 애플리케이션을 위해서.

읽기, 쓰기는 common

쓰기는 자주 (이러한) dependent X

Distributed File System (DFS)

- **Chunk servers**

- Files are divided into chunks, which are typically 16 – 64MB
- Chunks are replicated (usually 2x or 3x) at different nodes
 - Try to keep replicas in different racks *다른 rack에 replica*
- Chunk servers also serve as compute servers *// (result 저장할 수 있음)*

- **Master node** (or a name node) *manage all info (정보가 어디에 있나?)*

- Another small file storing metadata about where files are stored
- Master node itself is replicated
 - File system directory knows where to find the master node copies

Outline

1. Distributed Computing for Data Mining
2. **MapReduce: Basics**
3. MapReduce: Details

MapReduce

- **MapReduce** is a style of programming designed for:
 1. Easy parallel programming
 2. Invisible management of hardware and software failures
 3. Easy management of very-large-scale data
- It has several implementations, including Hadoop, Spark, Flink, and the original Google implementation just called “MapReduce”



Example: Word Counting

- Suppose we have a huge text document
- Count the number of times each distinct word appears in the file
 - E.g., (ball, 30), (basket, 10), (great, 5), and so on
- **Many applications:**
 - Analyze web server logs to find popular URLs
 - Statistical machine translation
 - Need to count the frequency of all 5-word sequences in documents

each machine have partial part
각 마스터기 계산할 때씩 일부분

MapReduce: Components

- **Map:**

- Apply a user-written *Map function* to each input element
- The output of the Map function is a set of key-value pairs

