

Data Engineering

MG-GY 8441

Big Data

Agenda

- Distributed Storage
- Parallel Computing
- Dependency Graphs
- Spark

References

- Leskovec, Rajaraman, Ullman, Mining of Massive Datasets (Chapter 2.1 - 2.3)
- (Optional) Dean, Ghemawat, MapReduce: Simplified Data Processing on Large Clusters

Examples

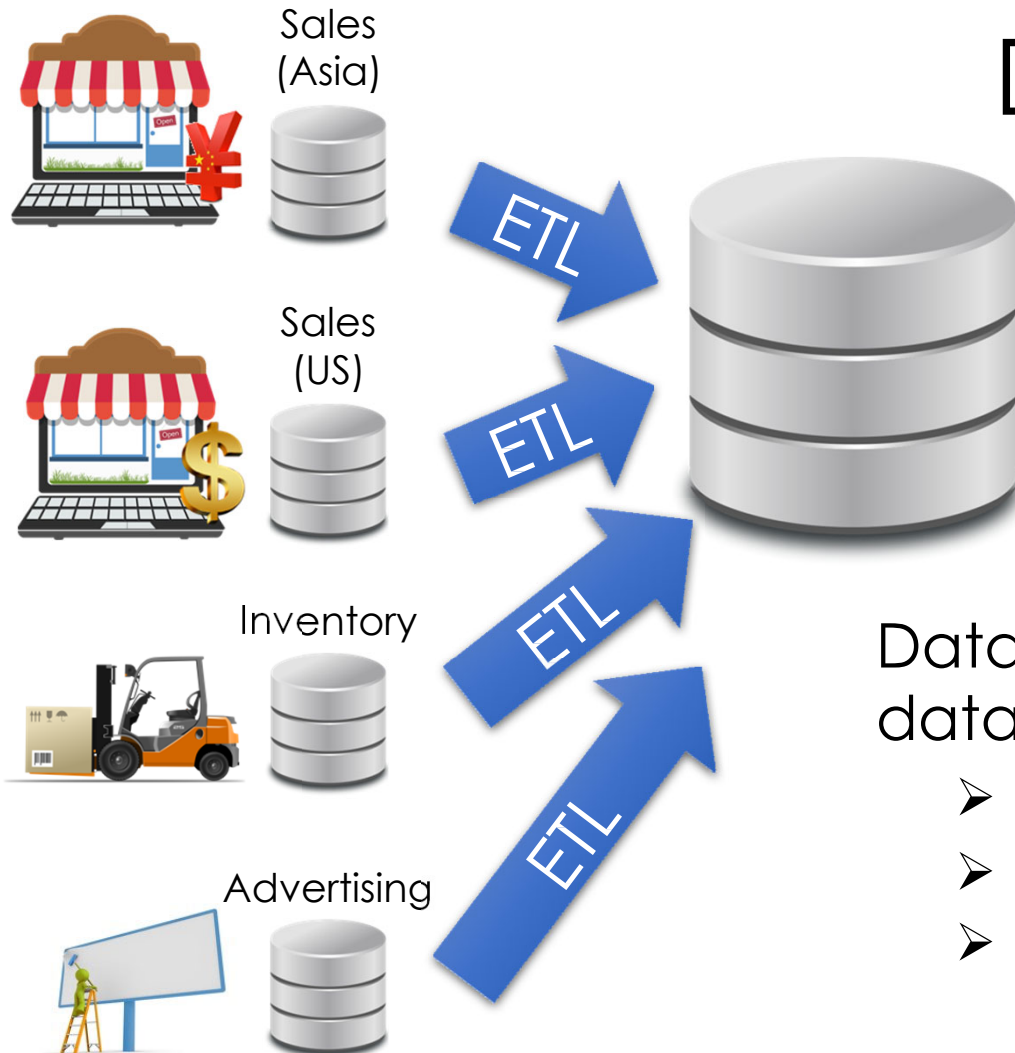
- (2019) Facebook generates 4 Petabytes
- (2019) YouTube collects 300 hours of HD video every minute
- (2017) Twitter stores >500PB of data
- (2017) CERN data center: > 200PB

Data Warehouse

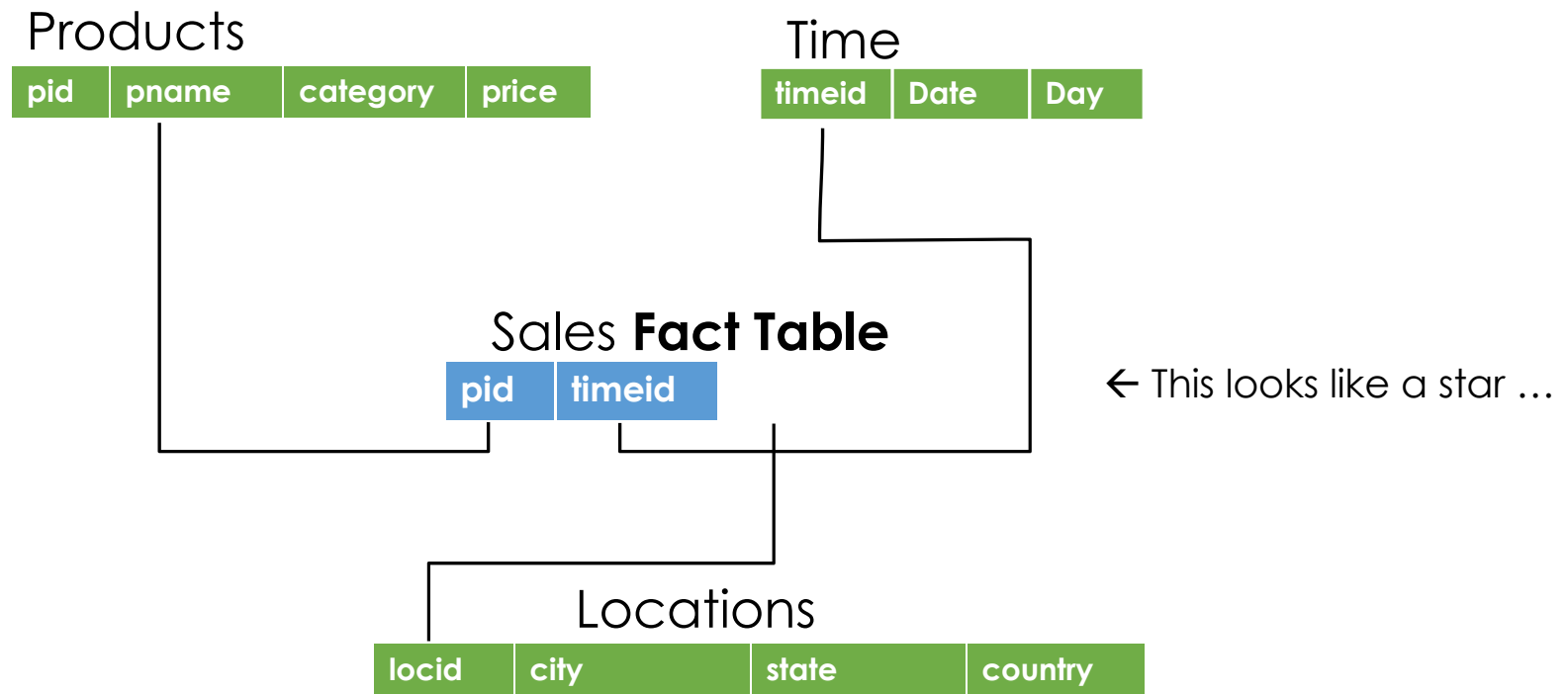
Collects and organizes historical data from multiple sources

Data is *periodically* **ETL**ed into the data warehouse:

- **Extracted** from remote sources
- **Transformed** to standard schemas
- **Loaded** into the (typically) relational (SQL) data system



Schema



Data Warehouse

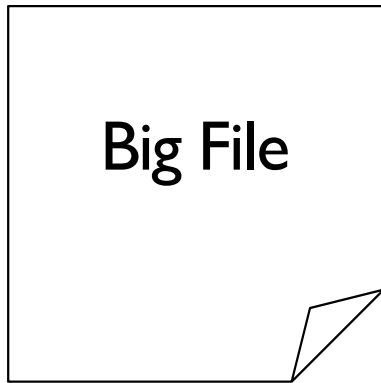
Collects and organizes historical data from multiple sources



- How do we deal with semi-structured and unstructured data?
- Do we really want to force a schema on load?

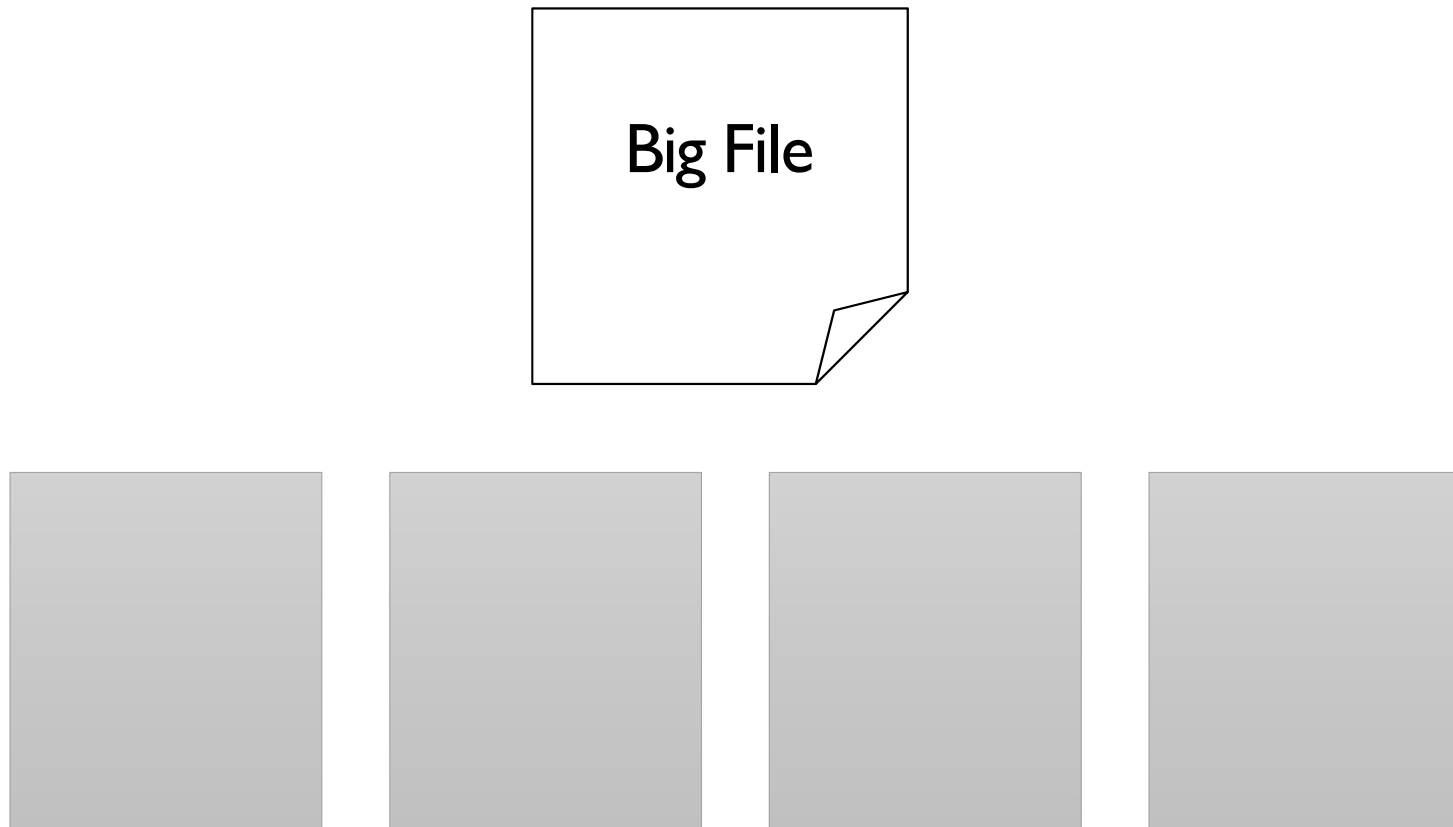
Distributed File Systems

Fault Tolerant Distributed File Systems (e.g., HDFS)



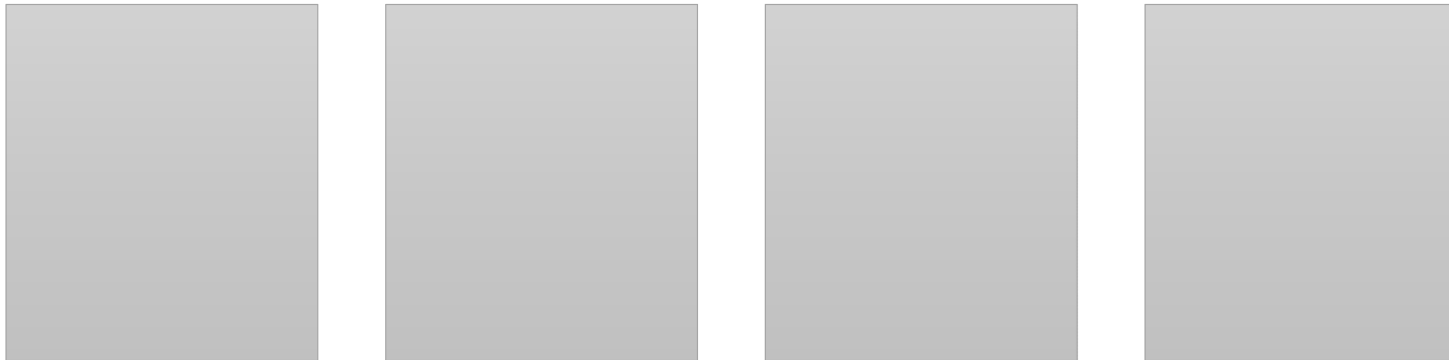
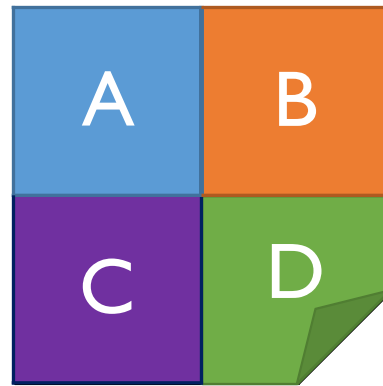
How do we **store** and **access** very **large files** across **cheap** commodity devices ?

