# An Introduction to Cluster Computing

**Apostolos N. Papadopoulos**

Assistant Professor
Data Engineering Lab,
Department of Informatics
Aristotle University of Thessaloniki
GREECE

# Outline

- Why one machine is not enough ?

- Parallel architectures

- Scaling out vs scaling up

- Important issues in cluster computing

- Hadoop MR recap

- Theoretical Issues

# Motivation

We need **more CPUs** because:

    we can run programs faster

We need **more disks** because:

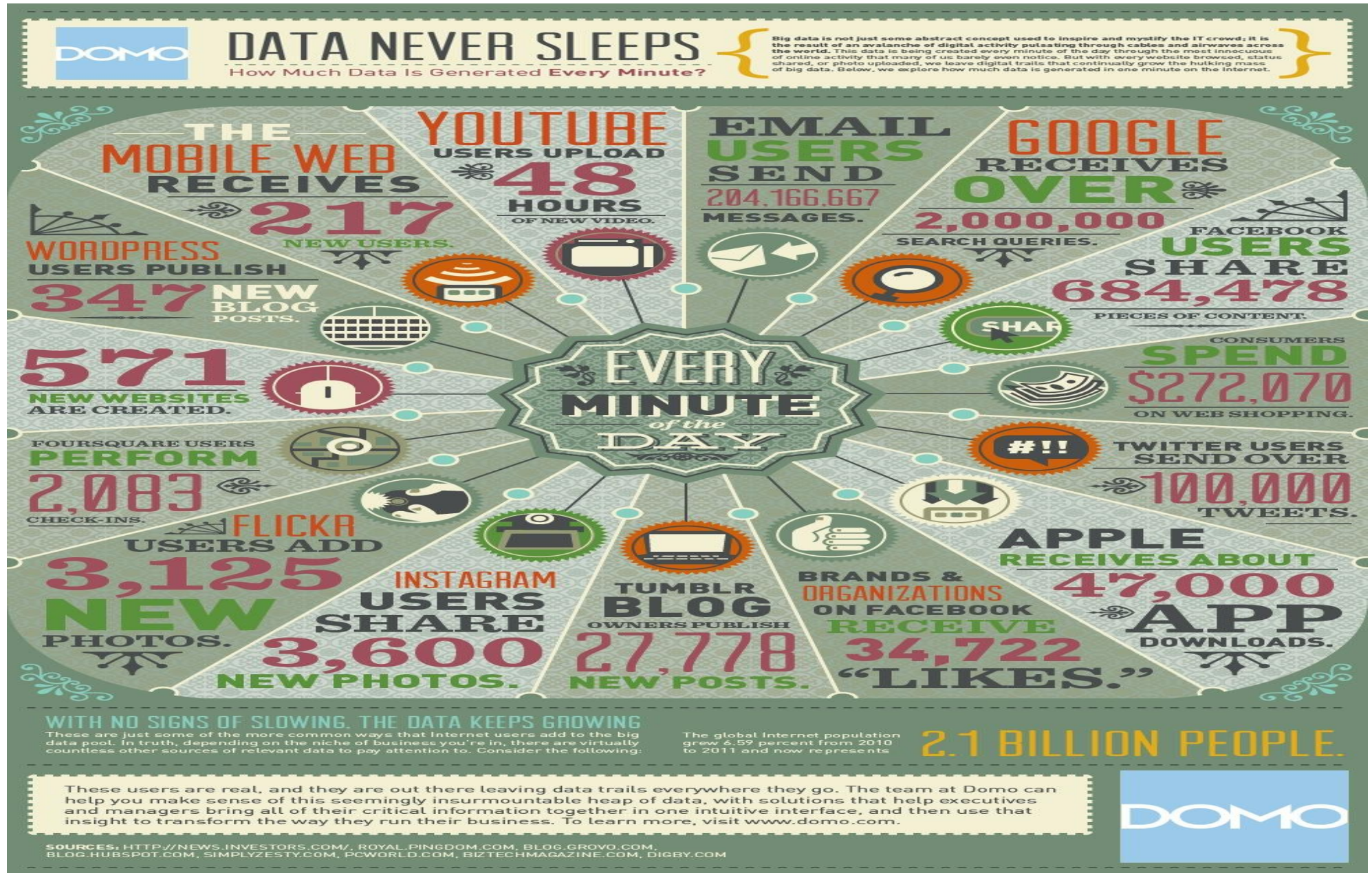    modern applications require huge amounts of data

    with many disks we can perform I/O in parallel

Assume that we are able to build a single disk with **500 TB** capacity. This is enough to store more than **20 billion webpages** (assuming an average size per page of **20KB**).