ex) document >url line 3

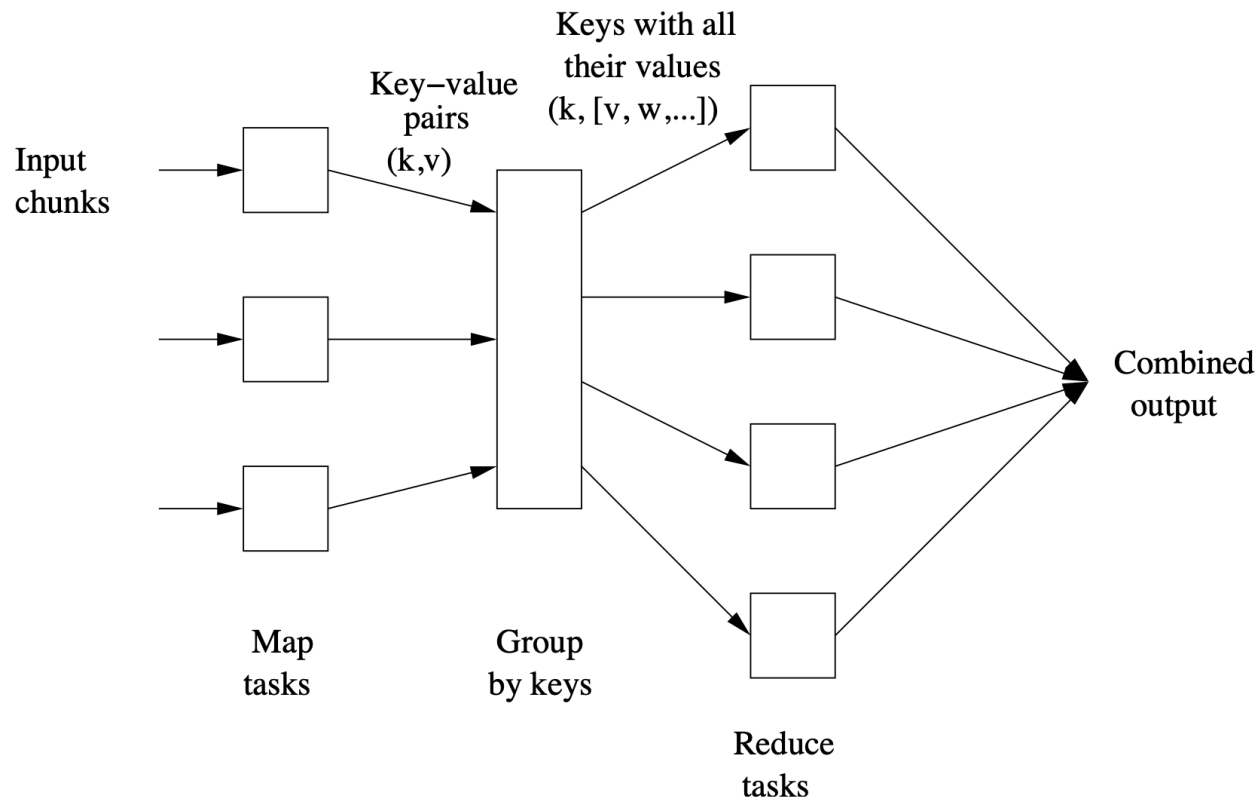
- **Group by key:** Sort and shuffle

- System sorts all the key-value pairs by key
- Outputs key-(list of values) pairs

- **Reduce:**

- User-written *Reduce function* is applied to each key-(list of values)

MapReduce Pattern



Source: Textbook

Jaemin Yoo

Map

- **Input**

- Chunks consisting of elements, which are key-value pairs
 - To allow composition of several MapReduce processes
 - Keys of the input (i.e., initial) elements can be ignored

- **Output**

- Zero or more key-value pairs

Group by Key

각 machine 의 다 키값을 키 bucket 으로 만든다.

• Input

- Key-value pairs from Map tasks

• Output

- The master controller knows the number of reducer tasks (r)
- Applies a hash function to the keys, producing bucket numbers 0 to $r - 1$
- The key value pairs are put into r local files
 - Each local file is destined for a specific Reduce task
- The master controller merges the files and returns $(k, [v_1, v_2, \dots, v_n])$

각 key 에
대해
hashing을
해서 bucket 으로
분류

Reduce

- **Input**

- One or more keys and their associated value lists
 - One Reducer can take several keys

- **Output**

- The Reduce function is applied to each key-(value list) pair
- The outputs from all the Reduce tasks are merged into a single file

Map and Reduce for Word Counting

```
// key: document name; value: text of the document
```

```
map(key, value):
```

```
    for each word w in value:
```

```
        emit(w, 1)
```

```
// key: a word; values: an iterator over counts
```

```
reduce(key, values):
```

```
    result = 0
```

