

Hadoop Training - III

Architecture & Theory

Credit

- Show up: 50 points
- Asking questions: 10 points
- Answering questions: 15 points
- Exam: 20 points

Hadoop Introduction

History

- Created by Doug Cutting
- Named after his son's toy elephant
- Developed to support distribution for the Nutch search engine project

www.nytimes.com/imagepages/2009/03/16/business/17cl ☆

The New York Times

March 16, 2009



Peter DaSilva for The New York Times

Doug Cutting with the stuffed elephant that inspired the name Hadoop, the software program he developed.

Close Window

http://www.nytimes.com/2009/03/17/technology/business-computing/17cloud.html?_r=1

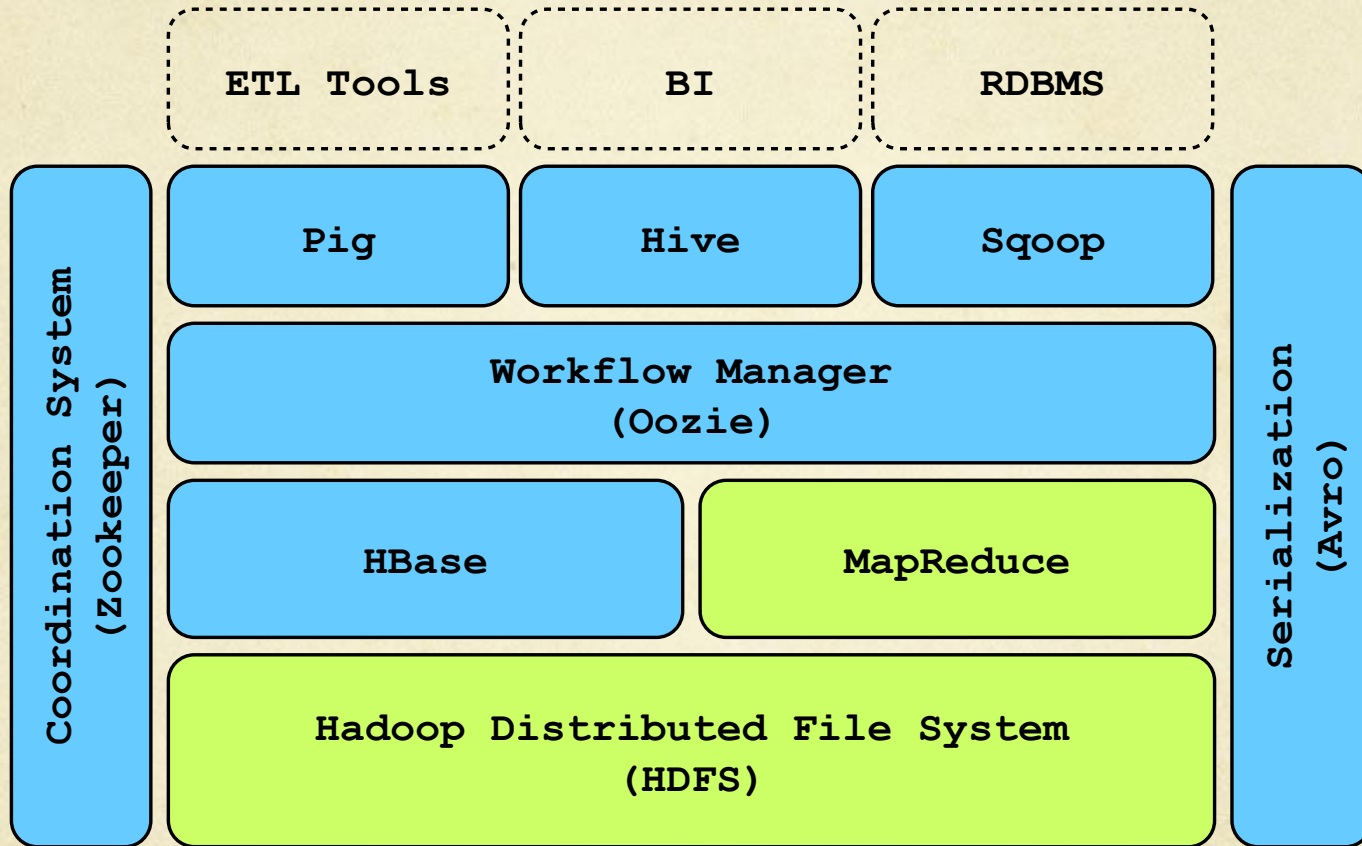
Assumptions and Goals

- Large Data Sets (GB, TB, PB, EB?, ZB?, YB?)
- Write-Once-Read-Many Access Model
- Streaming Data Access
- Moving Computation is Cheaper than Moving Data
- Hardware Failure
- Portability Across Heterogeneous Hardware and Software Platforms

What Hadoop is Not

- Hadoop is **not** a substitute for a database
- MapReduce is **not** always the best algorithm
- HDFS is **not** a substitute for a High Availability SAN-hosted File System
- HDFS is **not** a POSIX File System

