

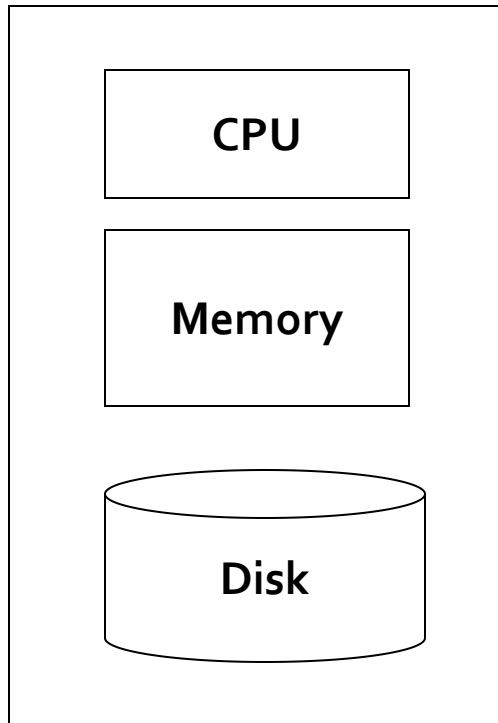
# Map-Reduce for large-scale data processing

from Mining of Massive Datasets (课程教参第二章)  
<http://www.mmds.org>

# MapReduce

- **Large scale computing for data mining**
- **Challenges:**
  - How to distribute computation?
  - Distributed/parallel programming is hard
- **Map-reduce** addresses all of the above
  - Google's computational/data manipulation model
  - Elegant way to work with big data

