# MapReduce

# Agenda

- Why distributed data processing?
- Simple data parallelization
- The MapReduce system
- How MapReduce works
- Google File System for distributed storage
- MapReduce examples
- Fault-tolerance
- Map reduce performance Issues
- Limitations

# Processing Data

- Data on file system on disk

- Simple processing flow:
  - 1. Read into memory
  - 2. Process data (apply some function)
  - 3. Write back to disk
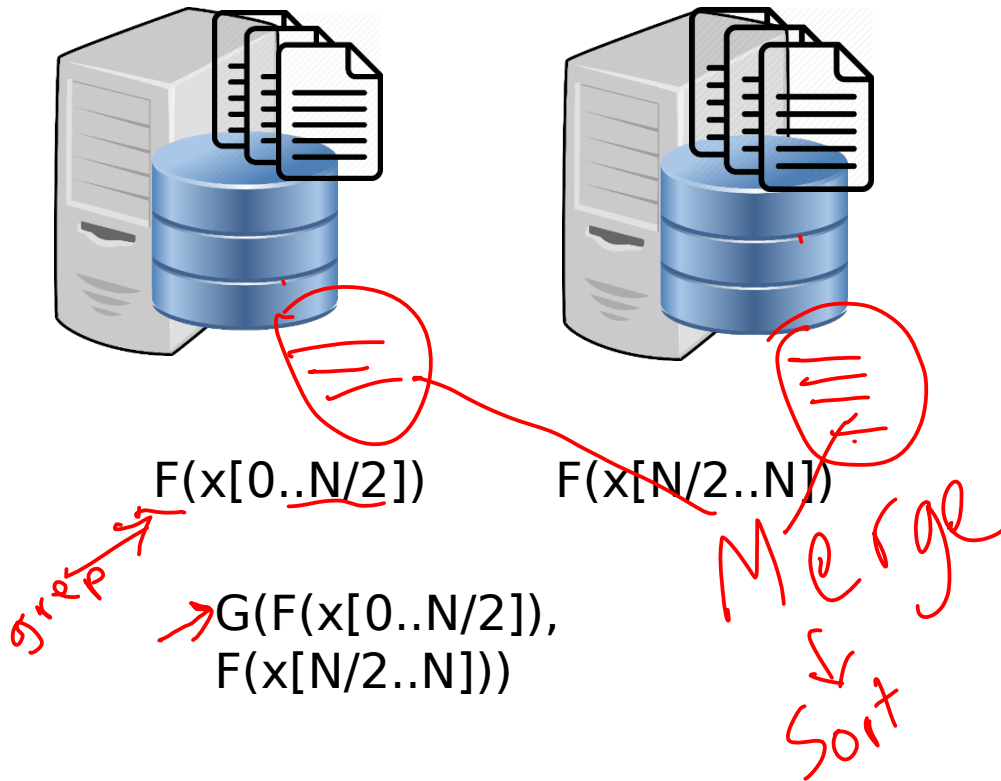
# Why Multiple Servers?

- Data size larger than disk capacity
- Sequentially processing data can be slow
  - Limited by disk read/write speeds
  - Parallel processing can significantly reduce time
- Single point of failure with a single server

# Multiple Servers: Divide and Conquer

*grep*

- Divide input data among multiple servers
- Each server processes a partition in parallel
- ???
- Profit!?

- Data processed by individual servers must be "aggregated" or collected
- May require significant communication

F(x[0..N/2])          F(x[N/2..N])

G(F(x[0..N/2]),
  F(x[N/2..N]))

*grep*

*Merge & Sort*

# MapReduce: Simplified Data Processing on Large Clusters

Jeffrey Dean and Sanjay Ghemawat

jeff@google.com, sanjay@google.com

*Google, Inc.*

## Abstract

MapReduce is a programming model and an associated implementation for processing and generating large data sets. Users specify a *map* function that processes a key/value pair to generate a set of intermediate key/value

given day, etc. Most such computations are conceptually straightforward. However, the input data is usually large and the computations have to be distributed across
h
a
al

Published at OSDI 2004.
>25,000 citations

# Why MapReduce?

- Automated parallelization and distribution of data and processing
- Clean, powerful, well-understood abstraction (Map and Reduce)
- Fault-tolerance
- Scalability: 1000s of servers, TBs of data, …
- Apache Hadoop: widely used open source Java implementation

# Some MapReduce Problems

- Word-count
- Creating an inverted index
  - Which documents does a word occur in?
  - Useful if you are building a search engine
- Log-processing
  - Filter log messages which match some condition
  - IP address == 156.54.61.*
- Matrix processing
  - PageRank, matrix multiplication, …

# (List (Processing (Like) It's) 1959)

- LISP: Everything is a List
- '(1 2 3 4 5 6)
- Map: Apply a function element-wise
- (map square '(1 2 3 4 5))
  - -> '(1 4 9 16 25)

- Reduce: Aggregate values in a list
  - Also called fold
- (reduce sum '(1 2 3 4 5))
  - -> 15
- Can pass any associative function

