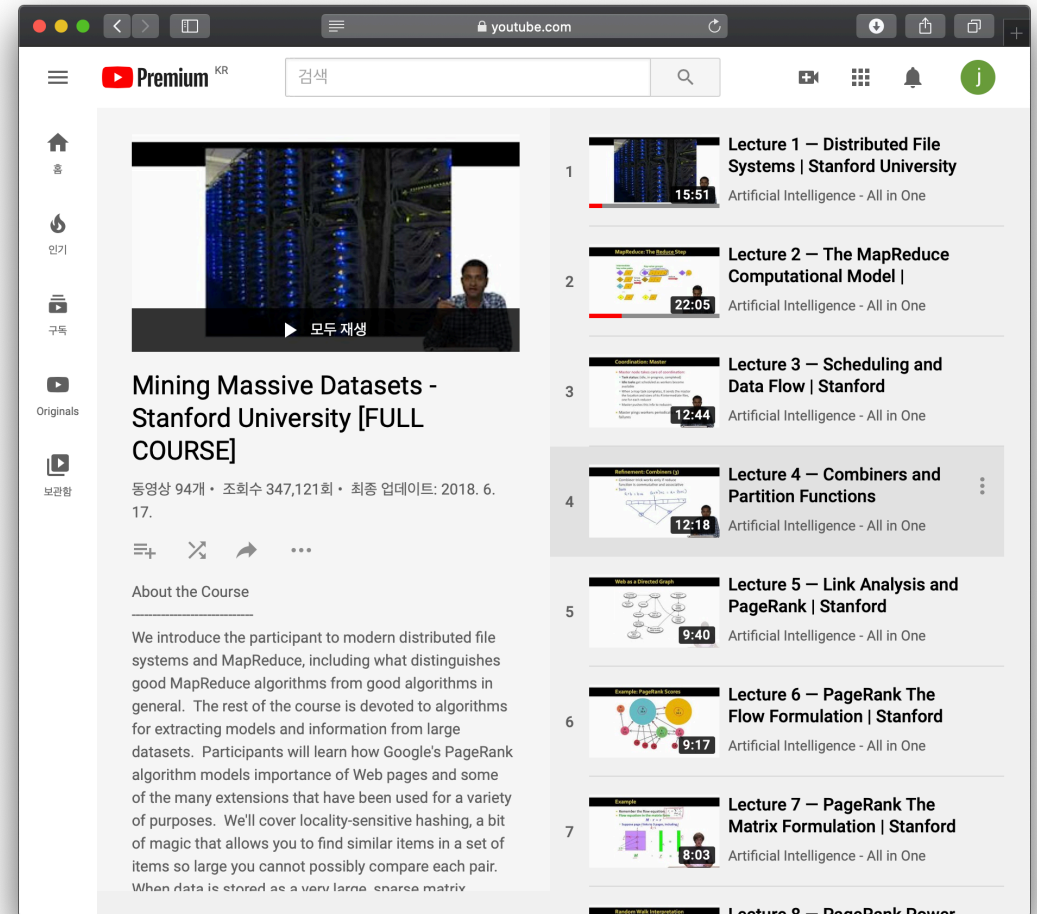


Mining Massive Datasets

Jisung Jeong / jisung0920.github.io

Reference

- Open Course :
Mining Massive Datasets (Stanford)
- Textbook :
Mining of Massive Datasets
(Jure Leskovec, Anand Rajaraman, Jeff Ullman)

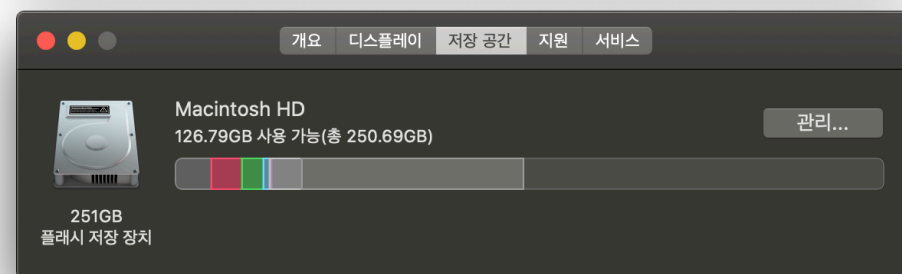
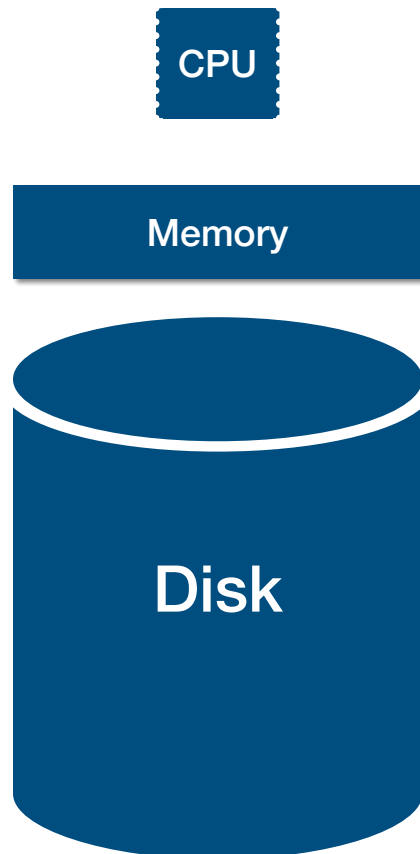


1. MapReduce
2. PageRank
3. Locality-Sensitive Hashing
4. Distance Measures, Nearest Neighbors
5. Frequent Itemsets
6. Communities in Social Networks
7. Stream Algorithms
8. Recommender Systems
9. Dimensionality Reduction
10. Clustering
11. Computational Advertising
12. Machine Learning
13. More About MapReduce
14. More About Locality-Sensitive Hashing
15. More About Link Analysis

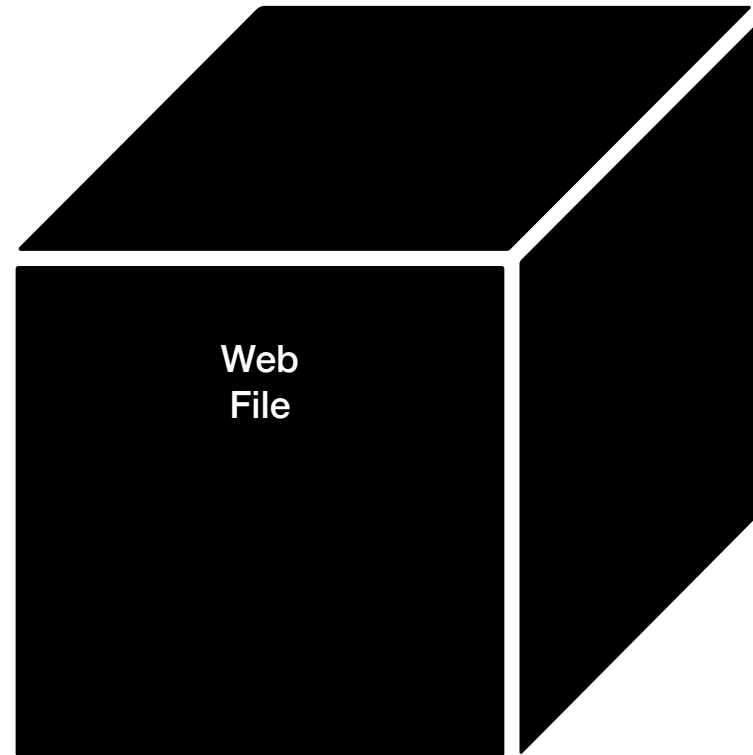
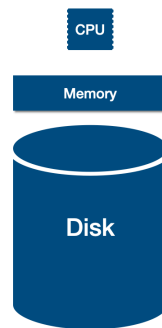
MapReduce