Hadoop Ecosystem



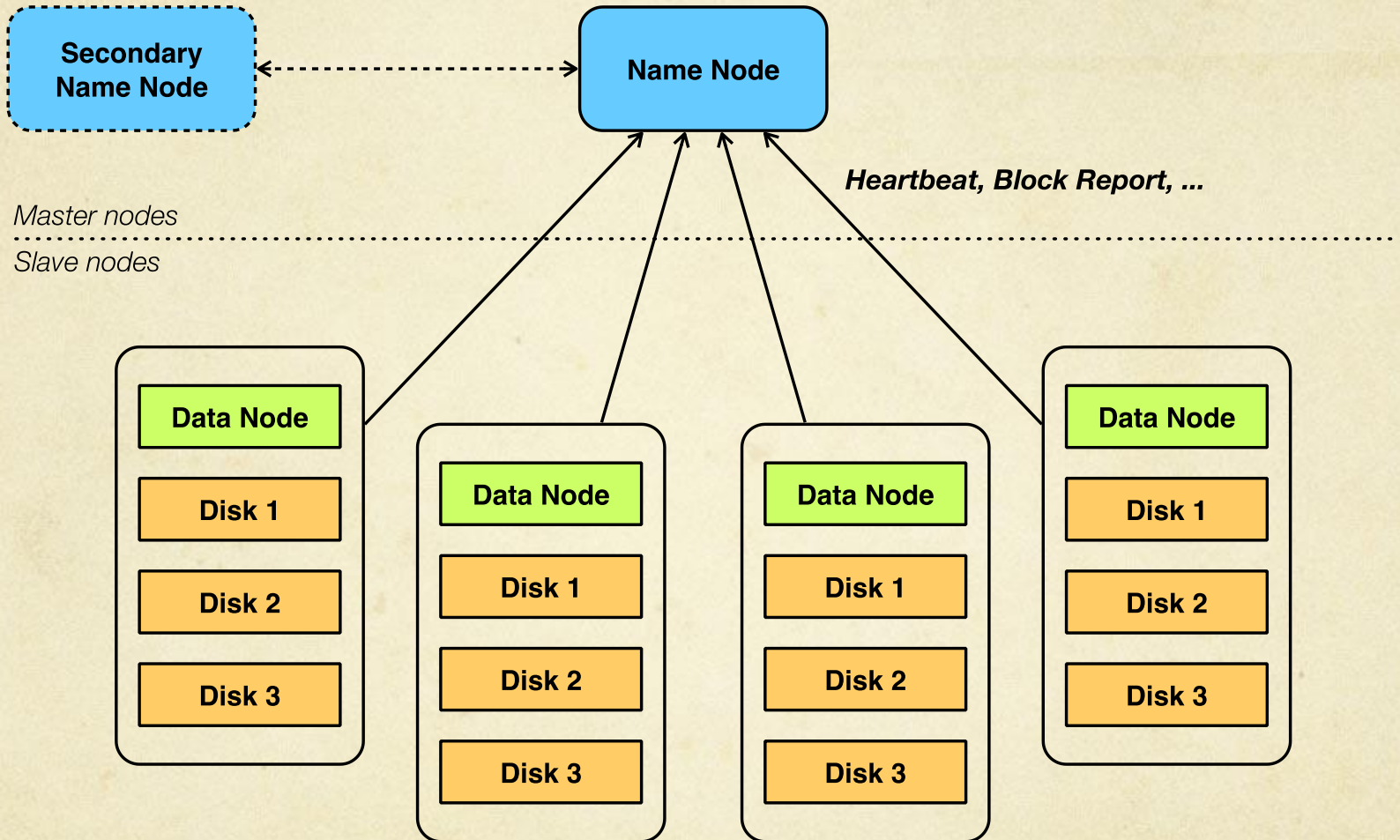
Hadoop File System

HDFS Concepts

- Files are broken into block-sized chunks
 - `dfs.block.size`, default is 64 MB
- File Replication
 - `dfs.replication`, default is 3
 - `dfs.replication.max`, default is 512
 - `dfs.replication.min`, default is 1

HDFS Architecture

- Master-Slave Architecture
- Name Node (Master)
- Data Node (Slave)
- Secondary Name Node



Name Node

- Manage File System Namespace
 - Maintain two critical tables
 1. File Name → Block Sequence
 2. Block ID → Machine List
- Single Point of Failure

File Name	Replicas	Block Sequence	Others
/data/part-0	2	B1, B2, B3	user, group, ...
/data/part-1	3	B4, B5	foo, bar, ...

Memory

Disk

fsimage

File Name	Replicas	Block Sequence	Others
/data/part-0	3	B1, B2, B3	user, group, ...
/data/part-1	3	B4, B5	user, group, ...

edits

OP Code	Operands
OP_SET_REPLICATION	"/data/part-0", 2
OP_SET_OWNER	"/data/part-1", "foo", "bar"

fsimage

edits



Merge



File Name	Replicas	Block Sequence	Others
/data/part-0	2	B1, B2, B3	user, group, ...
/data/part-1	3	B4, B5	foo, bar, ...

.....

Block ID	Machine (Data Node) List
B1	DN-1, DN-2
B2	DN-2, DN-3
B3	DN-3, DN-4
B4	DN-4, DN-1, DN-2
B5	DN-2, DN-3, DN-1

