# MapReduce and Hadoop

## Alan Sussman, Department of Computer Science

**With thanks to MMDS authors Leskovec, Rajaraman, Ullman, www.mmds.org**

UNIVERSITY OF
MARYLAND

# Announcements

- Assignment 4 posted Tuesday, due May 2 at 11:59 pm

- Quiz 2 done, 3rd quiz likely last week of class

# MapReduce

- **Challenges:**

  - How to distribute computation?

  - Distributed/parallel programming is hard

- **Map-reduce** addresses this problem for certain kinds of computations

  - Started as Google's computational/data manipulation model

  - Overall, an elegant way to work (in parallel) with big data
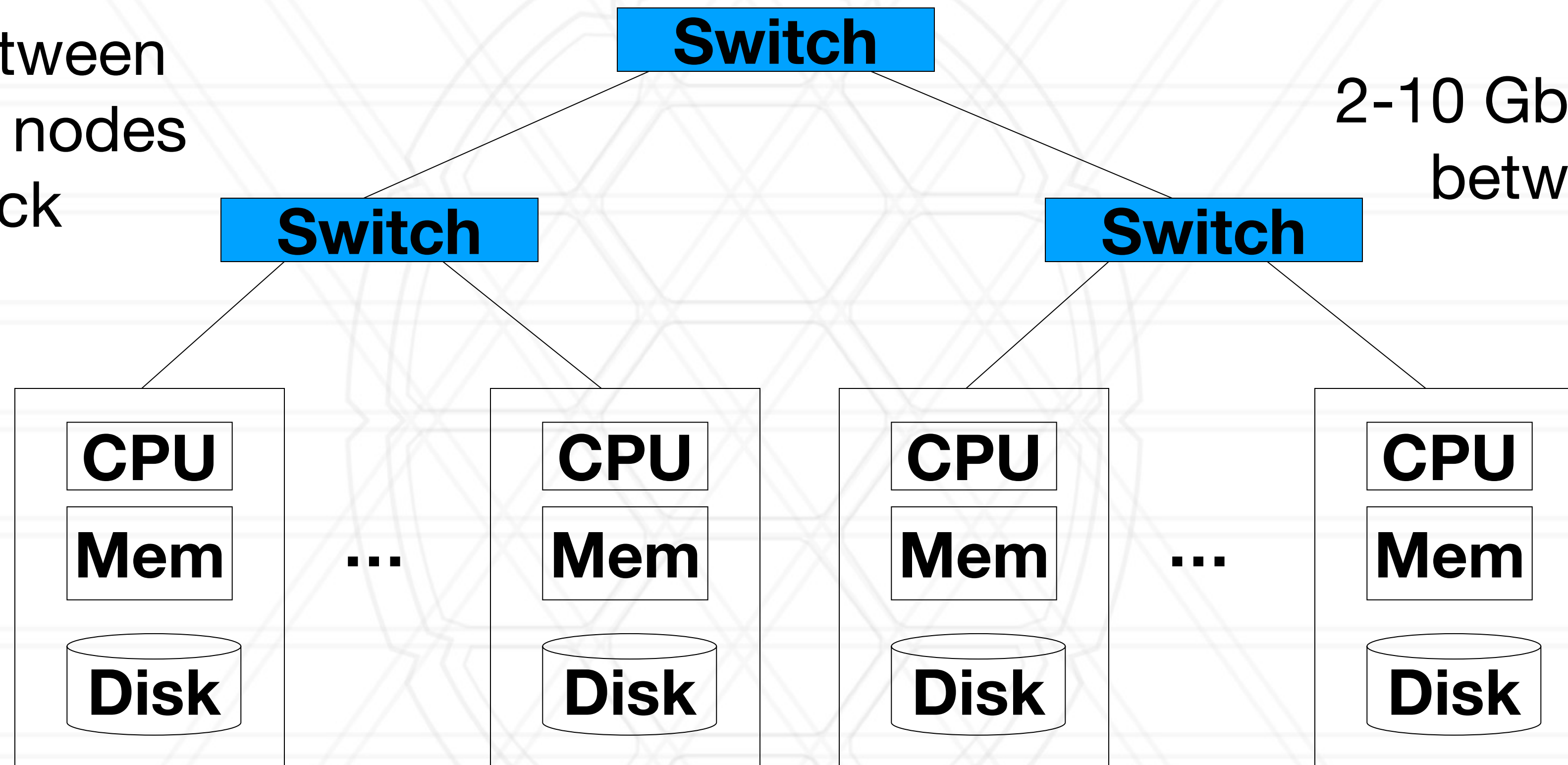
DEPARTMENT OF
COMPUTER SCIENCE

# Motivation: Google Example

- 20+ billion web pages x 20KB = 400+ TB

- If 1 computer reads 30-35 MB/sec from disk

  - ~4 months to read the web

- ~1,000 hard drives to store the web

- Takes even more to **do** something useful with the data!

- **A standard architecture for such problems emerged several years ago:**

  - Cluster of commodity Linux nodes

  - Commodity network (e.g., Ethernet) to connect them

# Cluster Architecture several years ago

1 Gbps between any pair of nodes in a rack

2-10 Gbps backbone between racks

**Switch**

**Switch**

**Switch**

| CPU |
| Mem |
| Disk |

...

| CPU |
| Mem |
| Disk |

| CPU |
| Mem |
| Disk |

...

| CPU |
| Mem |
| Disk |

Each rack contains 16-64 nodes

In 2011 it was guesstimated that Google had 1M machines, http://bit.ly/Shh0RO

DEPARTMENT OF
COMPUTER SCIENCE

# Large-scale Computing

- **Large-scale computing** for **data analytics** 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 lose 1/day

    - People estimated Google had ~1M machines in 2011

      - 1,000 machines fail every day!

DEPARTMENT OF
COMPUTER SCIENCE

# Idea and Solution

- **Issue:** **Copying data over a network takes time**

- **Idea:**

  - Bring computation close to the data

  - Store files multiple times in multiple locations 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:**

    - Provides global file namespace

    - Google GFS; Hadoop HDFS;

- **Typical usage pattern**

  - Huge files (100s of GB to TB)

  - Data is rarely updated in place

  - Reads and appends are common

DEPARTMENT OF
COMPUTER SCIENCE

# 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 in the cluster
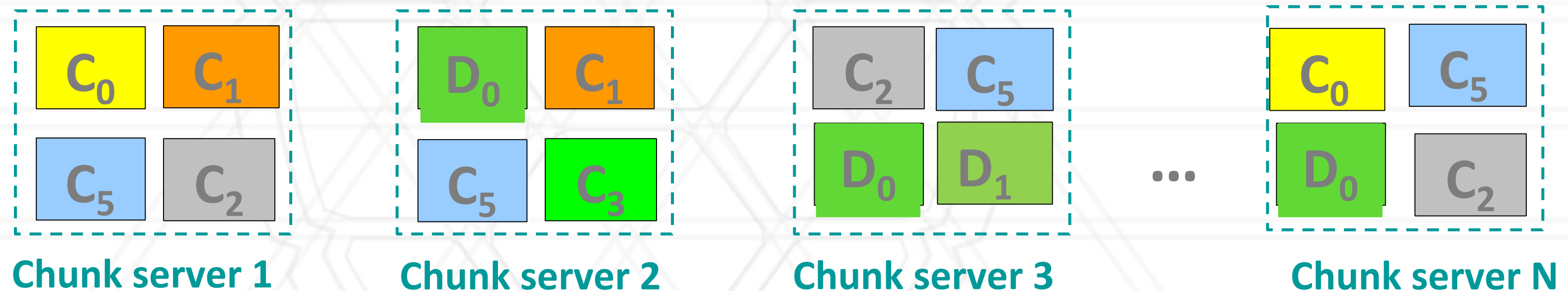
- **Primary node**

  - a.k.a. Name Node in Hadoop's HDFS

  - Stores metadata about where files are stored

  - Might be replicated

- **Client library for file access**

  - Talks to primary to find chunk servers

  - Connects directly to chunk servers to access data

DEPARTMENT OF
COMPUTER SCIENCE

# 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

| Chunk server 1 | Chunk server 2 | Chunk server 3 | Chunk server N |
|---|---|---|---|
| $C_0$ $C_1$ | $D_0$ $C_1$ | $C_2$ $C_5$ | $C_0$ $C_5$ |
| $C_5$ $C_2$ | $C_5$ $C_3$ | $D_0$ $D_1$ | $D_0$ $C_2$ |

...