However, **just to scan** these 500 TB we need **more than 4 months** if the disk can bring **40 MB/sec**. Imagine the time required to process the data !

# What is Happening Today

# In the Near Future

"*IBM Research* *and Dutch astronomy agency* *Astron* *work on new technology to handle*

***one exabyte of raw data per day***

*that will be gathered by the world largest radio telescope, the* ***Square Kilometer Array***, *when activated in* ***2024***."

# Some Challenges

- Scalability
- Load balancing
- Fault Tolerance
- Efficiency
- Data Stream processing
- Support for complex objects
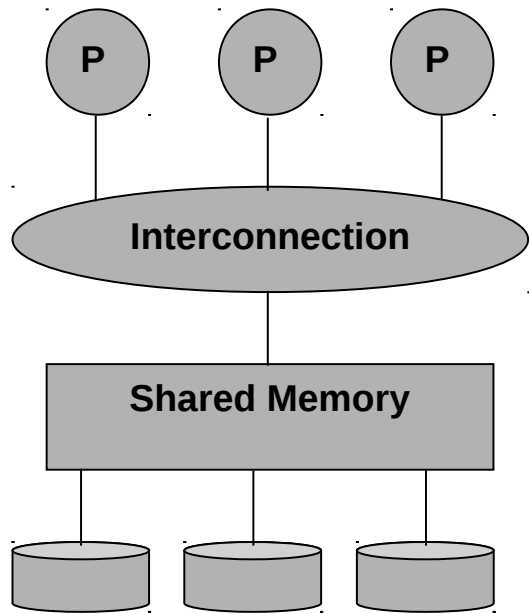- Accuracy/Speed tradeoffs (with performance guarantees)

# Parallel Architectures

**Shared Memory**: processors share a common main memory and also share secondary storage (e.g., disks)
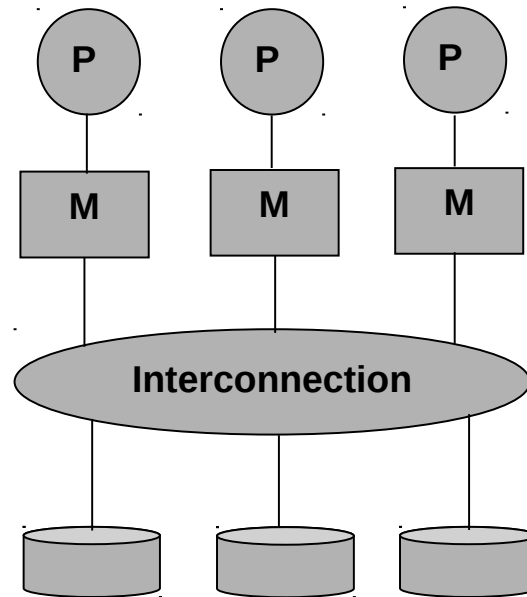
**Shared Disk**: processors share only secondary storage, whereas each processor has its own private memory