# Single Node Architecture



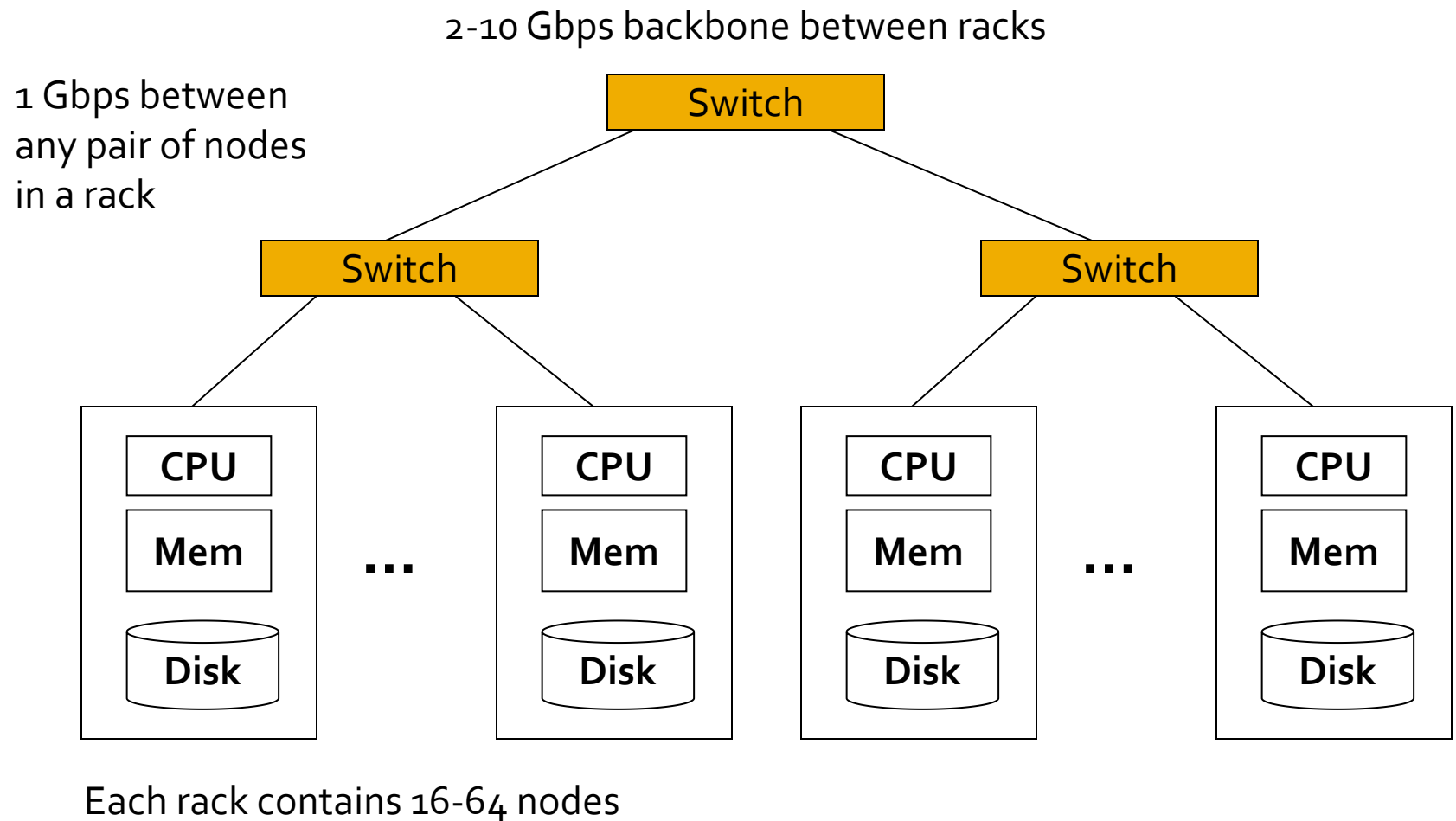
**Machine Learning, Statistics**

**“Classical” Data Mining**

# Motivation: Google Example

- 20+ billion web pages x 20KB = 400+ TB
- 1 computer reads 30-35 MB/sec from disk
  - ~4 months to read the web
- Takes even more to **do** something useful with the data!
- **A standard architecture for such problems:**
  - Cluster of commodity Linux nodes
  - Commodity network (ethernet) to connect them

# Cluster Architecture





# Large-scale Computing

- **Large-scale computing for data mining problems on commodity hardware**
- **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 (1,000 days)
    - If you have 1,000 servers, expect to loose 1/day
    - Estimated ~1,000 machines fail every day in Google

# Idea and Solution

- **Issue:** Copying data over a network takes time
- **Idea:**
  - Bring computation close to the data
  - Store files multiple times for reliability
- **Map-reduce** addresses these problems
  - Google's computational/data manipulation model
  - Elegant way to work with big data
  - **Storage Infrastructure – File system**
    - Google: GFS. Hadoop: HDFS
  - **Programming model**
    - Map-Reduce



# Storage Infrastructure

- **Problem:**

- If nodes fail, how to store data persistently?

- **Answer:**

- **Distributed File System:**

- Chunks stored on nodes with redundancy
    - Provides global file namespace

- **Typical usage pattern**

- Huge files (~TB)
  - Data is rarely updated in place
  - Reads and appends are common

# Distributed File System

## ■ Chunk servers

- File is split into contiguous chunks
- Typically each chunk is 16-64MB
- Each chunk replicated (usually 2x or 3x)
- Try to keep replicas in different racks

## ■ Master node

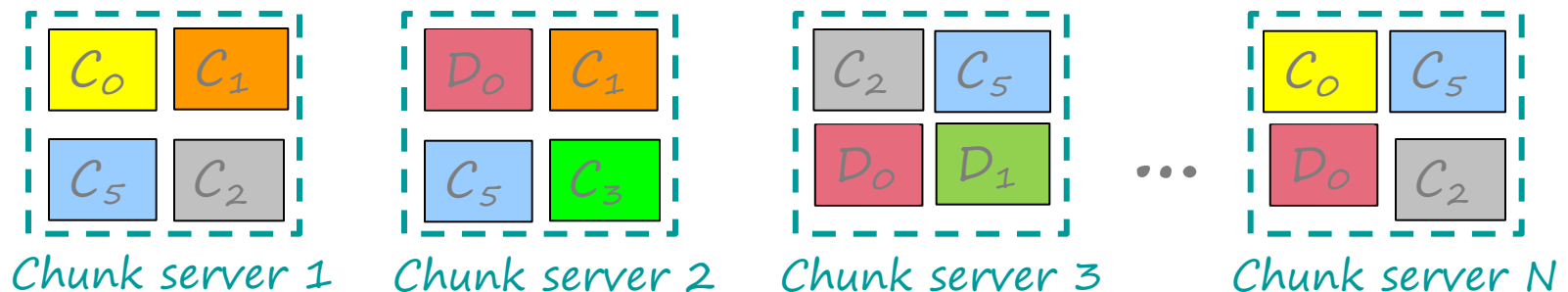
- i.e. Name Node in Hadoop's HDFS
- Stores metadata about where files are stored
- Might be replicated

## ■ Client library for file access

- Talks to master to find chunk servers
- Connects directly to chunk servers to access data

# Distributed File System

- **Reliable distributed file system**
- Data kept in “chunks” spread across machines
- Each chunk **replicated** on different machines
  - Seamless recovery from disk or machine failure



Bring computation directly to the data!