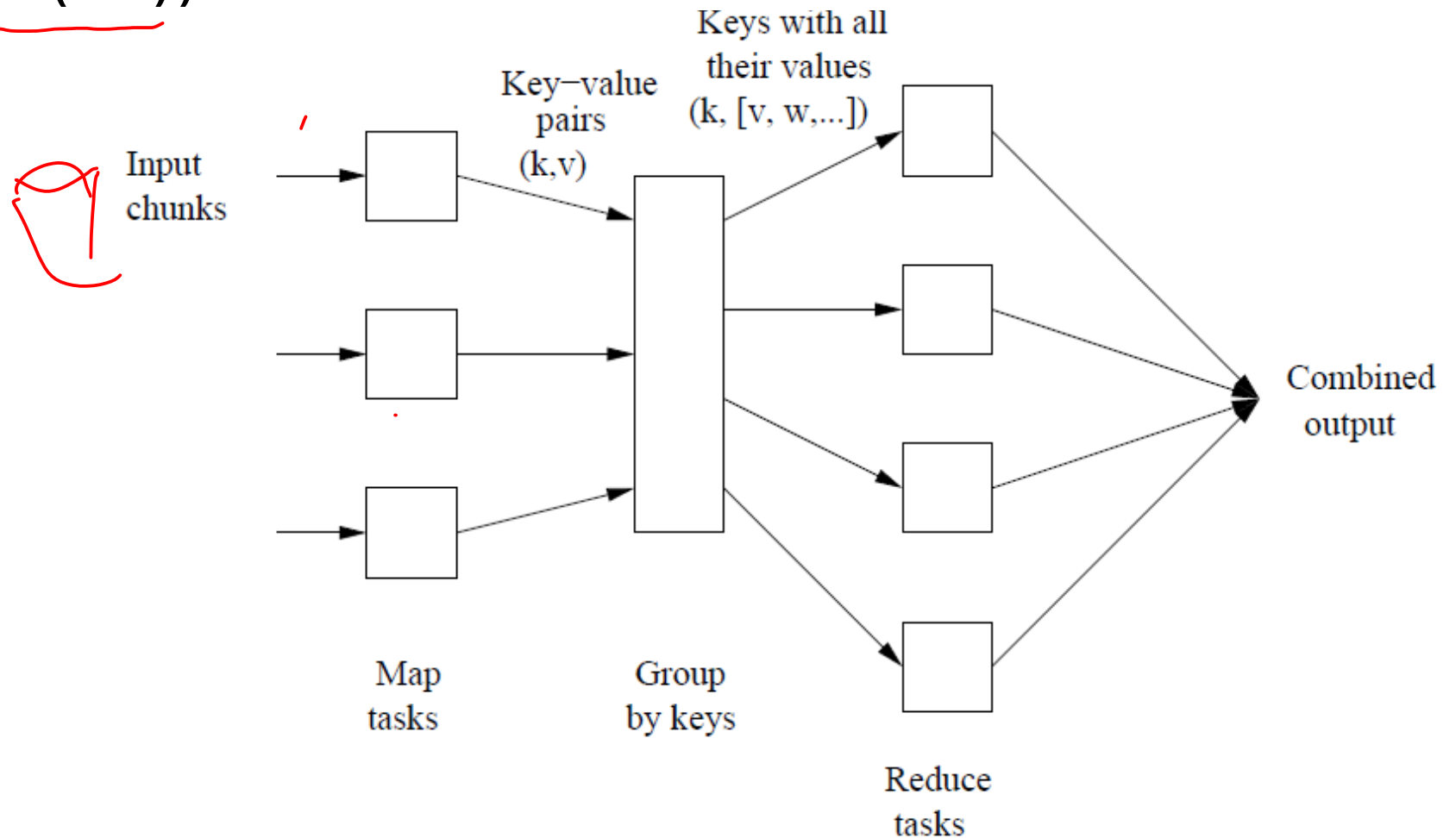(reduce sum (map square '(1 2 3 4 5)))

-> Sum of squares

*(handwritten annotations)*
$f(x, y) \to z$
$+1 (2\ 3\ 4\ 5)$
$+1 (+2 \cdot (3\ 4\ 5))$
$+\ x\ y\ z$
$\equiv (+\ x\ y) + z$
$\equiv x + (y + z)$

# Map Reduce Semantics

- Map: (k1, v1) -> list(k2, v2)
- Reduce: (k2, list(v2)) -> list(v3)