Fault Tolerant Distributed File Systems (e.g., HDFS)



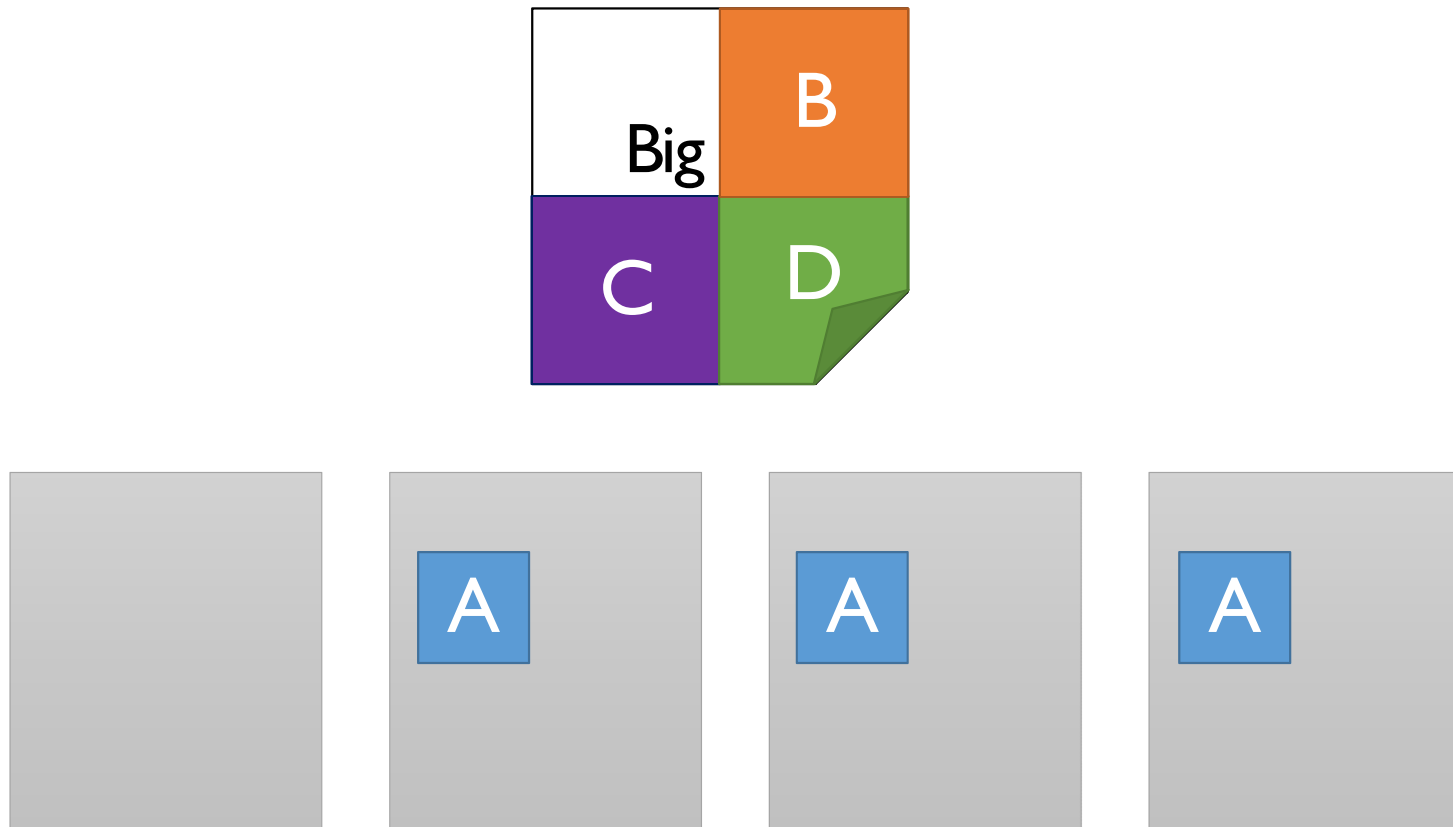
[Ghemawat et al., SOSP'03]

Fault Tolerant Distributed File Systems (e.g., HDFS)



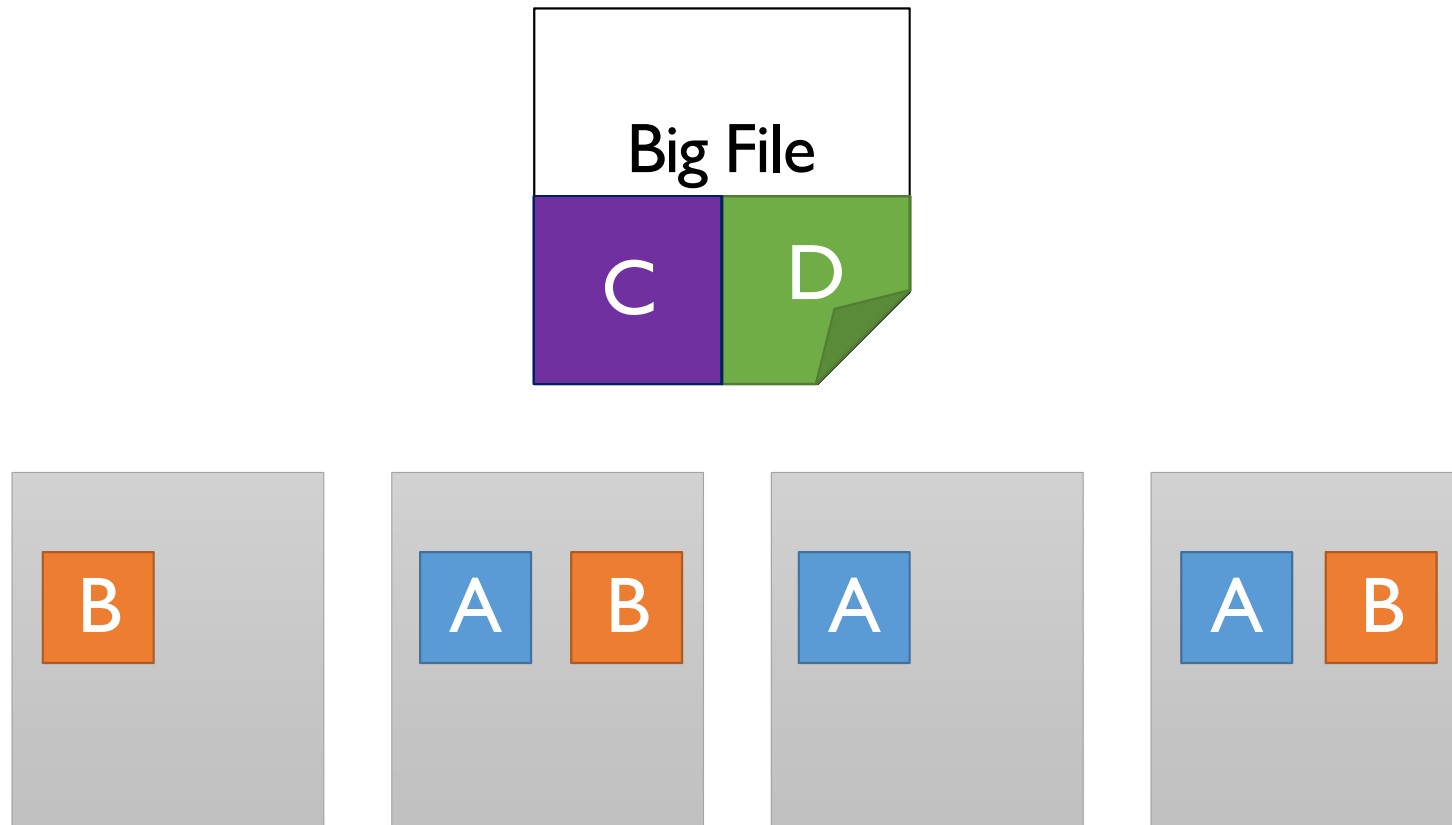
[Ghemawat et al., SOSP'03]

Fault Tolerant Distributed File Systems (e.g., HDFS)



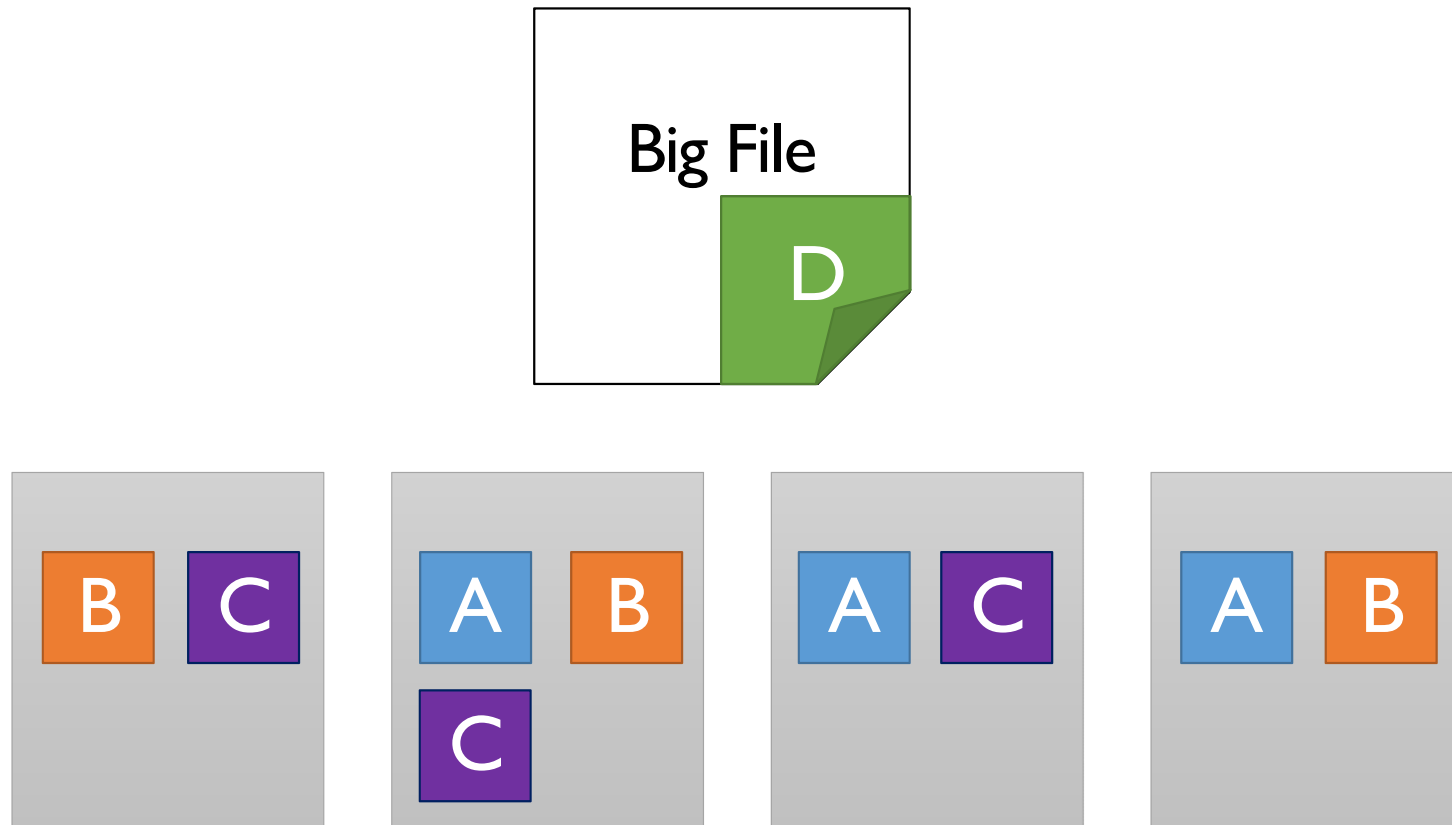
[Ghemawat et al., SOSP'03]

