# BigData and Map Reduce

**VITMAC03**

# Motivation

‣ Process lots of data
  - Google processed about 24 petabytes of data per day in 2009.
‣ **A single machine** cannot serve all the data
  - You need a distributed system to store and process **in parallel**
‣ Parallel programming?
  - **Threading** is hard!
  - How do you facilitate communication between nodes?
  - How do you **scale to more machines**?
  - How do you handle machine failures?

2

# MapReduce

‣ MapReduce [OSDI'04] provides
  • Automatic parallelization, distribution
  • I/O scheduling
    • Load balancing
    • Network and data transfer optimization
  • Fault tolerance
    • Handling of machine failures

‣ **Need more power: Scale out, not up!**
  • Large number of **commodity servers** as opposed to some high end specialized servers
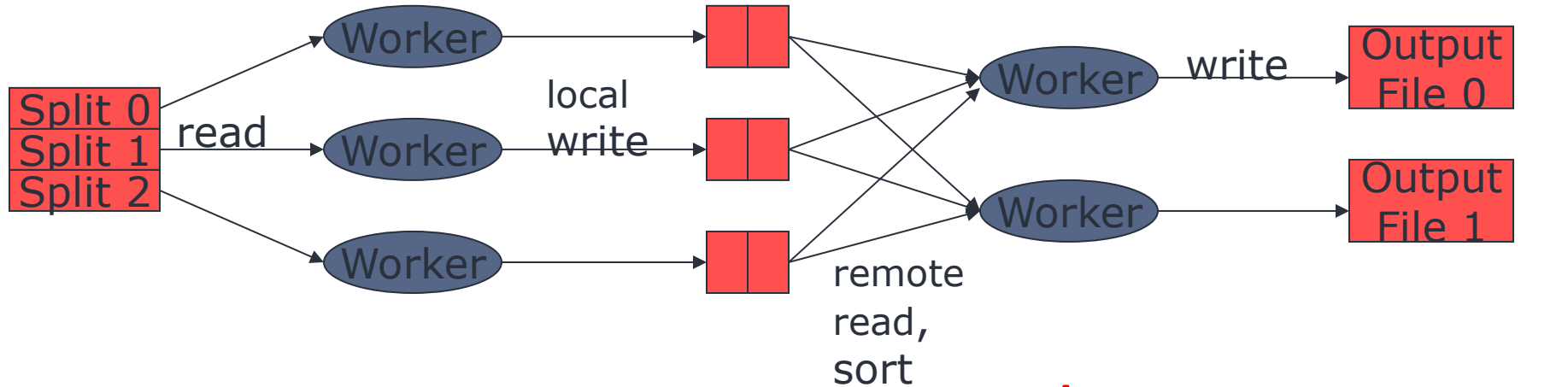
**Apache Hadoop:**
Open source implementation of MapReduce