Block Report



DN-1

DN-2

DN-3

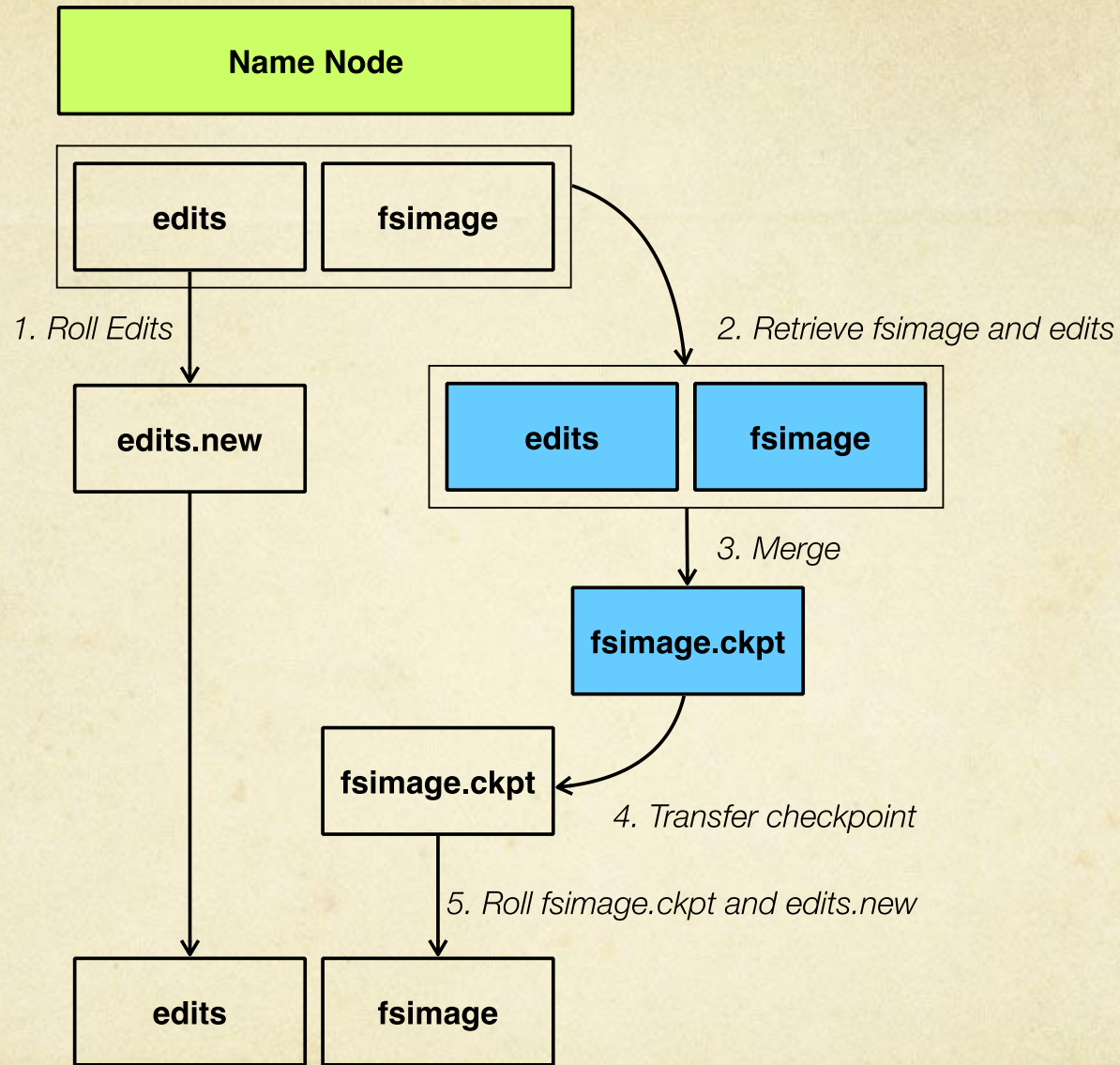
DN-4

Data Node

- Manage Storage
- Replicate Blocks
- Block Report
- Send Heartbeat to Name Node
 - `dfs.heartbeat.interval`, default 3 seconds