Chunk servers also serve as compute servers

# Programming Model: 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**, and the original Google implementation just called “MapReduce”

# Programming Model: MapReduce

## Warm-up task:

- We have a huge text document
- Count the number of times each distinct word appears in the file
- **Sample application:**
  - Analyze web server logs to find popular URLs

# Task: Word Count

## Case 1:

- File too large for memory, but all <word, count> pairs fit in memory

## Case 2:

- Count occurrences of words:
  - `words (doc.txt) | sort | uniq -c`
    - where **words** takes a file and outputs the words in it, one per a line
- Case 2 captures the essence of **MapReduce**
  - Great thing is that it is naturally parallelizable

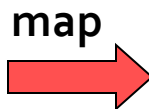
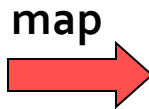
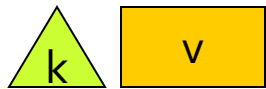
# MapReduce: Overview

- Sequentially read a lot of data
- **Map:**
  - Extract something you care about
- **Group by key:** Shuffle and Sort
- **Reduce:**
  - Aggregate, summarize, filter or transform
- Write the result

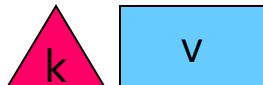
Outline stays the same, **Map** and **Reduce**  
change to fit the problem

# MapReduce: The Map Step

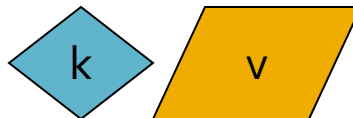
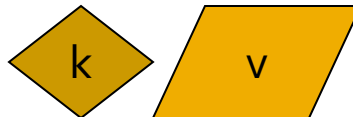
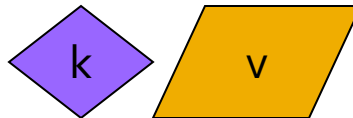
Input  
key-value pairs



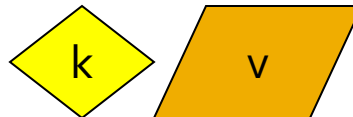
...



Intermediate  
key-value pairs



...