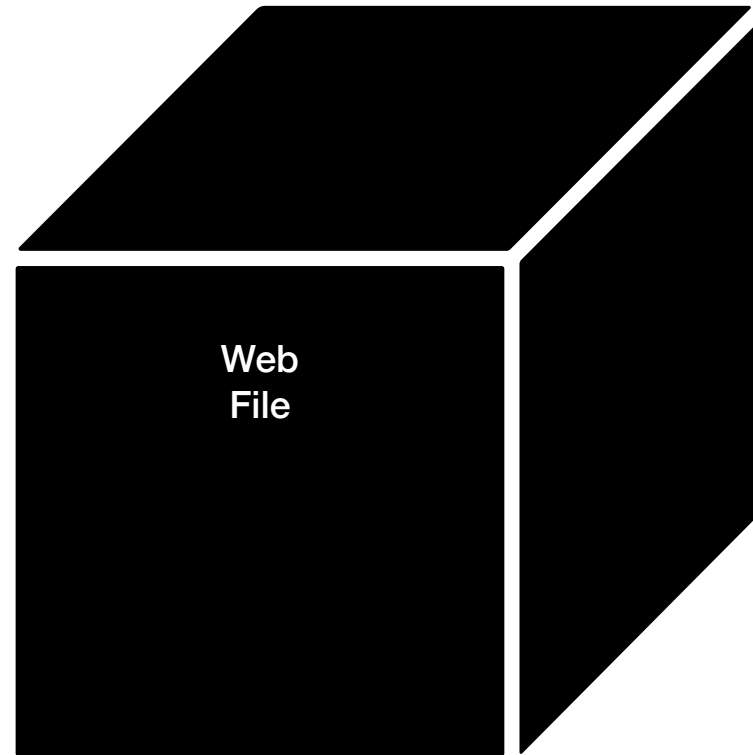
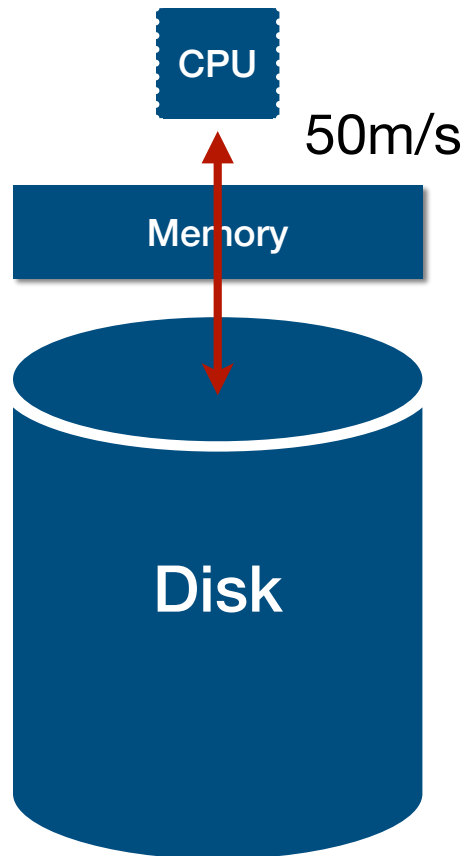
Physical Organization of Compute Nodes

Jisung Jeong / jisung0920.github.io



- 웹페이지 수 : 10Billion
 - 웹페이지 크기 : 20KB
- = 200TB

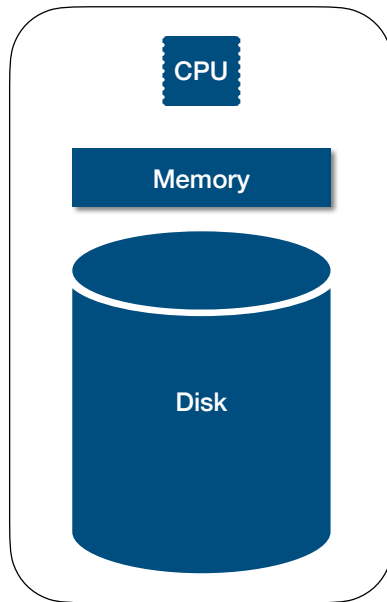




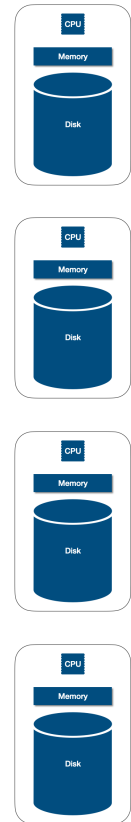
1. Physical Organization of Compute Nodes
2. Distributed File Systems
3. MapReduce