**Bring computation directly to the data!**

**Chunk servers also serve as compute servers**

# 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

## The Problem:

- Count occurrences of words in a document:

  - `words(doc.txt) | sort | uniq -c`

    - where `words` takes a file and outputs the words in it, one per line

- This pipeline captures the essence of **MapReduce**

  - Great thing is that it is naturally parallelizable

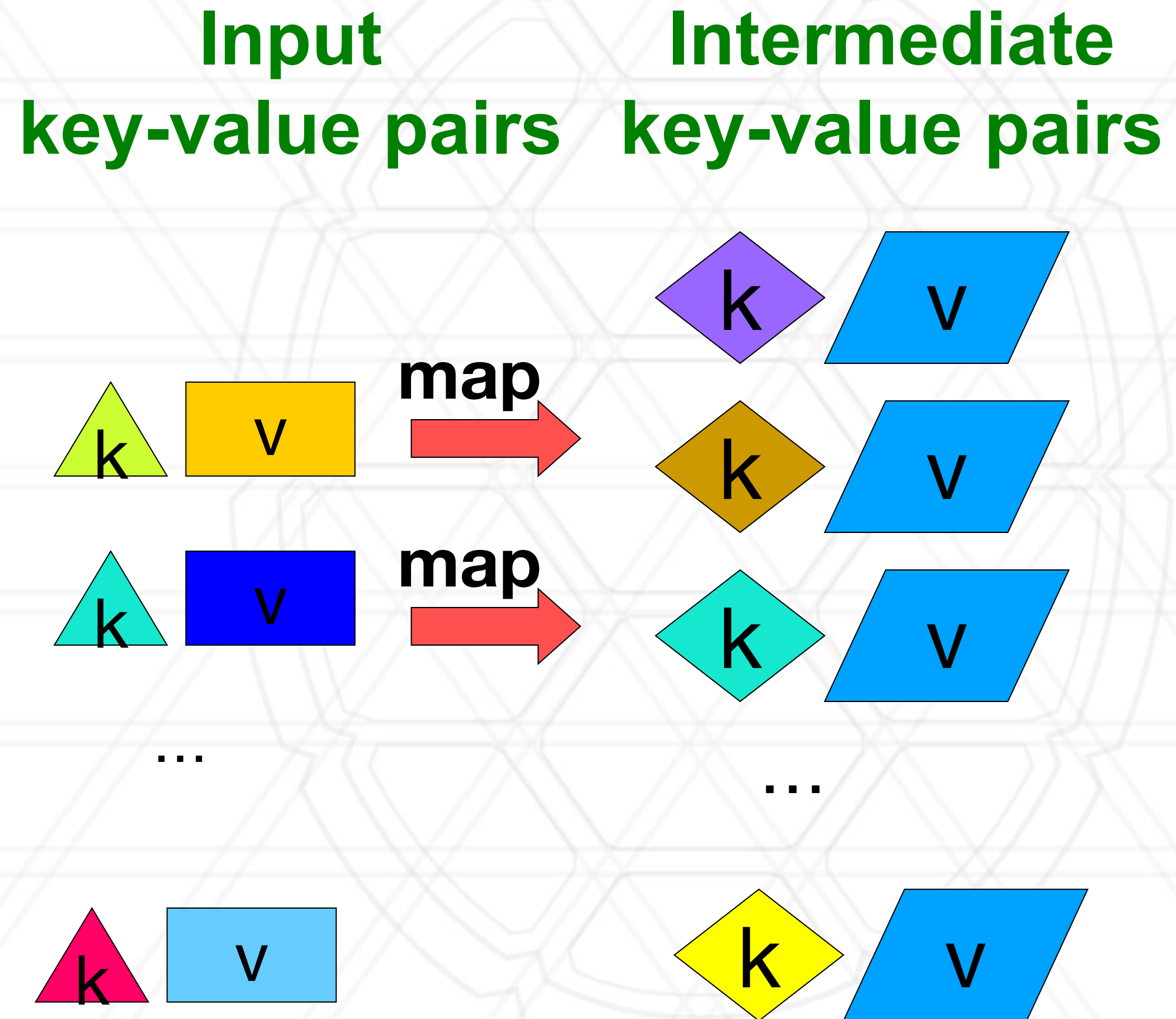# MapReduce: Overview

- Inspired by LISP

- Map(function, set of values)

  - Applies function to each value in the set
    <span style="color:red">(map 'length '(() (a) (a b) (a b c))) ⇒ (0 1 2 3)</span>

- Reduce(function, set of values)

  - Combines all the values using a binary function (e.g.,+)
    <span style="color:red">(reduce #'+ '(1 2 3 4 5)) ⇒ 15</span>
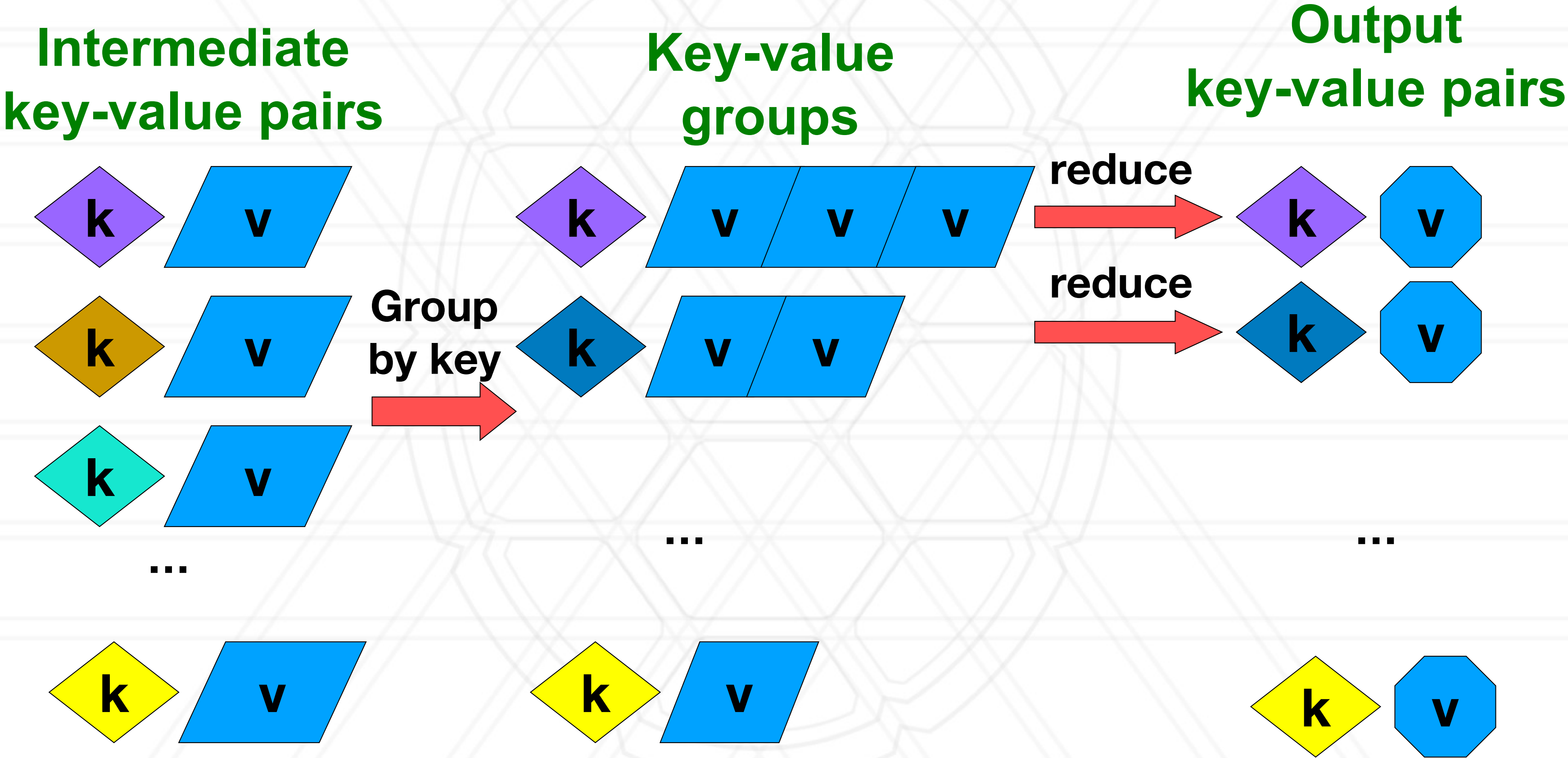
# MapReduce: Overview

- Sequentially read a lot of data

- **Map:**

  - Extract something you care about

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

- **Reduce:**

  - Aggregate, summarize, filter or transform

- Write the result

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

DEPARTMENT OF
COMPUTER SCIENCE

# MapReduce: The Map Step

**Input key-value pairs**

**Intermediate key-value pairs**

map

map

...