Secondary Name Node

- **Not** a backup of Primary Name Node
- Merge fsimage and edits log



Hadoop MapReduce

MapReduce Concepts

- Jobs are broken into tasks
- Parallelize
 - Computation
 - Disk IO

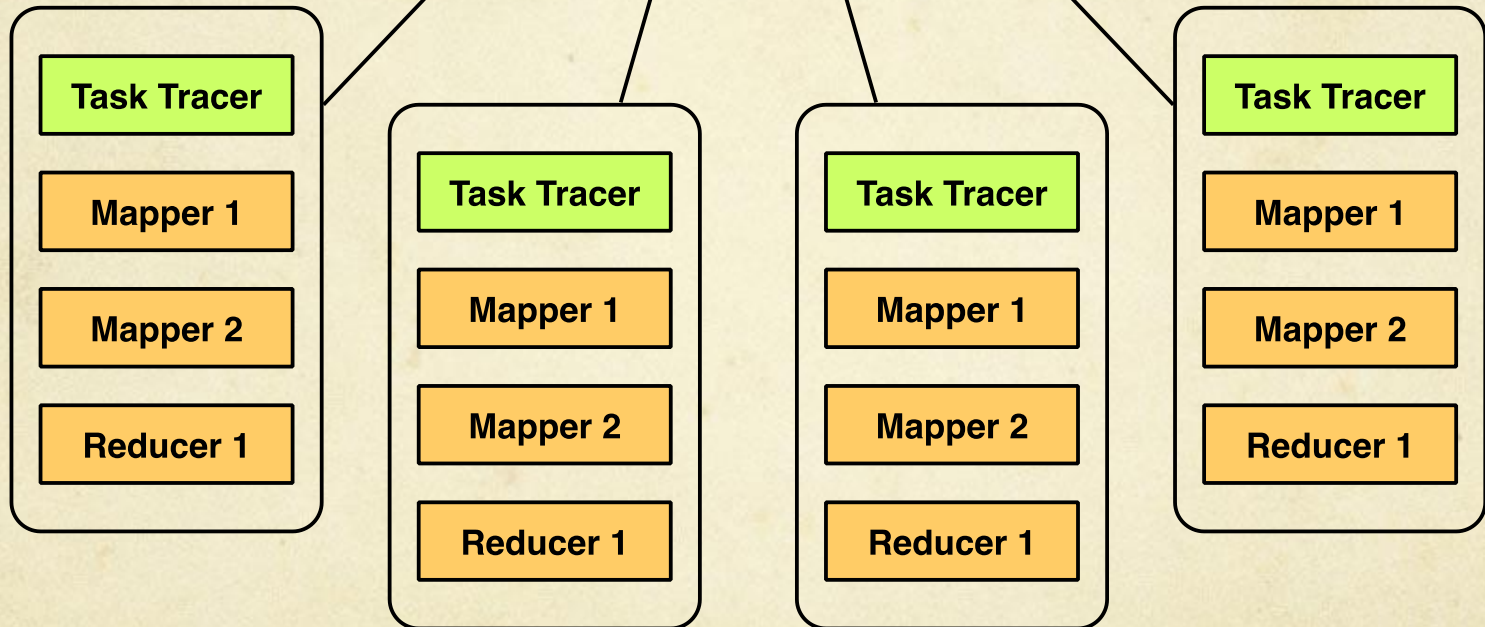
MapReduce Architecture

- Master-Slave Architecture
- Job Tracker (Master)
- Task Tracker (Slave)
 - Mapper
 - Reducer