Fault Tolerant Distributed File Systems (e.g., HDFS)



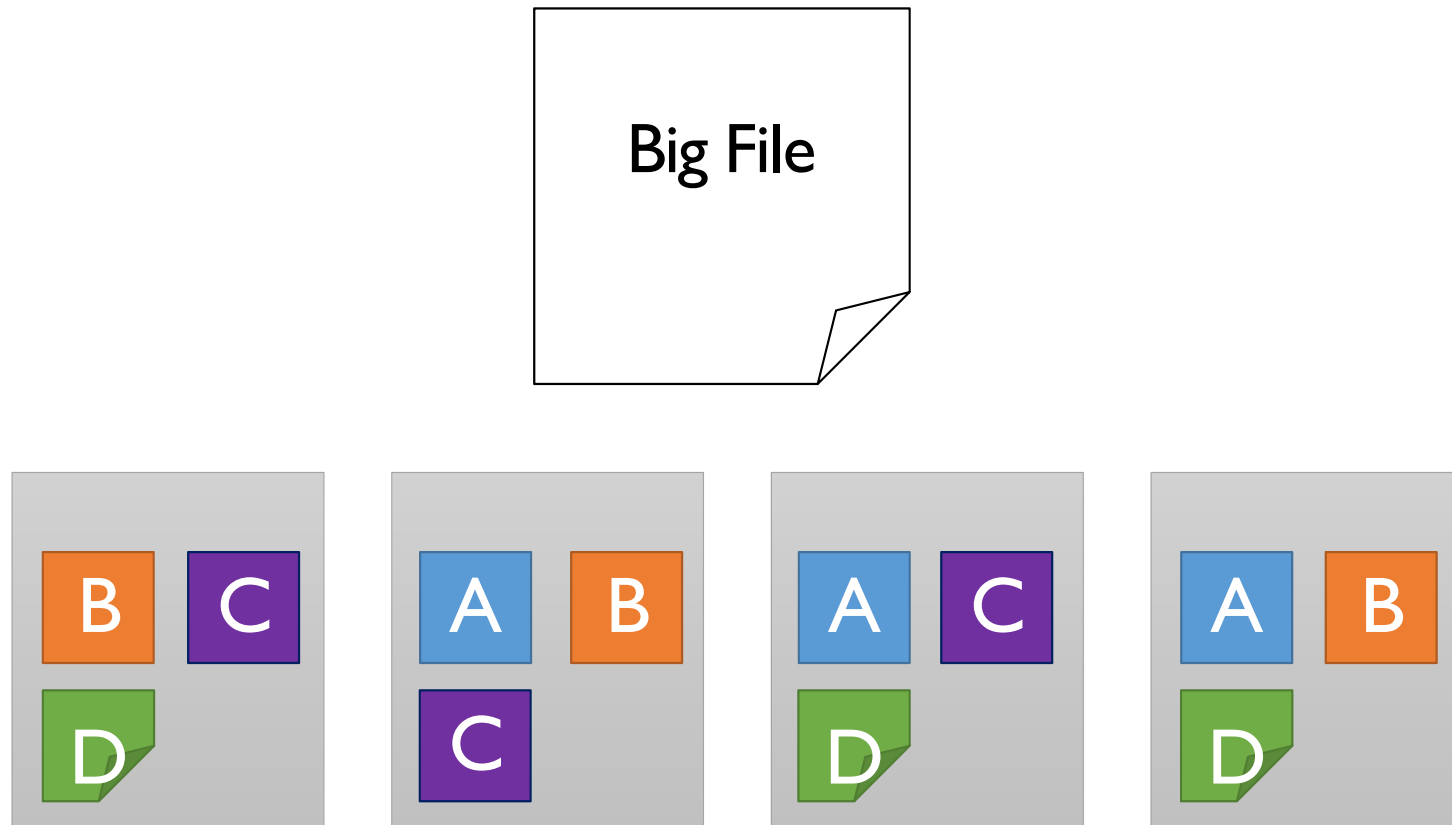
[Ghemawat et al., SOSP'03]

Fault Tolerant Distributed File Systems (e.g., HDFS)



[Ghemawat et al., SOSP'03]

Fault Tolerant Distributed File Systems (e.g., HDFS)

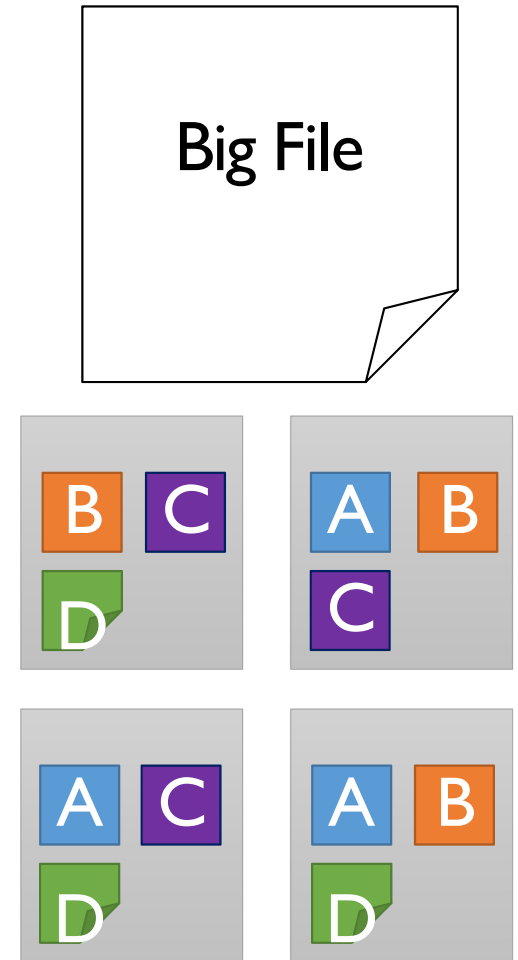


[Ghemawat et al., SOSP'03]

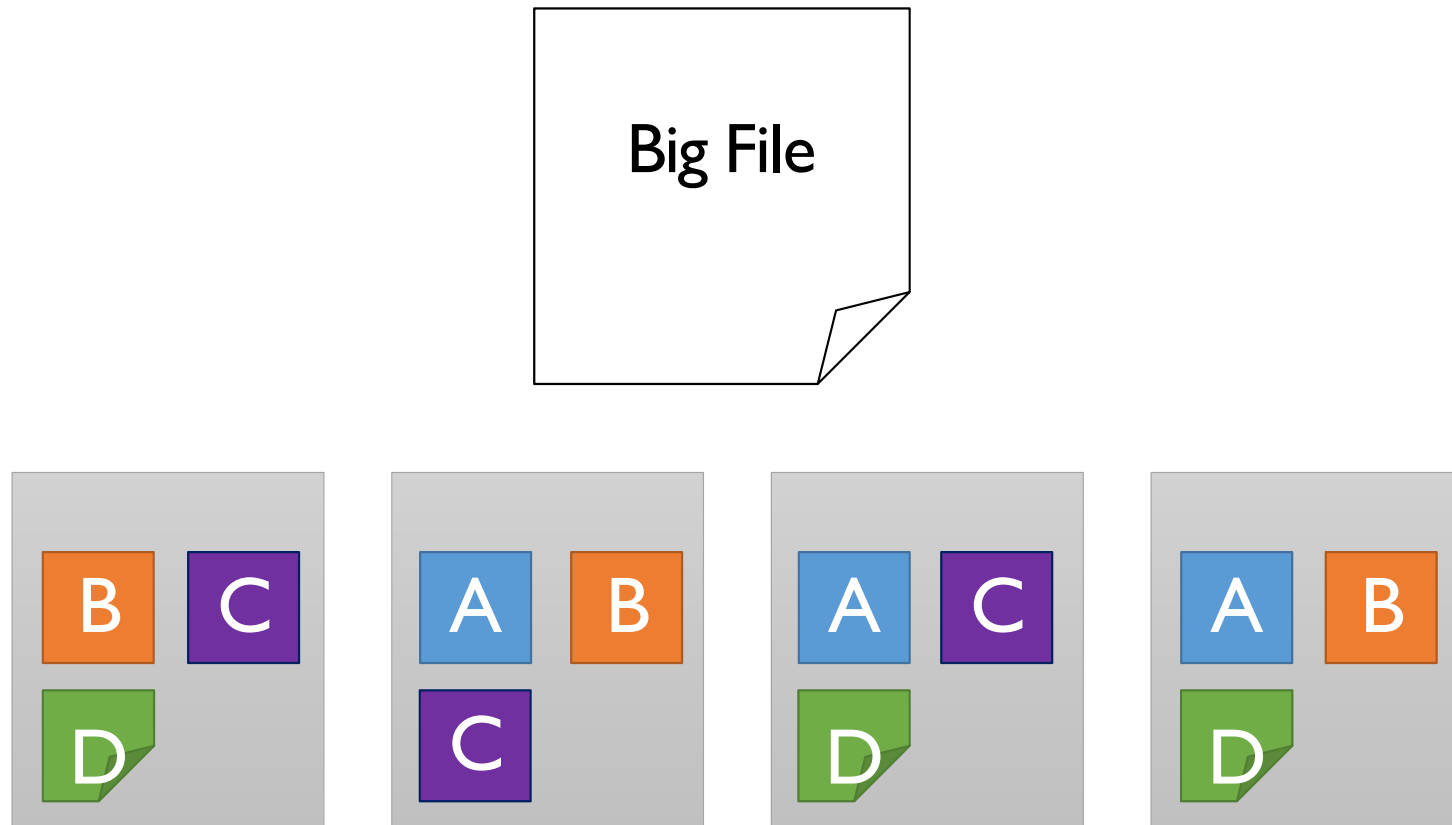
Fault Tolerant Distributed File Systems (e.g., HDFS)

- Split large files over multiple machines
 - Easily support very massive files spanning
- Read parts of file in parallel
 - Fast reads of large files
- Often built using cheap commodity storage devices

Cheap commodity storage devices will fail!