...

# MapReduce: The Reduce Step

Intermediate
key-value pairs

Key-value
groups

Output
key-value pairs

DEPARTMENT OF
COMPUTER SCIENCE

# More Specifically

- **Input:** a set of key-value pairs

- Programmer specifies two methods:

  - **Map(k, v)** $\rightarrow$ <k', v'>*

    - Takes a key-value pair and outputs a set of key-value pairs

      - E.g., key is the filename, value is a single line in the file

    - There is one Map call for every *(k,v)* pair

  - **Reduce(k', <v'>*)** $\rightarrow$ <k', v''>*

    - **All values *v'* with same key *k'* are reduced together and processed in *v'* order**

    - There is one Reduce function call per unique key *k'*

DEPARTMENT OF
COMPUTER SCIENCE

# MapReduce: Word Counting

**MAP:**
Read input and produces a set of key-value pairs

**Group by key:**
Collect all pairs with same key

**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/mache 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)
...

**Only sequential reads**

**Big document**        **(key, value)**        **(key, value)**        **(key, value)**

DEPARTMENT OF COMPUTER SCIENCE

Alan Sussman & Abhinav Bhatele (CMSC416)

# 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)
```

DEPARTMENT OF
COMPUTER SCIENCE

# Map-Reduce: Environment
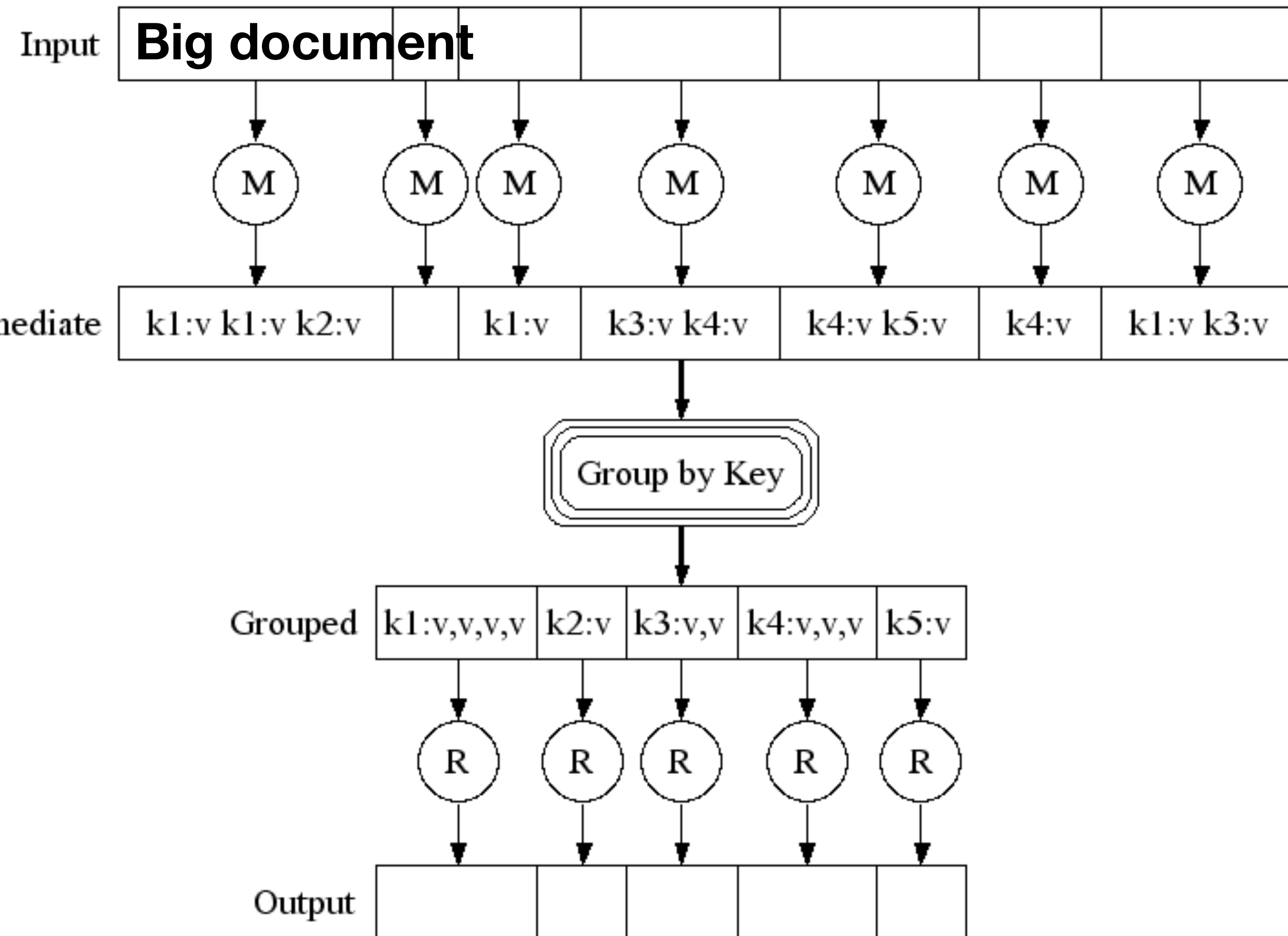
**Map-Reduce environment (runtime system) takes care of:**

- Partitioning the input data

- Scheduling the program's execution across a set of machines

- Performing the **group by key** step

- Handling machine failures

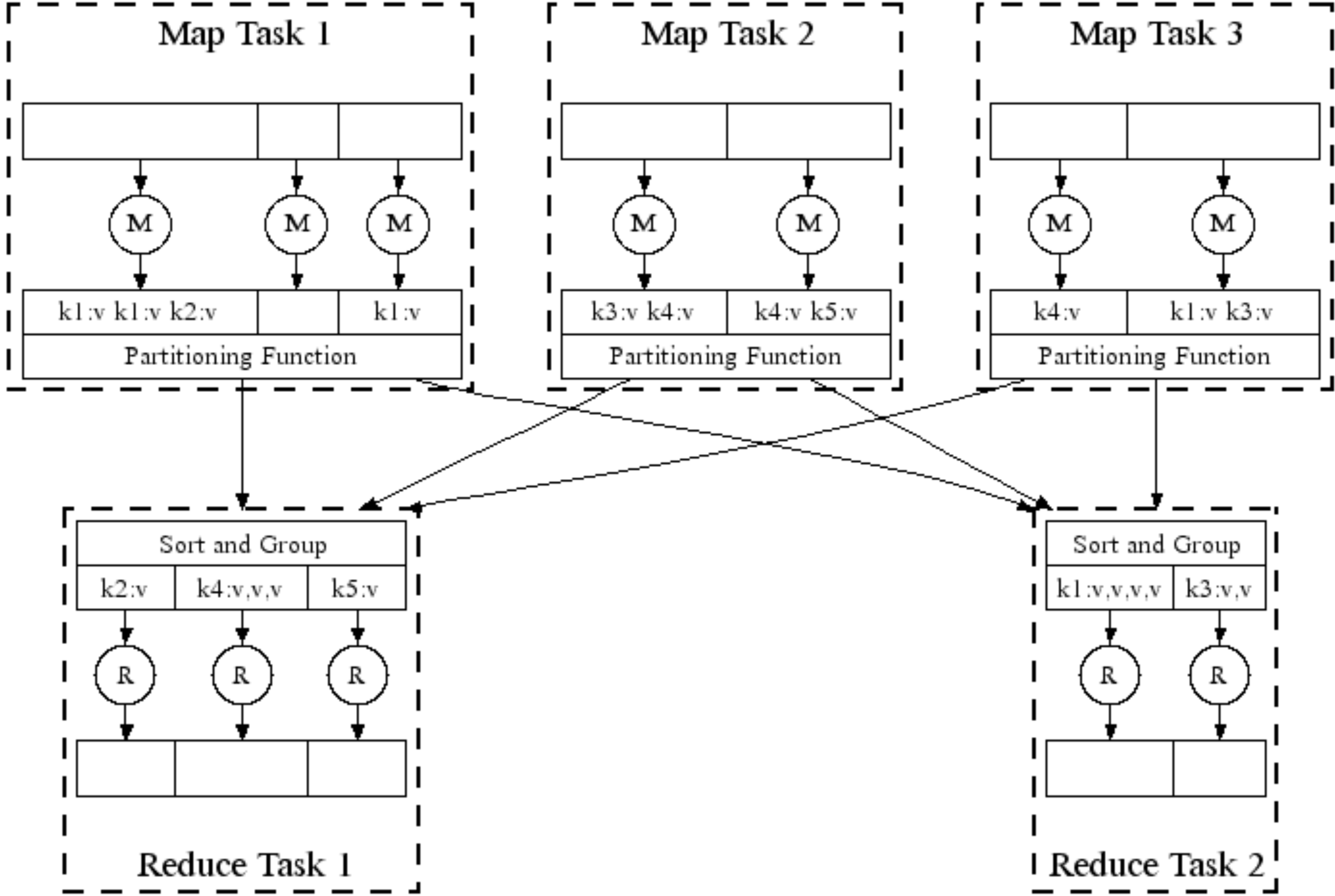- Managing required inter-machine communication

# Map-Reduce: A diagram

**MAP:**
Read input and produces a set of key-value pairs

**Group by key:**
Collect all pairs with same key
(Hash merge, Shuffle, Sort, Partition)

**Reduce:**
Collect all values belonging to the key and output

Input | **Big document**

M  M M  M  M  M  M

Intermediate | k1:v k1:v k2:v | | k1:v | k3:v k4:v | k4:v k5:v | k4:v | k1:v k3:v

Group by Key

Grouped | k1:v,v,v,v | k2:v | k3:v,v | k4:v,v,v | k5:v

R  R  R  R  R

Output

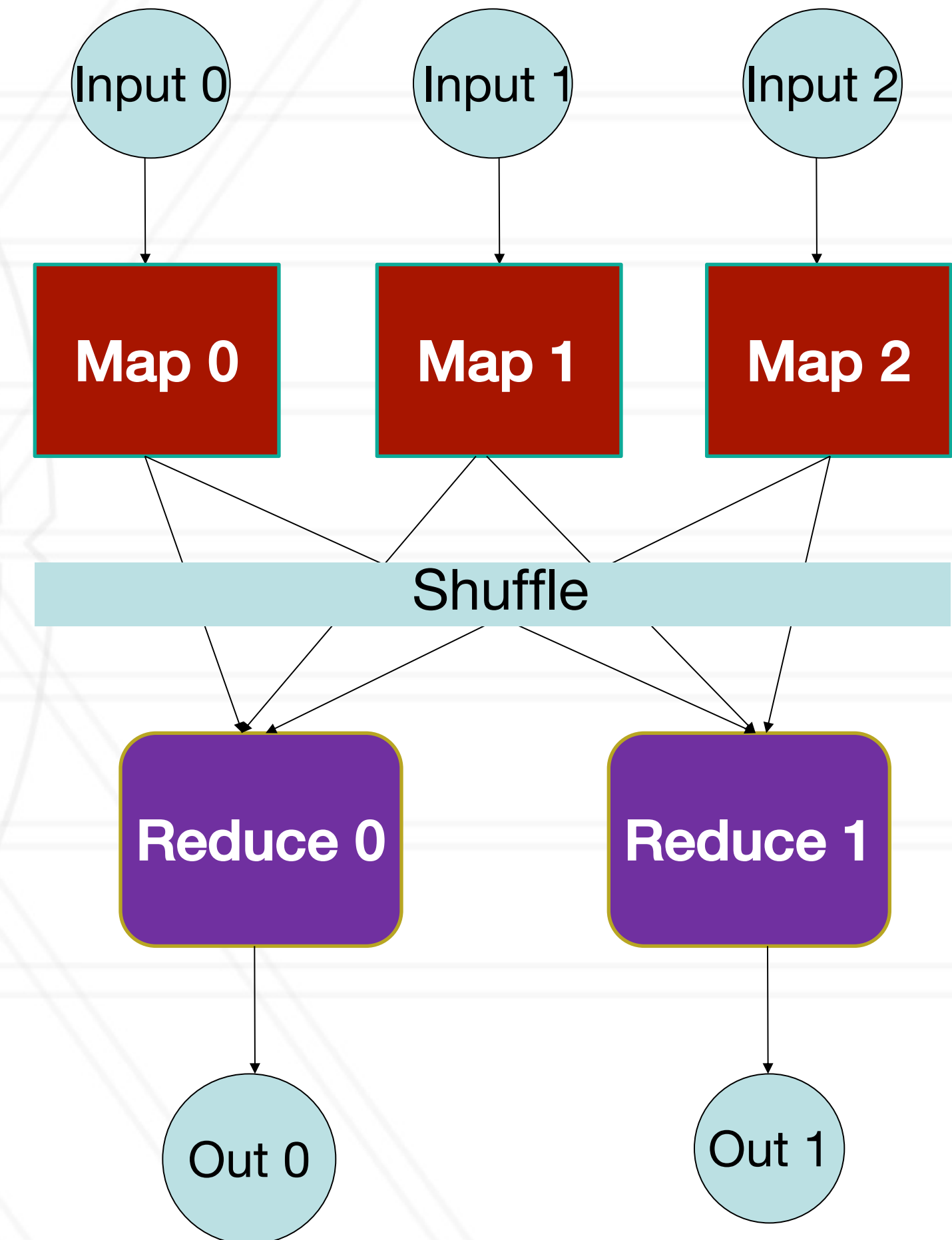DEPARTMENT OF COMPUTER SCIENCE

# 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

DEPARTMENT OF
COMPUTER SCIENCE