3

# Typical problem solved by MapReduce

‣ Read a lot of data

‣ Map: extract something you care about from each record

‣ Shuffle and Sort

‣ Reduce: aggregate, summarize, filter, or transform

‣ Write the results

4

# MapReduce workflow

Input Data

Output Data

Split 0
Split 1
Split 2

read

Worker

Worker

Worker

local
write

Worker

Worker

write

Output
File 0

Output
File 1

remote
read,
sort

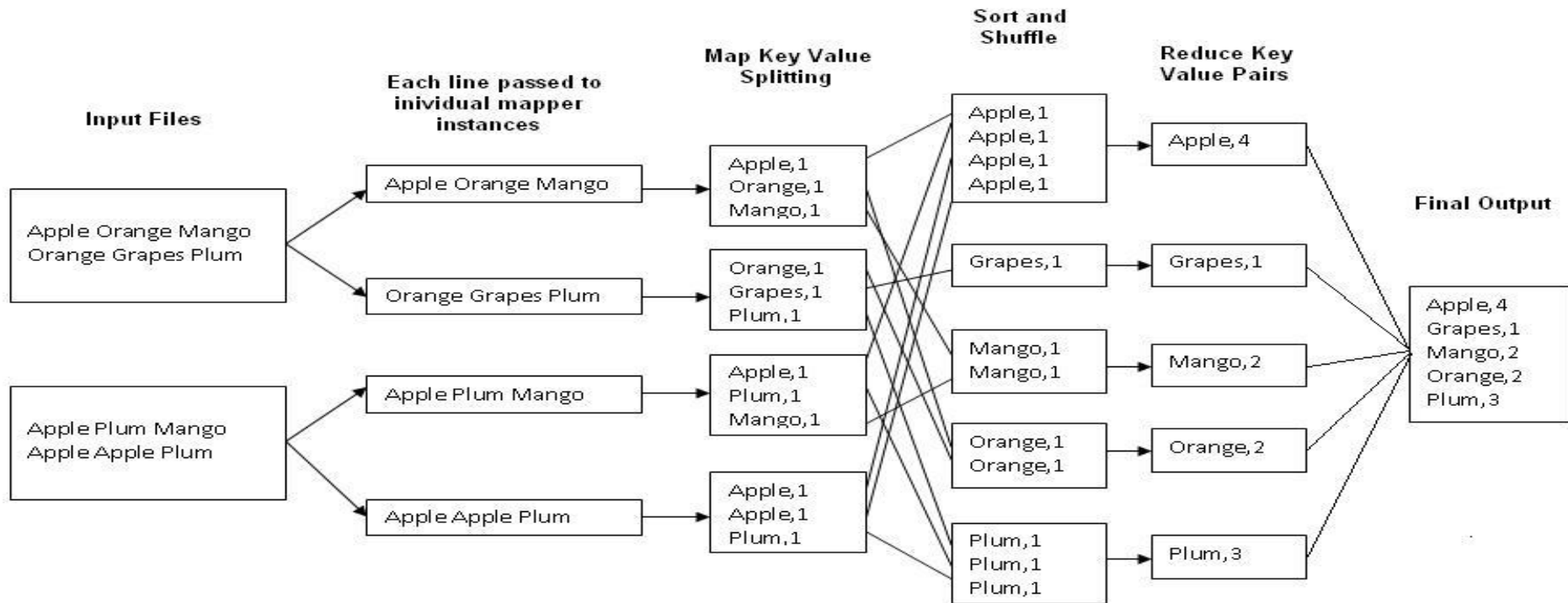**Map**

extract something you
care about from each
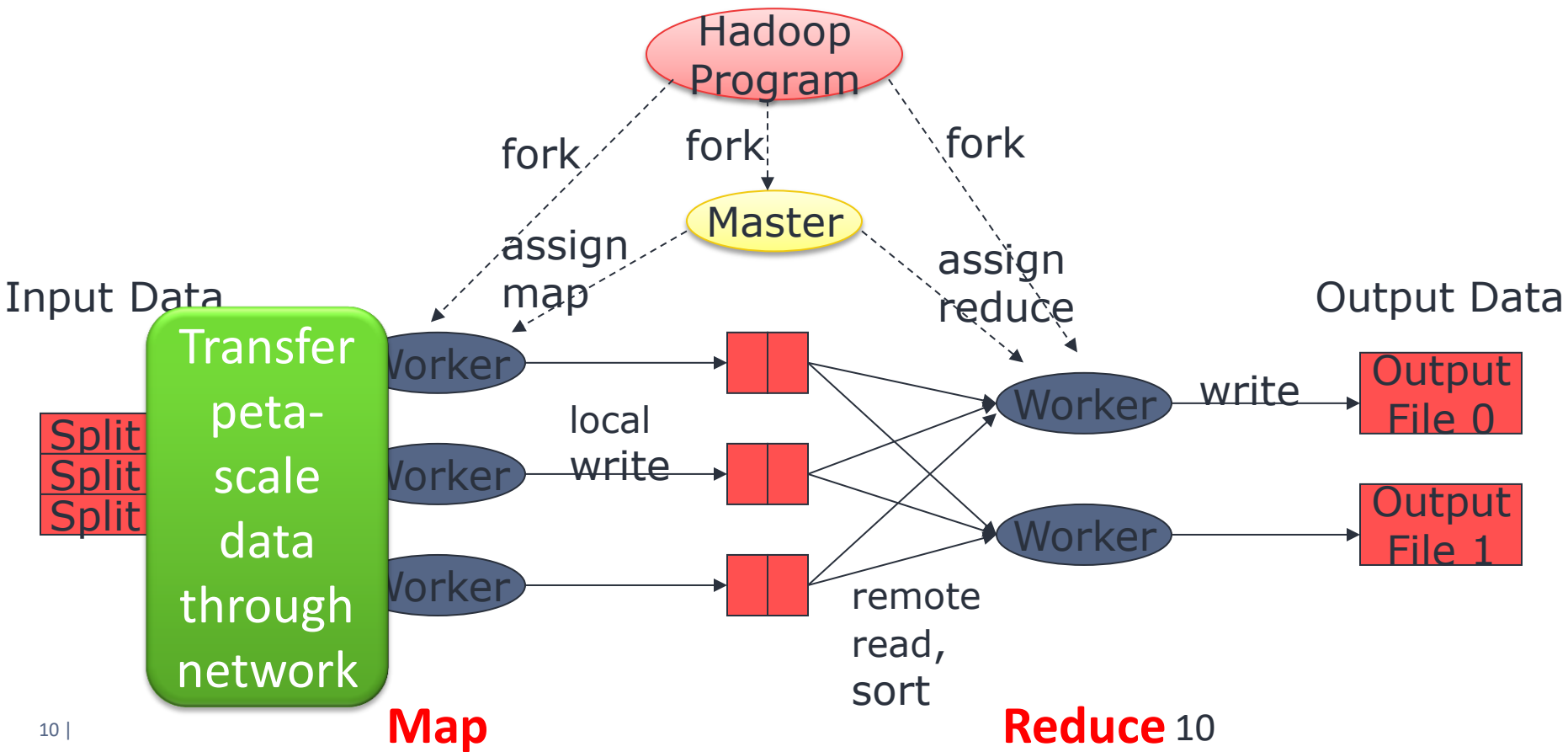record