1. Physical Organization of Compute Nodes

Compute Node



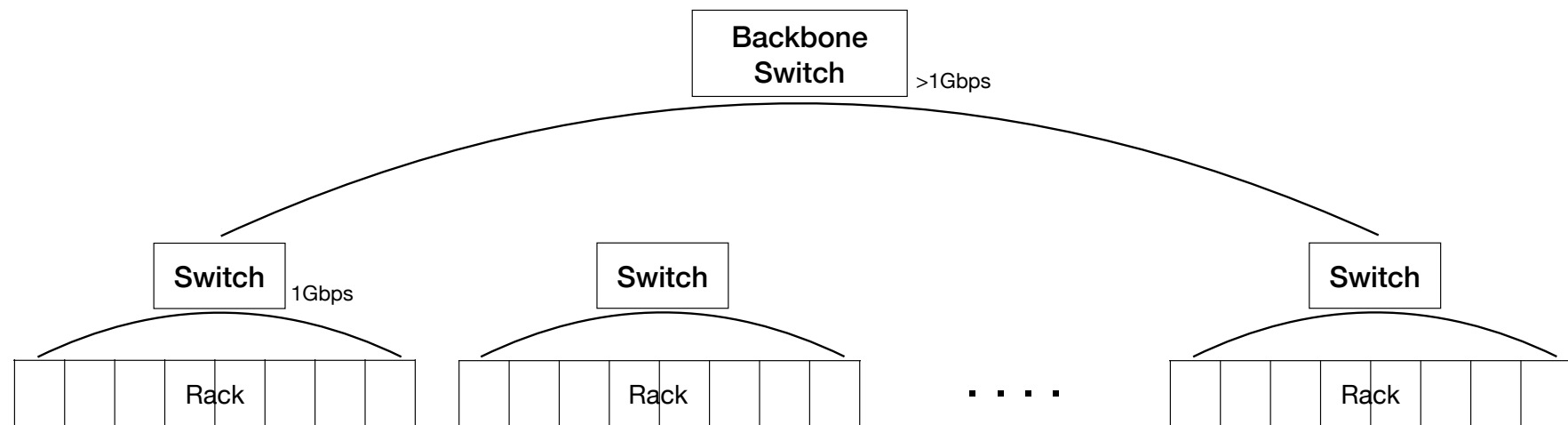
Rack



1. Physical Organization of Compute Nodes



1. Physical Organization of Compute Nodes



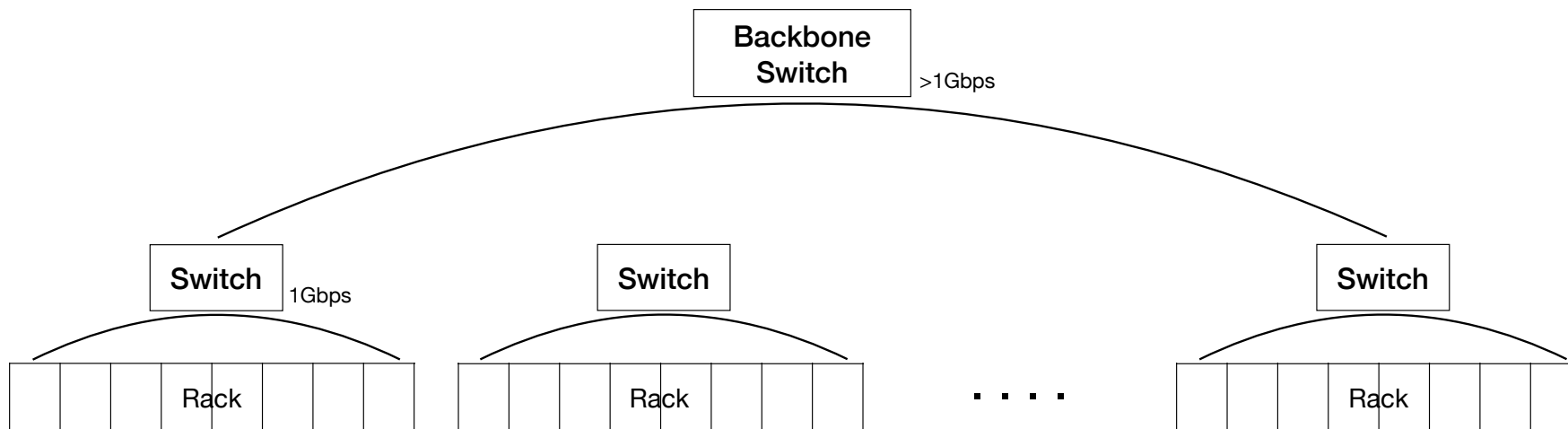
1. Physical Organization of Compute Nodes



MapReduce

Distributed File Systems

Jisung Jeong / jisung0920@gmail.com

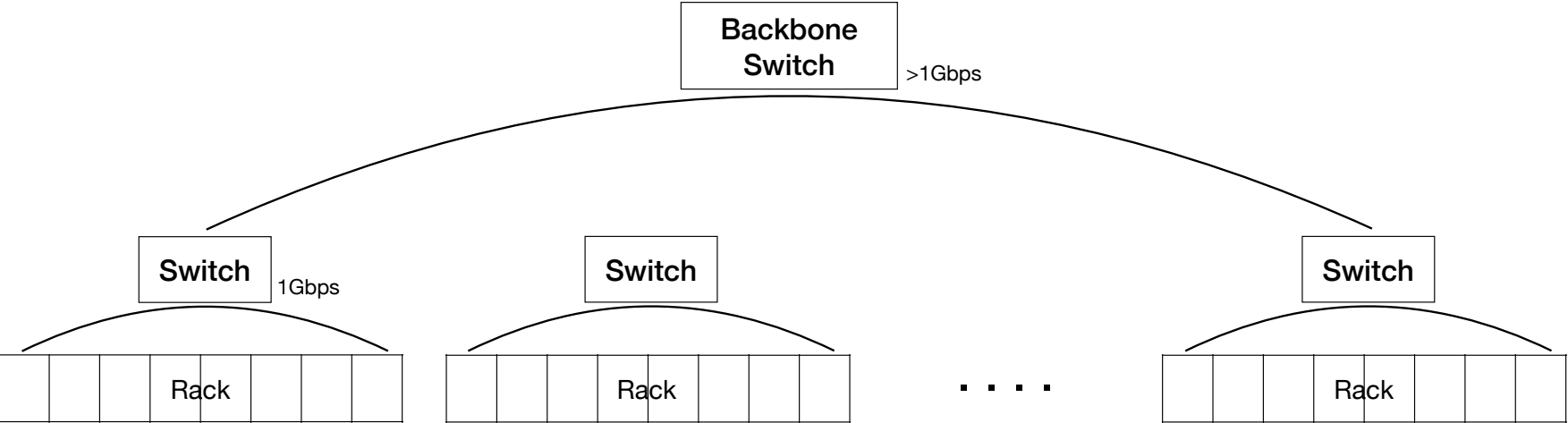


1. Node Failure

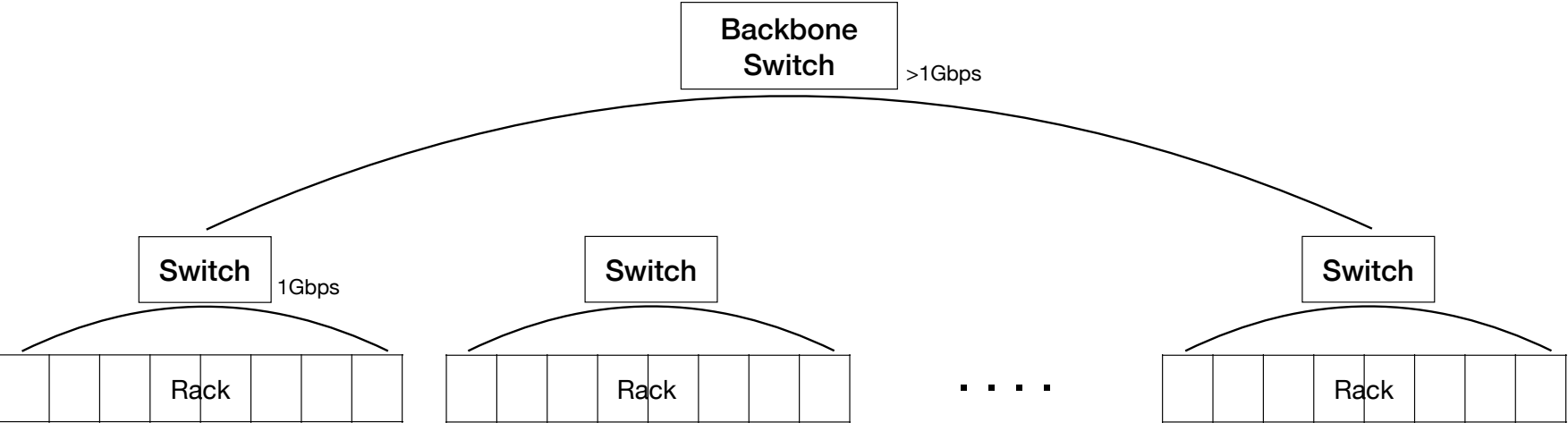
2. Network Bottleneck

3. Programming

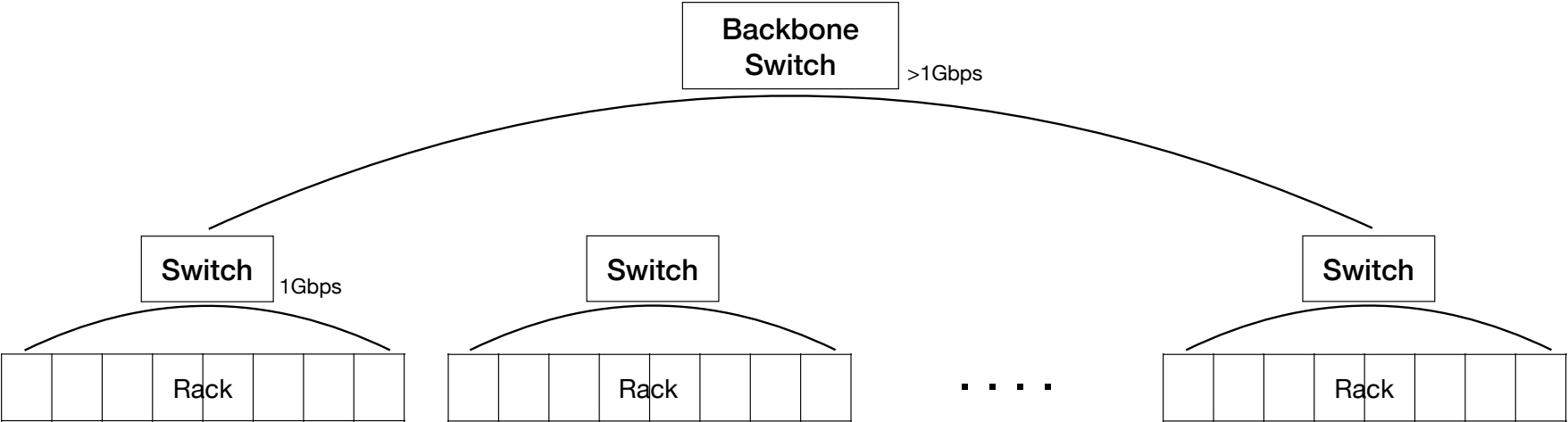
Node Failure



Network Bottleneck



Programming



Distributed File System

Distributed File System

file A

A1
A2
A3
.
.
.
.
.