**Shared Nothing**: processors do not share anything, each one has private secondary storage and memory
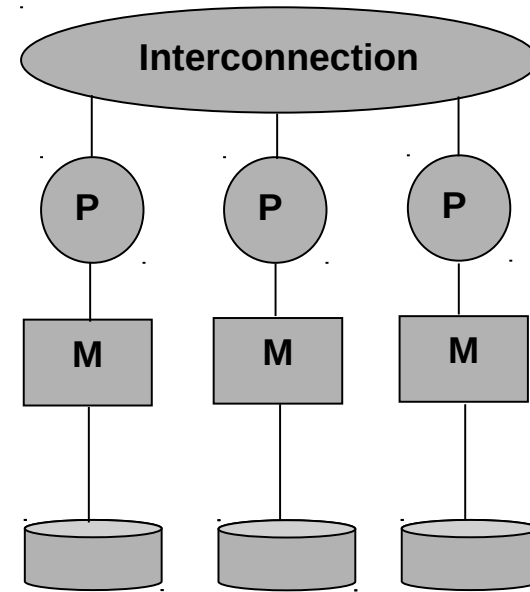
# Parallel Architectures



shared memory

shared disk

shared nothing

# Scalability

**Scale-Up**: put more resources into the system
to make it bigger and more powerful

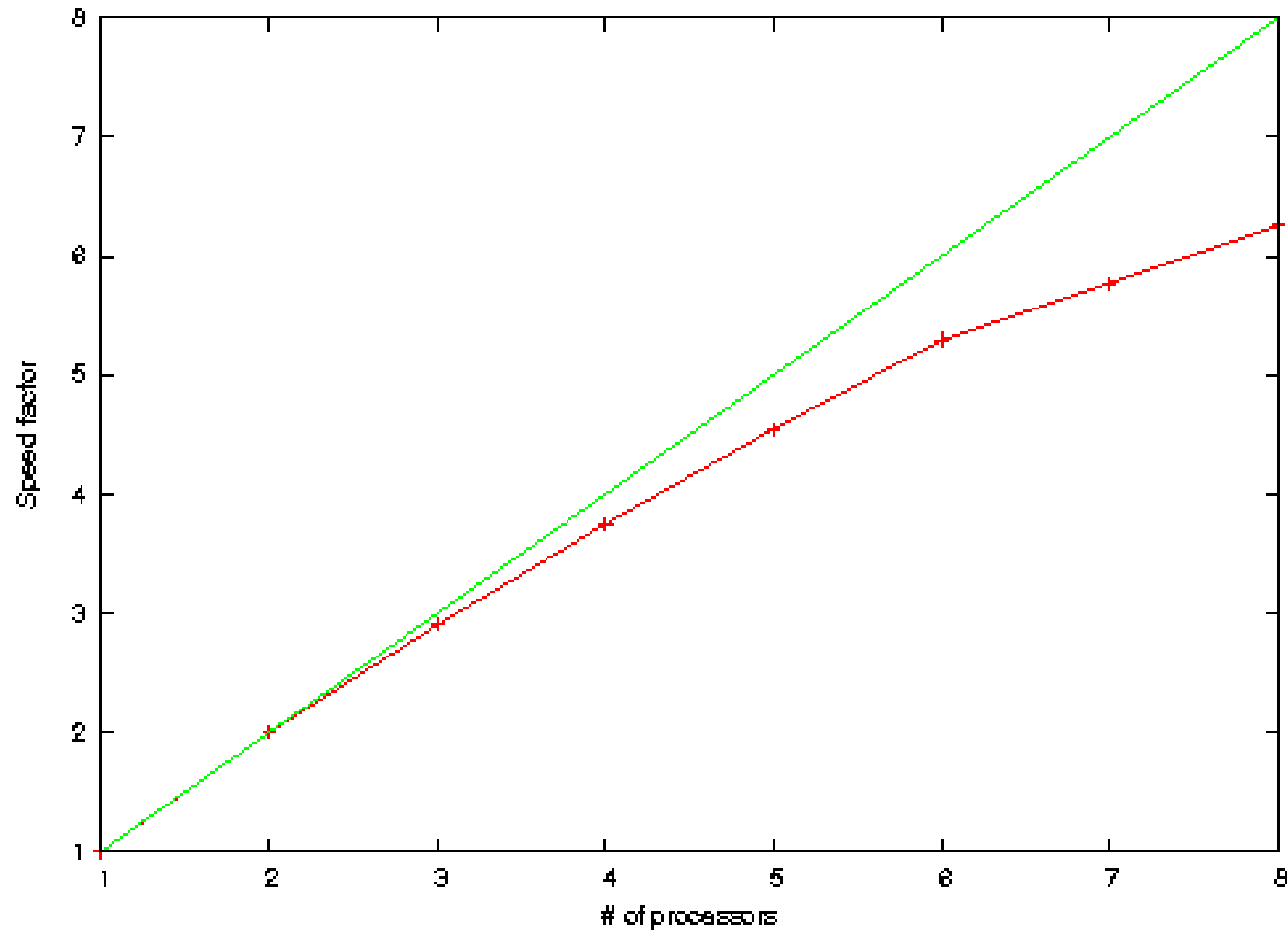**Scale-Out**: connect a large number of "ordinary" machines and create a cluster