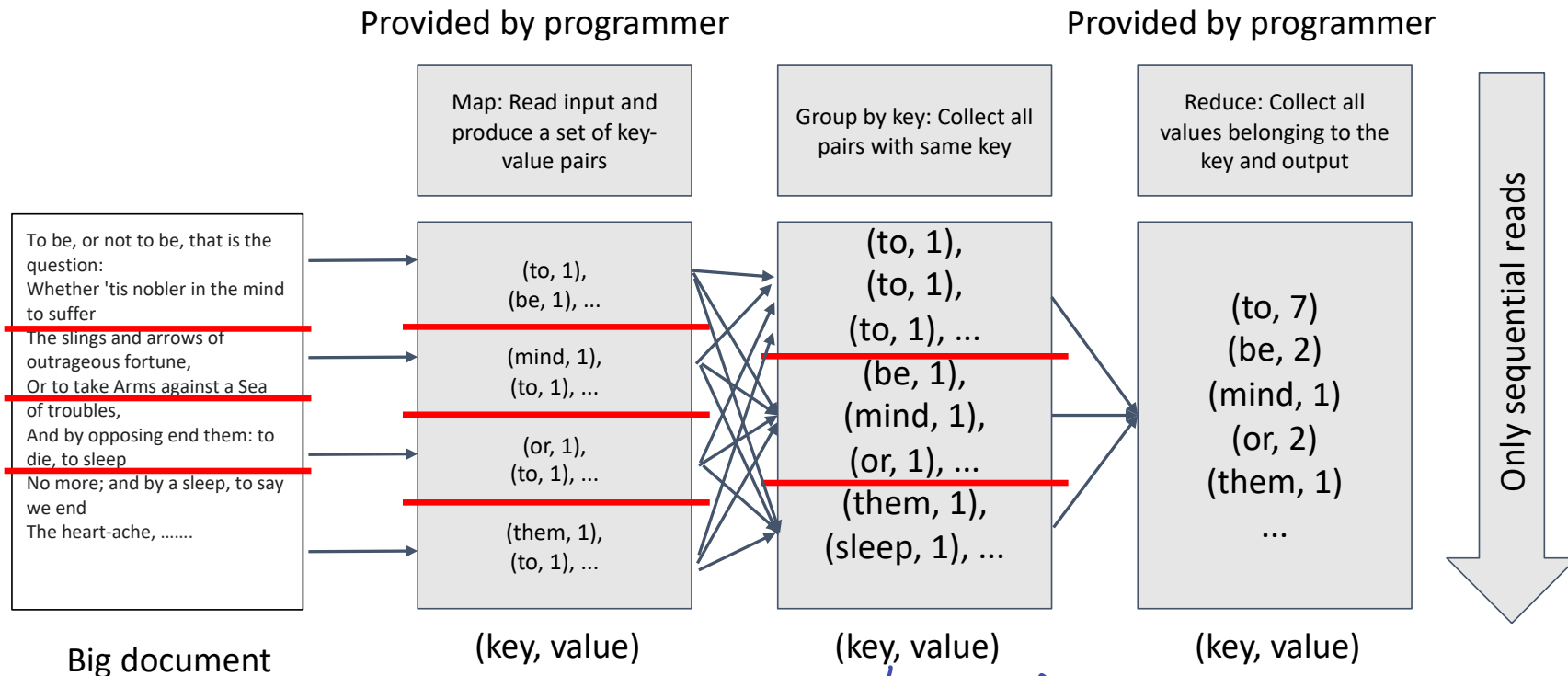
```
    for each count v in values:
```

```
        result += v
```

```
    emit(key, result)
```