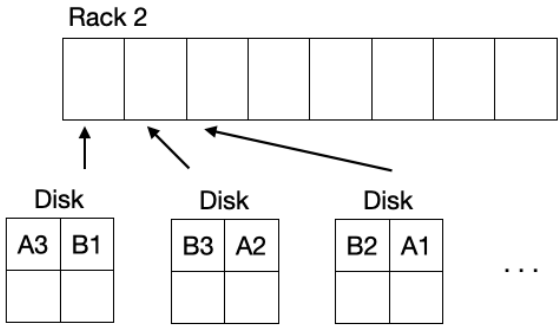
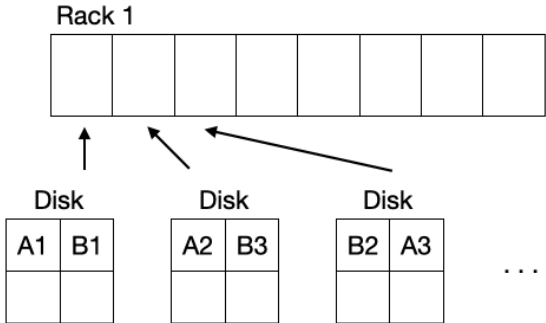
Distributed File System

file A

A1
A2
A3
.
.
.
.
.

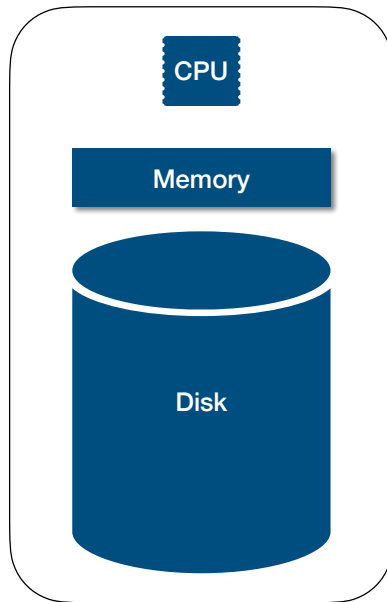
file B

B1
B2
B3
.
.
.
.
.

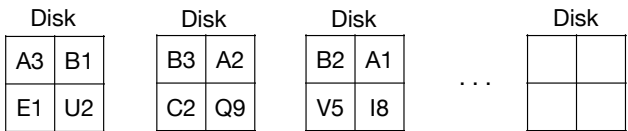
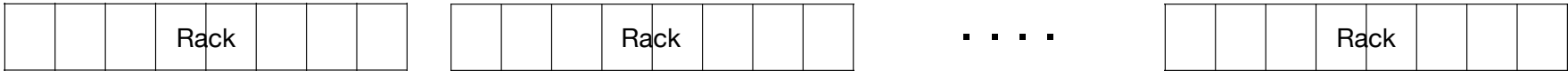


Distributed File System

Master Node



Distributed File System



MapReduce

MapReduce

Jisung Jeong / jisung0920@gmail.com



Cluster

File

A1
A2
A3
.
.
.
.
.