**Reduce**

aggregate,
summarize, filter,
or transform

5

# Mappers and Reducers

‣ Need to handle more data? Just add more Mappers/Reducers!

‣ No need to handle multithreaded code ☺

- Mappers and Reducers are typically single threaded and deterministic
  - Determinism allows for restarting of failed jobs
- Mappers/Reducers run entirely independent of each other
  - In Hadoop, they run in separate JVMs

# Example: Word Count



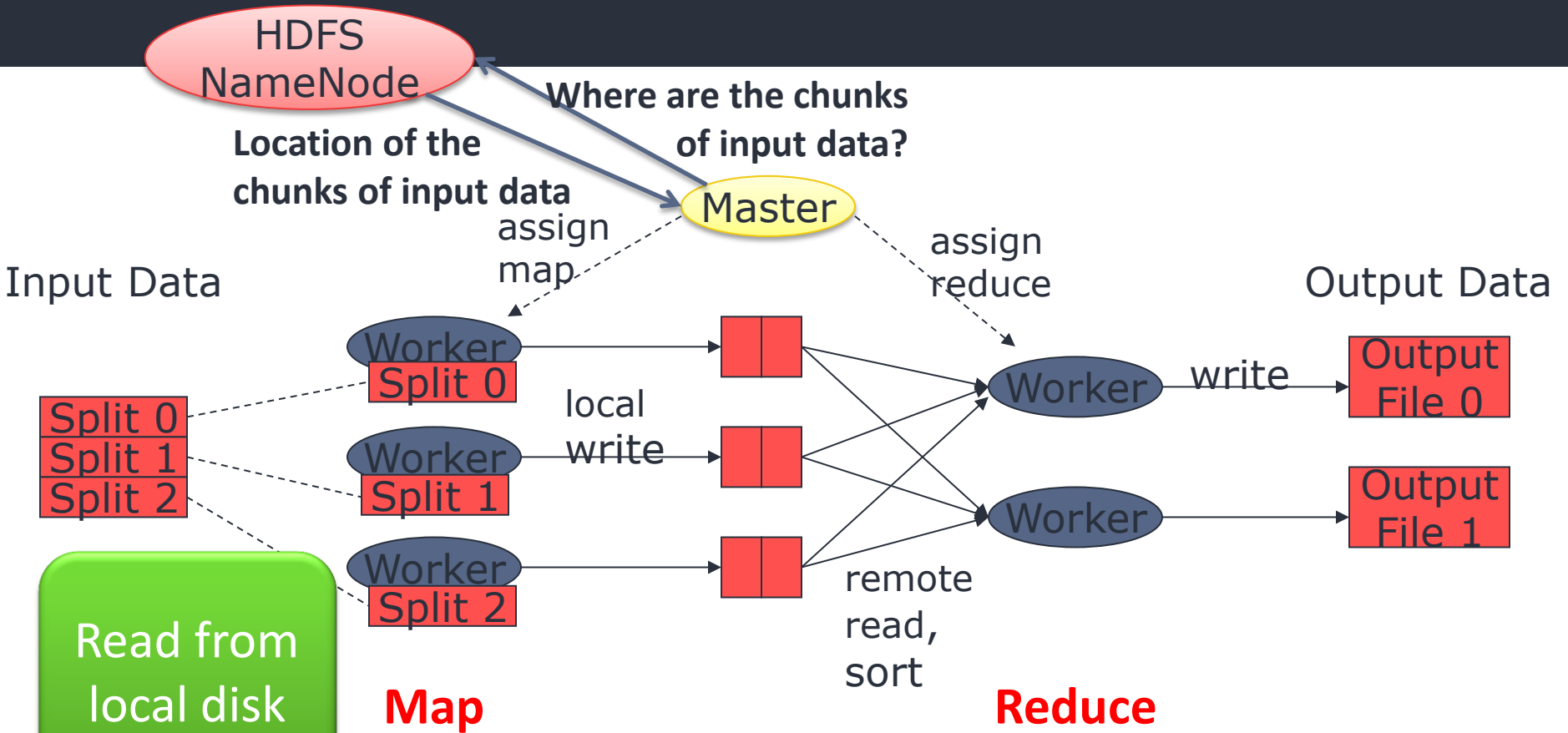http://kickstarthadoop.blogspot.ca/2011/04/word-count-hadoop-map-reduce-example.html