reduce시작 전 map의 결과물
분산된 데이터

MapReduce for Word Counting

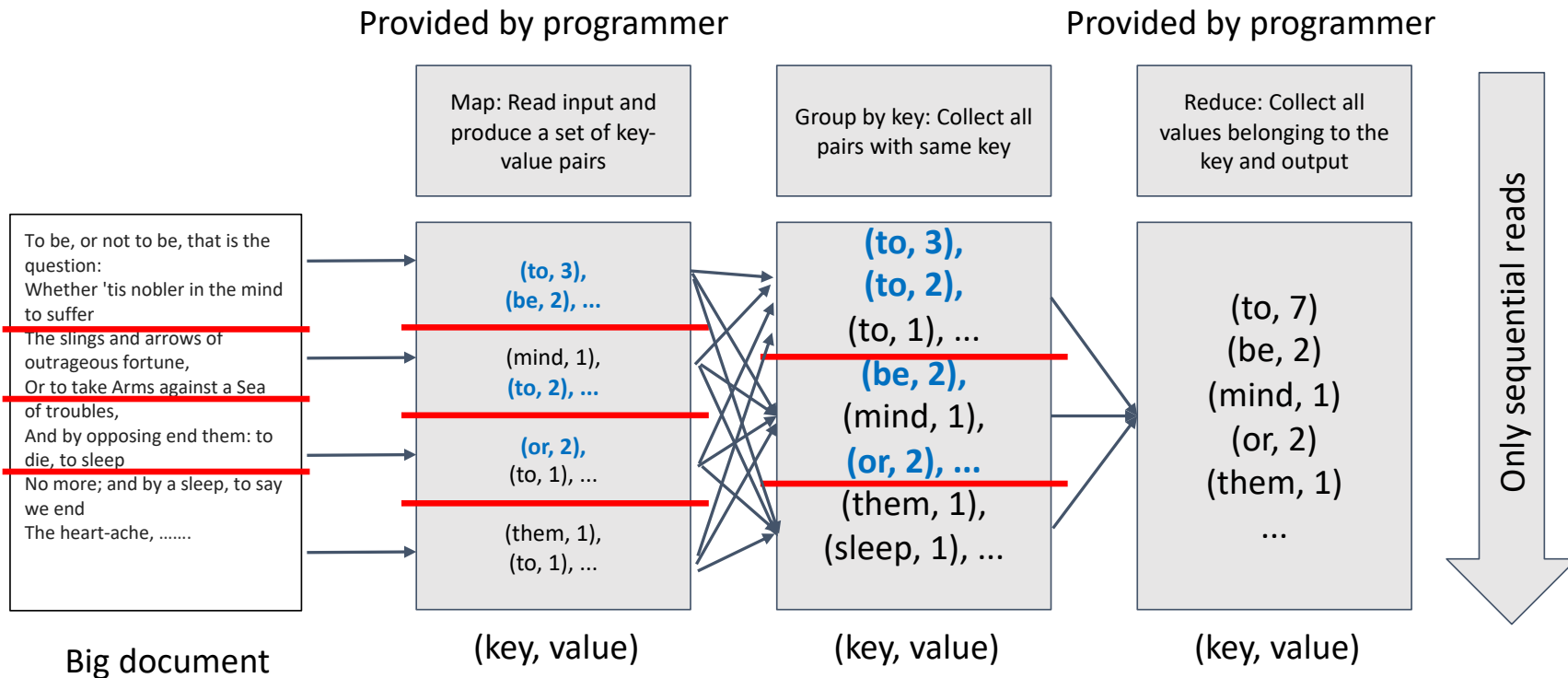


sorting
key value
binnise sorting을 위한

Combiners

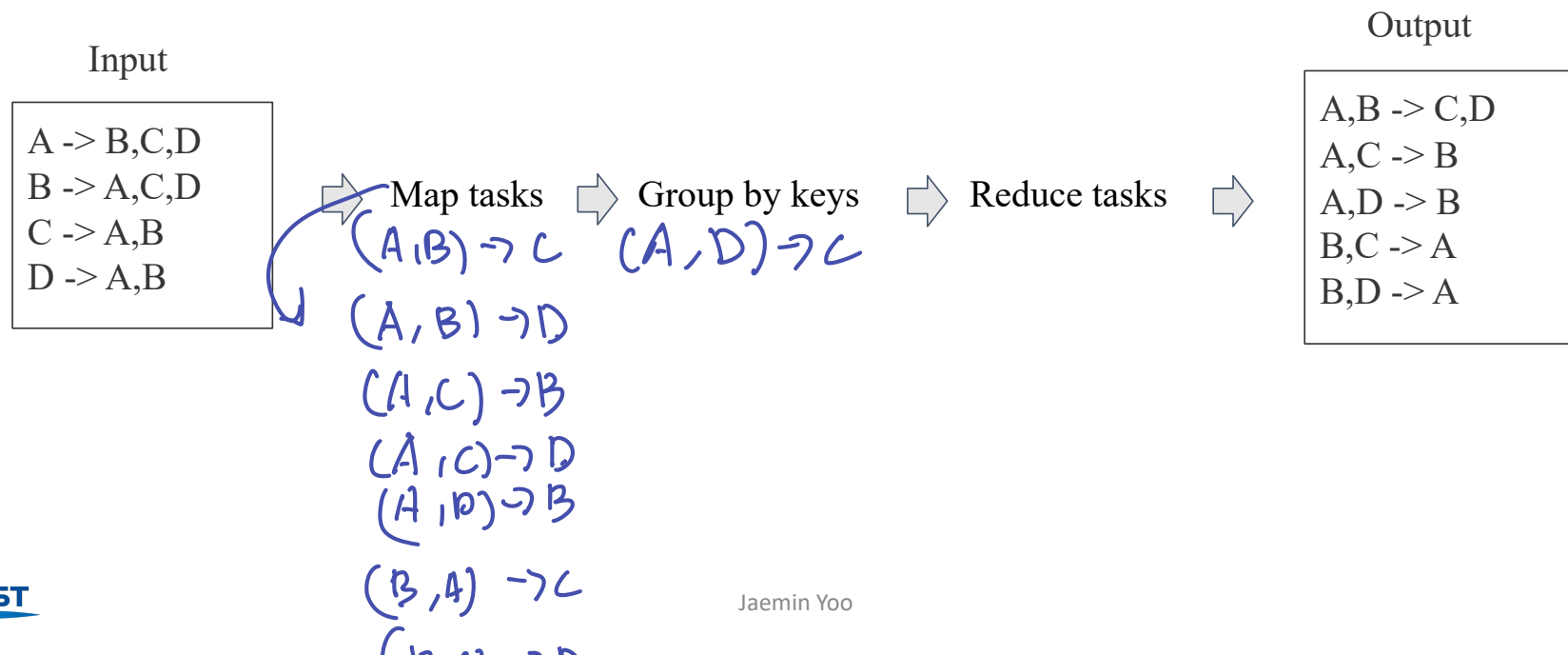
- If the Reduce function is ^{→ 교환}commutative and ^{→ 결합}associative, values can be combined in any order for the same result
 - E.g., word counting, matrix-vector multiplication
 - **Commutative:** $a + b = b + a$
 - **Associative:** $(a + b) + c = a + (b + c)$
- In this case, the Reducer's work can be pushed to the Map tasks ^{→ 여러 맵터들 map 단계에서 처리됨}
 - E.g., instead of producing $(w, 1), (w, 1), \dots$, produce (w, m)
- It reduces the amount of data shuffling for the group by key step