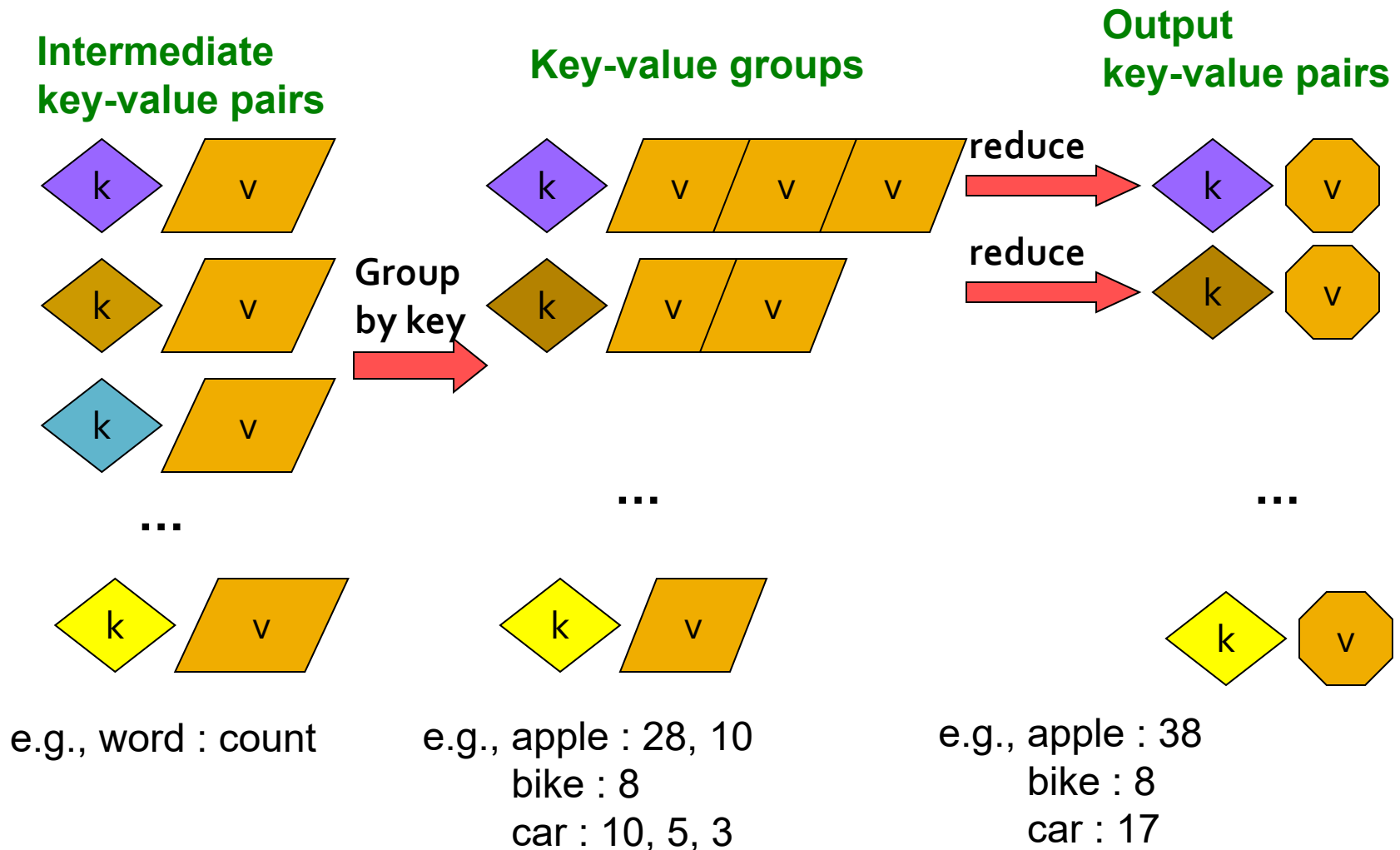


e.g., doc ID : doc text

e.g., word : count  
e.g., car : 10; apple : 28  
bike : 8; car : 5  
car : 2; apple : 10 ...



# MapReduce: The Reduce Step



# More Specifically

- **Input:** a set of key-value pairs
- Programmer specifies two methods:
  - **Map( $k, v$ )**  $\rightarrow \langle k', v' \rangle^*$ 
    - Takes a key-value pair and outputs a set of key-value pairs
      - E.g., input key is the filename, value is a single line in the file;
      - $k'$  word,  $v'$  count
    - There is one Map call for every  $(k, v)$  pair
  - **Reduce( $k', \langle v' \rangle^*$ )**  $\rightarrow \langle k', v'' \rangle^*$ 
    - **All values  $v'$  with same key  $k'$  are reduced together and processed in  $k'$  order**
    - There is one Reduce function call per unique key  $k'$

# MapReduce: Word Counting

Provided by the  
programmer

## MAP:

Read input and  
produces a set of  
key-value pairs

## Group by key:

Collect all pairs  
with same key

Provided by the  
programmer

## Reduce:

Collect all values  
belonging to the  
key and output

The crew of the space  
shuttle Endeavor recently  
returned to Earth as  
ambassadors, harbingers of  
a new era of space  
exploration. Scientists at  
NASA are saying that the  
recent assembly of the  
Dextre bot is the first step in  
a long-term space-based  
man/machine partnership.  
"The work we're doing now  
-- the robotics we're doing -  
is what we're going to  
need .....

(The, 1)  
(crew, 1)  
(of, 1)  
(the, 1)  
(space, 1)  
(shuttle, 1)  
(Endeavor, 1)  
(recently, 1)  
....

(crew, 1)  
(crew, 1)  
(space, 1)  
(the, 1)  
(the, 1)  
(the, 1)  
(shuttle, 1)  
(recently, 1)  
...

(crew, 2)  
(space, 1)  
(the, 3)  
(shuttle, 1)  
(recently, 1)  
...