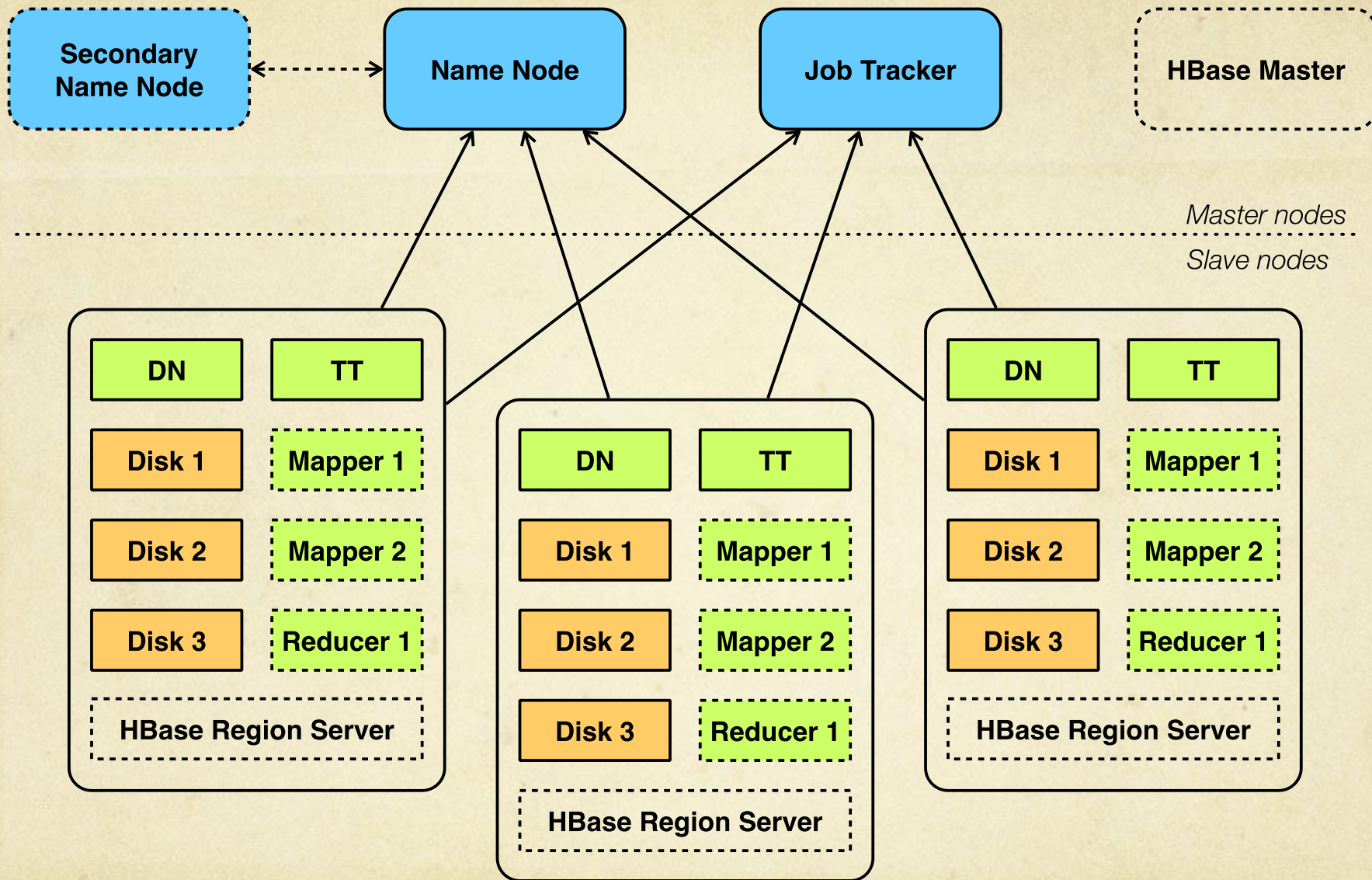


Master nodes

Slave nodes



Hadoop Cluster



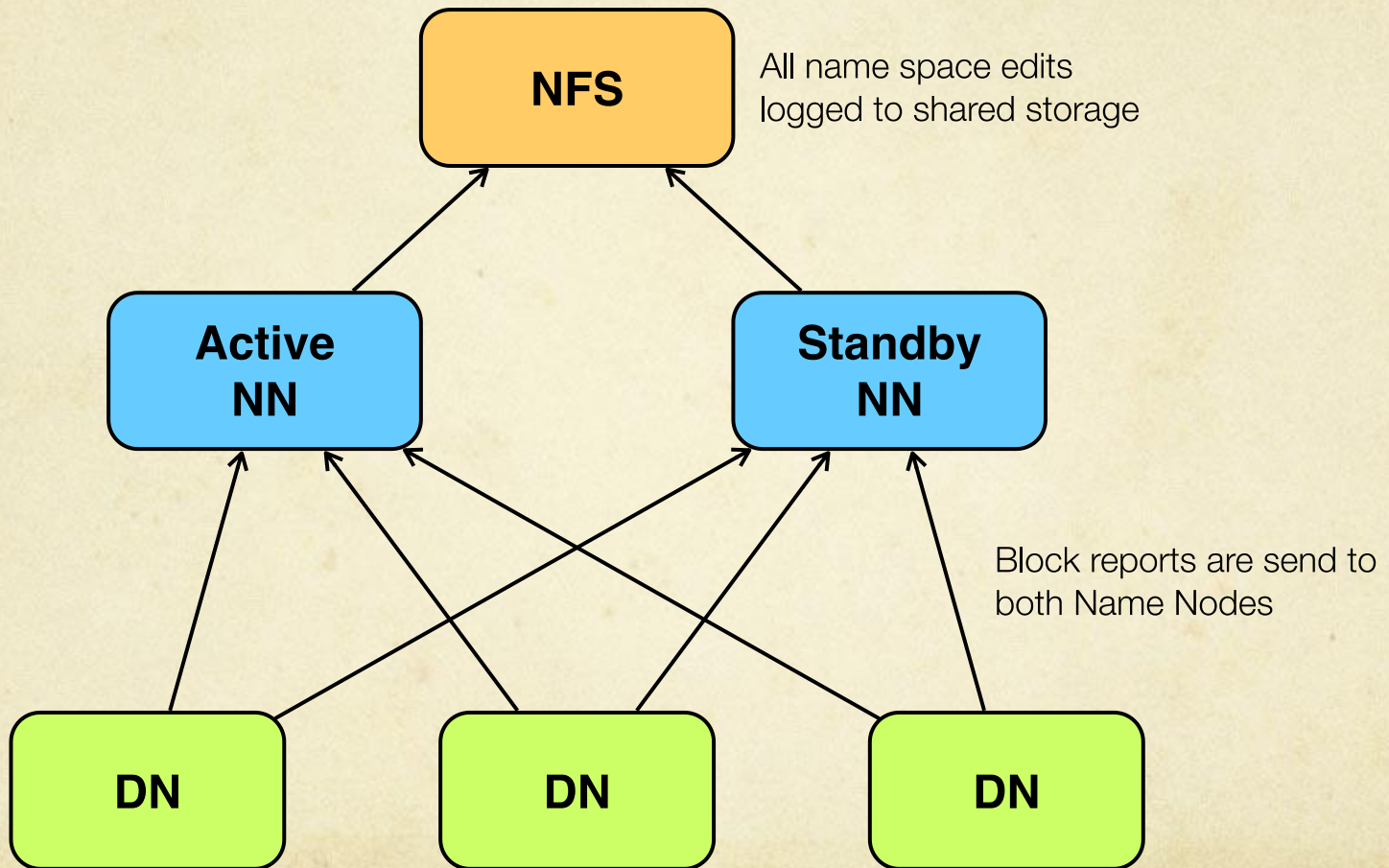
Slave Node

- Data Node
- Task Tracker
 - Mapper
 - Reducer
- HBase Region Server

High Availability for the HDFS (Name Node)

- Storage
 - RAID 1
 - Shared Storage
- Host
 - DRDB-LinuxHA
 - HDFS HA
 - Avatar Node