Scale-Out is **more powerful** than Scale-Up, and also **less expensive**

# Scalability: measures

Among the three parallel architectures, shared-nothing is the one that **scales best.** This is the main reason for being adopted for building massively parallel systems (thousands of processors)

- **Speedup**: monitor performance by increasing the number of processors

- **Sizeup**: monitor performance by increasing only the dataset size

- **Scaleup**: monitor performance by increase both the number of processors and the dataset size

# The Speedup Curve
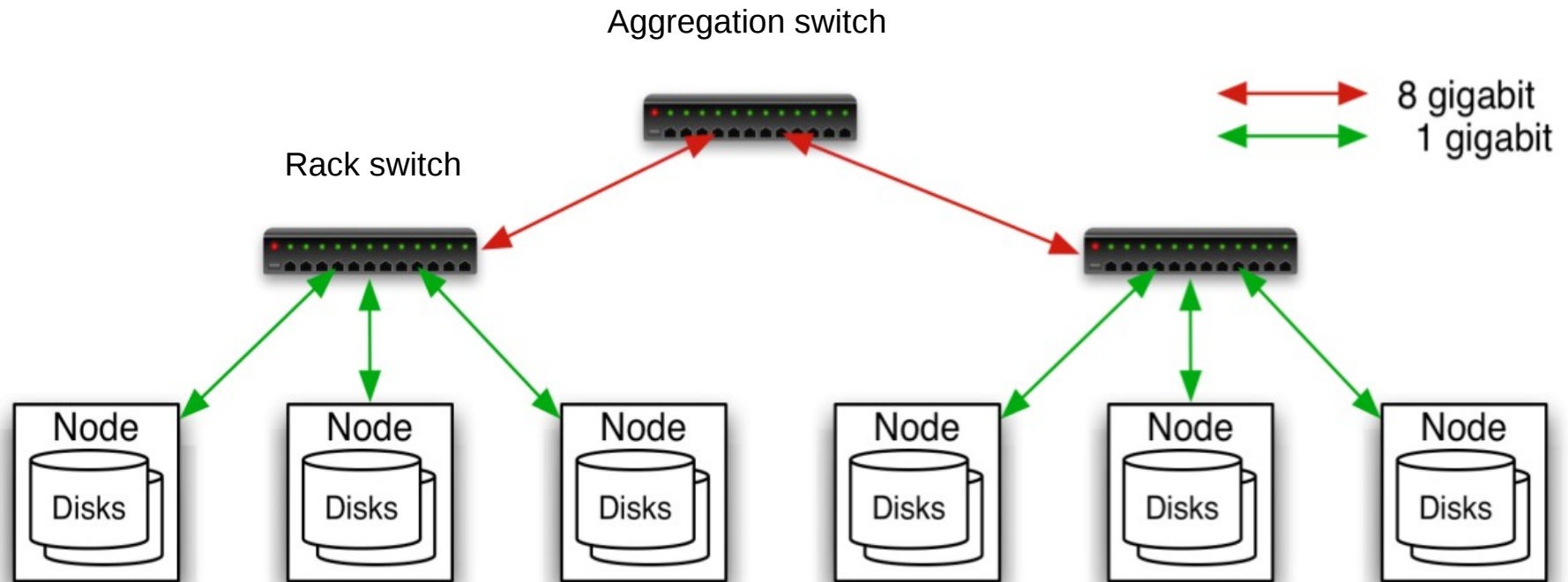
# Real Curves are Non-Linear

Why ?