Big document

(key, value)

(key, value)

(key, value)

# Word Count Using MapReduce

**map(key, value) :**

```
// key: document name; value: text of the document
  for each word w in value:
    emit(w, 1)
```

**reduce(key, values) :**

```
// key: a word; value: an iterator over counts
  result = 0
  for each count v in values:
    result += v
  emit(key, result)
```

# Map-Reduce: Environment

## Map-Reduce environment takes care of:

- Partitioning the input data
- Scheduling the program's execution across a set of machines
- Performing the **group by key** step
- Handling machine failures
- Managing required inter-machine communication

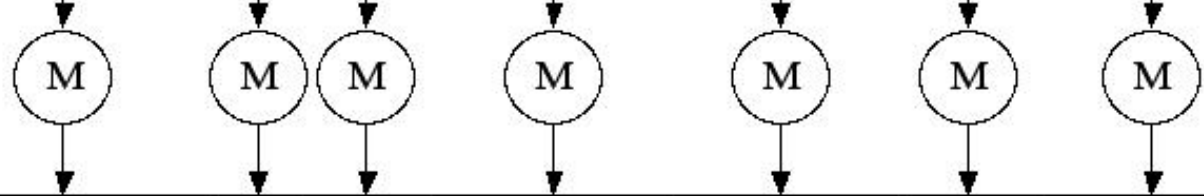
# Map-Reduce: A diagram

Input

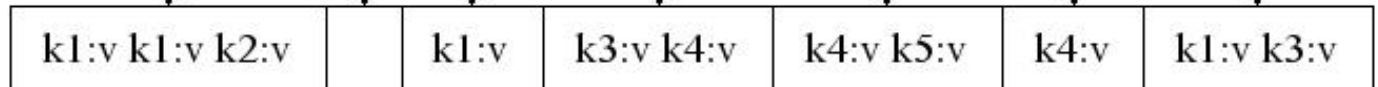


## MAP:

Read input and produces a set of key-value pairs

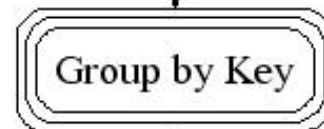


Intermediate

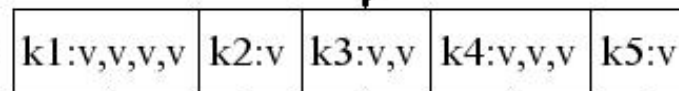


## Group by key:

Collect all pairs with same key  
(Hash merge, Shuffle, Sort, Partition)