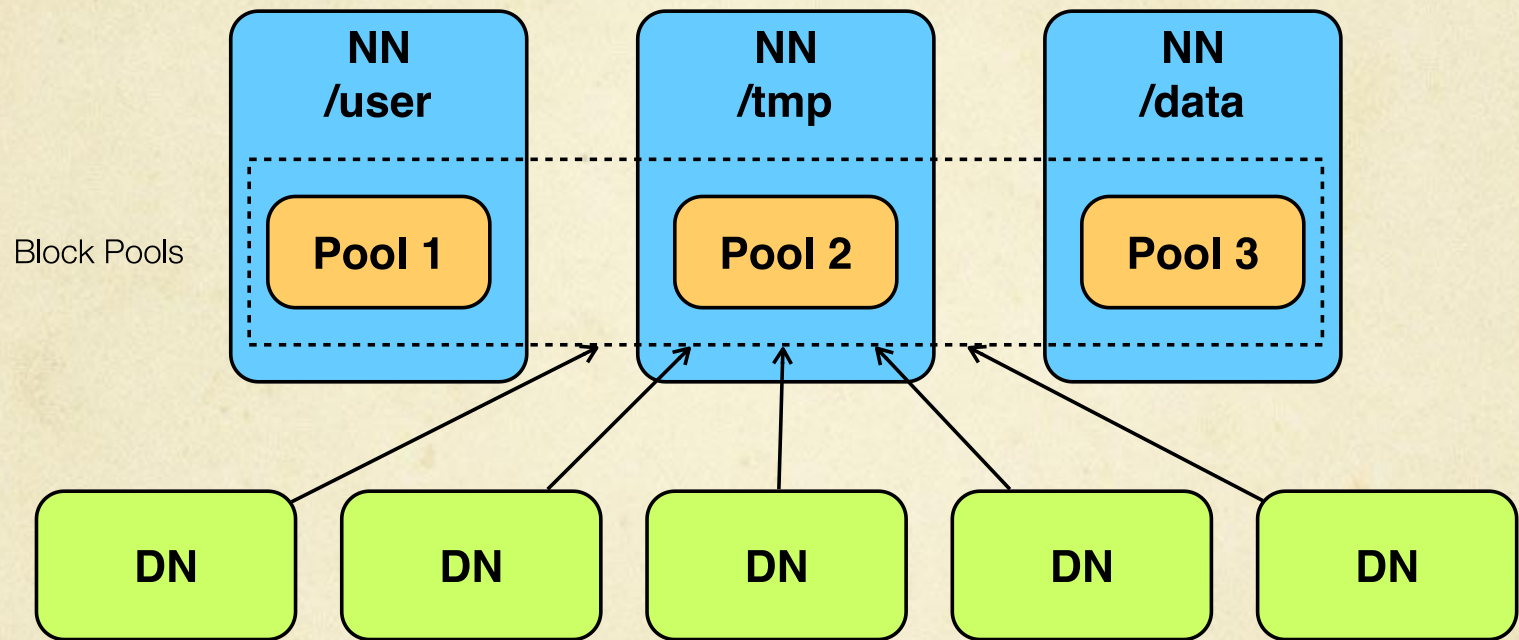
HDFS HA



Avatar Node

- Developed by Facebook
- Roles
 - Primary Avatar Node
 - Standby Avatar Node (Hot Standby)
 - Avatar Data Node

HDFS Federation



MapReduce

Program Model

- Mapper

- $M(K_{in}, V_{in}) \rightarrow (K_{tmp}, V_{tmp})$

- Reducer

- $R(K_{tmp}, [V_{tmp1}, V_{tmp2}, \dots]) \rightarrow (K_{out}, V_{out})$

Example: Word Count

Hello Hadoop
Goodbye Hadoop



K_{in}	V_{in}
1	Hello Hadoop
2	Goodbye Hadoop



Mapper



K_{tmp}	V_{tmp}
Hello	1
Hadoop	1
Goodbye	1
Hadoop	1

K_{tmp}	$[V_{tmp}]$
Goodbye	[1]
Hadoop	[1, 1]
Hello	[1]