# Hadoop

# Announcements

- Assignment 4 due May 2 at 11:59 pm

  - Questions?

- Quiz 3 last week of class

DEPARTMENT OF
COMPUTER SCIENCE

# Data Flow – Hadoop architecture

- **Input and final output** are stored in 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**
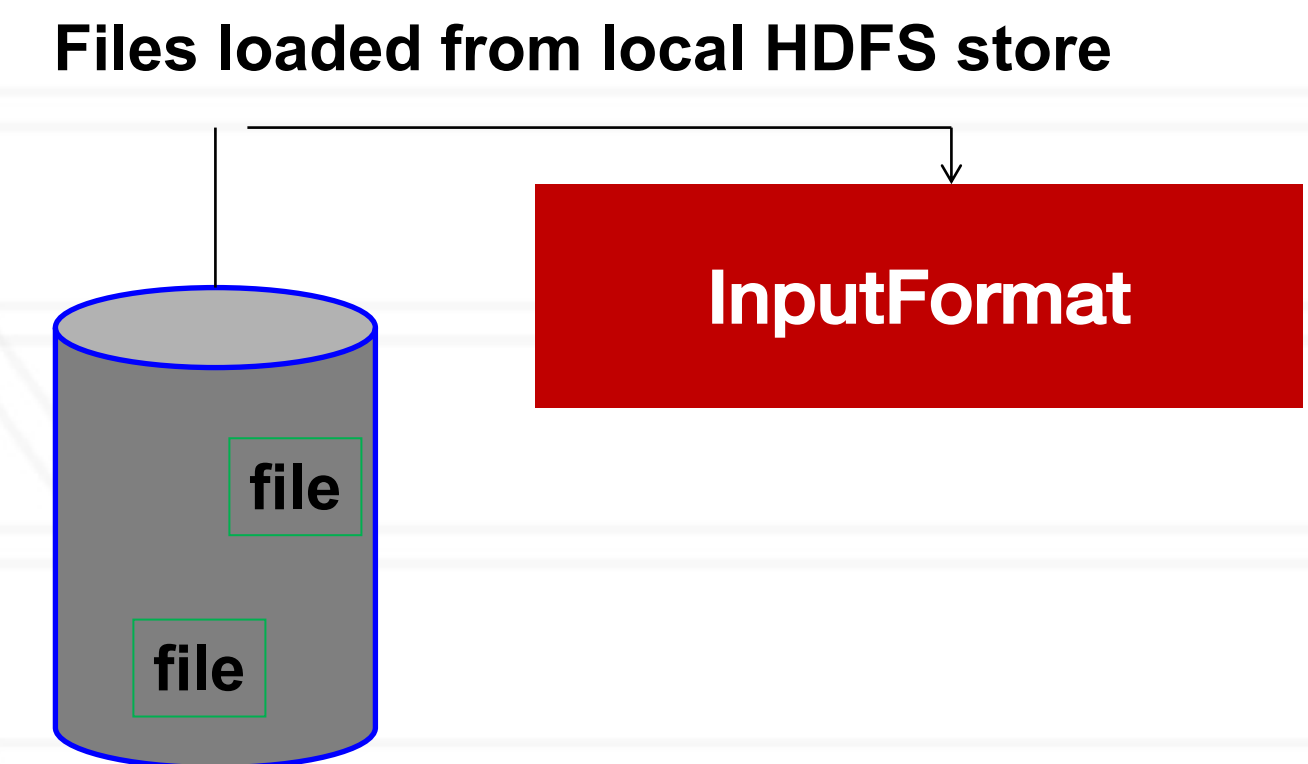
DEPARTMENT OF
COMPUTER SCIENCE

# Input Files

- *Input files* are where the data for a MapReduce task is initially stored

- The input files typically reside in a distributed file system (e.g., HDFS)

- The format of input files is arbitrary

  - Line-based log files

  - Binary files

  - Multi-line input records

  - Or something else entirely, e.g., a database

file

file

Gregory Kesden, Carnegie Mellon U.

DEPARTMENT OF
COMPUTER SCIENCE

# InputFormat

- How the input files are split up and read is defined by the *InputFormat*

- InputFormat is a class that does the following:

  - Selects the files that should be used for input

  - Defines the *InputSplits* that break a file

  - Provides a factory for *RecordReader* objects that read the file

**Files loaded from local HDFS store**

**InputFormat**

**file**

**file**

Gregory Kesden, Carnegie Mellon U.

DEPARTMENT OF
COMPUTER SCIENCE

# InputFormat Types

- Several InputFormats are provided with Hadoop:

| InputFormat | Description | Key | Value |
|---|---|---|---|
| TextInputFormat | Default format; reads lines of text files | The byte offset of the line | The line contents |
| KeyValueInputFormat | Parses lines into (K, V) pairs | Everything up to the first tab character | The remainder of the line |
| SequenceFileInputFormat | A Hadoop-specific high-performance binary format | user-defined | user-defined |

Gregory Kesden, Carnegie Mellon U.
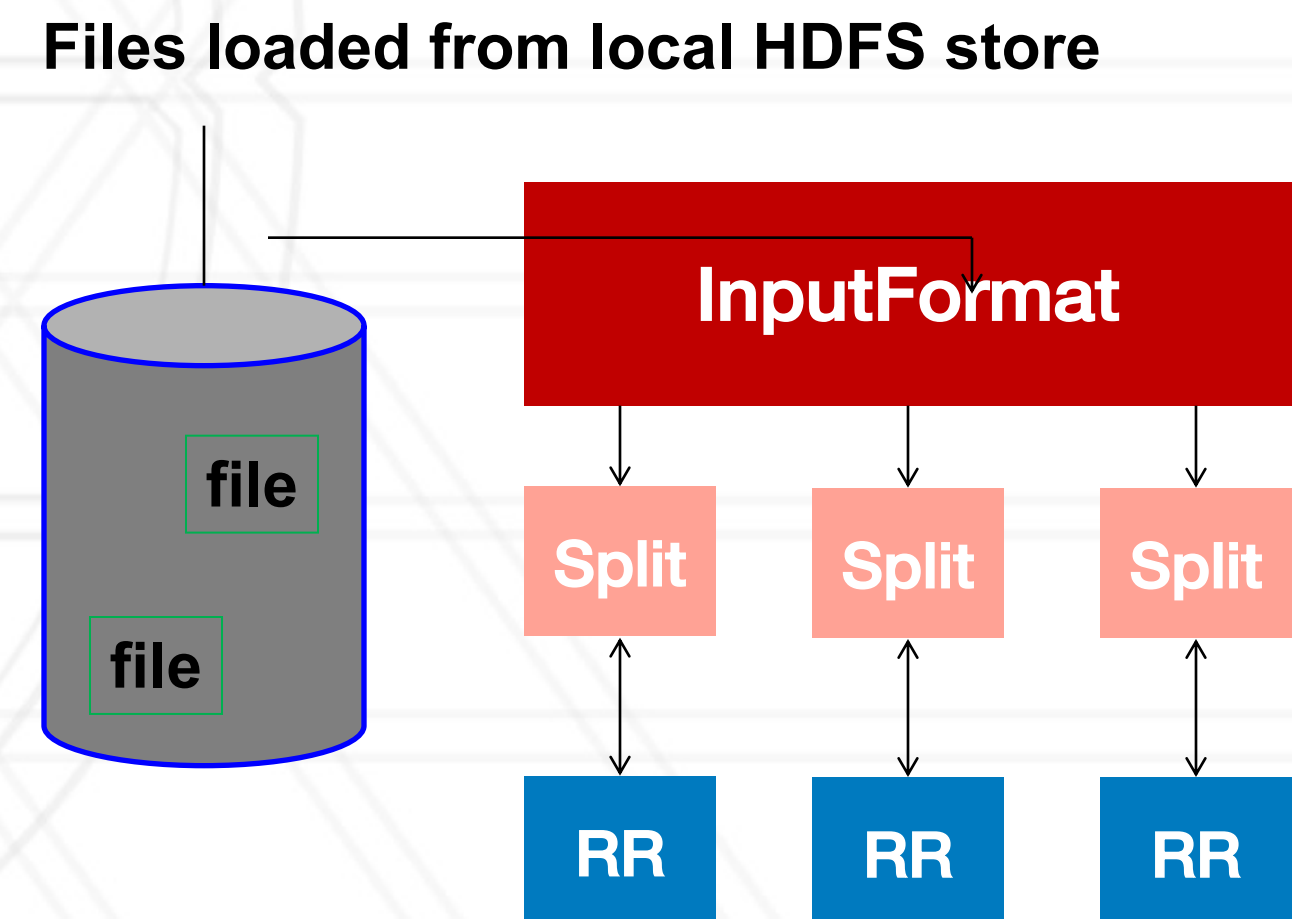
DEPARTMENT OF
COMPUTER SCIENCE

# Input Splits

- An *input split* describes a unit of work that comprises a single map task in a MapReduce program

- By default, the InputFormat breaks a file up into 64MB splits

- By dividing the file into splits, allow several map tasks to operate on a single file in parallel

- If the file is very large, this can improve performance significantly through parallelism

- Each map task corresponds to a *single* input split