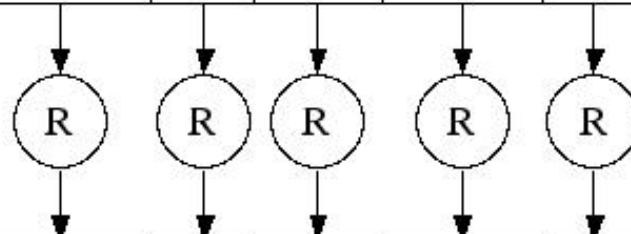


Grouped

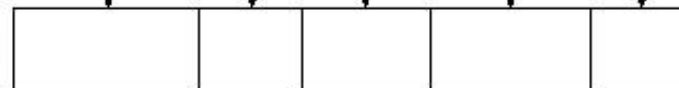


## Reduce:

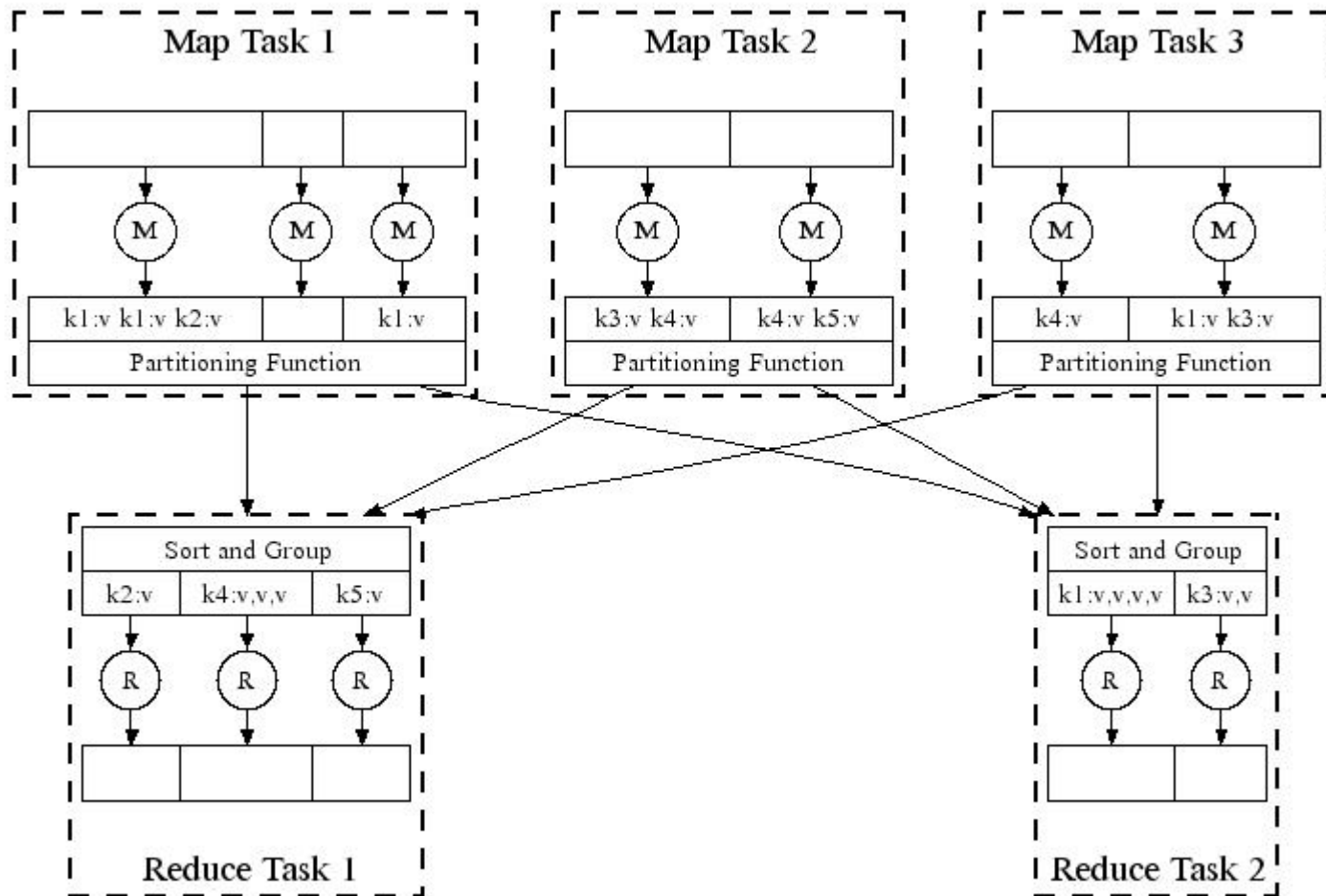
Collect all values belonging to the key and output