B2	A1
V5	I8

B3	A2
C2	Q9

A3	B1
E1	U2

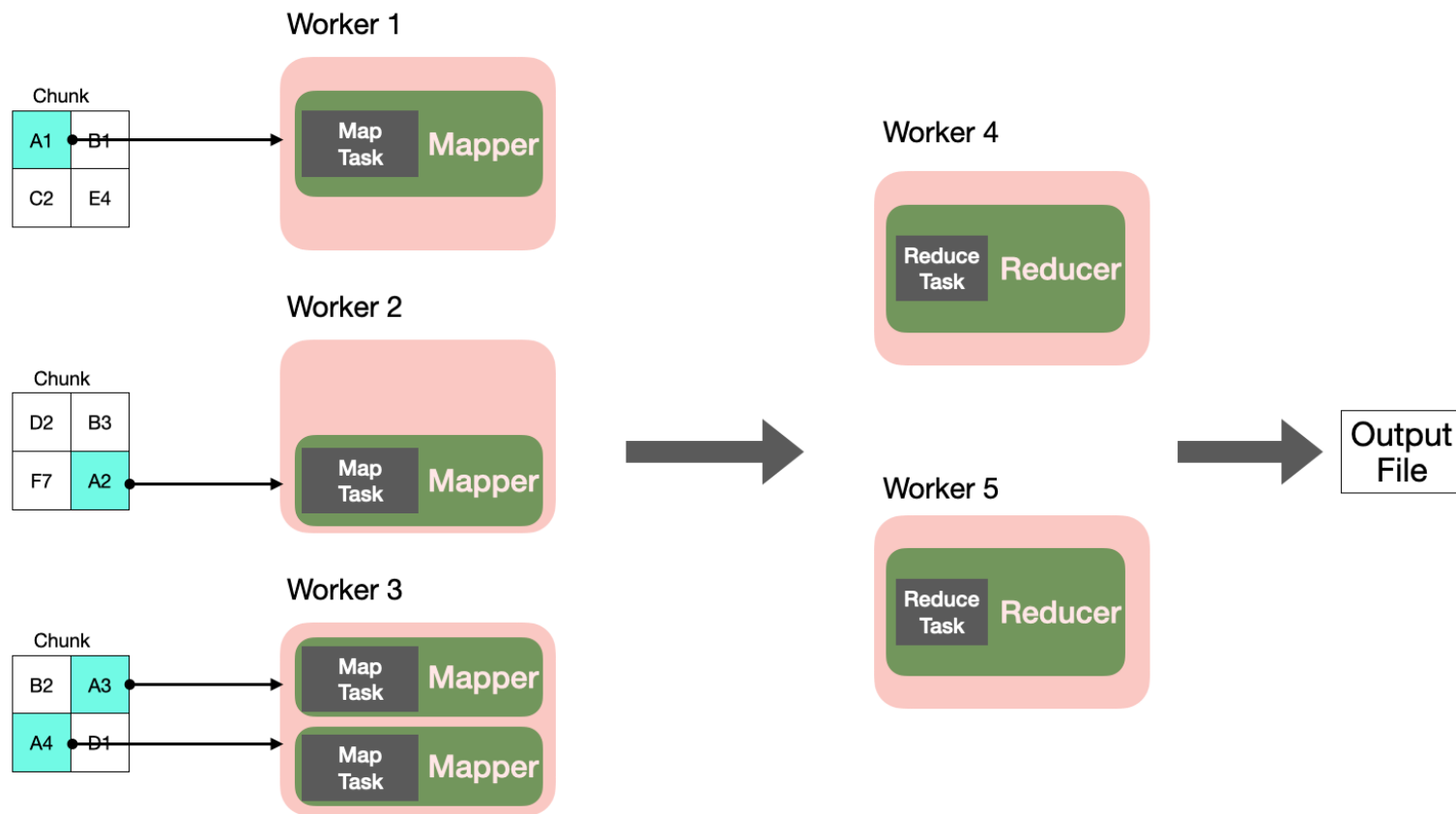
DFS

1. Map

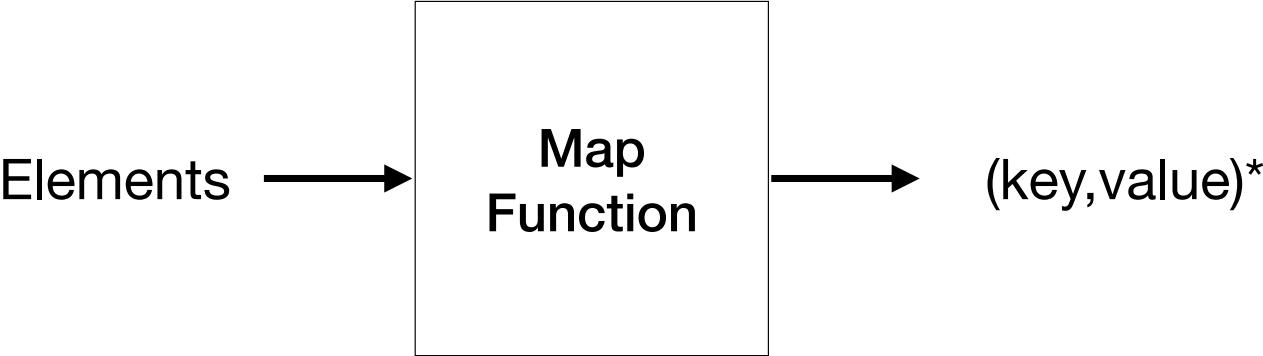
2. Group by Key

3. Reduce

- Worker
- Task
- Map Task
- Reduce Task
- Mapper
- Reducer



Map



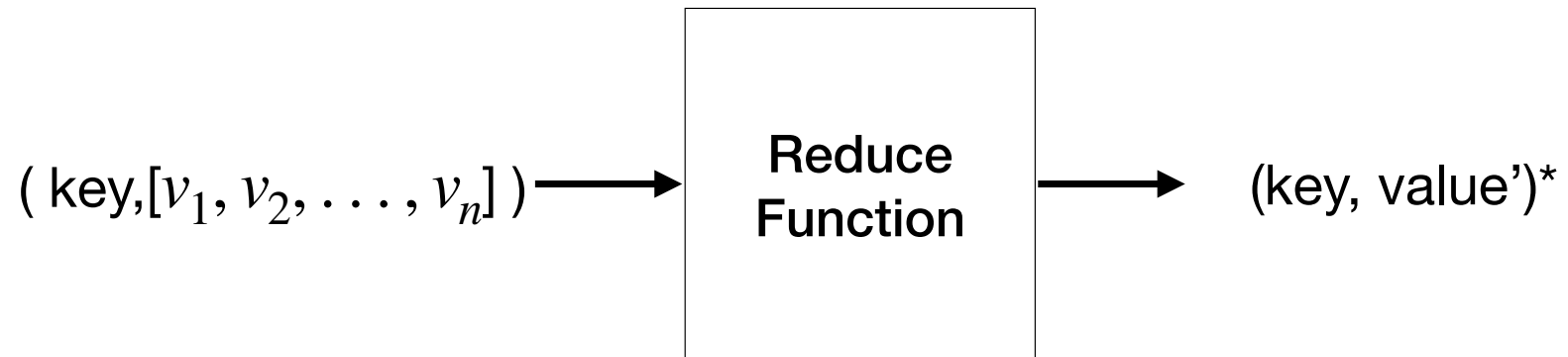
Group by key

$(k_1, v_a), (k_2, v_a), (k_3, v_b), (k_1, v_b), (k_1, v_c) \dots$