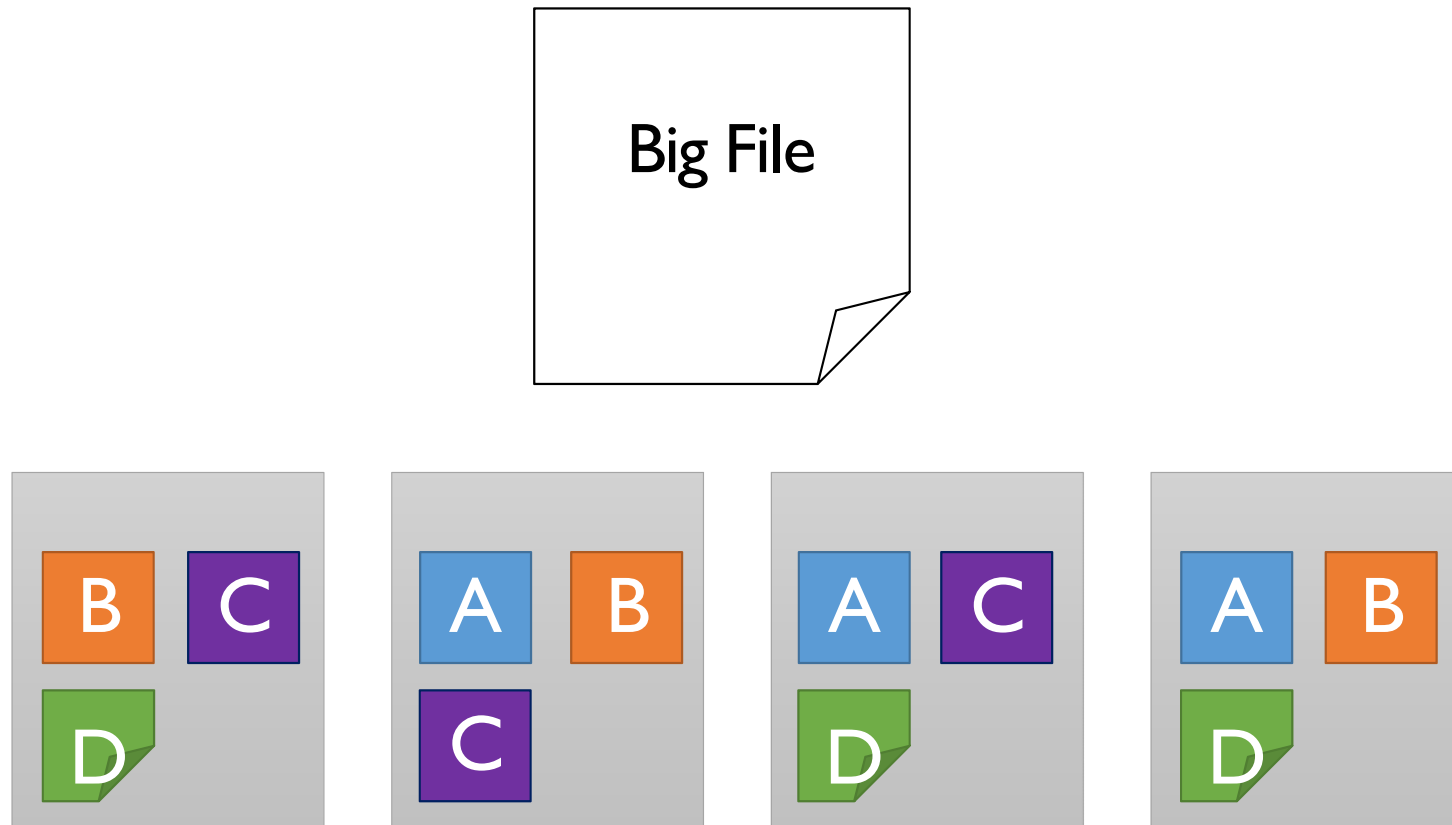
Fault Tolerant Distributed File Systems (e.g., HDFS)



[Ghemawat et al., SOSP'03]

Fault Tolerant Distributed File Systems

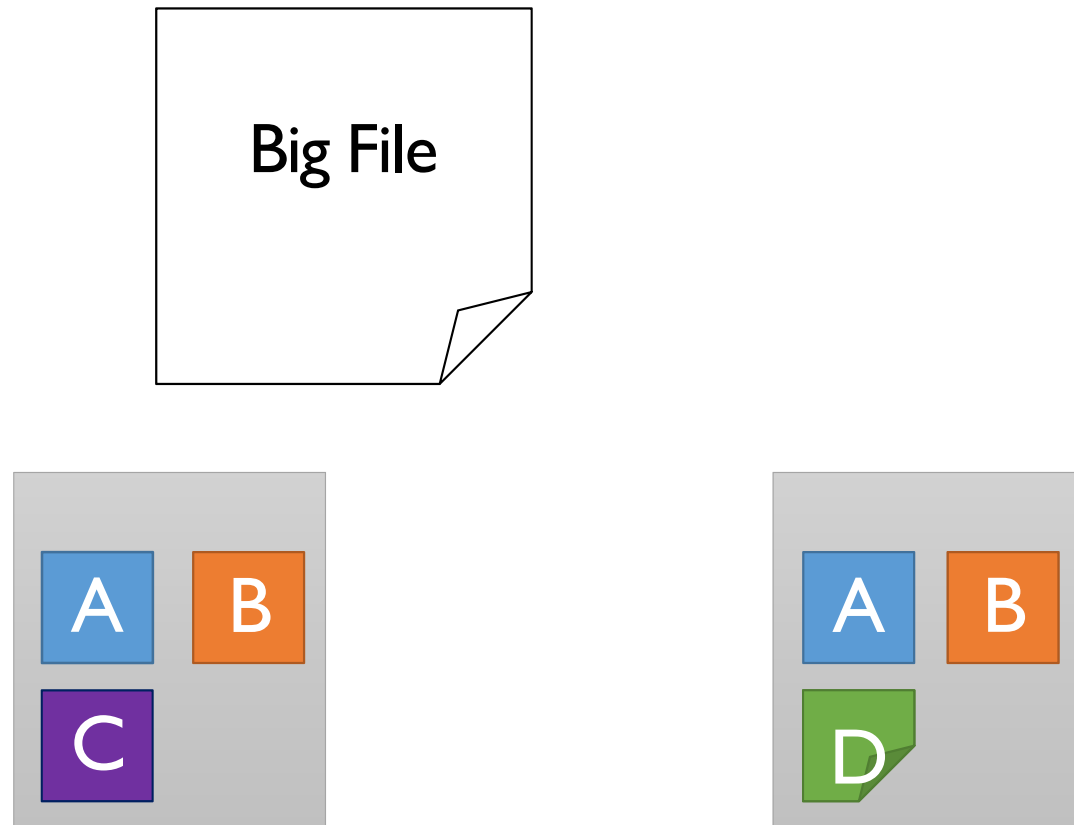
Failure Event



[Ghemawat et al., SOSP'03]

Fault Tolerant Distributed File Systems

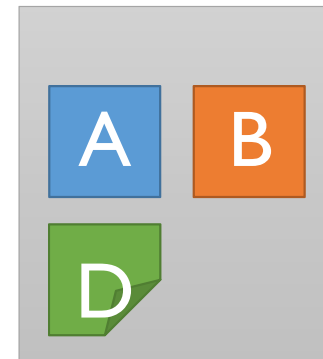
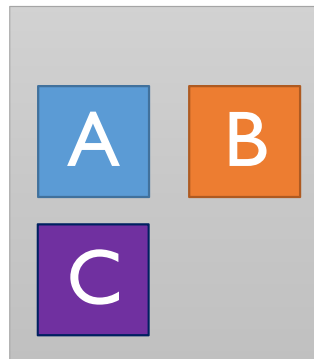
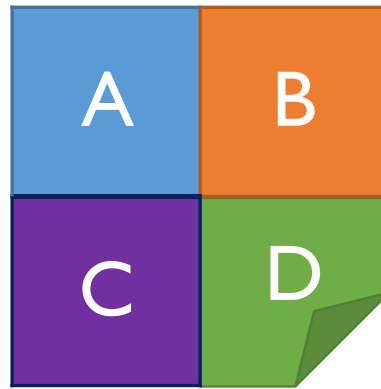
Failure Event



[Ghemawat et al., SOSP'03]

Fault Tolerant Distributed File Systems

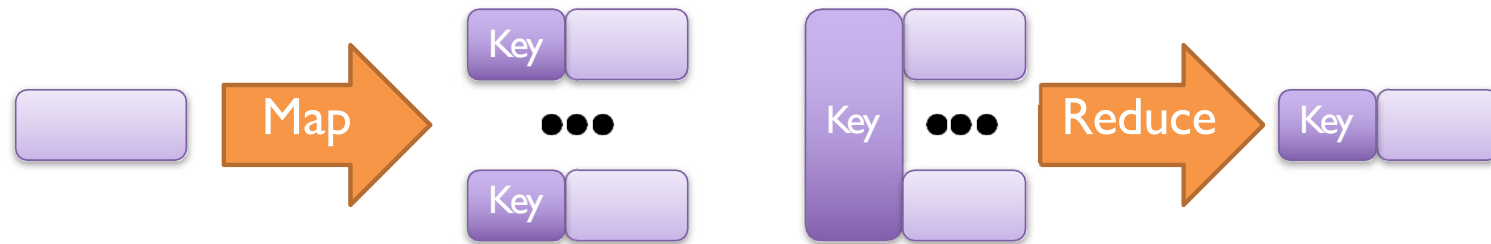
Failure Event



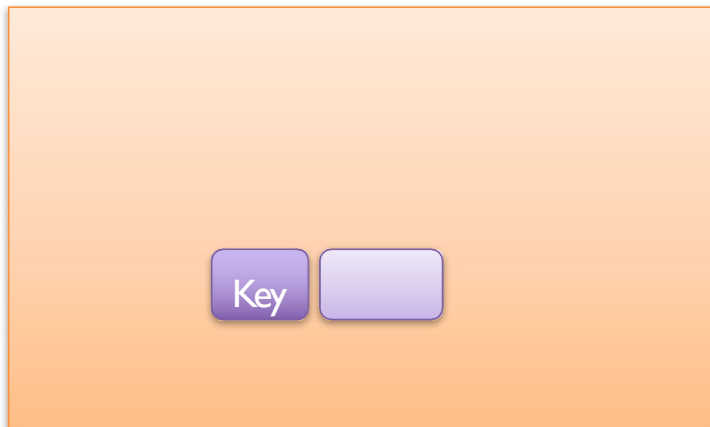
[Ghemawat et al., SOSP'03]

Map-Reduce

The Map Reduce Abstraction



Example: *Word-Count*



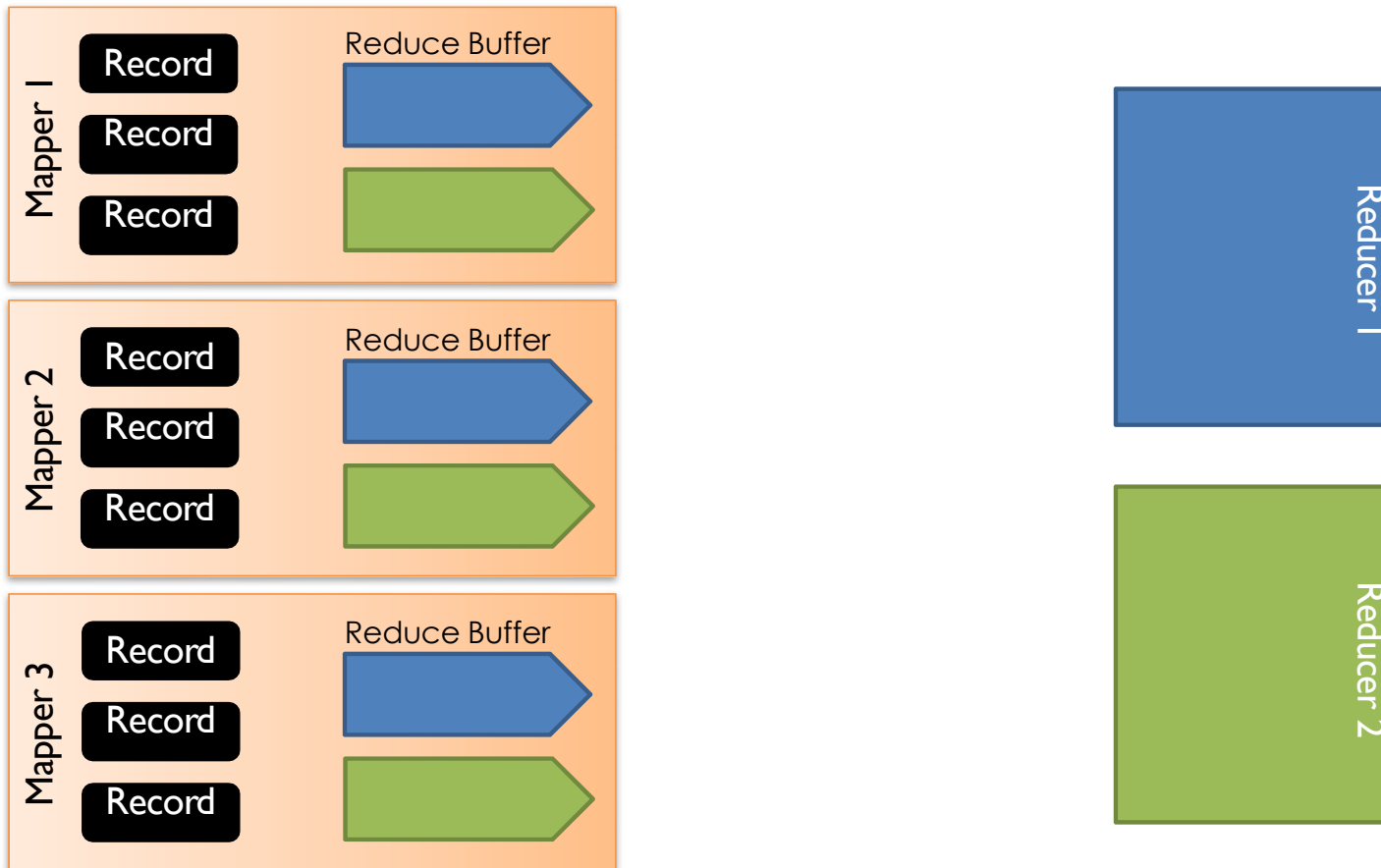
```
Reduce(word, counts) {  
  emit (word, SUM(counts))  
}
```