**Files loaded from local HDFS store**

file

file

**InputFormat**

Split  Split  Split

Gregory Kesden, Carnegie Mellon U.

# RecordReader

- The input split defines a slice of work but does not describe how to access it

- The *RecordReader* class actually loads data from its source and converts it into (K, V) pairs suitable for reading by Mappers

- The RecordReader is invoked repeatedly on the input until the entire split is consumed

- Each invocation of the RecordReader lead to another call of the map function defined by the programmer
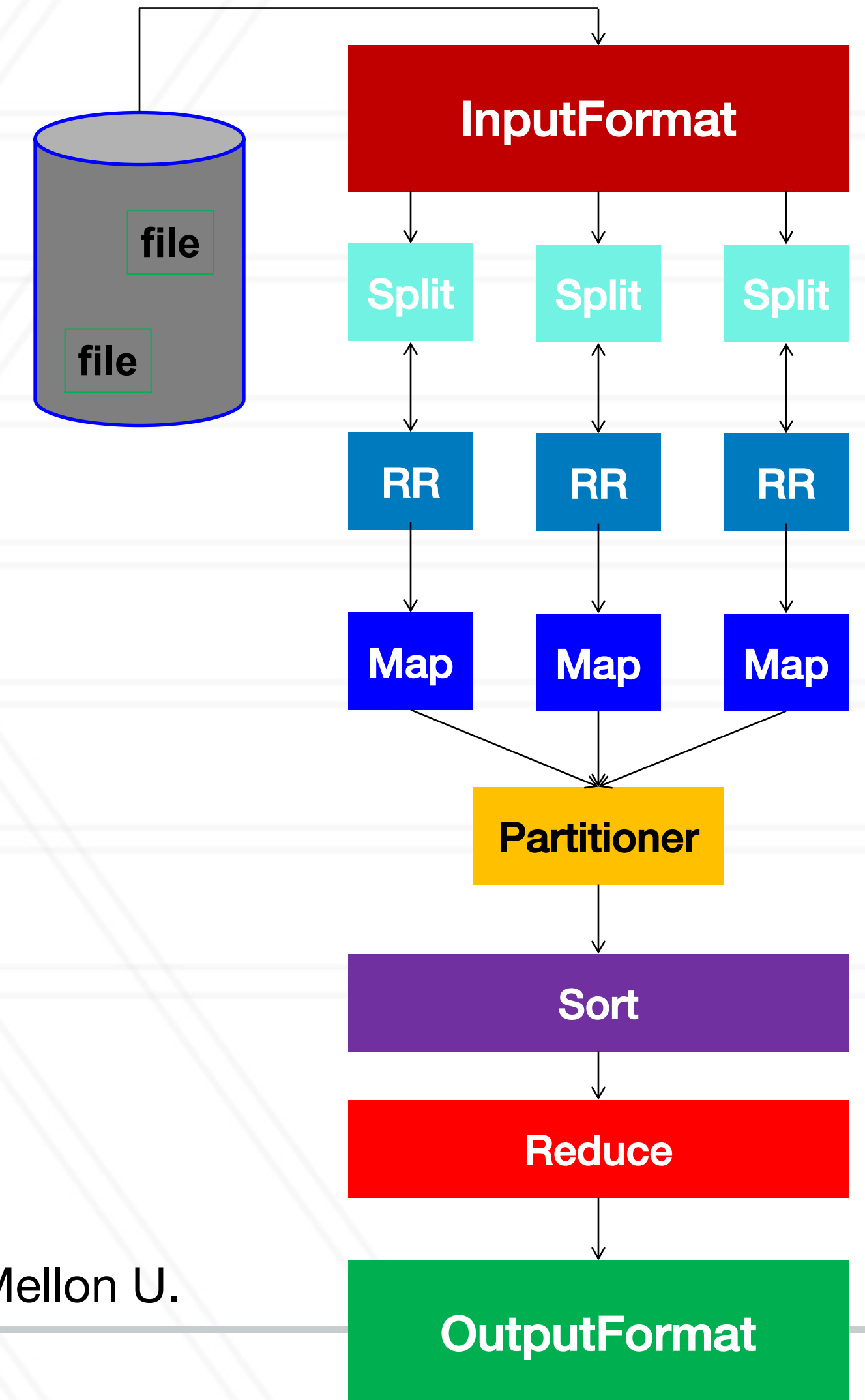
**Files loaded from local HDFS store**

**InputFormat**

file

file

Split   Split   Split

RR   RR   RR

Gregory Kesden, Carnegie Mellon U.

Alan Sussman & Abhinav Bhatele (CMSC416)

# OutputFormat

- The *OutputFormat* class defines the way (K,V) pairs produced by Reducers are written to output files

- The instances of OutputFormat provided by Hadoop write to files on the local disk or in HDFS

- Several OutputFormats are provided by Hadoop:

| OutputFormat | Description |
|---|---|
| TextOutputFormat | Default; writes lines in "key \t value" format |
| SequenceFileOutputFormat | Writes binary files suitable for reading into subsequent MapReduce jobs |
| NullOutputFormat | Generates no output files |

Gregory Kesden, Carnegie Mellon U.

Alan Sussman & Abhinav Bhatele (CMSC416)

# 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, one for each reducer

  - Master pushes this info to reducers

- Master pings workers periodically to detect failures

# Dealing with Failures

- ## Map worker failure

  - Map tasks completed or in-progress at worker are reset to idle

  - Reduce workers are notified when task is rescheduled on another worker

- ## Reduce worker failure

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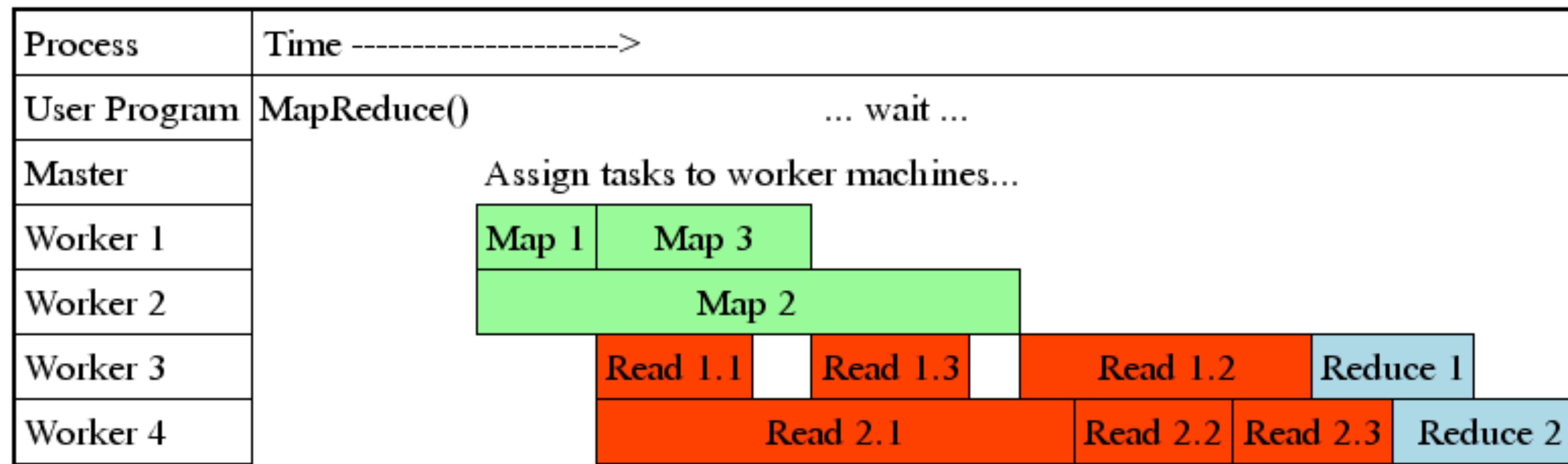
  - Reduce task is restarted

- ## Master failure

  - MapReduce task is aborted and client is notified

DEPARTMENT OF
COMPUTER SCIENCE

# How many Map and Reduce jobs?

- *M* map tasks, *R* reduce tasks

- **Rule of a thumb:**

  - Make *M* much larger than the number of nodes in the cluster

  - One DFS chunk per map is common

  - Improves dynamic load balancing and speeds up recovery from worker failures