Output



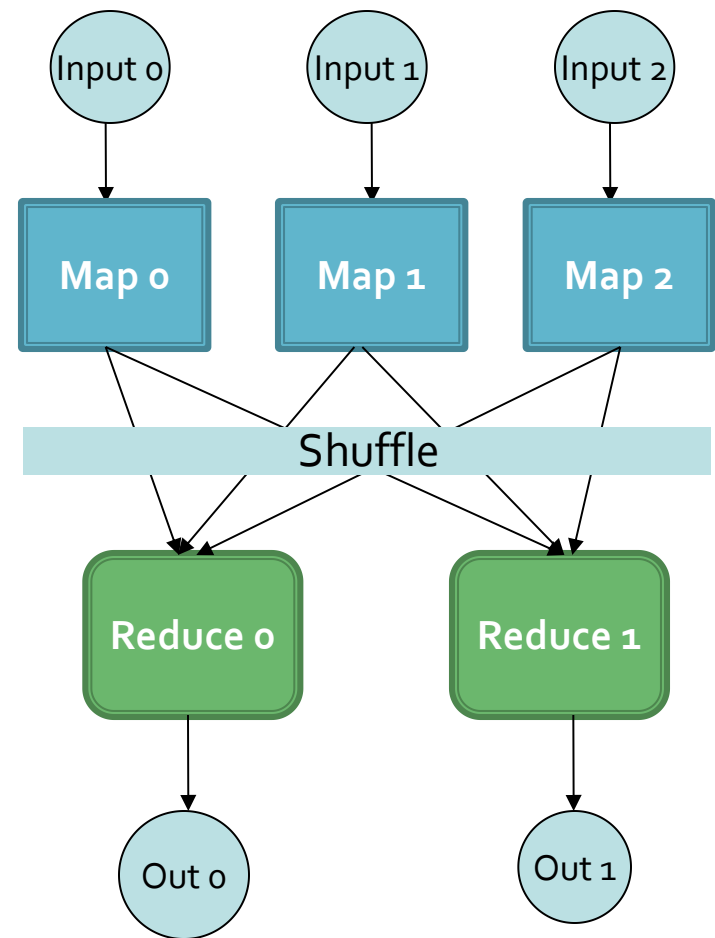
# Map-Reduce: In Parallel



All phases are distributed with many tasks doing the work

# Map-Reduce

- Programmer specifies:
  - Map and Reduce and input files
- **Workflow:**
  - Read inputs as a set of key-value-pairs
  - **Map** transforms input kv-pairs into a new set of k'v'-pairs
  - Sorts & Shuffles the k'v'-pairs to output nodes
  - All k'v'-pairs with a given k' are sent to the same **reduce**
  - **Reduce** processes all k'v'-pairs grouped by key into new k''v''-pairs
  - Write the resulting pairs to files
- All phases are distributed with many tasks doing the work





# Data Flow

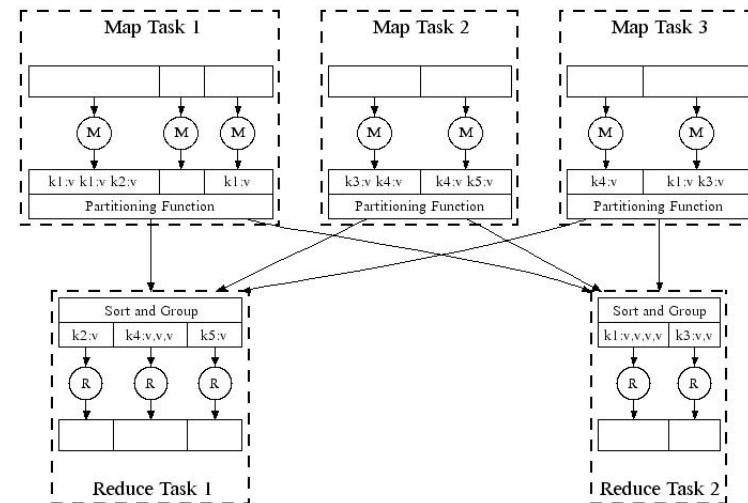
- **Input and final output are stored on a distributed file system (FS):**
  - Scheduler tries to schedule map tasks “close” to physical storage location of input data
- **Intermediate results are stored on local FS of Map and Reduce workers**
- **Output is often input to another MapReduce task**

# Coordination: Master

- **Master node takes care of coordination:**
  - **Task status:** (idle, in-progress, completed)
  - **Idle tasks** get scheduled as workers become available
  - When a map task completes, it sends the master the location and sizes of its  $R$  intermediate files
  - Master pushes this info to reducers
- Master pings workers periodically to detect failures