The Map Reduce System



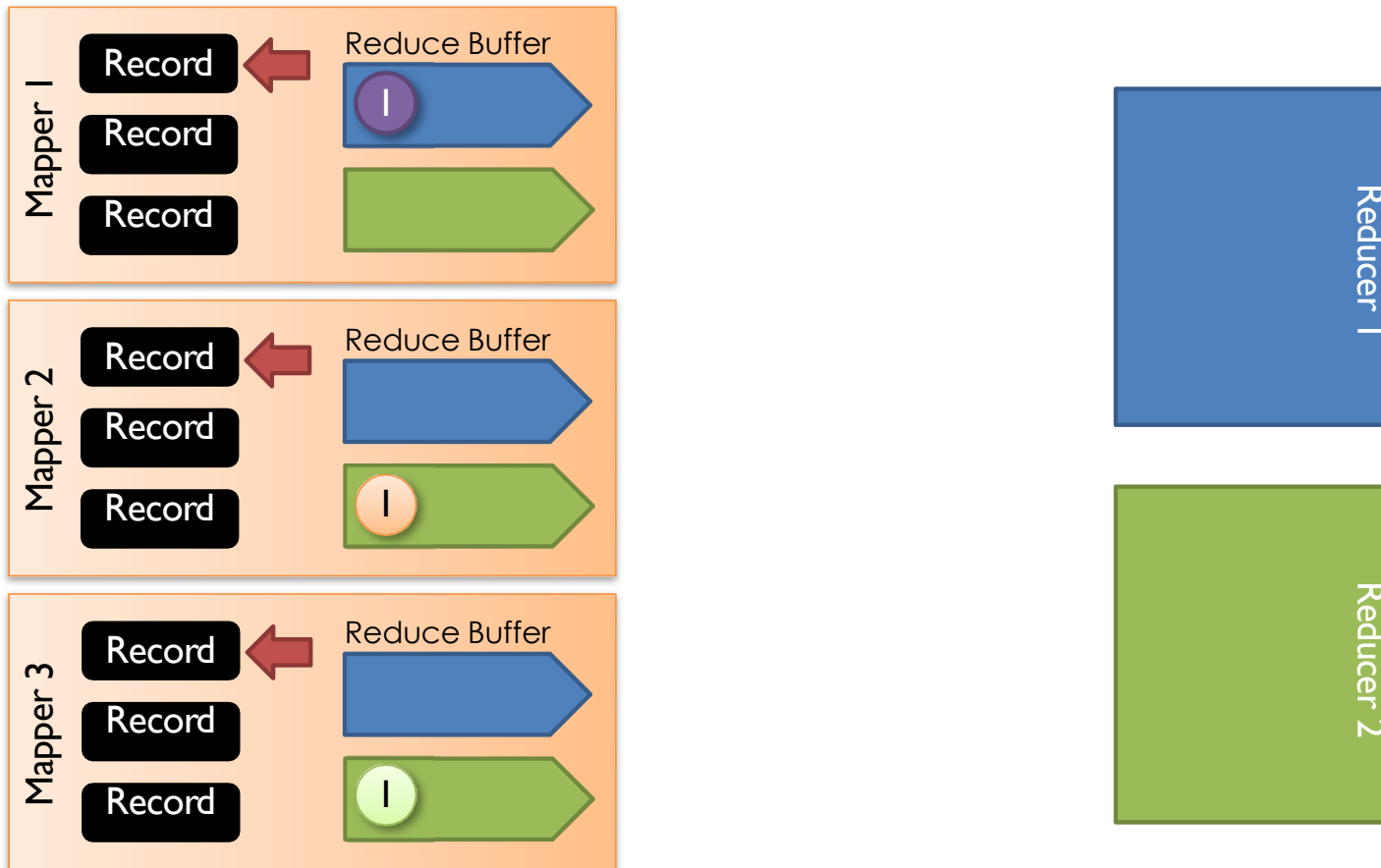
[Dean & Ghemawat, OSDI'04]

The Map Reduce System



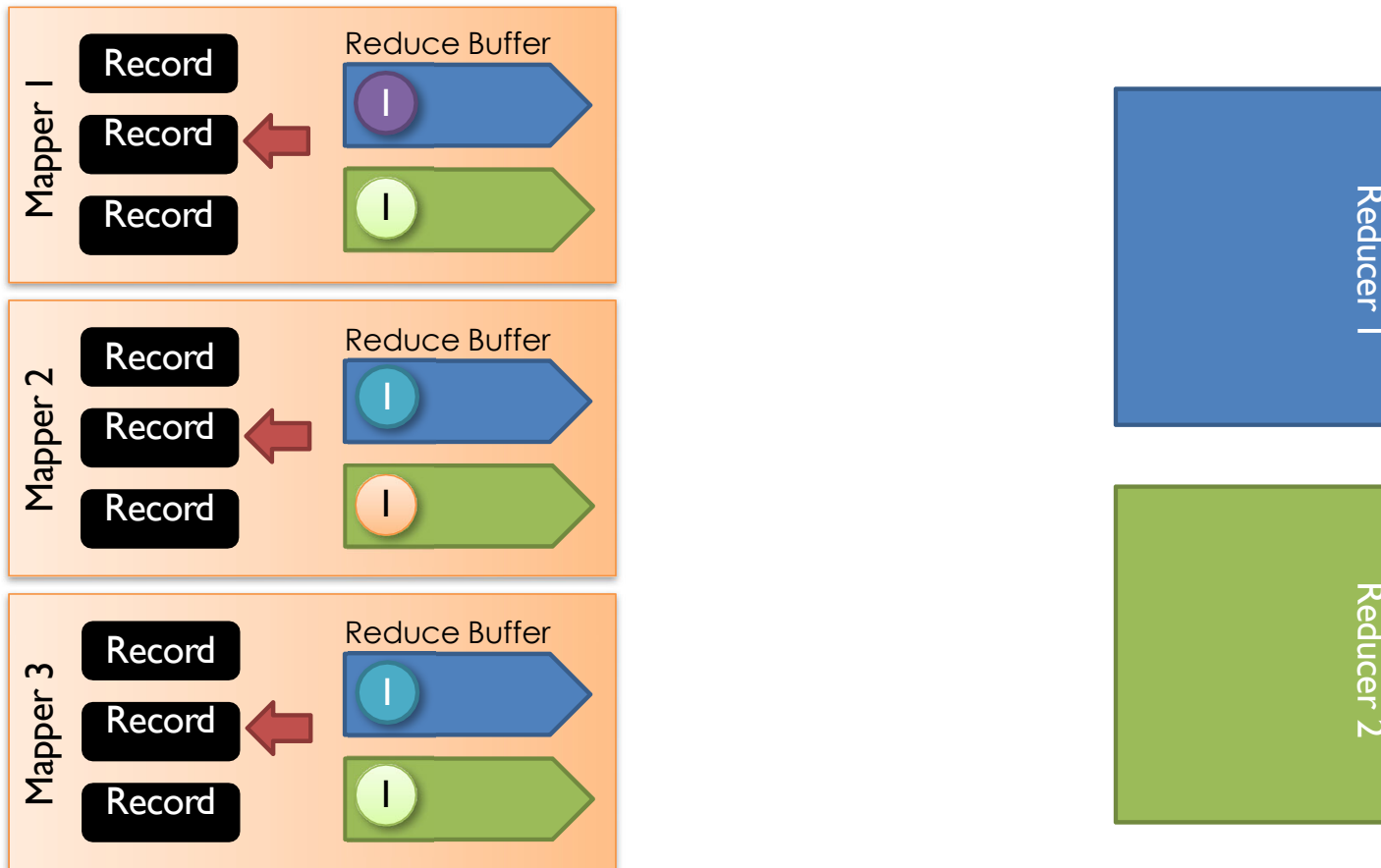
[Dean & Ghemawat, OSDI'04]

The Map Reduce System



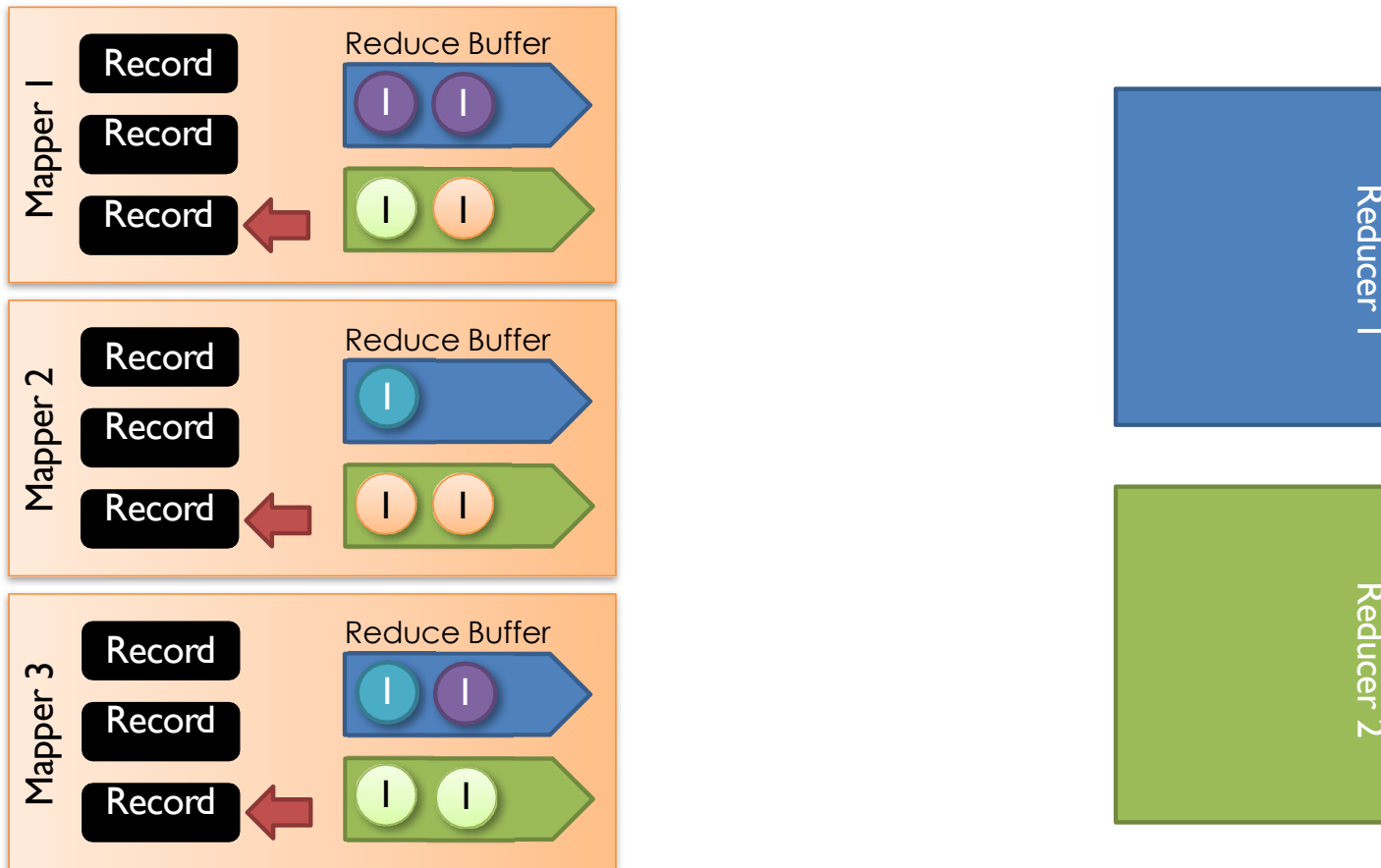
[Dean & Ghemawat, OSDI'04]