# MapReduce System Architecture
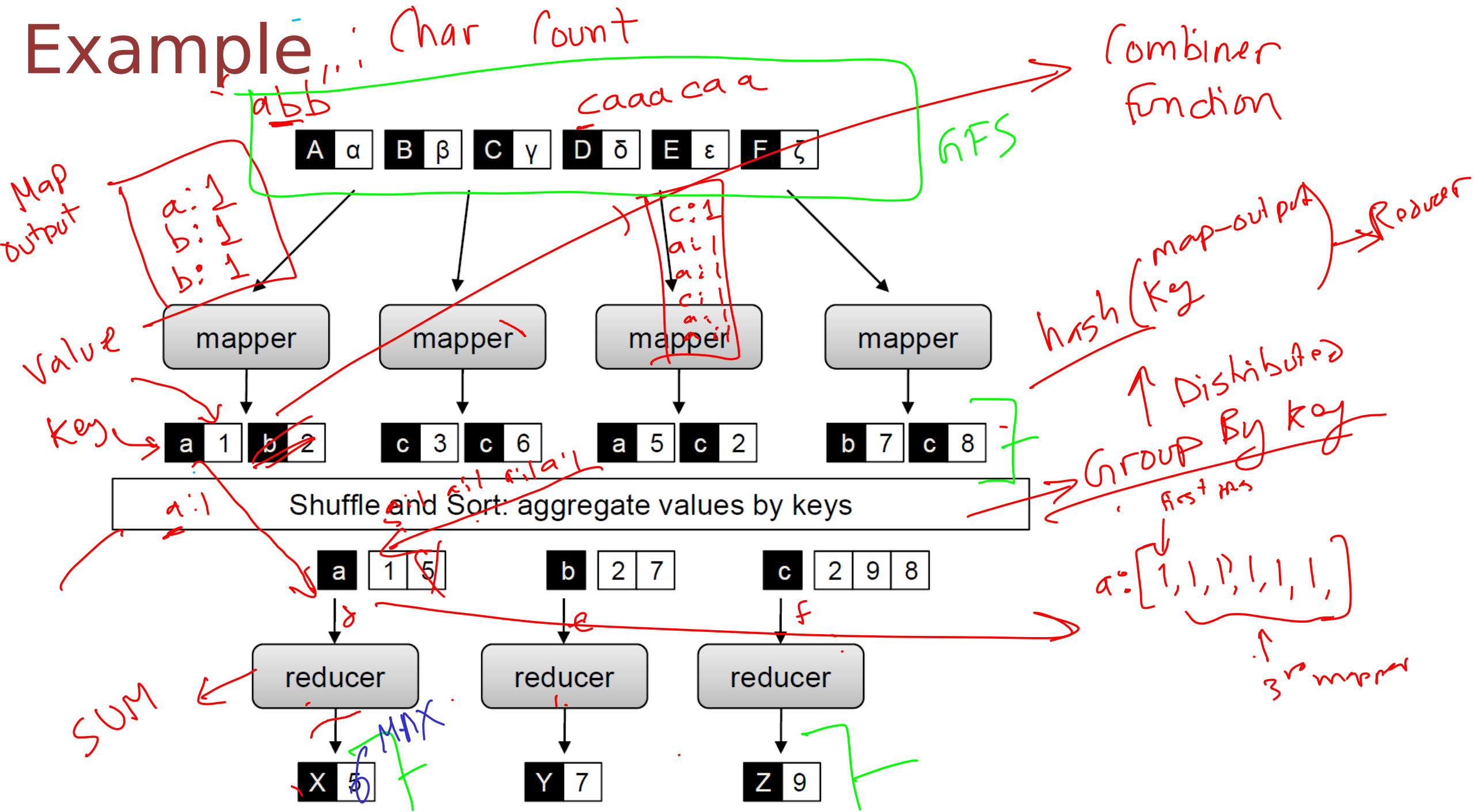
# Computation Flow

- Map tasks turn chunks into series of key-value pairs based on user-provided fn

- Key-value pairs from all map tasks are collected and sorted/grouped by key

- Each reducer task gets a subset of keys
  - Master uses a hash function to assign a key to one of R reducers
  - All key-val pairs with the same key are processed by the same reducer task

- Reducer tasks process one key at a time, and combine all values associated with that key based on reduce function

# Example



Char count

Combiner function

abb

caaa ca a

| A α | B β | C γ | D δ | E ε | F ζ |

GFS

Map Output
a: 1
b: 1
b: 1

Value

Key

c: 1
a: 1
a: 1
c: 1
a: 1

hash(key

hash(map-output) → Reducer

Distributed
Group By key

mapper | mapper | mapper | mapper

| a 1 | b 2 | | c 3 | c 6 | | a 5 | c 2 | | b 7 | c 8 |

a: 1

a: 1 c: 1 a: 1 a: 1

Shuffle and Sort: aggregate values by keys

first map

| a | 1 | 5 | | b | 2 | 7 | | c | 2 | 9 | 8 |

a: [1, 1, 1, 1, 1, 1,

γ | e | f

SUM | reducer | reducer | reducer

3rd mapper

MAX

| X 5 | | Y 7 | | Z 9 |

# Distributed File System

- MapReduce distributes **computation**
- Distributing and managing **data** is as important, if not more
- Storing data centrally:
  - I/O bottleneck
  - Processing data in parallel not very useful if limited by single disk

# The Google File System

Sanjay Ghemawat, Howard Gobioff, and Shun-Tak Leung

Google*

## ABSTRACT

We have designed and implemented the Google File System, a scalable distributed file system for large distributed data-intensive applications. It provides fault tolerance while running on inexpensive commodity hardware, and it delivers high aggregate performance to a large number of clients.

While sharing many of the same goals as previous distributed file systems, our design has been driven by observations of our application workloads and technological environment, both current and anticipated, that reflect a marked departure from some earlier file system assumptions. This
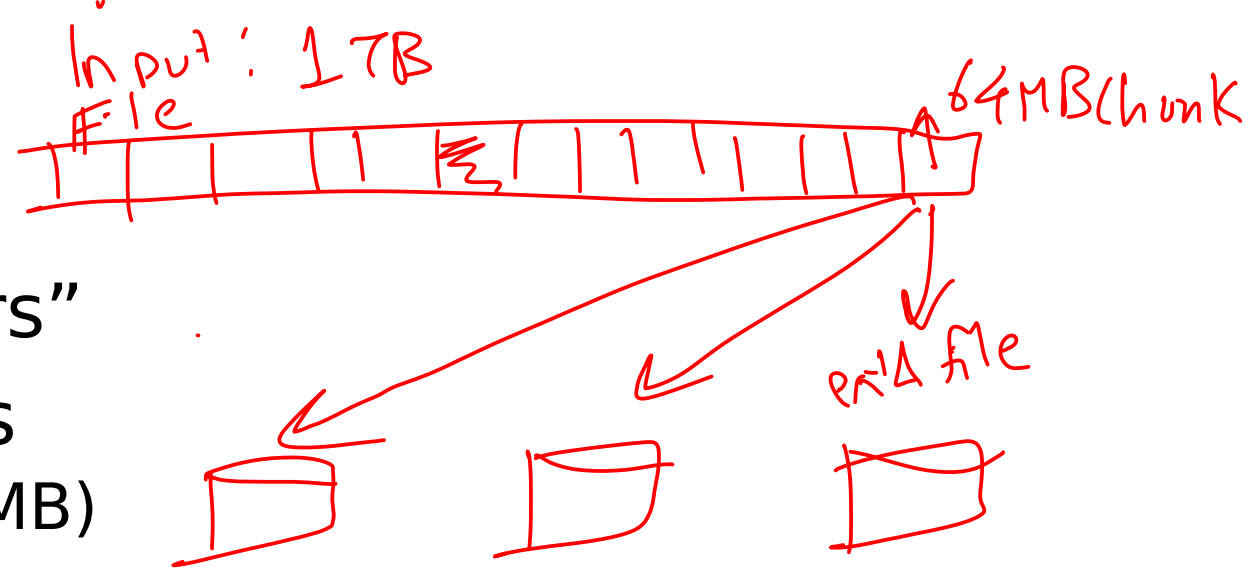
## 1. INTRODUCTION