**Start-up costs**: cost for starting an operation in a processor

**Interference**: cost for communication among processors and resource congestion

**Skew**: either in data or tasks → the slowest processor becomes the bottleneck

**Result formation**: partial results from each processor must be combined to provide the final result.

# Cluster Configuration Example

Aggregation switch

Rack switch

8 gigabit
1 gigabit

Node
Disks

Node
Disks

Node
Disks

Node
Disks

Node
Disks

Node
Disks

- 40 nodes/rack, 1000-4000 nodes in cluster
- 1 Gbps bandwidth within rack, 8 Gbps out of rack
- 8 x 2GHz cores, 8 GB RAM, 4 disks (= 4 TB?)

*Source: Matei Zaharia*

13

# Fault Tolerance

**Failures are very common in massively parallel systems**

Let $P$ the probability that a disk will fail in the next month. If we have $D$ disks in total, the probability that at least one disk will fail is given by:

Prob {*at least one disk failure*} = $1 - (1-P)^D$

e.g., $D = 10000$, $P = 0.0001$

Prob {*at least one disk failure*} = 0.63

# Fault Tolerance

Failures may happen because of:

**Hardware** not working properly

    Disk failure

    Memory failure (8% of DIMMs have problems)

    Inadequate cooling (CPU overheating)

**Resource** unavailability

    Due to overload

We must provide fault tolerance in the cluster!

# Fault Tolerance

Simplest protocol: if there is a failure, restart the job.

Assume a job that requires 1 week of processing. If there is a failure once per week, the job will never finish!

# Fault Tolerance

A better protocol:

Replicate the data and also split the job in parts and replicate them as well.

# Hadoop

A very successful model and platform to run jobs in massively parallel systems (thousands of processors and disks)
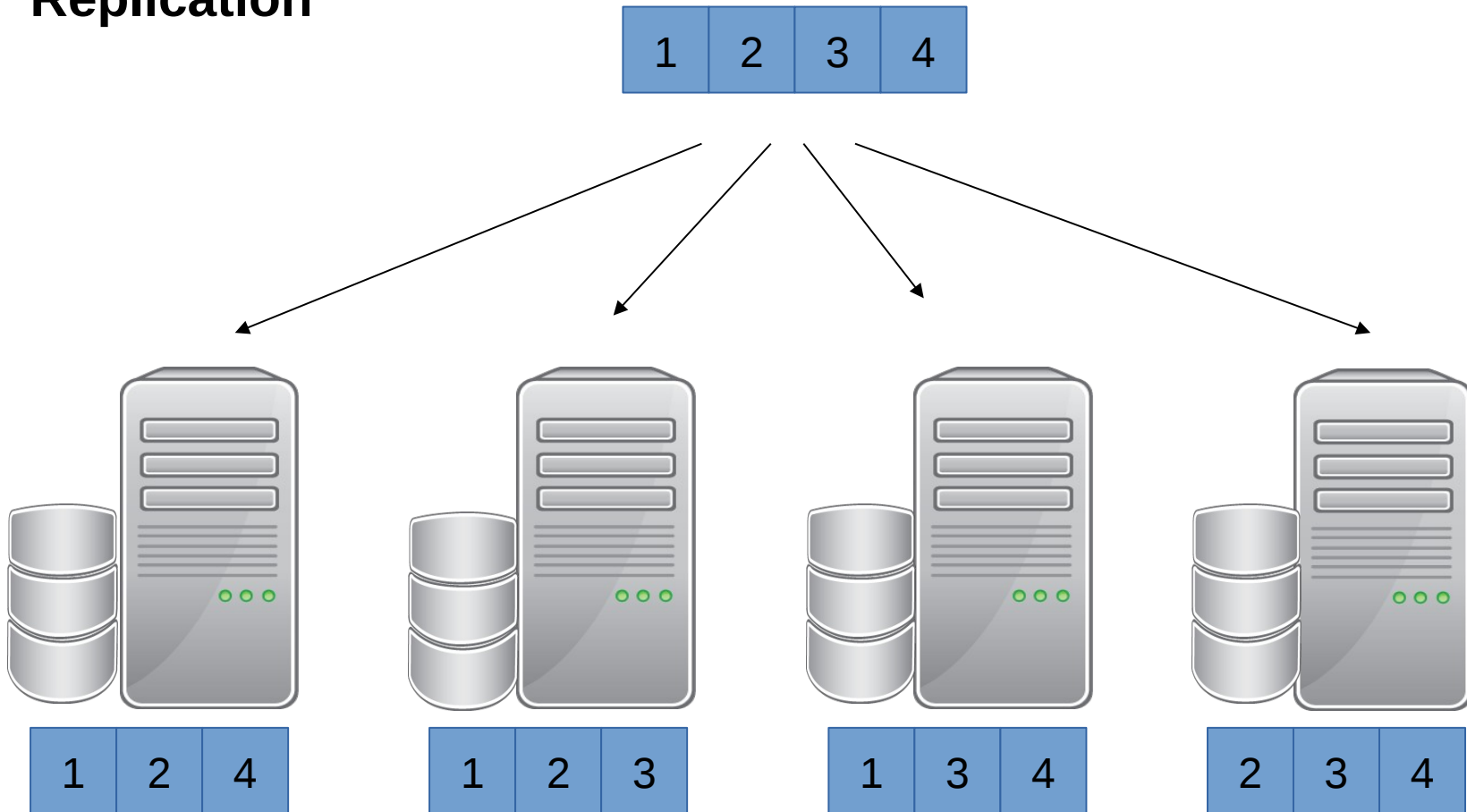
It contains two parts:

- the Hadoop MapReduce layer

- the Hadoop Distributed File System (HDFS)

Hadoop is the open-source alternative of MapReduce and Google File System (GFS) invented by Google. It has been used in Google's data centers mainly for: 1) constructing and maintaining the **Inverted Index** and 2) executing the **PageRank** algorithm.

# Hadoop

**Replication**



The the file is split in chunks. Each is replicated three times in this example.

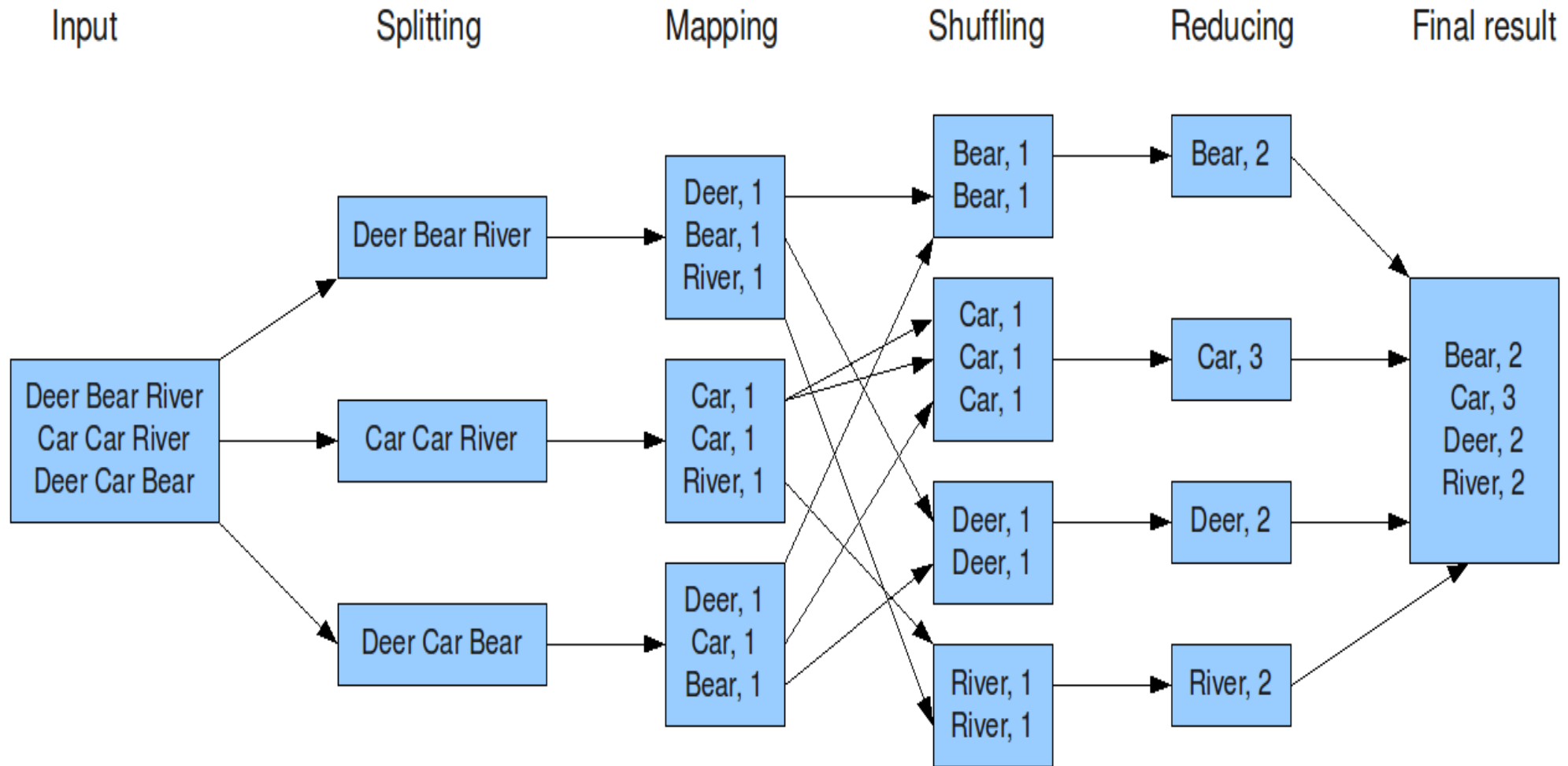# Processing in Hadoop

Based on key-value pairs

Each job is composed of one or more MR stages

Each MR stage comprises:

- the **map** phase
- the **shuffle-and-sort** phase
- the **reduce** phase

The programmer **focuses on the problem**. Replication, fault tolerance, scheduling, re-scheduling and other low level processes are handled by Hadoop.

# WordCount: the "hello world" of Hadoop

# Hadoop MR API

The programmer must implement the following functions:

**map()**: accepts a set of key-value pairs and generates another list of key-value pairs.

**combine()**: performs an aggregation before sending the data to reducers (reduces network traffic).
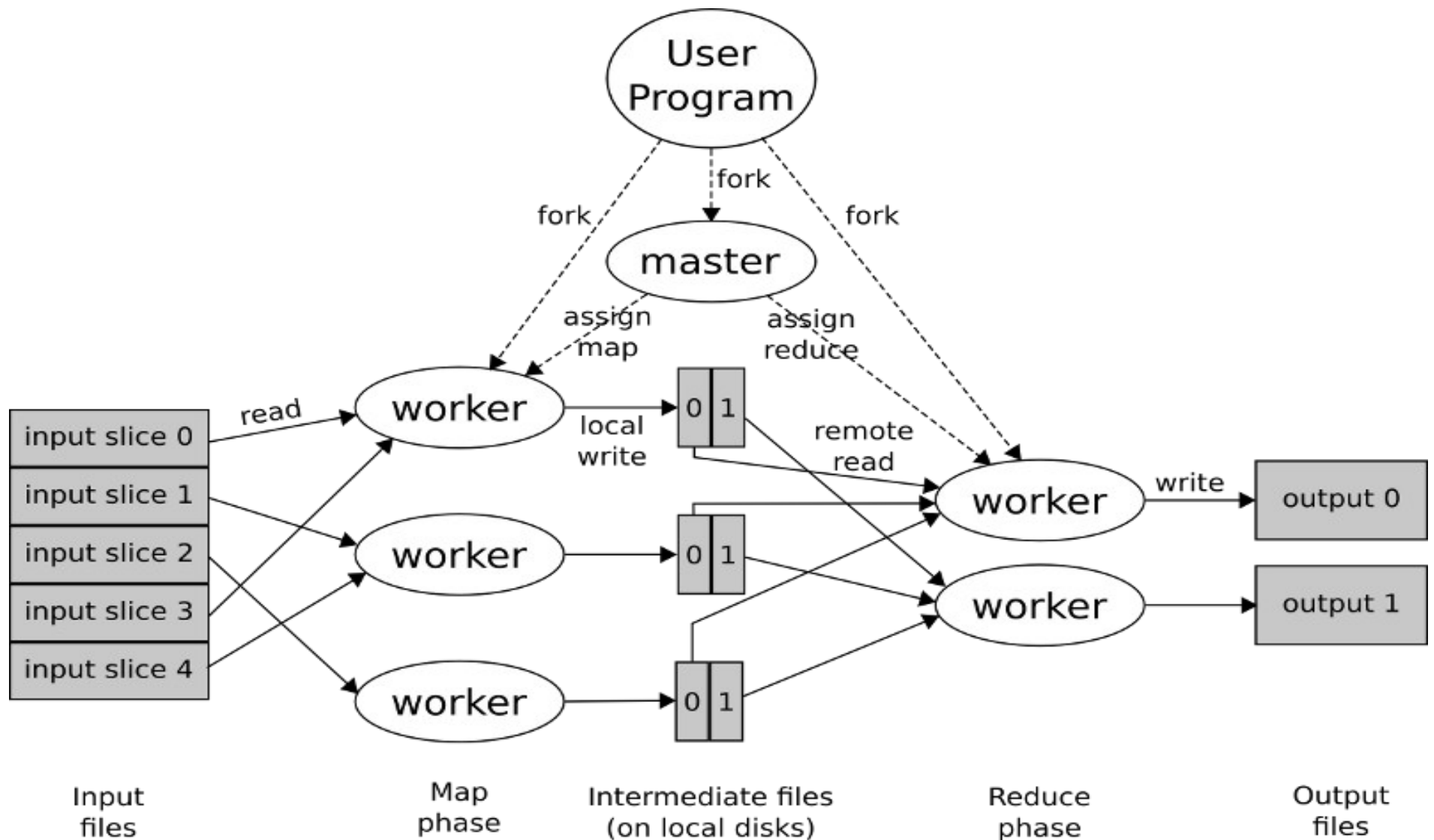
**partition()**: uses a hash function to distribute data to reducers (load balancing, avoids hotspots).

**reduce()**: accepts a key and a list of values for this specific key and performs an aggregation.

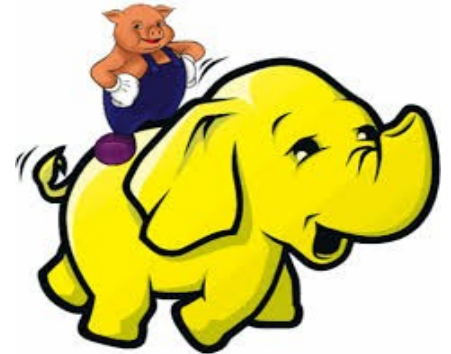*Note: combine() and partition() are optional*

# Workflow in Hadoop



23

# Hadoop Ecosystem

# Theoretical Issues in MR

Paper titles:

- ➢ "On the Computational Complexity of MapReduce"

- ➢ "A new Computation Model for Cluster Computing"

- ➢ "Fast Greedy Algorithms in MapReduce and Streaming"

- ➢ "Minimal MapReduce Algorithms"

- ➢ "Filtering: A Method for Solving Graph Problems in MapReduce"

- ➢ "A Model of Computation for MapReduce" (**SODA 2010**)

# MR Limitations

Difficult to design efficient/optimal algorithms (everything must be expressed in key-value pairs)

A lot of disk I/Os (mappers reading HDFS and writing local data)

A lot of network traffic (shuffling is expensive)

Difficult to handle data skew (the curse of the last reducer!)

Not very good for iterative processing (requires many MR stages)

Not very good for **streaming applications**

# Is There an Alternative ?