Combiner for Word Counting



Pop Quiz

- Let's find common friends of friends on a social network
- What are the map and reduce functions?

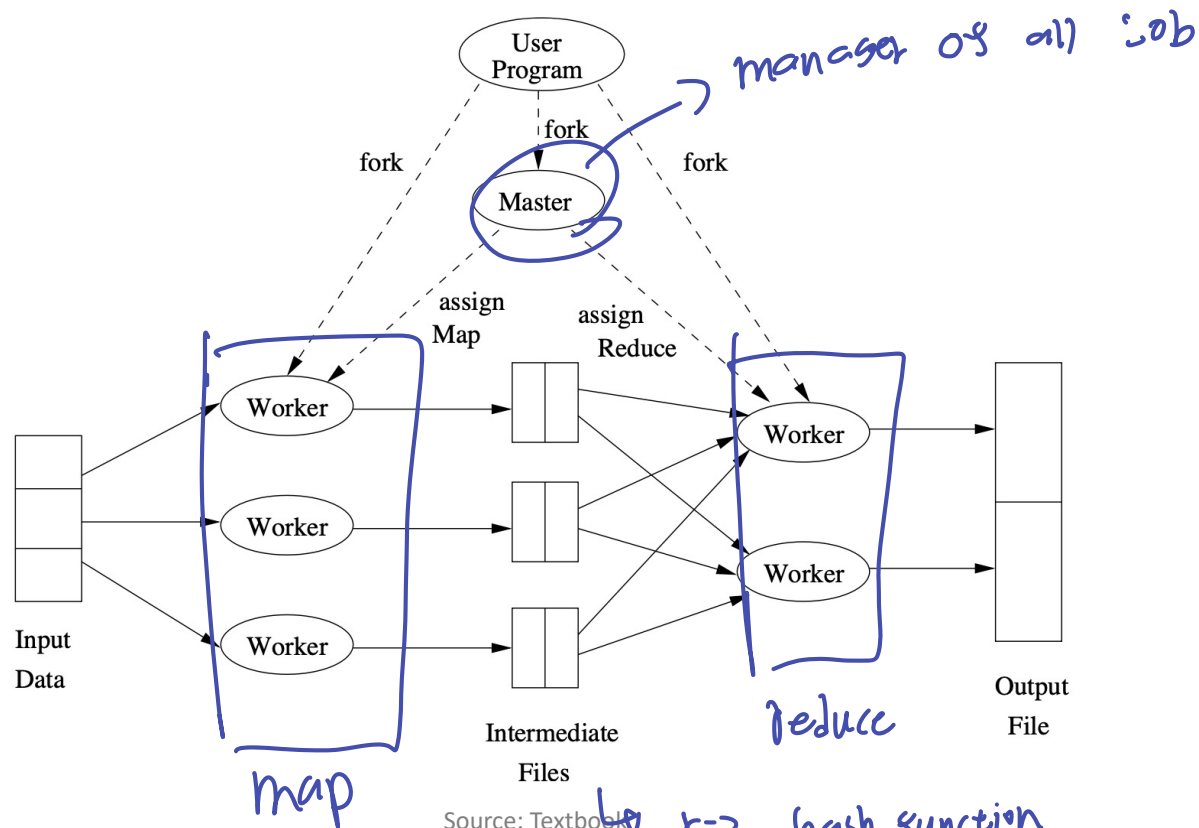


$$(B, A) \rightarrow B$$
$$(C, A) \rightarrow B$$

Outline

1. Distributed Computing for Data Mining
2. MapReduce: Basics
3. **MapReduce: Details**

MapReduce Program Execution



Source: Textbook

by $r=2$, hash function //

MapReduce Program Execution

- User program forks a Master controller and Worker processes
 - A worker handles either Map or Reduce tasks, but not both
- **The Master controller**
 - Schedules Map and Reduce tasks and assigns them to Worker processes
 - One Map task per chunk is reasonable
 - Reduce tasks should be fewer than Map tasks
 - Otherwise, too many intermediate files per Map task
 - Keeps track of task status (idle, executing at a Worker, or completed)
 - Worker process reports to Master when it finishes a task

MapReduce Program Execution

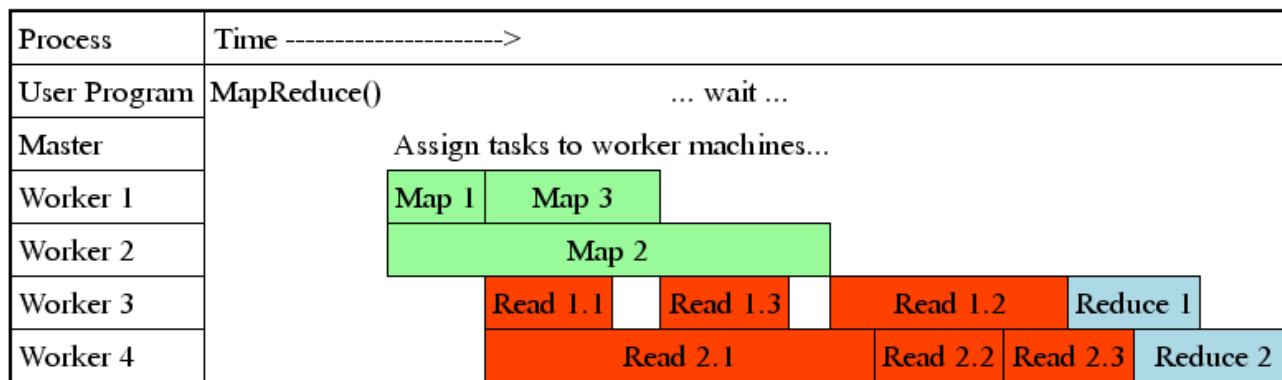
- **Map task**

- Creates a file for each Reduce task on the Map Worker's local disk
- The Master is informed of the files and the destination Reduce tasks

- **Reduce task**

- When assigned to a Worker, it is given all the files that form its input
- Executes code written by the user and writes its output to a file

MapReduce Pipelining



Source: Textbook Slides

Coping with Node Failures

need to return
to algorithm way

- **Master failure**

- Entire MapReduce job must be restarted

- **Map worker failure**

- Detected by Master because it periodically pings Worker processes
- 2.1. [• All Map tasks assigned this Worker have to be redone, even if completed
 - Output for Reduce tasks reside in compute node and are unavailable
 - Master informs each Reduce task that the location of its input has changed

- 2. [• **Reduce worker failure**

- Only in-progress tasks are reset to idle and the Reduce task is restarted

Number of Reduce Tasks

- Parallelism can be maximized by using one Reduce task per key
- However, this plan is not usually the best
 - Overhead associated with each task created
 - Often there are far more keys than there are compute nodes
- Significant variation in the # of values for different keys
 - Reduce tasks will exhibit skew in processing time
 - **Rule of thumb:** # of nodes < # of Reduce tasks < # of keys !!