We have designed and implemented the Google File System (GFS) to meet the rapidly growing demands of Google's data processing needs. GFS shares many of the same goals as previous distributed file systems such as performance, scalability, reliability, and availability. However, its design has been driven by key observations of our application workloads and technological environment, both current and anticipated, that reflect a marked departure from some earlier file system design assumptions. We have reexamined traditional choices and explored radically different points in the

# Google File System

- Distributed FS for large datasets
- Supports read, write, open, close, record-append operations
- Can run on off-the-shelf computing clusters
  - No specialized hardware required
- Provides fault tolerance via replicating each piece of data three times
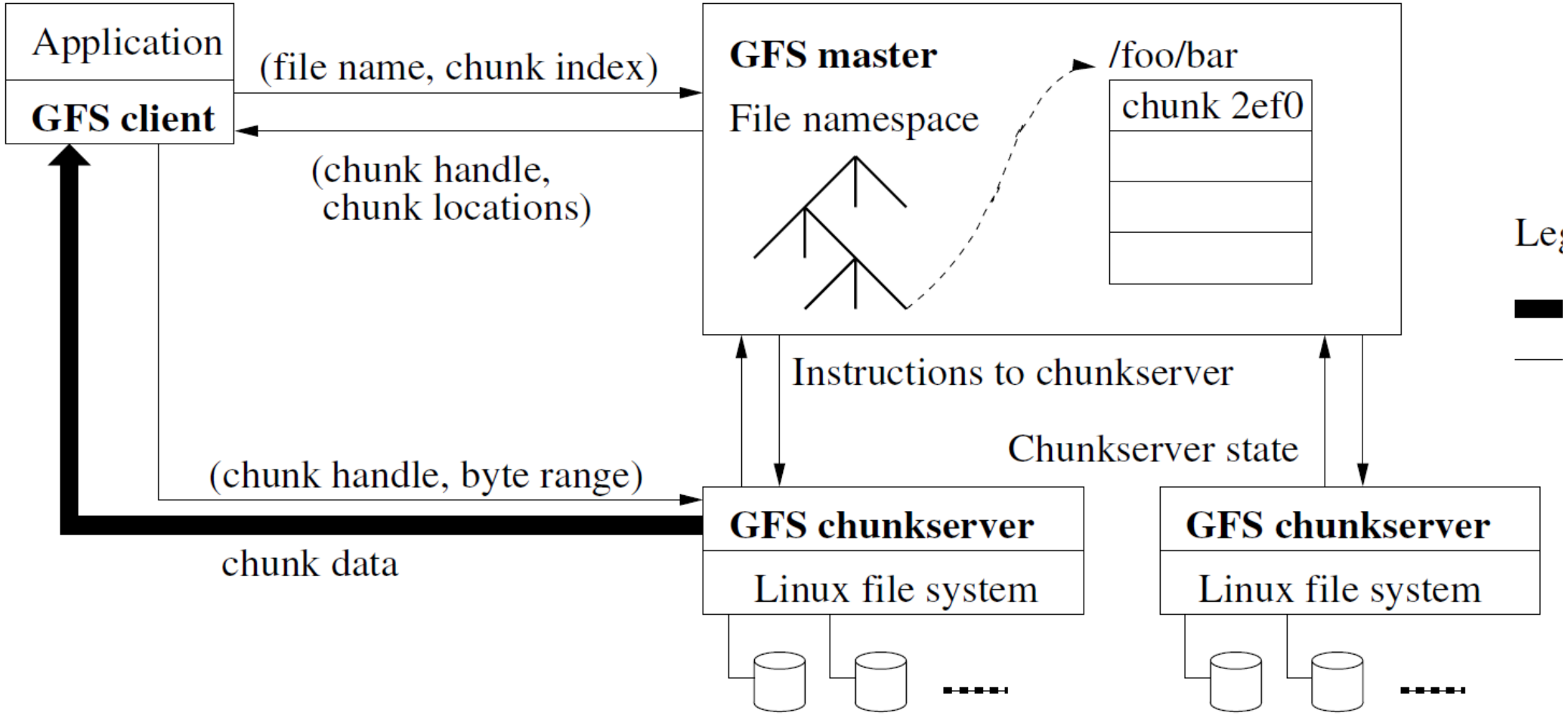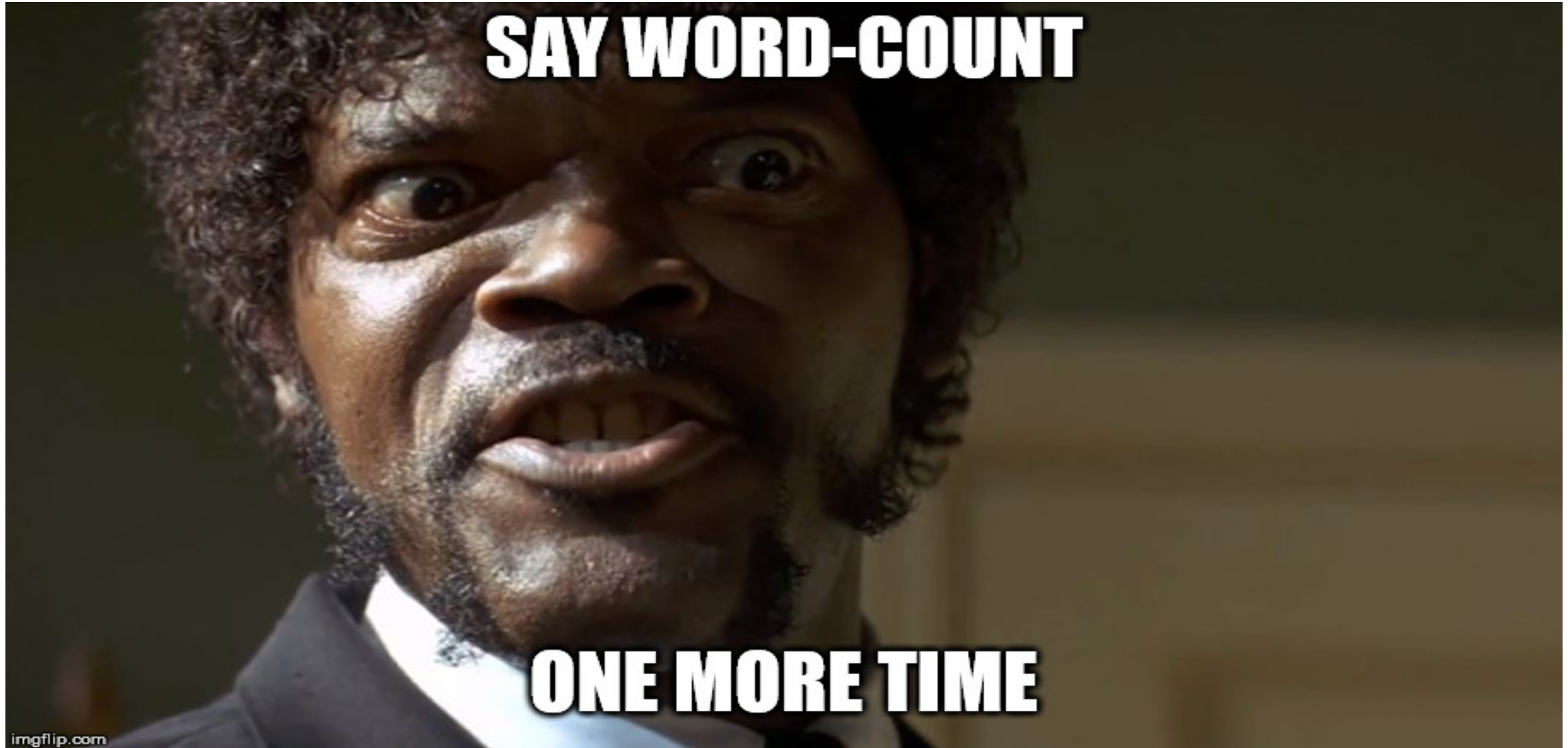  - Replicas can also be used for load balancing