- **Usually *R* is smaller than *M***

  - Because output is spread across *R* files

DEPARTMENT OF
COMPUTER SCIENCE

# Task Granularity & Pipelining

- **Fine granularity tasks:** map tasks >> machines

    - Minimizes time for fault recovery

    - Can do pipeline shuffling with map execution

    - Better dynamic load balancing

| Process | Time --------------------> | | | | | |
|---|---|---|---|---|---|---|
| User Program | MapReduce() | | | ... wait ... | | |
| Master | | Assign tasks to worker machines... | | | | |
| Worker 1 | | Map 1 | Map 3 | | | |
| Worker 2 | | Map 2 | | | | |
| Worker 3 | | Read 1.1 | Read 1.3 | Read 1.2 | Reduce 1 | |
| Worker 4 | | Read 2.1 | | Read 2.2 | Read 2.3 | Reduce 2 |

DEPARTMENT OF COMPUTER SCIENCE

# Refinements: Backup Tasks

- **Problem**

  - Slow workers significantly lengthen the job completion time:

    - Other jobs on the machine

    - Bad disks

    - Weird things

- **Solution**

  - Near end of phase, spawn backup copies of tasks

    - Whichever one finishes first "wins"

- **Effect**

  - Dramatically shortens job completion time

# Refinement: Combiners

- Often a Map task will produce many pairs of the form $(k, v_1)$, $(k, v_2)$, ... for the same key $k$
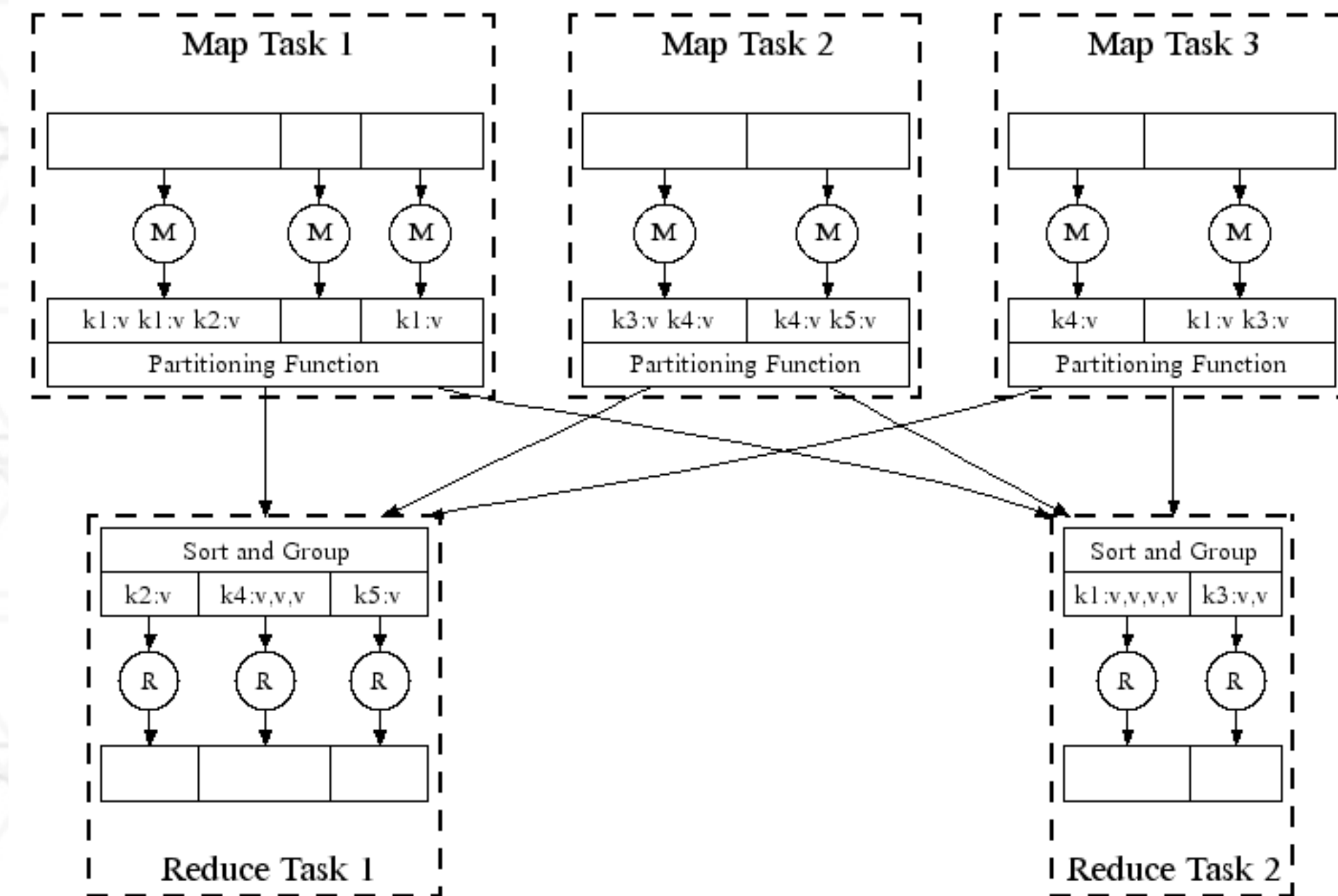
  - E.g., popular words in the word count example

- **Can save network time by** <span style="color:green">**pre-aggregating values in the mapper:**</span>

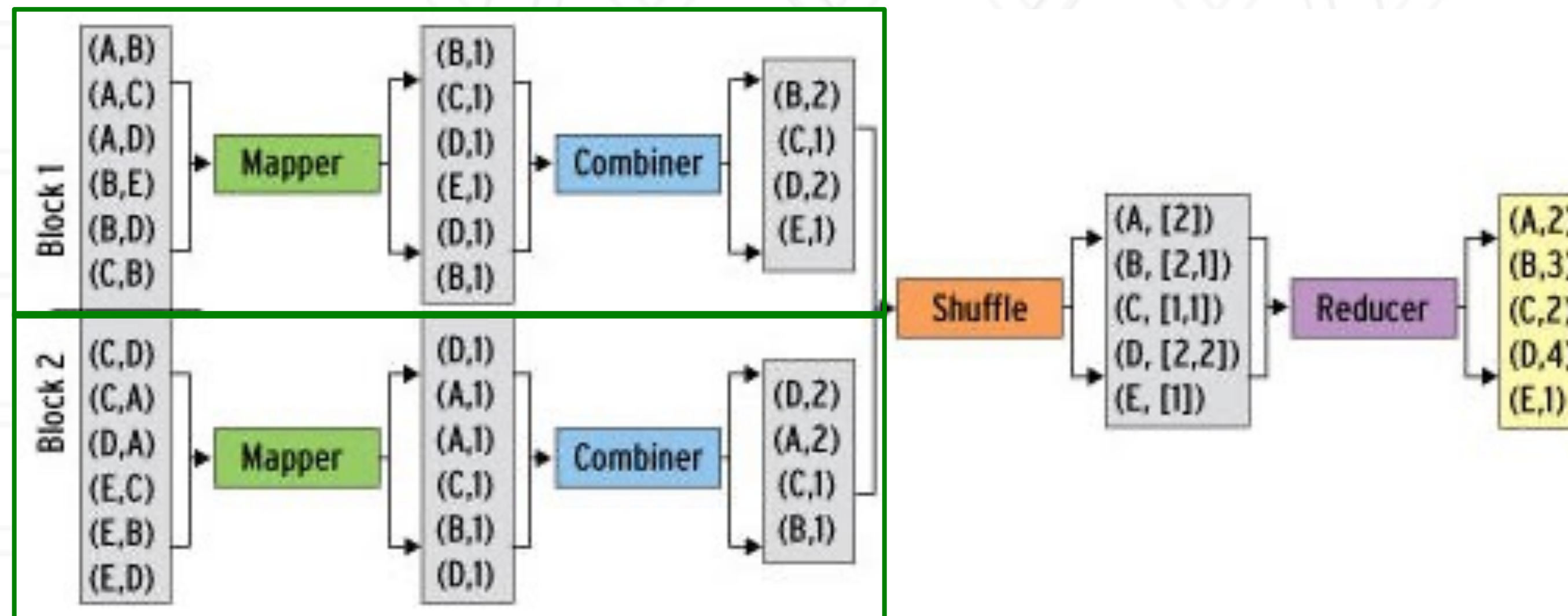  - combine(k, list($v_1$)) $\rightarrow$ $v_2$

  - Combiner is usually same as the reduce function

- Works only if reduce function is commutative and associative

# Refinement: Combiners

- **Back to the 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!

DEPARTMENT OF
COMPUTER SCIENCE

# Refinement: Partition Function

- **Want to control how keys get partitioned**

  - Inputs to map tasks are created by contiguous splits of input file

  - Reduce needs to ensure that records with the same intermediate key end up at the same worker

- **System uses a default partition function:**

  - **hash(key) mod $R$**

- **Sometimes useful to override the hash function:**

  - E.g., **hash(hostname(URL)) mod $R$** ensures URLs from a host end up in the same output file