# MapReduce



Hadoop Program

fork · fork · fork

Master

Input Data

Transfer peta-scale data through network

assign map

assign reduce

Output Data

Split
Split
Split

Worker

Worker

local write

Worker

Worker

write

Output File 0

Worker

Output File 1

remote read, sort

**Map**

**Reduce** 10

## Hadoop Distributed File System (HDFS)

‣ Split data and store 3 replica on commodity servers

# MapReduce

‣ Failures are norm in commodity hardware

‣ **Worker** failure

- Detect failure via periodic heartbeats
- Re-execute in-progress map/reduce tasks

‣ **Master** failure

- Single point of failure; Resume from Execution Log

‣ **Robust**

- Google's experience: lost 1600 of 1800 machines once!, but finished fine

14

# Summary

‣ MapReduce

- Programming paradigm for data-intensive computing
- Distributed & parallel execution model
- Simple to program
  - The framework automates many tedious tasks (machine selection, failure handling, etc.)

# Zoom in: GFS in more detail

# Motivation: Large Scale Data Storage

‣ Manipulate large (Peta Scale) sets of data

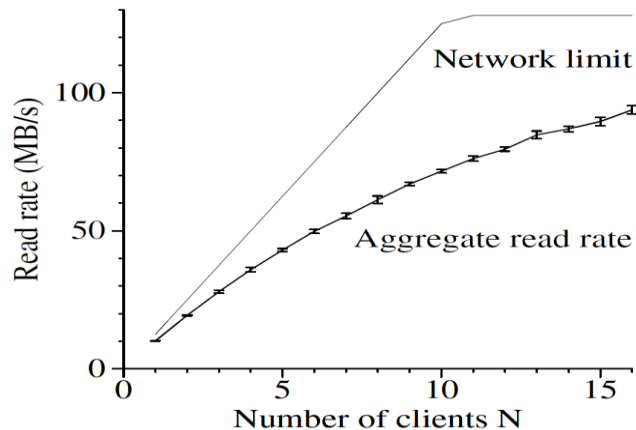‣ Large number of machines with commodity hardware

‣ Component failure is the norm

‣ Goal: **Scalable**, **high performance**, **fault tolerant** distributed file system

20

# Why a new file system?

‣ None designed for their failure model

‣ Few scale as highly or dynamically and easily

‣ Lack of special primitives for large distributed computation

▸ Designed for Google's application
- Control of both file system and application
- Applications use a few specific access patterns
  - Append to larges files
  - Large streaming reads
- Not a good fit for
  - low-latency data access
  - lots of small files, multiple writers, arbitrary file modifications