# Overlay File Systems

- Data is stored in "chunk servers"
- HDFS files: sequence of chunks
  - Chunks are just large blocks (64MB)
- Chunk-servers store chunks as **files in a local file-system (such as ext4)**
- HDFS files are indexed by a central "namenode"
- (HDFS-file-name, chunk-number) -> Chunk-servers, chunk-handle

# GFS Architecture

# MapReduce Examples

# Word-count

- Mapper: Emit (word,1) pairs
- Grouping by key :
  - (the,1), (the,1),....      : Sent to Reducer 0
  - (apple, 1), (apple,1),....    : Sent to Reducer 1
  - (car, 1), (car,1).....        : Sent to Reducer 2
  - (dog, 1), (dog,1),....     : Sent to Reducer 0
  - ....
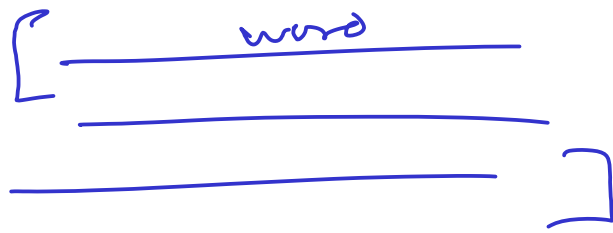- Reducer function is simple addition, and each reducer outputs:
  - (the, 102)
  - (apple, 4)
  - ...

**Handwritten annotations:**

Input: [list of words]

Input$_2$ : [Words are lame]

[(list,1), (of,1), (words,1)]   [(words,1), (are,1)]
(lame,1)

hash(list)/3

list = 1
words = 2

# Grep

- Mapper: Emit (word, line-number) pairs
- Grouping by key :
  - (word1, [line-1, line-4, line-3, line-2, line-100])
  - (word2, [line-3, line-4, line-6, line-7])
- Map output is usually sorted by key:
  - (word1, [1,2,3,4,100])
  - (word2, [3,4,6,7]
- Reducer:
  - Do nothing / Identity / Null