$(k_1, [v_a, v_b, v_c]), (k_2, [v_a]), \dots$

Reduce



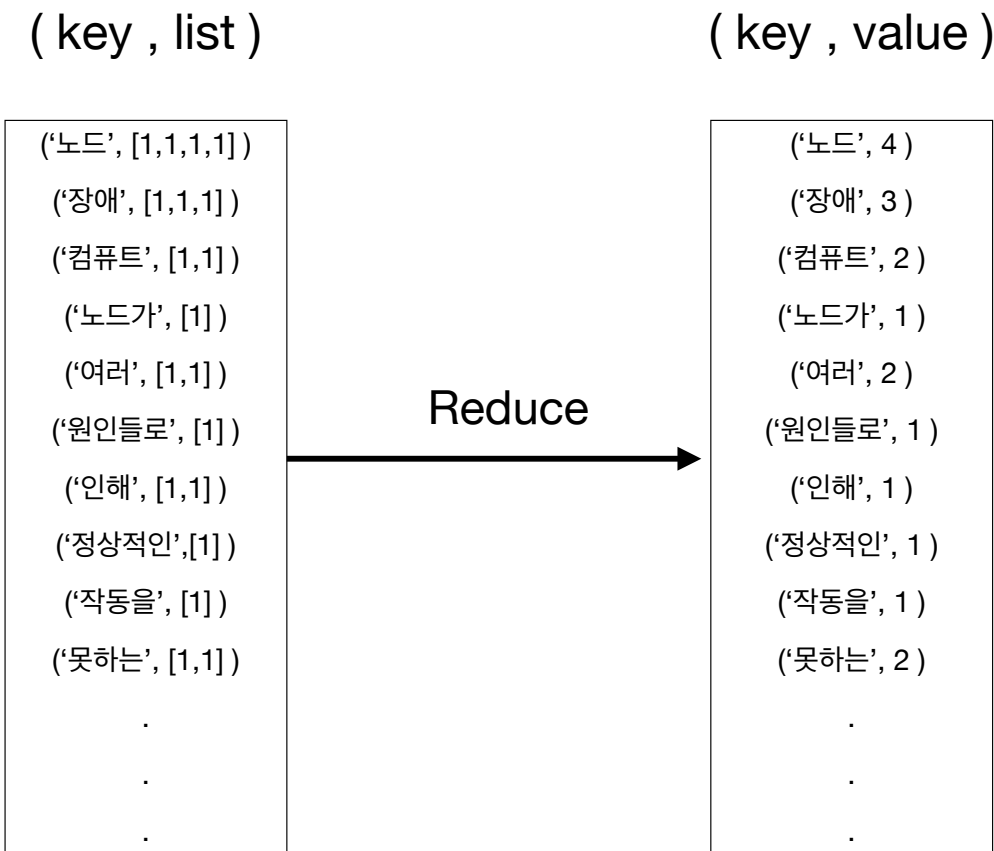
Example - Word Counter



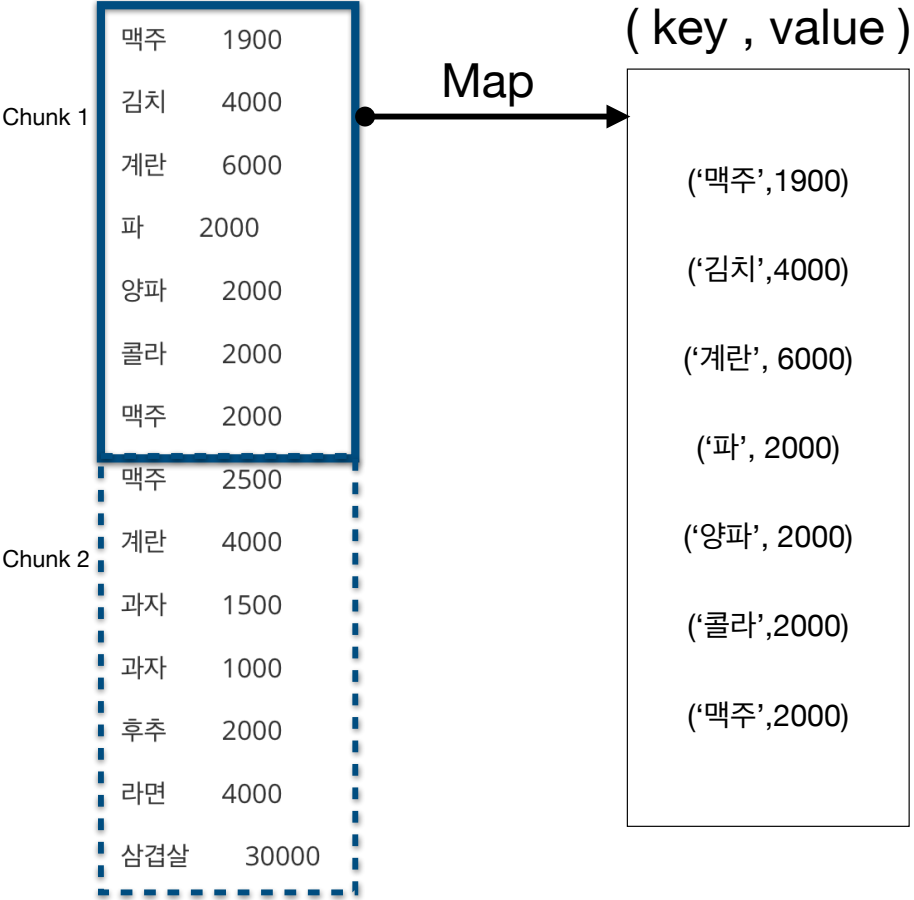
Example - Word Counter



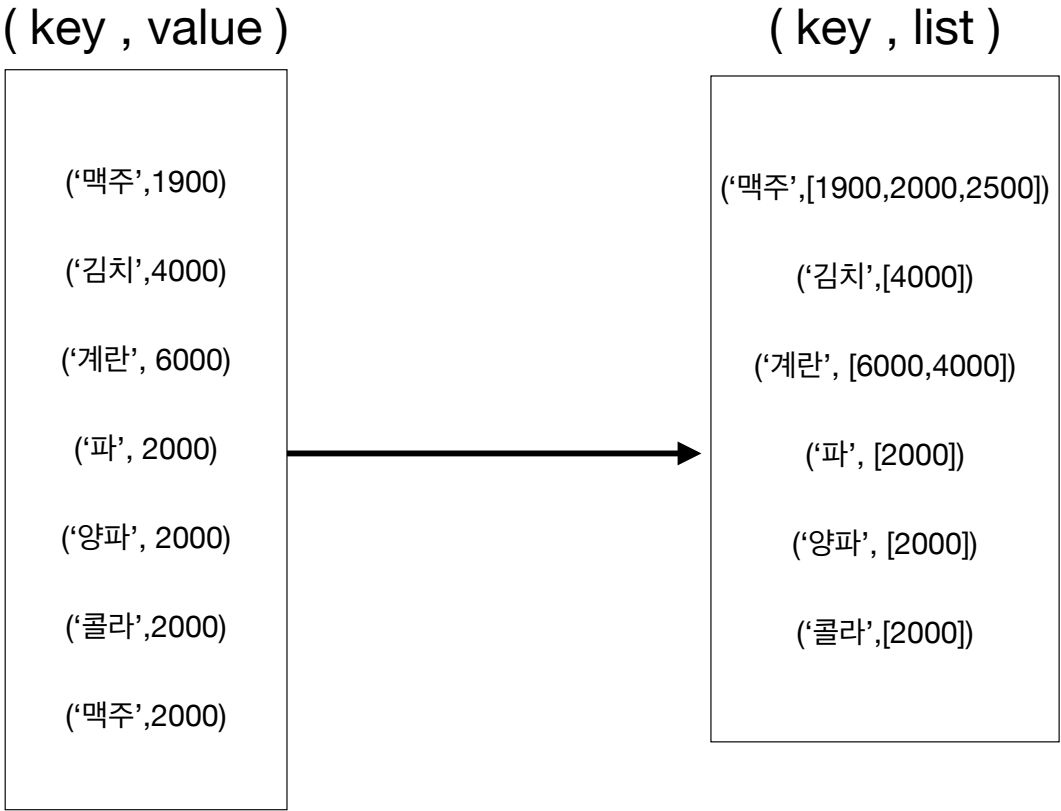
Example - Word Counter



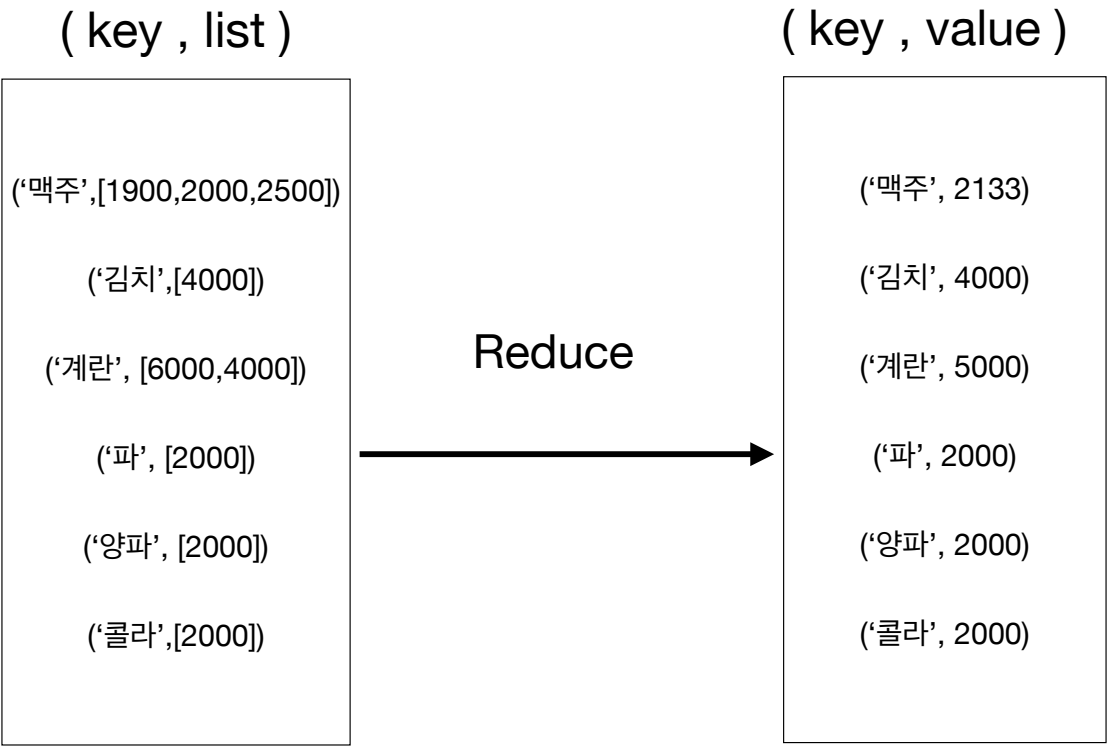
Example - Average Price