▸ Not POSIX, although mostly traditional
- Specific operations: RecordAppend

# Contents

- ‣ Motivation

- ‣ **Design overview**
  - Write Example
  - Record Append

- ‣ Fault Tolerance & Replica Management

- ‣ Conclusions
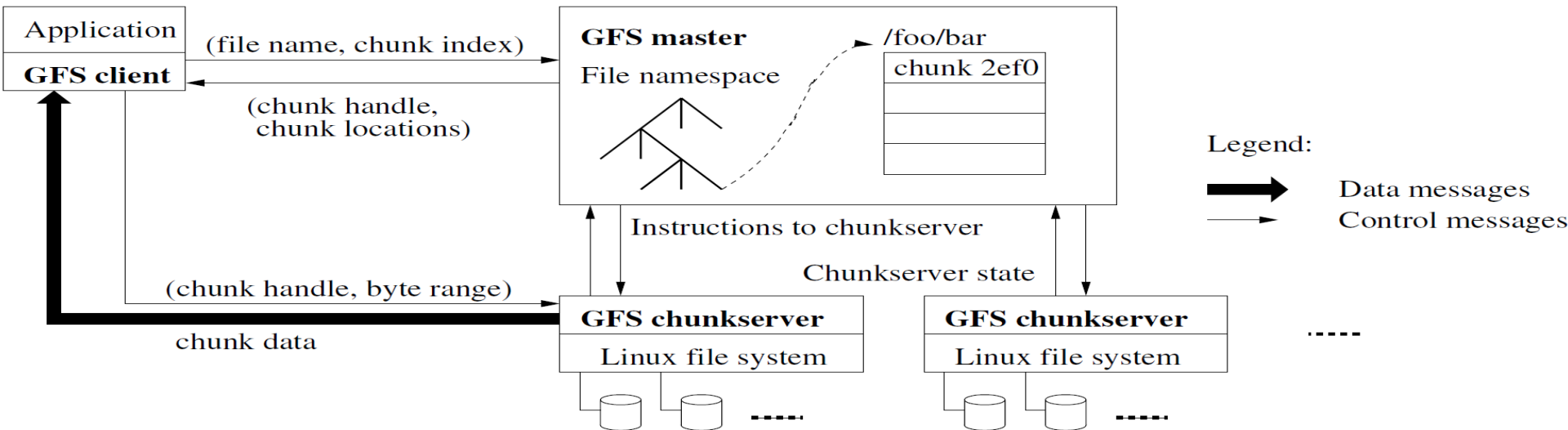
‣ **Master (NameNode)**

- Manages metadata (namespace)
- Not involved in data transfer
- Controls allocation, placement, replication
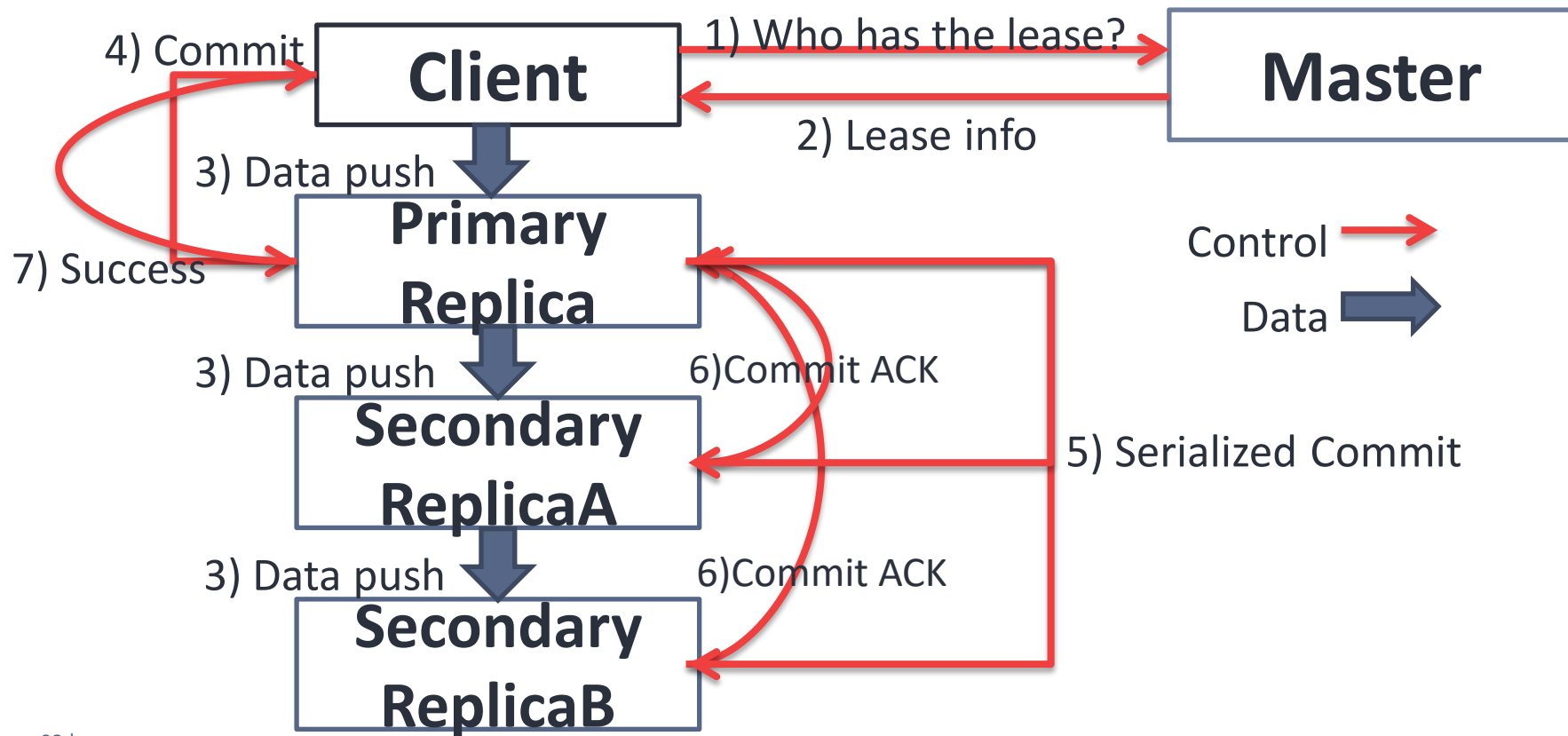
‣ **Chunkserver (DataNode)**