2/26 강의 내용, 알고리즘을 공부하기

Algorithms using MapReduce

- MapReduce is not a solution to every problem
 - DFS only makes sense when files are very large and rarely updated in place
 - MapReduce is not suitable for frequent database changes
 - E.g., managing online retails sales (say by Amazon)
- **Original purpose of Google:**
 - To execute very large matrix-vector multiplications for PageRank
 - Matrix-vector and matrix-matrix calculations fit nicely into MapReduce
 - Relational algebra (i.e., SQL) can also use MapReduce effectively

update를 하는 과정은 이따기
→ 정렬이 복잡함

PageRank

Matrix-Vector Multiplication

- Suppose we have an $n \times n$ matrix M and a vector \mathbf{v} of length n
- Matrix-vector product is vector \mathbf{x} of length n whose i -th element is:

$$x_i = \sum_{j=1}^n m_{ij} v_j$$

- If $n = 100$, we do not need DFS or MapReduce
- But in large search engines, e.g., Google, n can be tens of billions

Matrix-Vector Multiplication

- Assume v fits in main memory:

Map function

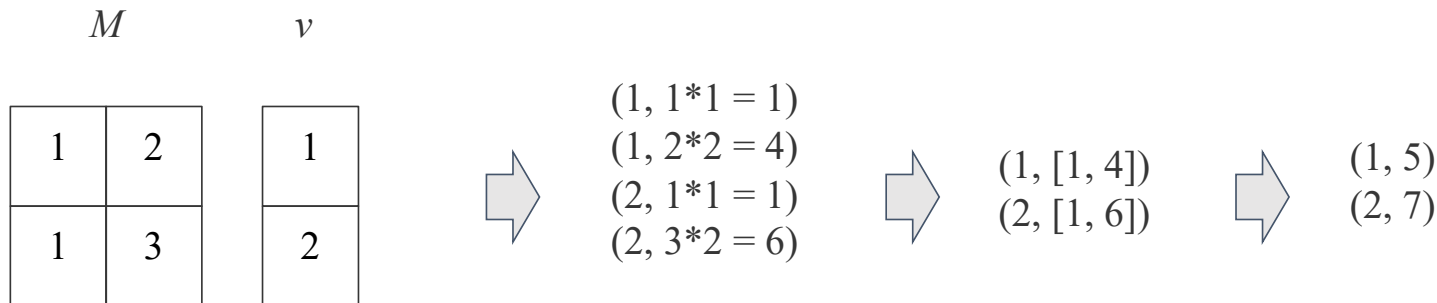
Input: m_{ij}

Output: $(i, m_{ij}v_j)$

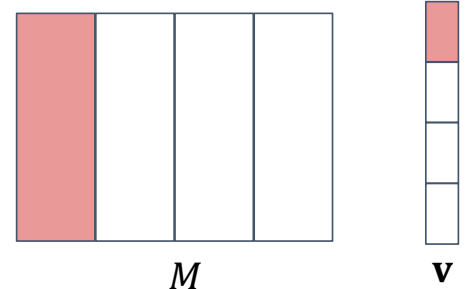
Reduce function

Input: $(i, [mi_1v_1, mi_2v_2, \dots, m_{in}v_n])$

Output: $(i, \sum_j m_{ij}v_j)$



Matrix-Vector Multiplication



- What if the vector \mathbf{v} cannot fit in main memory?
 - Divide M into vertical stripes and divide \mathbf{v} into horizontal stripes
 - The width of vertical strips = the height of horizontal stripes
 - Use enough stripes so that one stripe of \mathbf{v} fits in main memory
 - Each Map task is assigned a chunk from one of M 's stripes and the corresponding stripe in \mathbf{v}
 - The Map and Reduce tasks then run exactly the same as before

Summary

1. Distributed Computing for Data Mining

- Distributed File System

2. MapReduce: Basics

- Map, group by key, reduce, and combiners
- Example: Word counting

3. MapReduce: Details

- Coping with node failure
- Algorithms using MapReduce