*word*
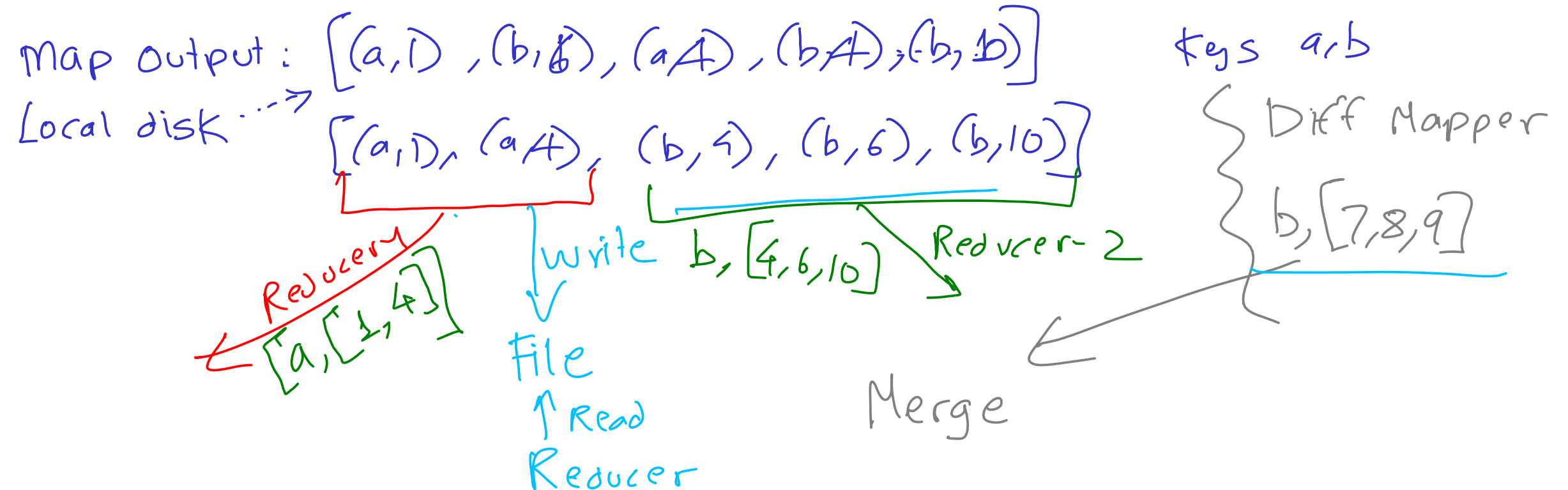
Intermediate Output

Reducer → Merge

Distributed Sort

# Combiners

- Sometimes, Reduce function is associative and commutative
  - f(f(a, b), c) =f(a, f(b, c))
  - f(a, b) = f(b, a)
  - Example: Addition
- In such cases, can combine map-output before sending to reducer
- In case of word-count example:
  - Map-output: [(word-1,1), (word-1, 1),…,, (word-2,1)… ]
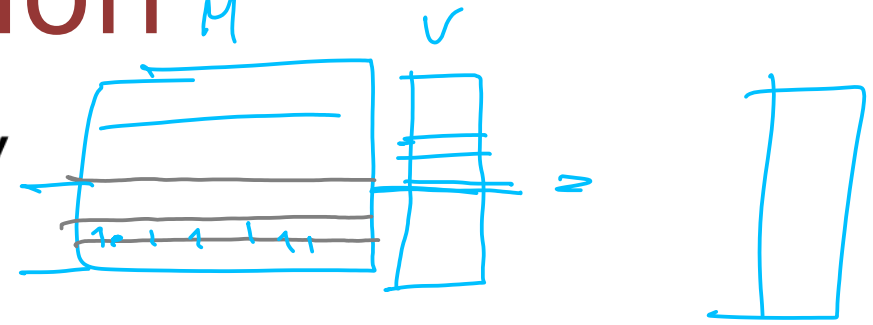  - With combiner: [(word-1, 20), (word-2, 12)]

User specified

# Map Output (Intermediate Data)

- Each mapper groups its output by key
    - Usually done by sorting entire local map output
- Creates an intermediate file for each reducer
- If M mappers and R reducers: up to M*R total intermediate files

Map Output: $[(a,1), (b,6), (a,4), (b,4), (b,10)]$

Local disk ⋯→ $[(a,1), (a,4), (b,4), (b,6), (b,10)]$

Reducer-1

$[a, [1,4]]$

Write

File

↑ Read

Reducer

$b, [4,6,10]$     Reducer-2

Merge

Keys a, b

Diff Mapper

$b, [7,8,9]$

# Matrix-Vector Multiplication

- Matrix, M and vector, v. Want to multiply : Mv
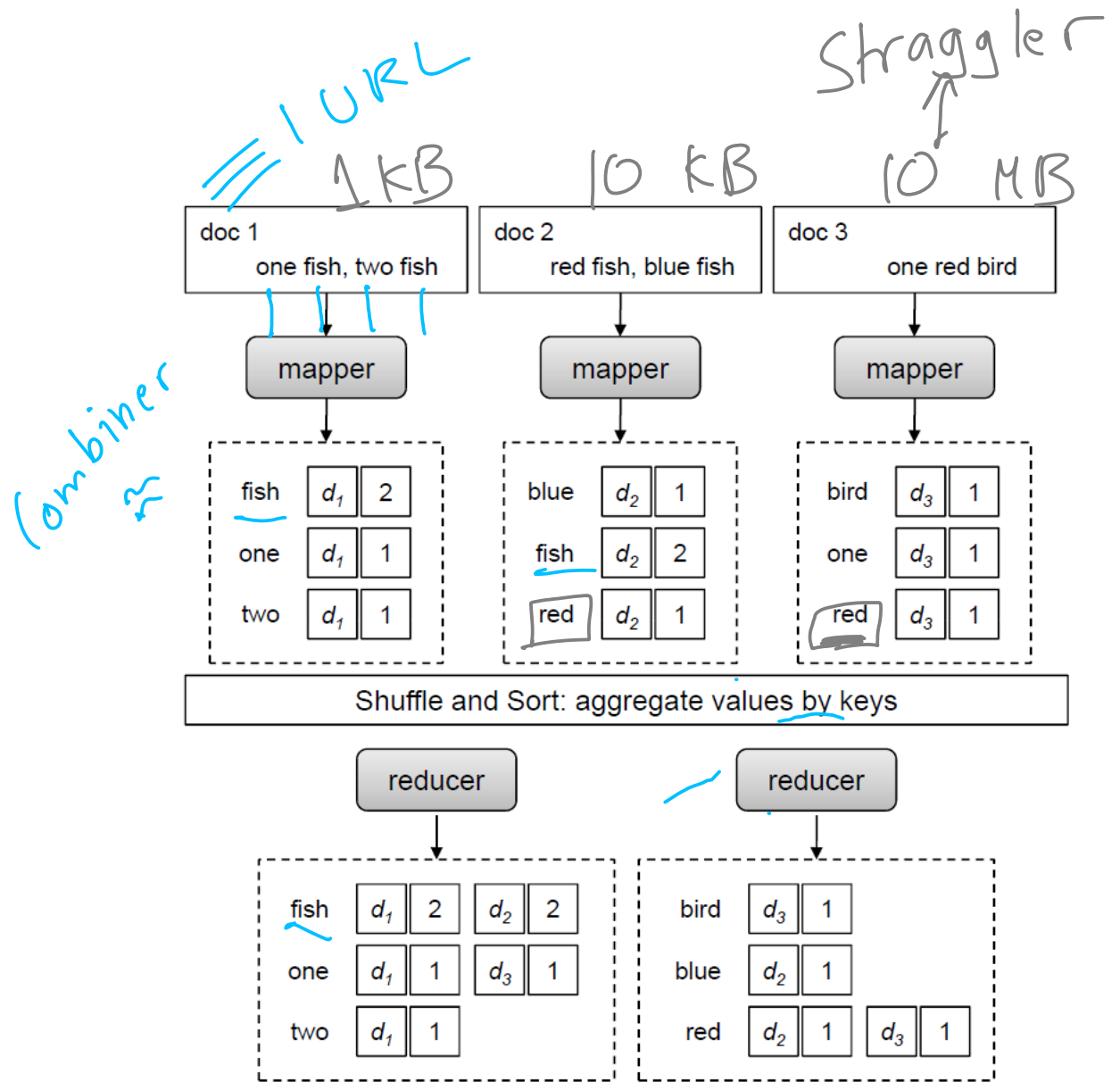  - V=[v1,…vn]. M is nXn matrix
- Matrix partitioned by row
- Map: outputs $(i, m_{ij} * v_j)$  $\forall j$
- Reduce: sum up all values for each key to output $(i, x_j)$