# Applications

- Three major classes:

  - Text tokenization, indexing, and search

  - Creation of other kinds of data structures (e.g., graphs)

  - Data mining and machine learning

- See list at https://cwiki.apache.org/confluence/display/HADOOP2/poweredby

- For Machine Learning algorithms, see MAHOUT at http://mahout.apache.org/

  - Default backend is now Spark

DEPARTMENT OF
COMPUTER SCIENCE

# Example: Host size
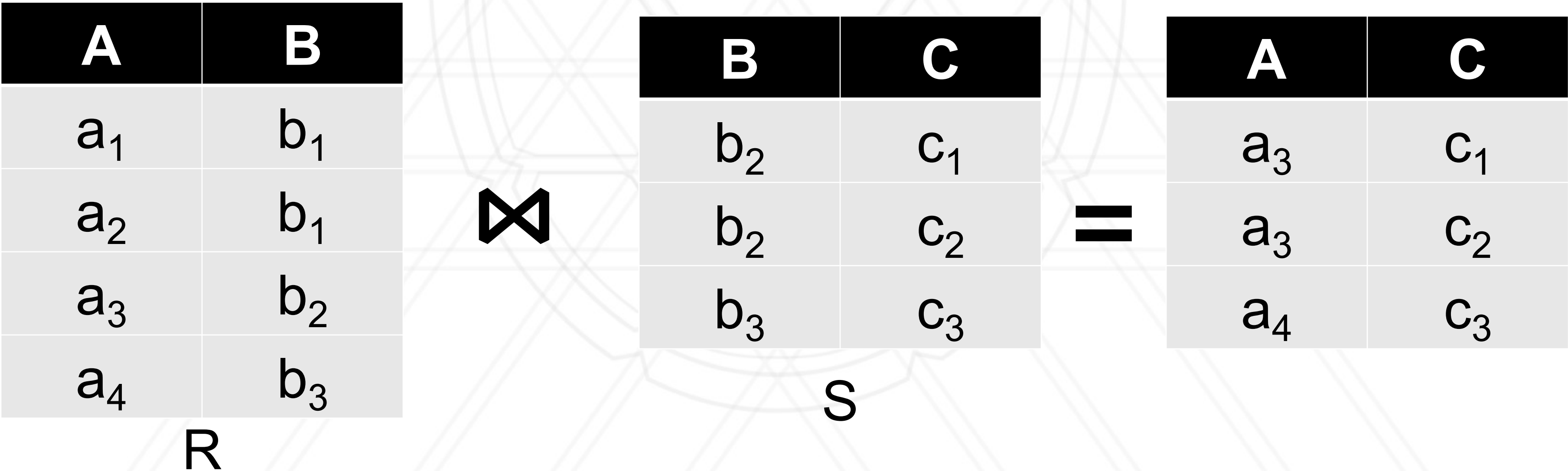
- **Suppose we have a large web corpus**

- 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

# Example: Language Model

- **Statistical machine translation:**

  - Need to count number of times every 5-word sequence occurs in a large corpus of documents

- **Very easy with MapReduce:**

  - **Map:**

    - Extract (5-word sequence, count) from document

  - **Reduce:**

    - Combine the counts

# Example: Database Join By Map-Reduce

- **Compute the natural join $R(A,B) \bowtie S(B,C)$**

- $R$ and $S$ are each stored in files

- Tuples are pairs *(a,b)* or *(b,c)*

| A | B |
|---|---|
| $a_1$ | $b_1$ |
| $a_2$ | $b_1$ |
| $a_3$ | $b_2$ |
| $a_4$ | $b_3$ |

R

$\bowtie$

| B | C |
|---|---|
| $b_2$ | $c_1$ |
| $b_2$ | $c_2$ |
| $b_3$ | $c_3$ |

S

=

| A | C |
|---|---|
| $a_3$ | $c_1$ |
| $a_3$ | $c_2$ |
| $a_4$ | $c_3$ |

DEPARTMENT OF
COMPUTER SCIENCE

# Map-Reduce Join

- **Use a hash function $h$ from B-values to $1...k$**

- **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 process $h(b)$

  - Hadoop does this automatically; just tell it what $k$ is.

- Each **Reduce process** matches all the pairs $(b,(a,R))$ with all $(b,(c,S))$ and outputs $(a,c)$.

# Cost Measures for Algorithms

- **In MapReduce we quantify the cost of an algorithm using**

1. *Communication cost* = total I/O of all processes

2. *Elapsed communication cost* = max of I/O along any path

3. (*Elapsed*) *computation cost* analogous, but count only running time of processes

Note that here the big-O notation is not the most useful
(adding more machines is always an option)

# Example: Cost Measures

- **For a map-reduce algorithm:**

  - **Communication cost =** input file size + 2 × (sum of the sizes of all files passed from Map processes to Reduce processes) + the sum of the output sizes of the Reduce processes.