The Map Reduce System



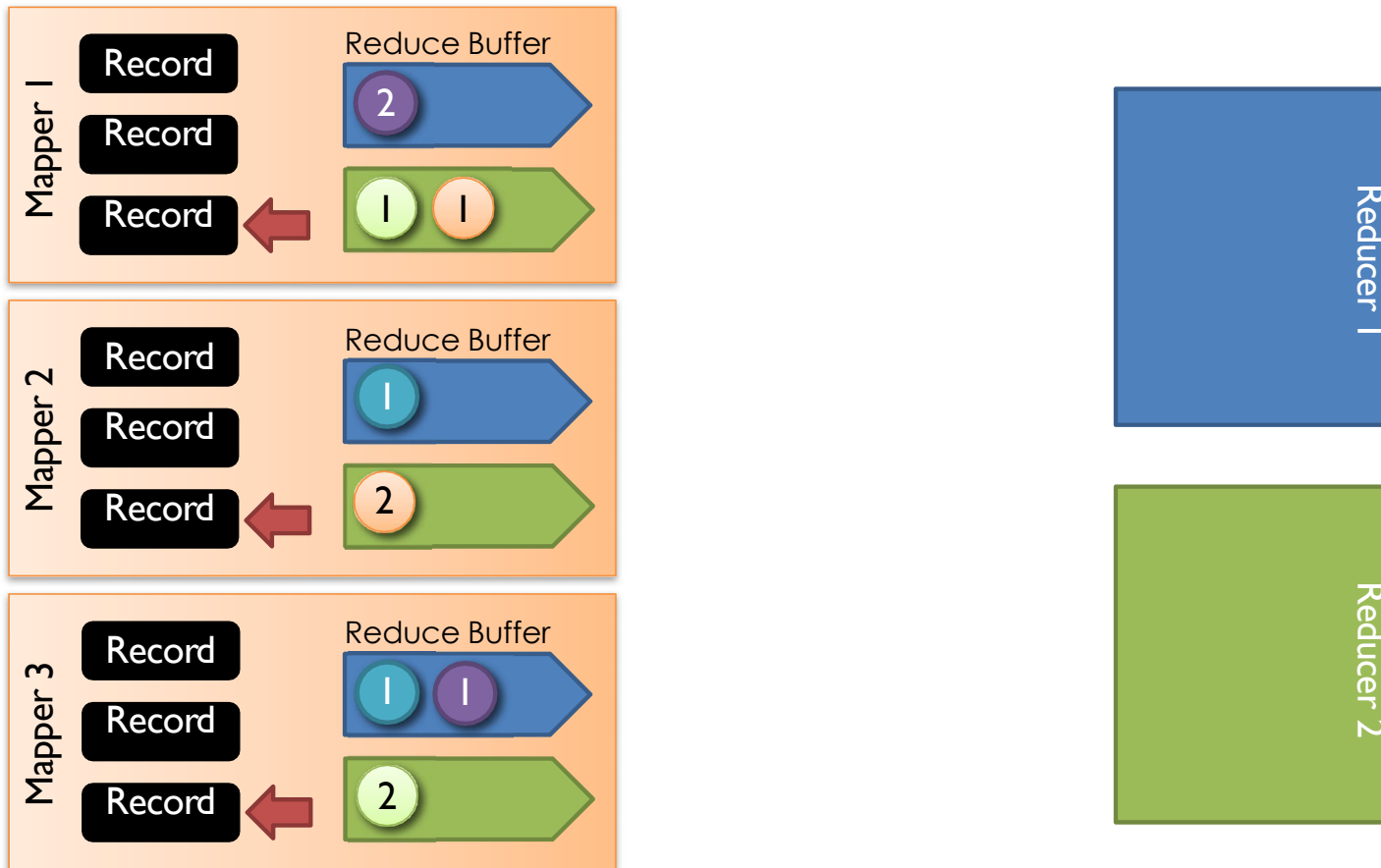
[Dean & Ghemawat, OSDI'04]

The Map Reduce System



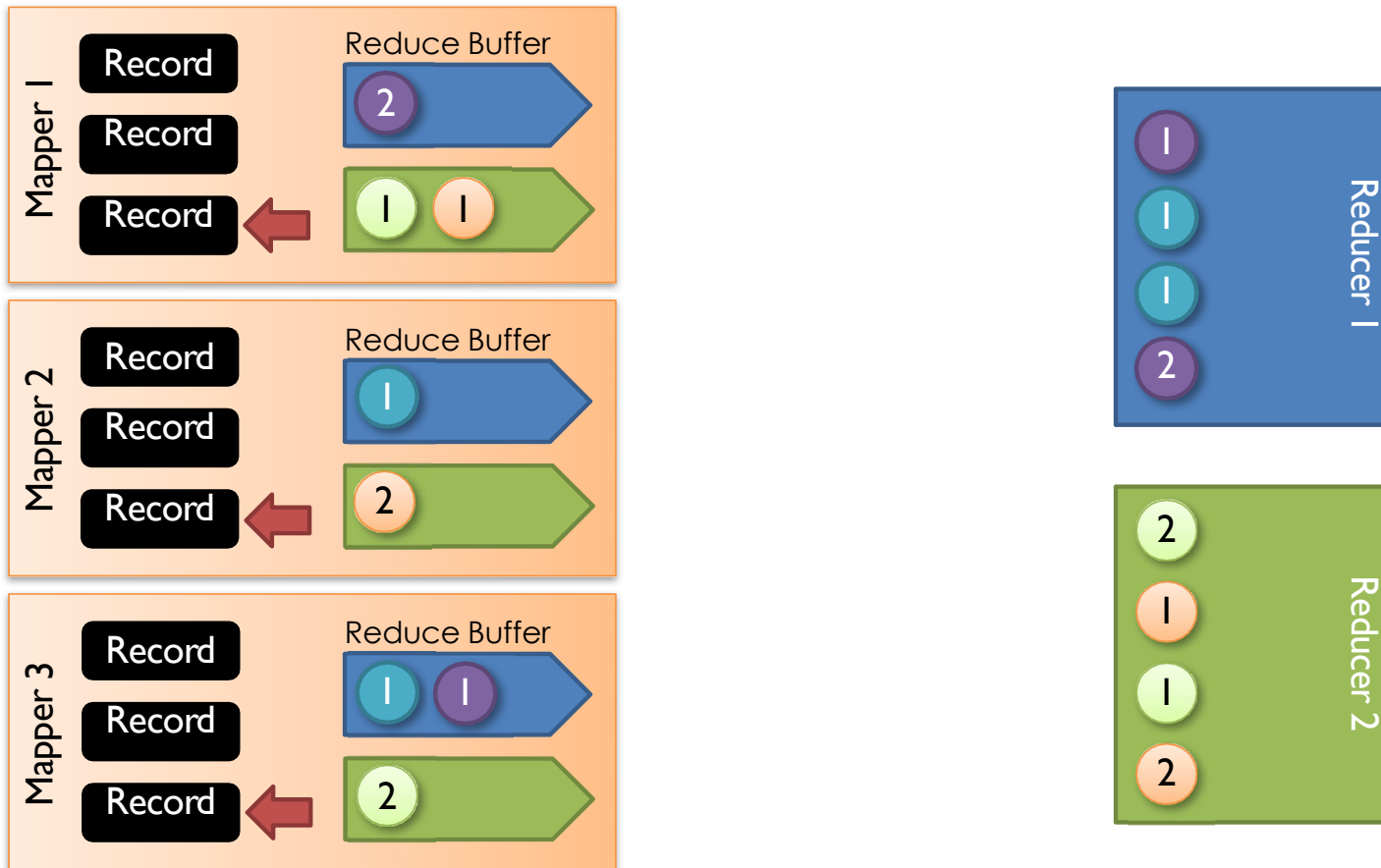
[Dean & Ghemawat, OSDI'04]

The Map Reduce System



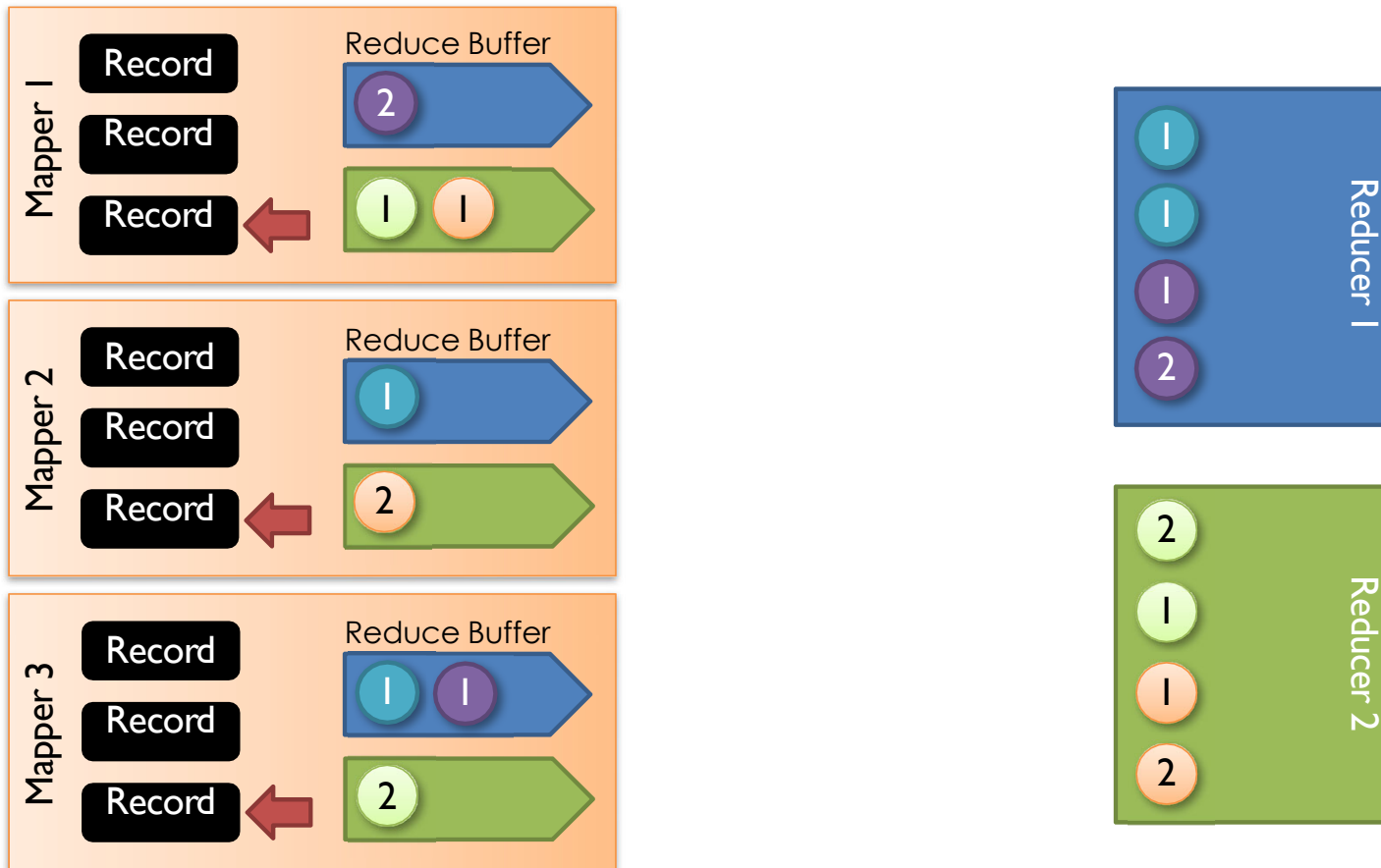
[Dean & Ghemawat, OSDI'04]