- QS: What if we can't fit v in memory?
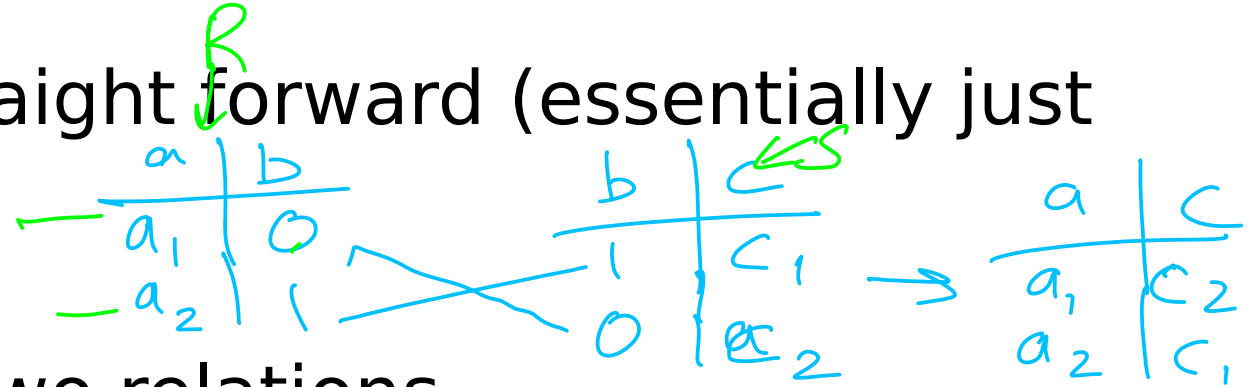- QS: Co-occurrence matrix of words in a corpus

# Inverted Index

- How many times does a word occur in each document?

# Relational Operations

- MapReduce can also be used to run relational operations
  - Select, Project, Joins, …
- Selection and projection are straight forward (essentially just filtering using map)
- Natural joins with MapReduce
- Assume: R(a,b) and S(b,c) are two relations
- Map: Emit (b, (R, a)) Where R is just the relation name and a "tag"
- Reduce: For each key (b), output all tuples in the values list (a,c)
- Apache Hive translates SQL queries into a MapReduce program