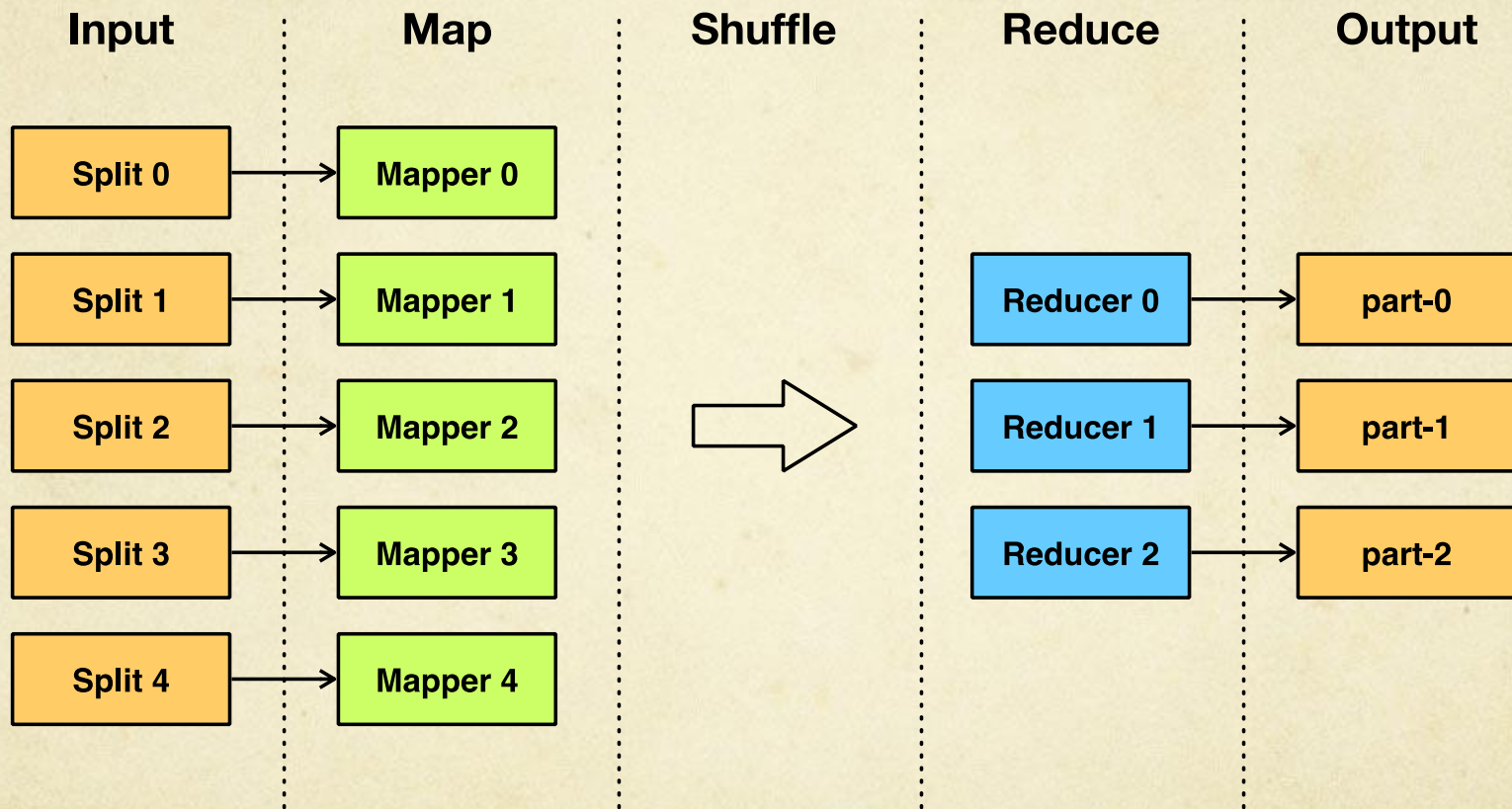


Reducer



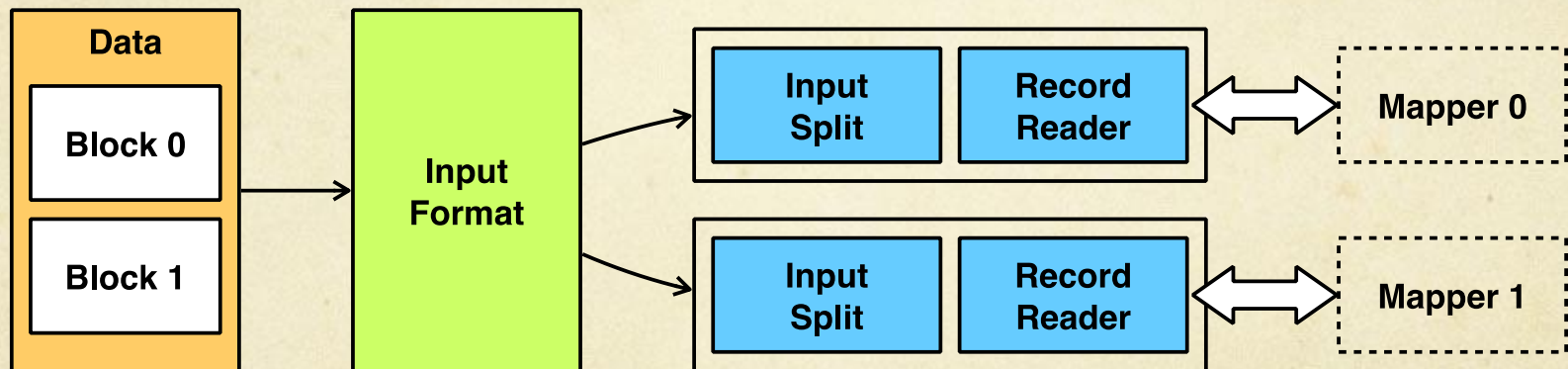
K_{out}	V_{out}
Goodbye	1
Hadoop	2
Hello	1