The Map Reduce System



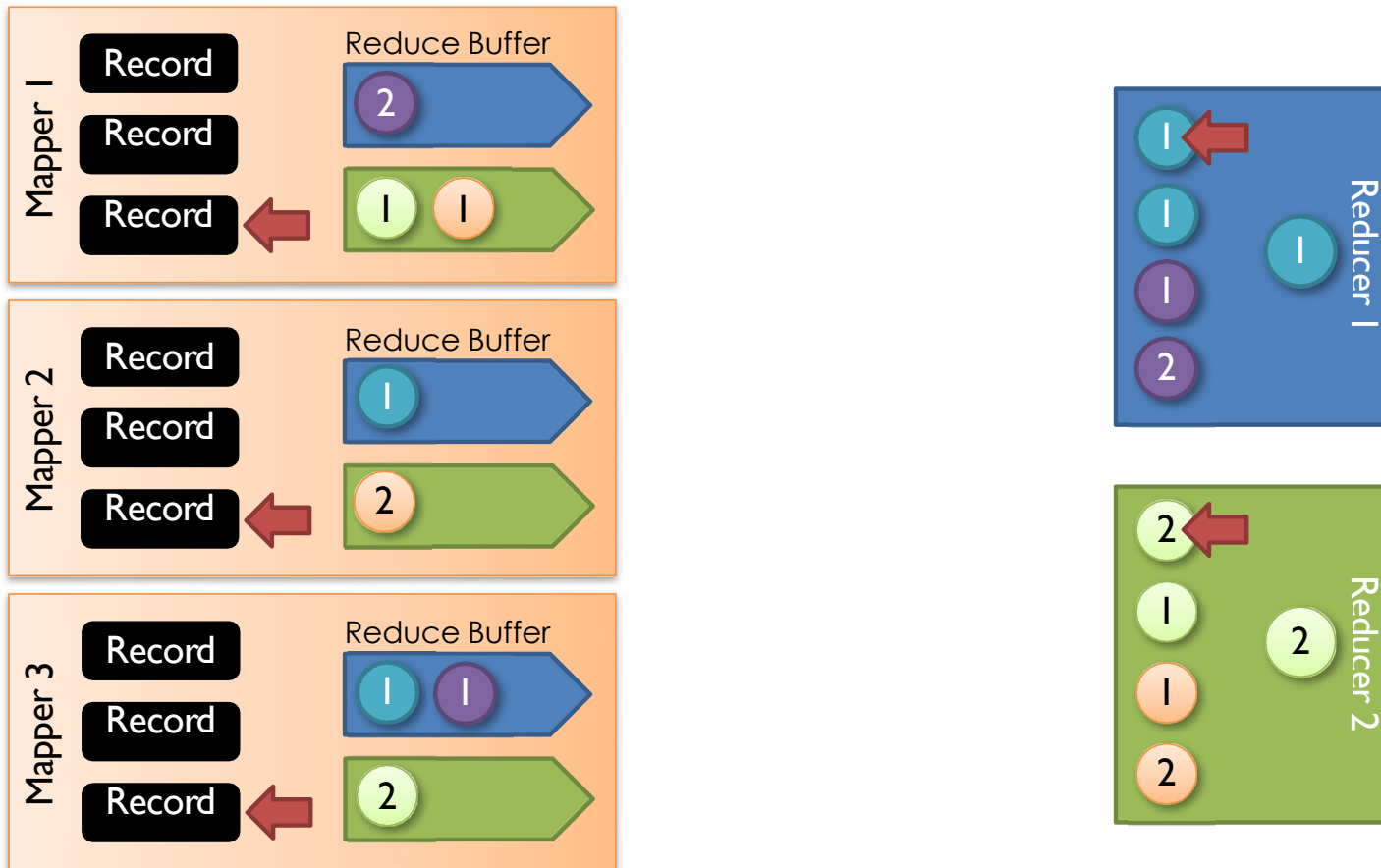
[Dean & Ghemawat, OSDI'04]

The Map Reduce System



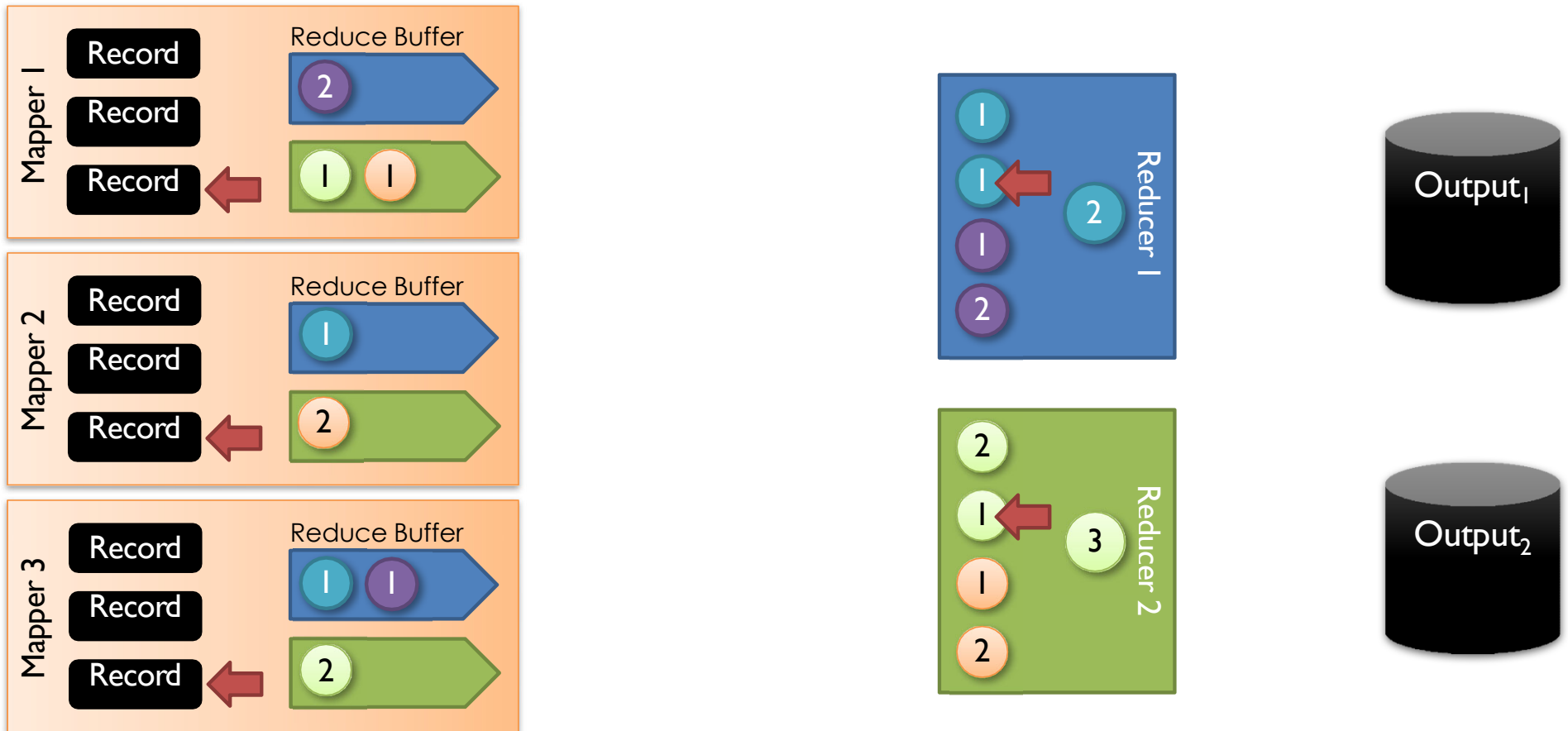
[Dean & Ghemawat, OSDI'04]

The Map Reduce System



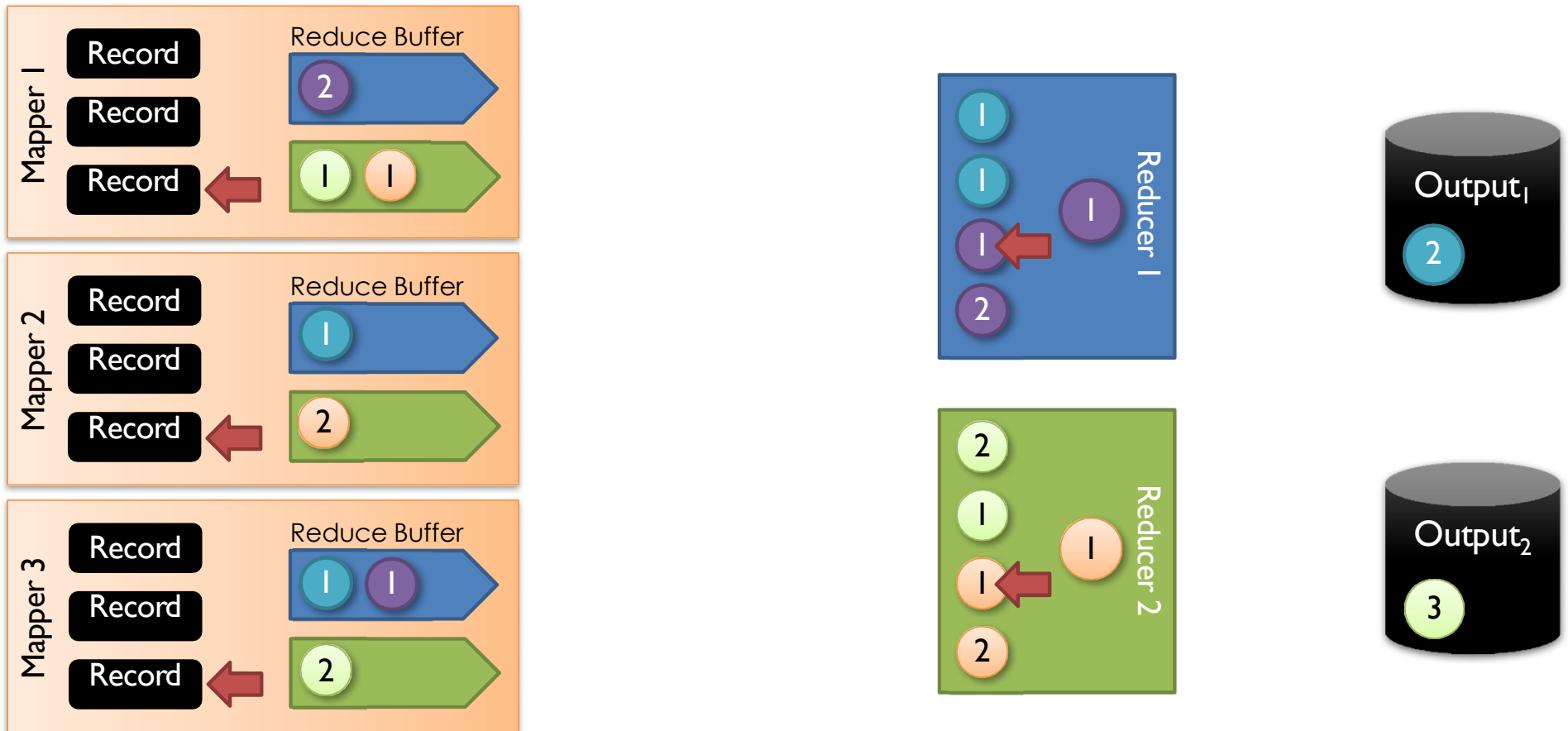
[Dean & Ghemawat, OSDI'04]

The Map Reduce System



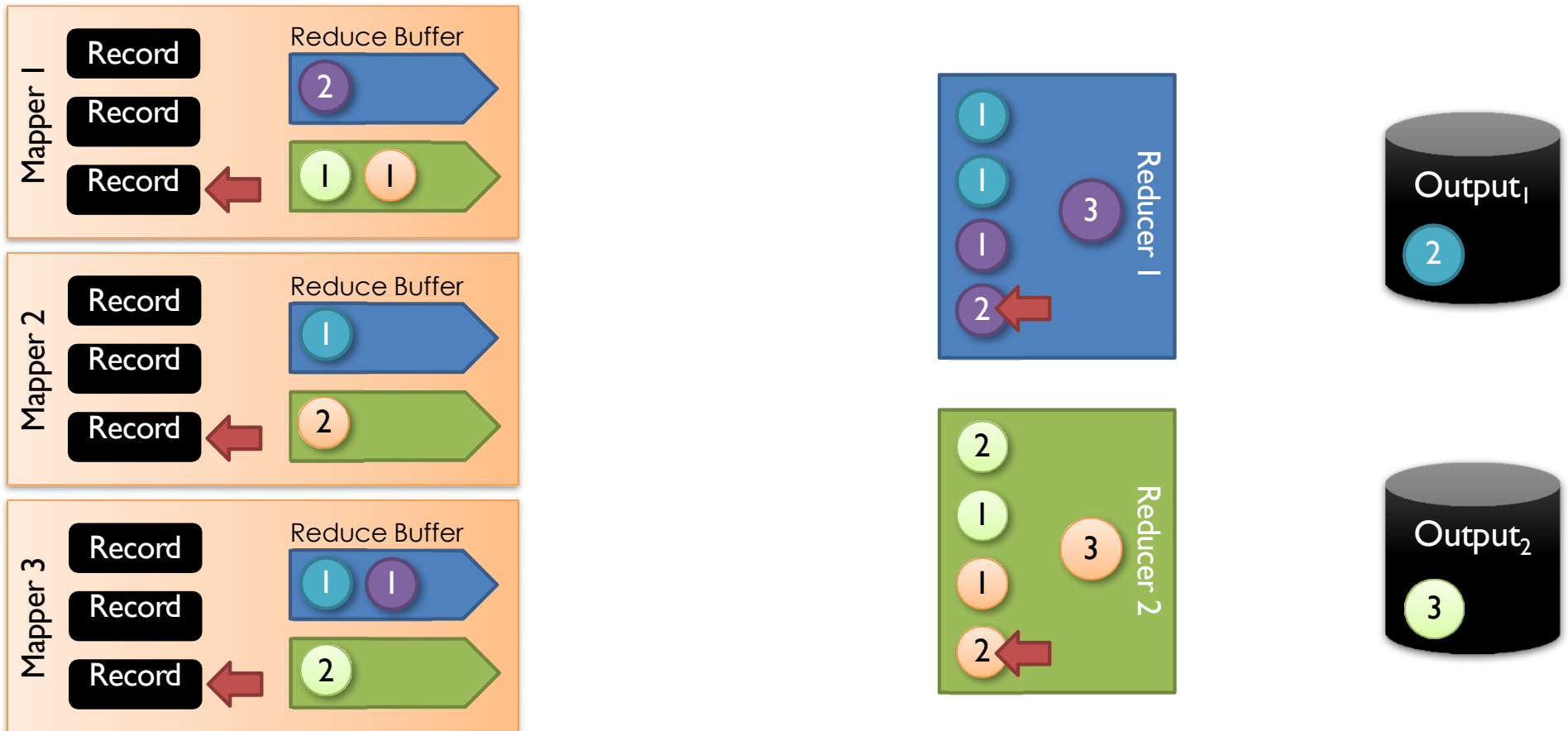
[Dean & Ghemawat, OSDI'04]