# Fault-Tolerance

- Master fails
  - Switch to secondary master
  - Restart entire job
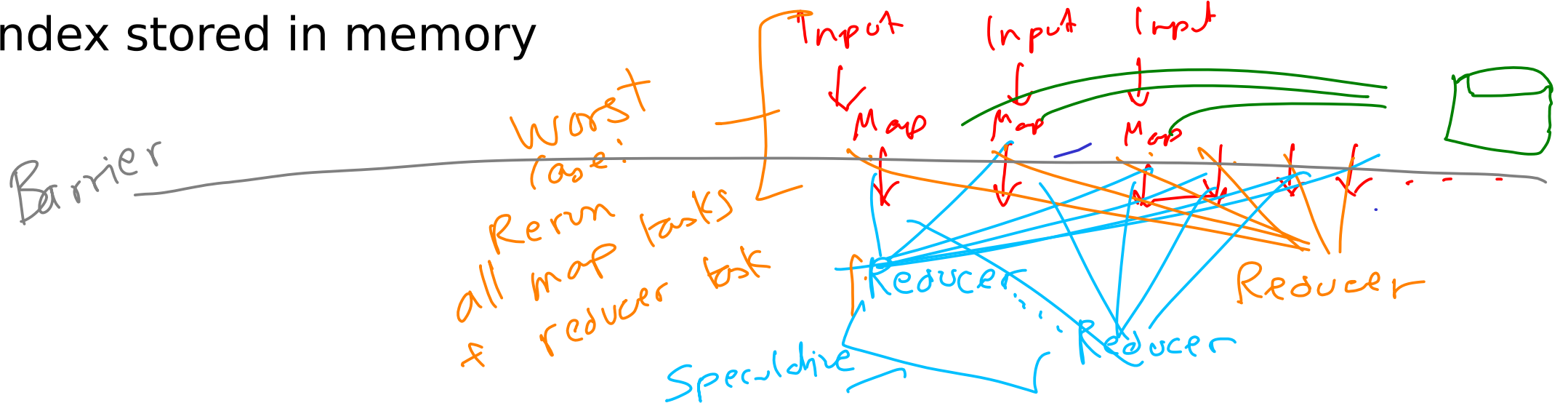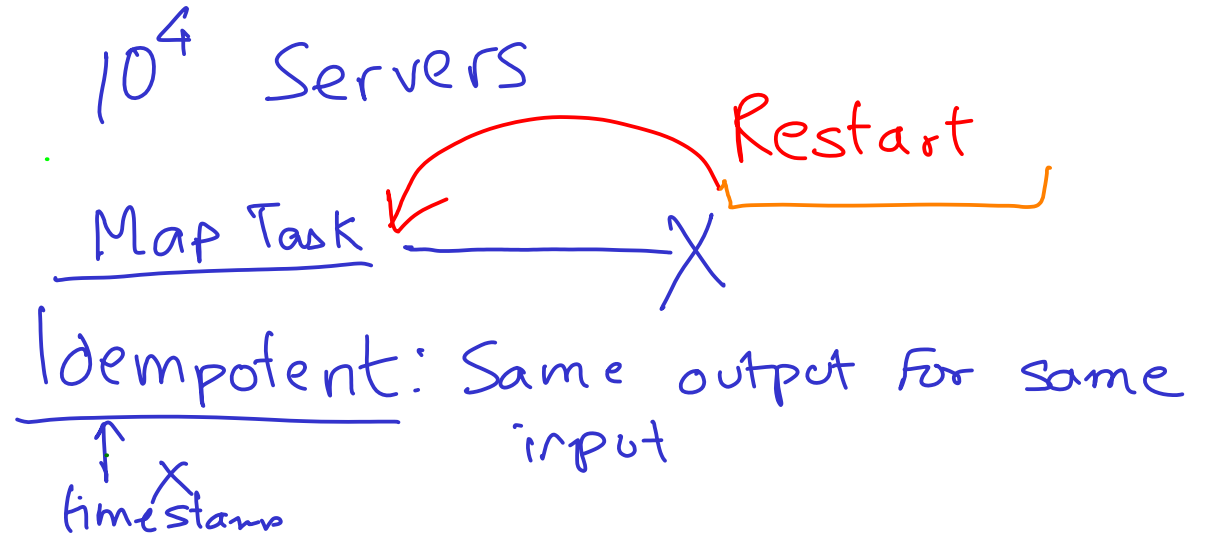
- Worker fails
  - If running map: restart map tasks on available/free worker nodes
  - If running reduce: restart reducer tasks

- File-system (HDFS/GFS) has a single master node
  - FS index stored in memory

$10^4$ Servers

Primary
Secondary
Replication

Restart

Map Task    X

Idempotent: Same output for same input

↑
X
timestamp

if no replicated master

Worst case:
Rerun all map tasks + reducer task

Barrier

Input   Input   Input

Map    Map    Map

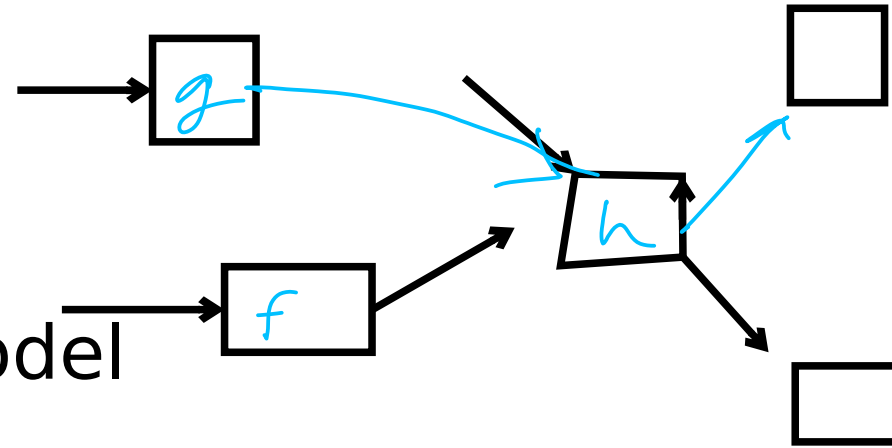Reducer       Reducer

Reducer

Speculative

# Performance Issues

- Map output <u>barrier</u> : Even partitioning of mappers workload required
  - Usually achieved by evenly splitting the input
  - Assumes that element-wise map function has uniform cost
- Speculative execution (backup tasks)
  - Run the **same** task on multiple workers
  - If some workers are slow (stragglers) → Hardware Slowdowns
    - Shared H/w
- Reducer skew: A major problem
- Recall that all values of a specific key must be handled by same reducer

Perf dictated by
time(slowest Mapper) + Grp By+
time (slowest Reducer)

- What if there's a really popular key?

"the" — $10^6$

# MapReduce Limitations

Idempotent

- Static data
  - Append-only usually OK
- Restrictive programming model
  - No support for dataflows

- How many mappers and reducers? → CPU, Mem
- How much memory to allocate?
- Purely batch processing
  - Jobs can take hours to complete
  - No streaming, interactive analytics

Stream of Input data

# When NOT To Use MapReduce

- Modern computing hardware has plenty of computing power:
  - A typical laptop: 8 cores, 16 GB RAM, 512 GB SSD (usually NVMe)
  - A typical server: 64 cores, 256 GB RAM, 2 TB SSD, ...
- Is your dataset really that large?
  - Wikipedia: 30 GB
- In many cases, an optimized non-distributed implementation beats a large cluster (!)
  - Frank McSherry's blog

+ Laptop