- Stores chunks of data
- No knowledge of GFS file system structure
- Built on local linux file system

www.cse.buffalo.edu/~okennedy/courses/cse704fa2012/2.2-HDFS.pptx

# Write(filename, offset, data)

# RecordAppend(filename, data)

- Significant use in distributed apps. For example at  Google production cluster:
  - 21% of bytes written
  - 28% of write operations
- Guaranteed: All data appended at least once as a single consecutive byte range
- Same basic structure as write

  - Client obtains information from master

  - Client sends data to data nodes (chunkservers)

  - Client sends "append-commit"

  - Lease holder serializes append

- **Advantage:** Large number of concurrent writers with minimal coordination

29

# RecordAppend (2)

‣ Record size is limited by chunk size

‣ When a record does not fit into available space,
  - chunk is padded to end
  - and client retries request.

# Contents

- ‣ Motivation

- ‣ Design overview
  - Write Example
  - Record Append

- ‣ **Fault Tolerance & Replica Management**

- ‣ Conclusions

# Fault tolerance

‣ # Replication

- High availability for reads

- User controllable, default 3 (non-RAID)

- Provides read/seek bandwidth

- Master is responsible for directing re-replication if a data node dies

‣ # Online checksumming in data nodes

- Verified on reads

# Replica Management

‣ Bias towards topological spreading

- Rack, data center

‣ Rebalancing

- Move chunks around to balance disk fullness

- Gently fixes imbalances due to:

  - Adding/removing data nodes

# Replica Management (Cloning)

- Chunk replica lost or corrupt
- Goal: minimize app disruption and data loss
  - Approximately in priority order
    - More replica missing-> priority boost
    - Deleted file-> priority decrease
    - Client blocking on a write-> large priority boost
  - Master directs copying of data

- Performance on a production cluster
  - Single failure, full recovery (600GB): 23.2 min
  - Double failure, restored 2x replication: 2min

34

# Garbage Collection

‣ Master does **not** need to have a strong knowledge of what is stored on each data node
- Master regularly scans namespace
- After GC interval, deleted files are removed from the namespace
- Data node periodically polls Master about each chunk it knows of.
- If a chunk is forgotten, the master tells data node to delete it.

# Limitations

‣ Master is a central point of failure

‣ Master can be a scalability bottleneck

‣ Latency when opening/stating thousands of files

‣ Security model is weak

# Conclusion

‣ Inexpensive commodity components can be the basis of a large scale reliable system

‣ Adjusting the API, e.g. RecordAppend, can enable large distributed apps

‣ Fault tolerant

‣ Useful for many similar apps

37