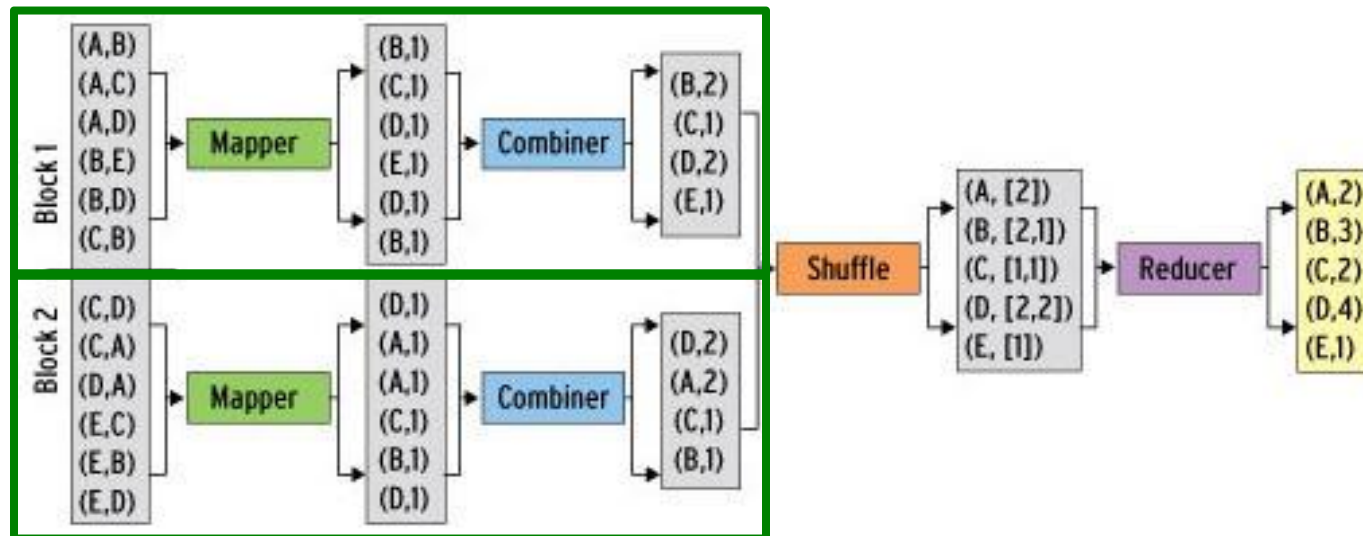
# Refinement: Combiners

- Often a Map task will produce many pairs of the form  $(k, v_1), (k, v_2), \dots$  for the same key  $k$ 
  - E.g., popular words in the word count example
- **Can save network time by pre-aggregating values in the mapper:**
  - $\text{combine}(k, \text{list}(v_1)) \rightarrow v_2$
  - Combiner is usually same as the reduce function



# Refinement: Combiners

- **Back to our word counting example:**
  - Combiner combines the values of all keys of a single mapper (single machine):



- Much less data needs to be copied and shuffled!

# Problems Suited for Map-Reduce

# Example: Host size

- **Suppose we have a large web corpus**
- **Look at the metadata file**
  - Lines of the form: (URL, size, date, ...)
- **For each host, find the total number of bytes**
  - That is, the sum of the page sizes for all URLs from that particular host
- **Other examples:**
  - Link analysis and graph processing
  - Machine Learning algorithms [1]

[1] Map-Reduce for Machine Learning on Multicore, CT Chu et al.

# Example: Join By Map-Reduce

- Compute the natural join  $R(A,B) \bowtie S(B,C)$  (on  $R.B=S.B$ )
- $R$  and  $S$  are each stored in files
- Tuples are pairs  $(a,b)$  or  $(b,c)$

A	B
a <sub>1</sub>	b <sub>1</sub>
a <sub>2</sub>	b <sub>1</sub>
a <sub>3</sub>	b <sub>2</sub>
a <sub>4</sub>	b <sub>3</sub>

R

$\bowtie$

B	C
b <sub>2</sub>	c <sub>1</sub>
b <sub>2</sub>	c <sub>2</sub>
b <sub>3</sub>	c <sub>3</sub>

S

=

A	C
a <sub>3</sub>	c <sub>1</sub>
a <sub>3</sub>	c <sub>2</sub>
a <sub>4</sub>	c <sub>3</sub>

# Map-Reduce Join

- **A Map process turns:**
  - Each input tuple  $R(a,b)$  into key-value pair  $(b,(a,R))$
  - Each input tuple  $S(b,c)$  into  $(b,(c,S))$
- **Map processes** send each key-value pair with key  $b$  to Reduce
- Each **Reduce process** matches all the pairs  $(b,(a,R))$  with all  $(b,(c,S))$  and outputs  $(a,b,c)$ .



# Implementations

## ■ Hadoop

- An open-source implementation in Java
- Uses HDFS for stable storage
- Download: <https://hadoop.apache.org/>

## ■ Spark

- Spark can process data in-memory
- Generally outperforms Hadoop, works well for smaller data sets that can all fit into a server's RAM.
- Has MLlib for machine learning
- Download: <https://spark.apache.org/>

They all have **standard-alone versions**, that you can try on a single computer.