Example - Average Price



Example - Average Price



MapReduce

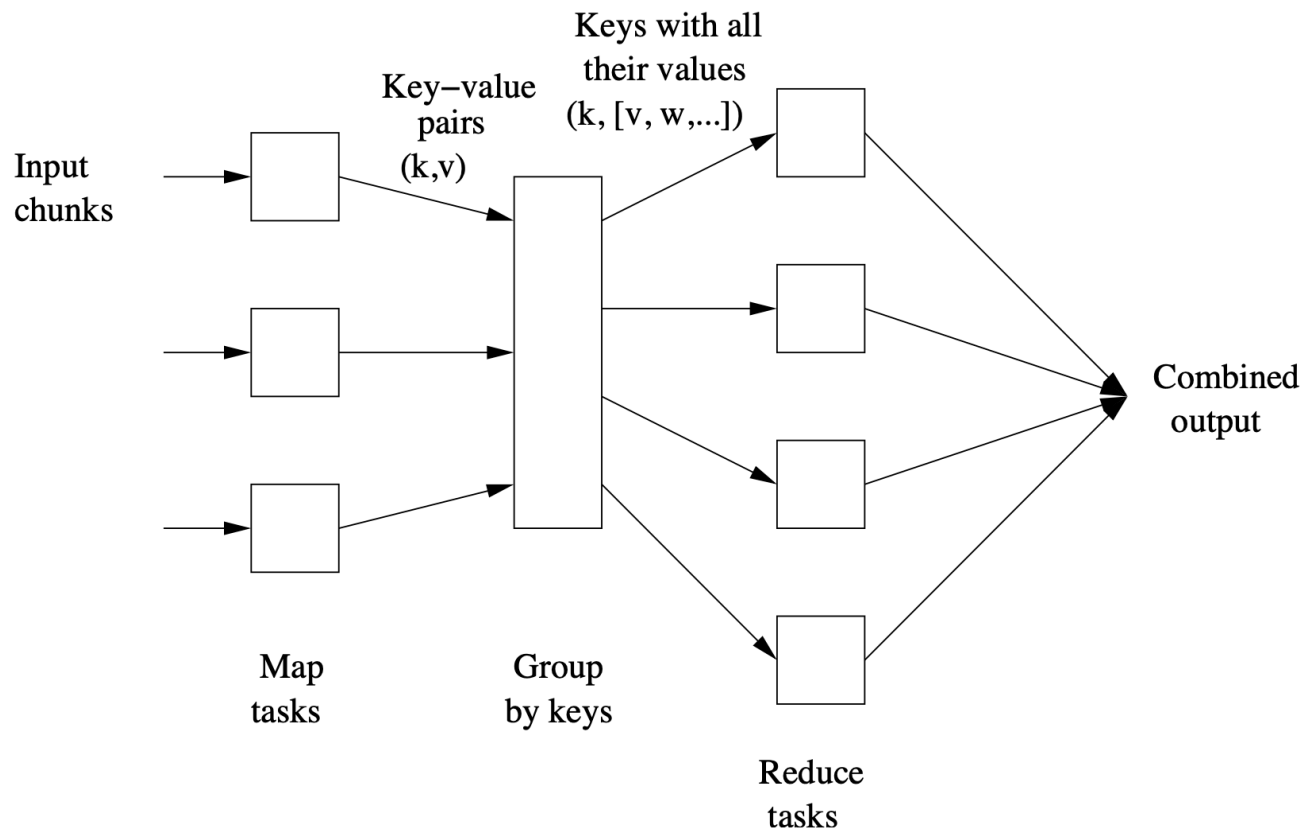


Figure 2.2: Schematic of a MapReduce computation

Mining of Massive Datasets - Figure 2.2

MapReduce

Advanced

Jisung Jeong / jisung0920@gmail.com

MapReduce

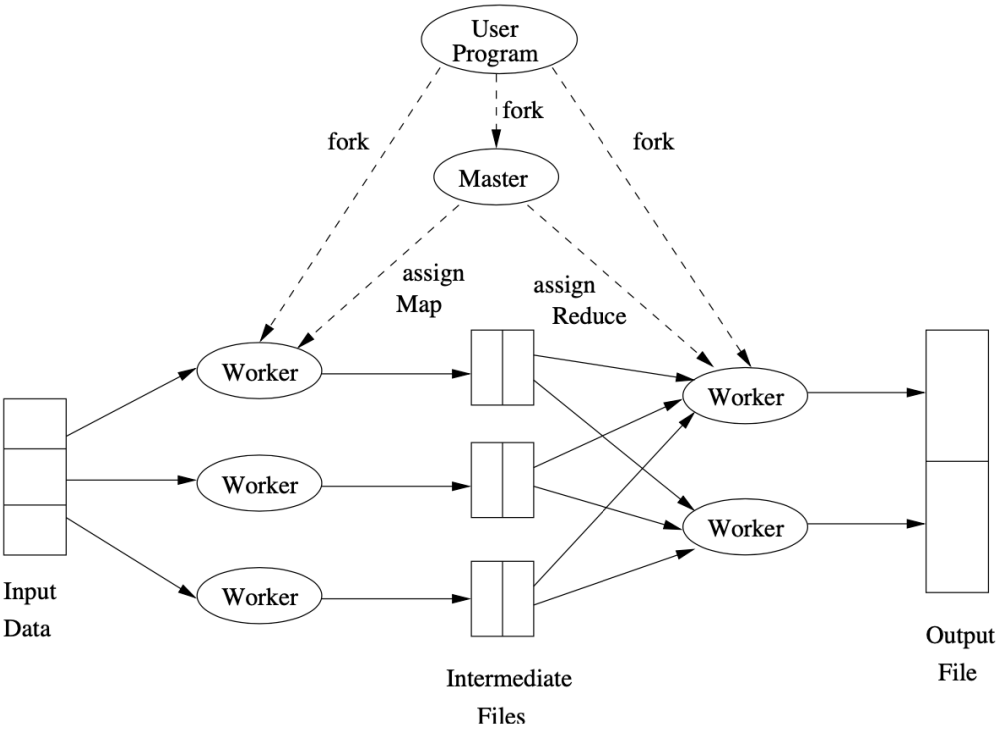
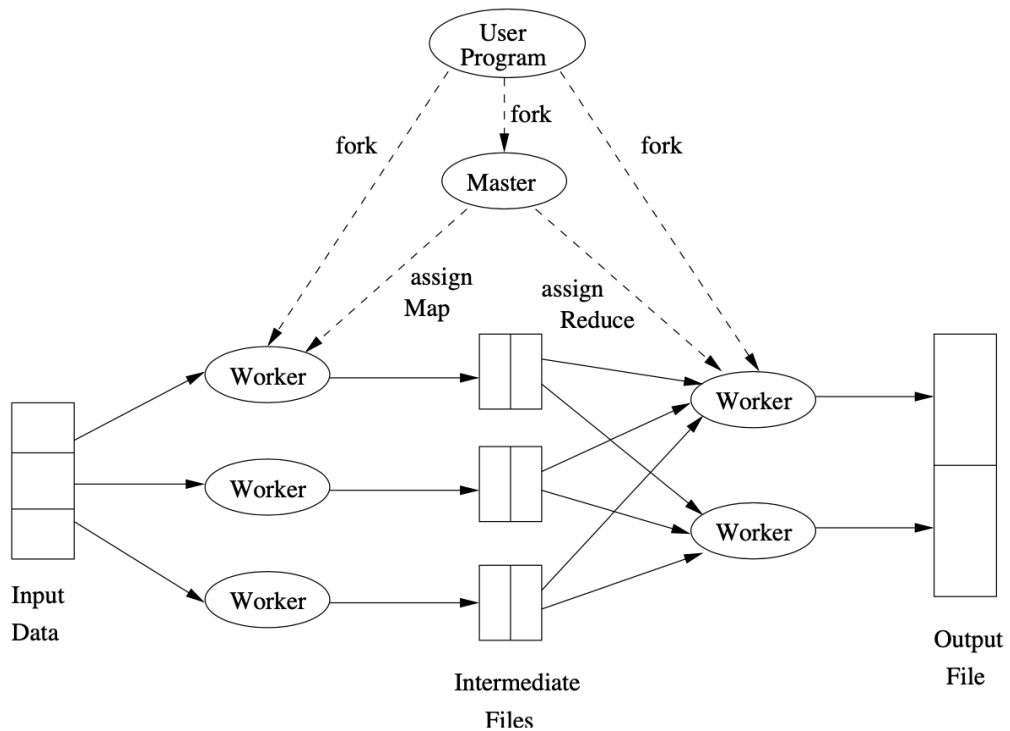