Phases



Input Phase

- Input Format
- Input Split
- Record Reader



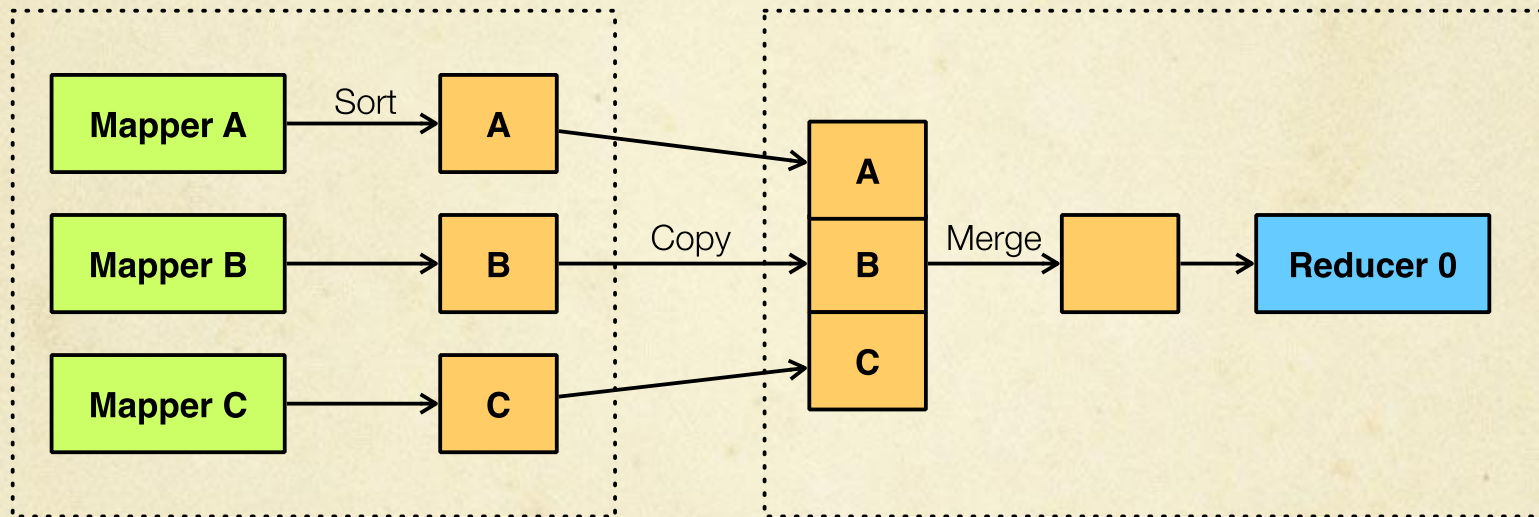
File Formats

- Text File
- Binary File
 - Sequence File
 - Map File
 - TFile
 - HAR

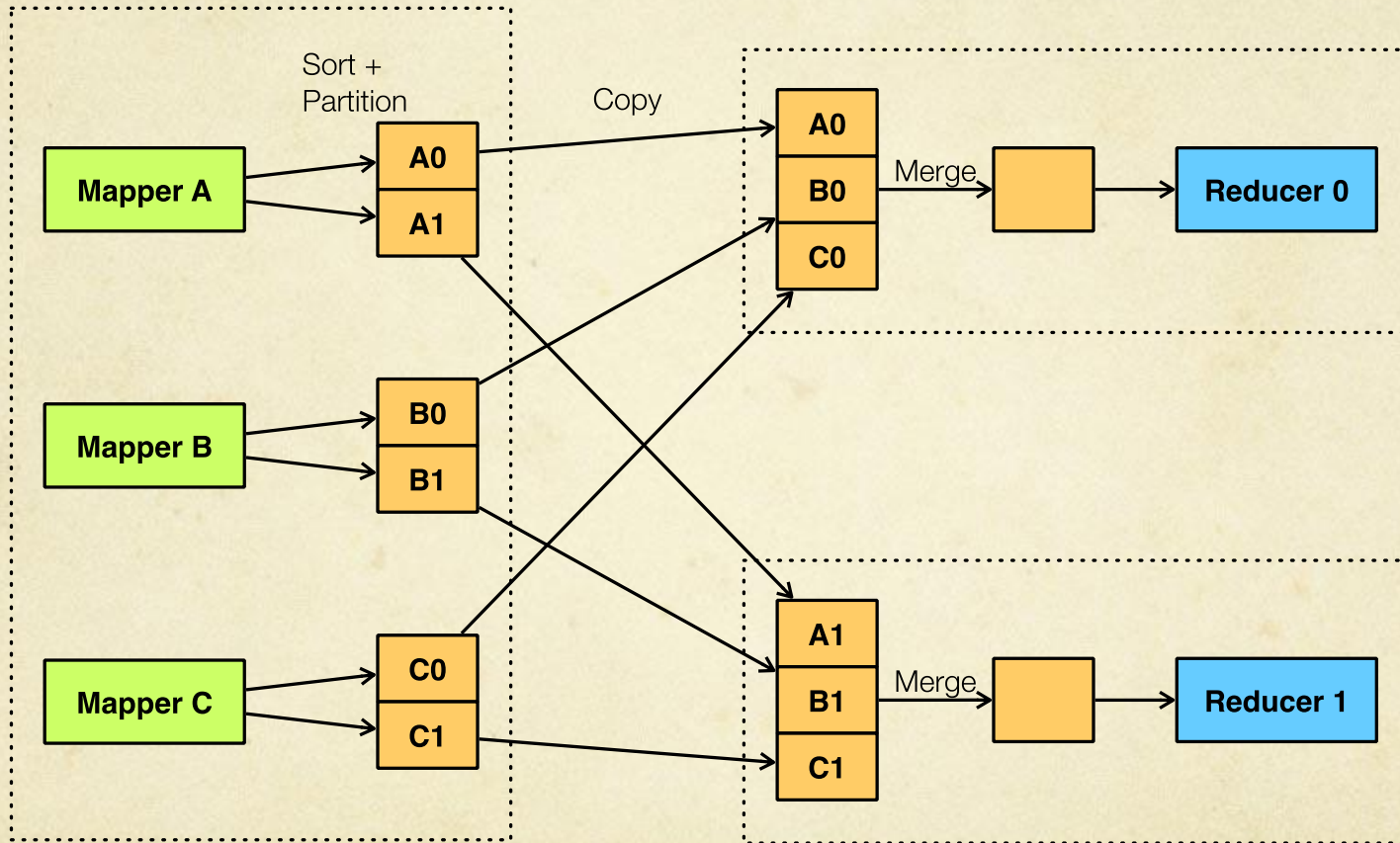
Compression and Split

- Why Compression?
 - Reduce Data Transfer Time (Disk IO, Network)
 - Save Storage
- Not all compression formats support splitting

Map, Shuffle and Reduce Phase

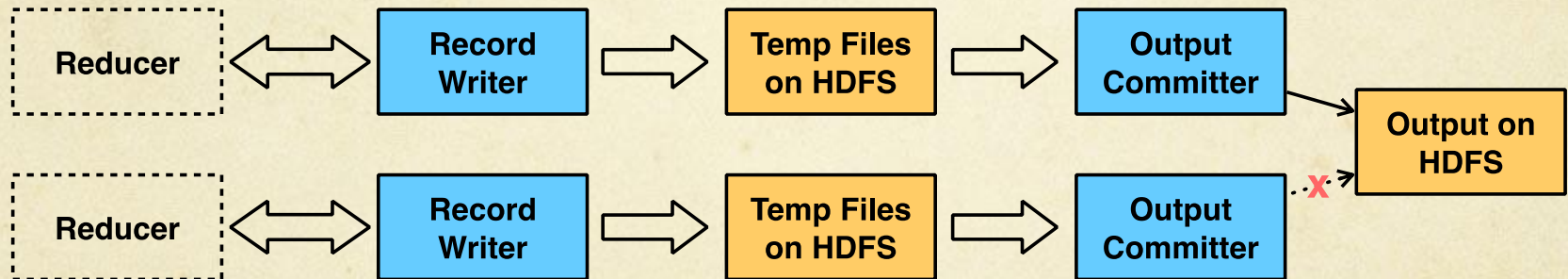


Multiple Reducers



Output Phase

- OutputFormat
- RecordWriter
- OutputCommitter



Attempts and Speculative Execution

- Attempts
 - `mapred.map.max.attempts`, default is 4
 - `mapred.reduce.max.attempts`, default 4
- Speculative Execution
 - `mapred.map.tasks.speculative.execution`, default is true
 - `mapred.reduce.tasks.speculative.execution`, default is true

- Combiner
 - **Not** Map-Side Reducer
 - No guarantee of how many times will be invoked
- Counters
- Number of Reducers
- Hadoop Streaming
- Hadoop Pipes

- Job Scheduler
 - Job Queue Task Scheduler (default)
 - Fair Scheduler
 - Capacity Scheduler
- MRv2 (YARN)
 - Apache Hadoop NextGen MapReduce
 - Included in hadoop-0.23, CDH4
 - Not stable yet