  - **Elapsed communication cost** is the sum of the largest input + output for any map process, plus the same for any reduce process

DEPARTMENT OF
COMPUTER SCIENCE

# What Cost Measures Mean

- Either the I/O (communication) or processing (computation) cost dominates

  - Ignore one or the other

- Total cost tells what you pay in rent from your friendly neighborhood cloud provider

- Elapsed cost is wall-clock time using parallelism

DEPARTMENT OF
COMPUTER SCIENCE

# Cost of Map-Reduce Join

- **Total communication cost**
  = $O(|R|+|S|+|R \bowtie S|)$

- **Elapsed communication cost** = $O(s)$

  - We're going to pick $k$ and the number of Map processes so that the I/O limit $s$ is respected

  - We put a limit $s$ on the amount of input or output that any one process can have. **$s$ could be:**

    - What fits in main memory

    - What fits on local disk

- With proper indexes, computation cost is linear in the input + output size

  - So computation cost is like comm. cost

DEPARTMENT OF
COMPUTER SCIENCE

Pointers and Further Reading

# Implementations

- Google

  - Not available outside Google

- **Hadoop**

  - An open-source implementation in Java

  - Uses HDFS for stable storage

  - Download: http://hadoop.apache.org/

- Amazon Elastic MapReduce (EMR)

  - Hadoop MapReduce running on Amazon EC2

  - Can also run Spark, HBase, Hive, Presto …

# Reading

- Jeffrey Dean and Sanjay Ghemawat: MapReduce: Simplified Data Processing on Large Clusters

  - https://research.google.com/archive/mapreduce-osdi04.pdf


- Sanjay Ghemawat, Howard Gobioff, and Shun-Tak Leung: The Google File System

  - https://research.google.com/archive/gfs-sosp2003.pdf

DEPARTMENT OF COMPUTER SCIENCE

# Resources

- Hadoop Resources

  - Introduction

    - https://hadoop.apache.org/

    - https://hadoop.apache.org/docs/stable/hadoop-project-dist/hadoop-common/SingleCluster.html

  - Map/Reduce Overview

    - https://hadoop.apache.org/docs/r1.2.1/mapred_tutorial.html

  - Eclipse Environment

    - https://people.apache.org/~srimanth/hadoop-eclipse/

  - Javadoc

    - http://hadoop.apache.org/docs/stable/api/

DEPARTMENT OF
COMPUTER SCIENCE