Figure 2.3: Overview of the execution of a MapReduce program

Mining of Massive Datasets - Figure 2.3

Node Failure

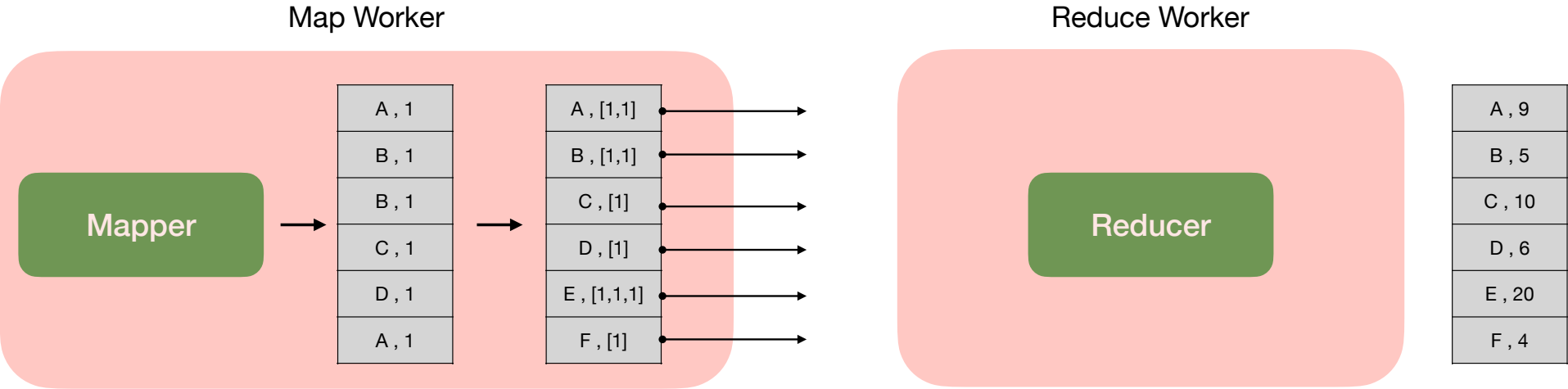


- Master
- Map Worker
- Reduce Worker

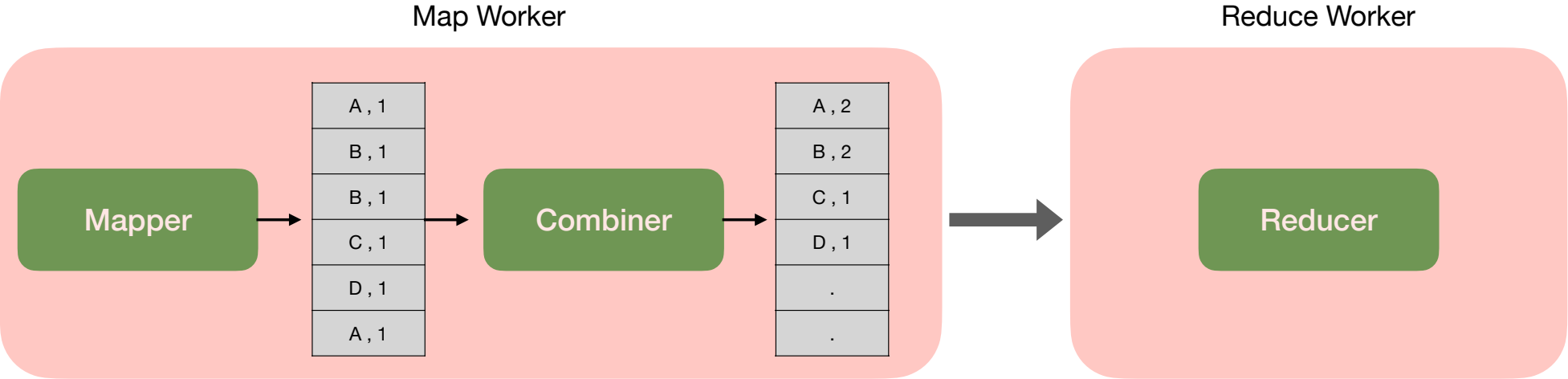
Figure 2.3: Overview of the execution of a MapReduce program

Mining of Massive Datasets - Figure 2.3

Combiner



Combiner



Partition Function

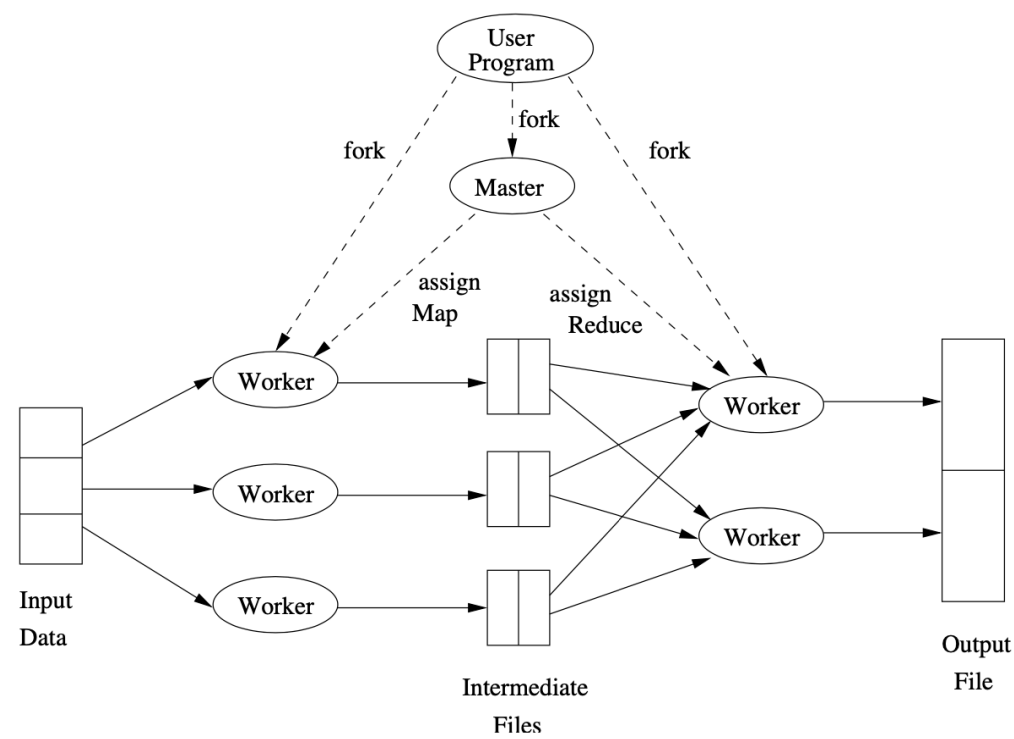


Figure 2.3: Overview of the execution of a MapReduce program

Mining of Massive Datasets - Figure 2.3

Partition Function