The Map Reduce System



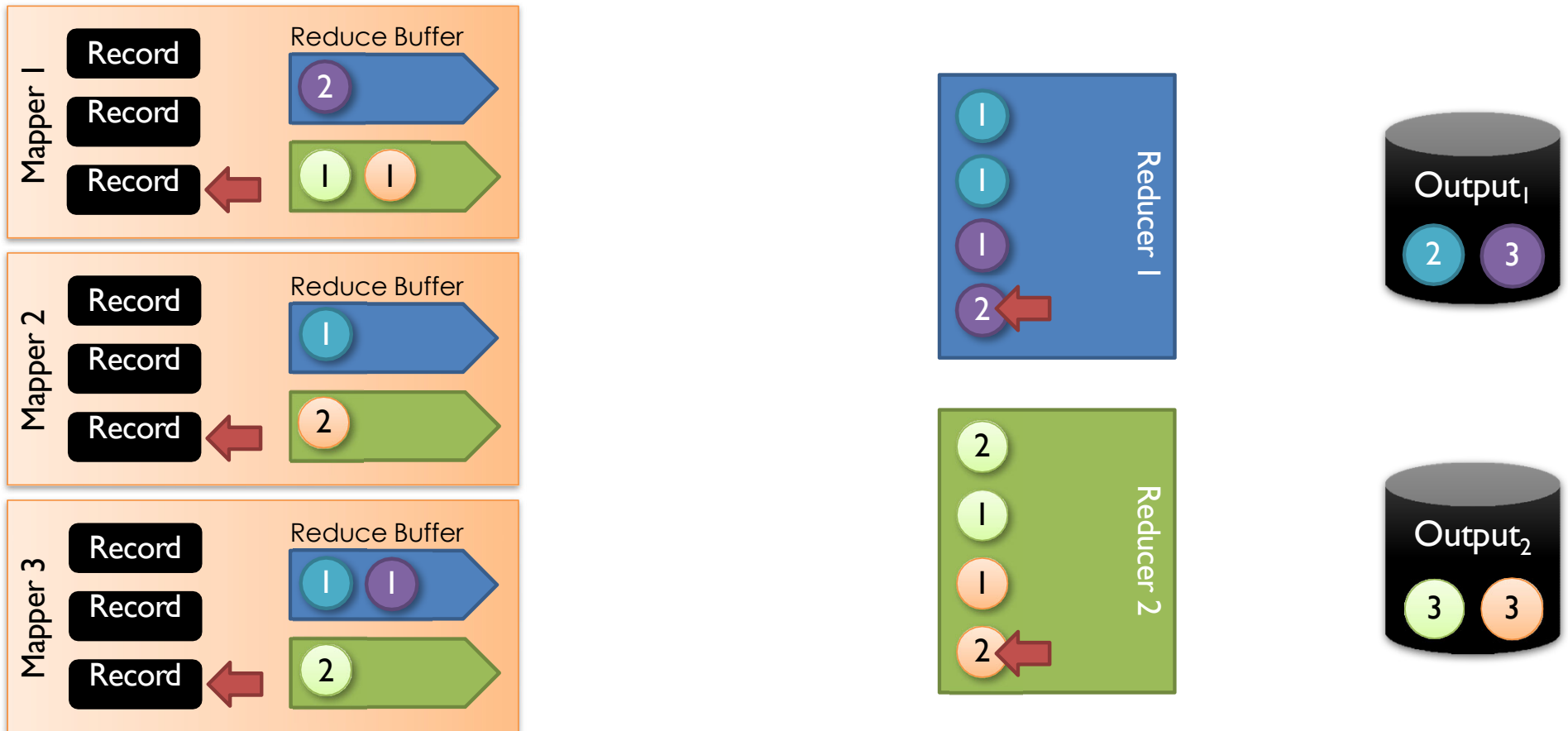
[Dean & Ghemawat, OSDI'04]

The Map Reduce System



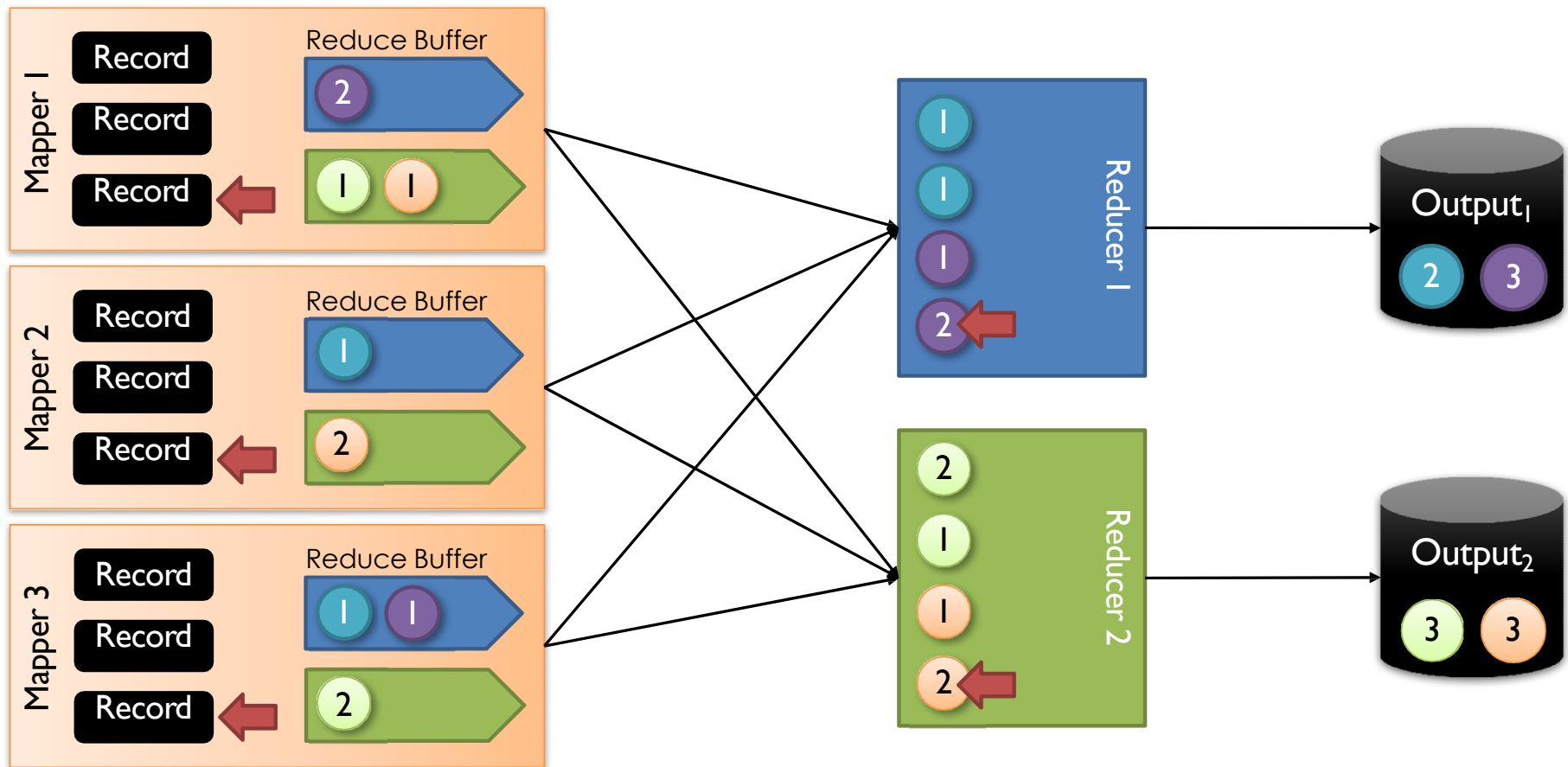
[Dean & Ghemawat, OSDI'04]

The Map Reduce System



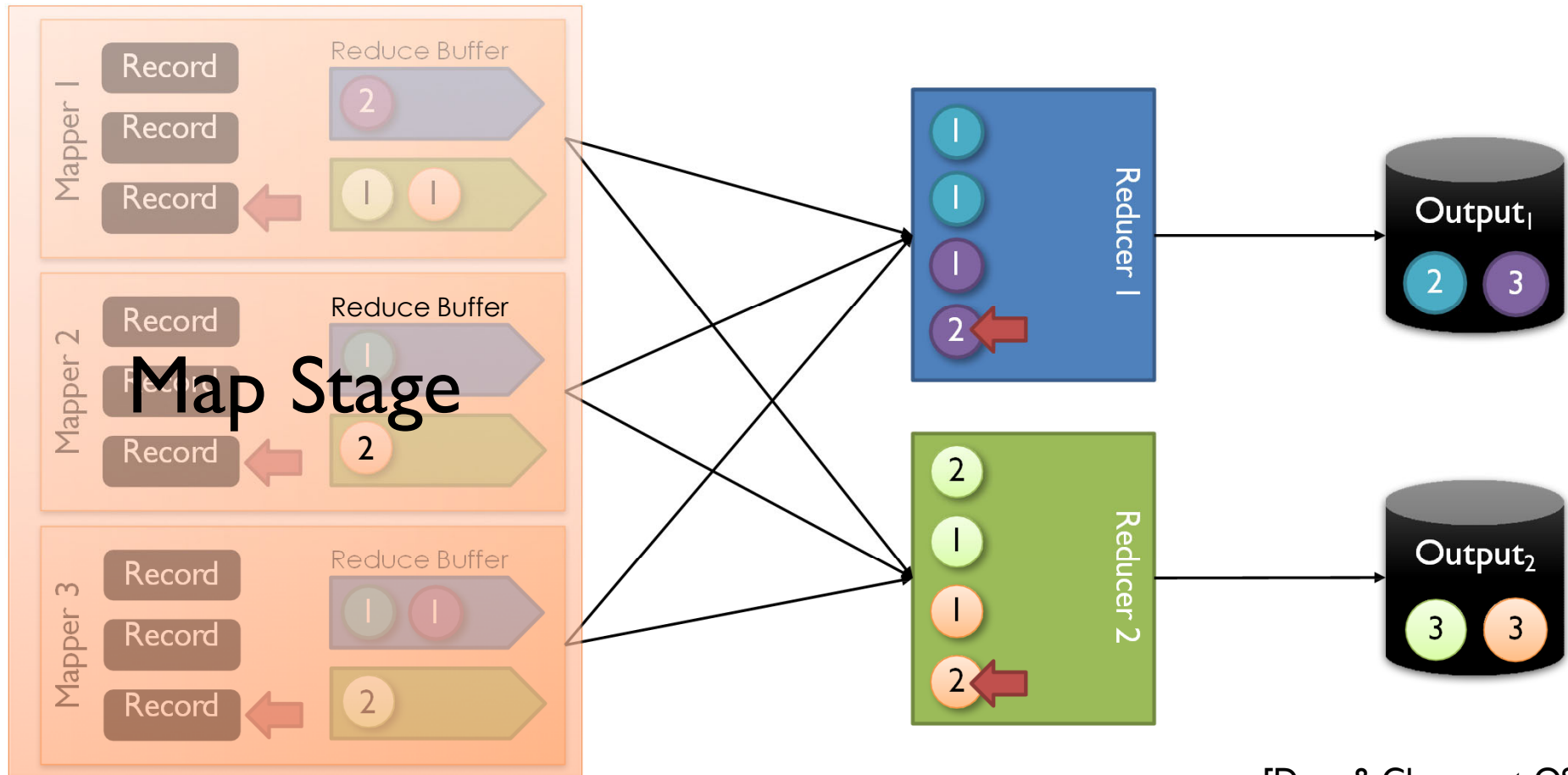
[Dean & Ghemawat, OSDI'04]

The Map Reduce System



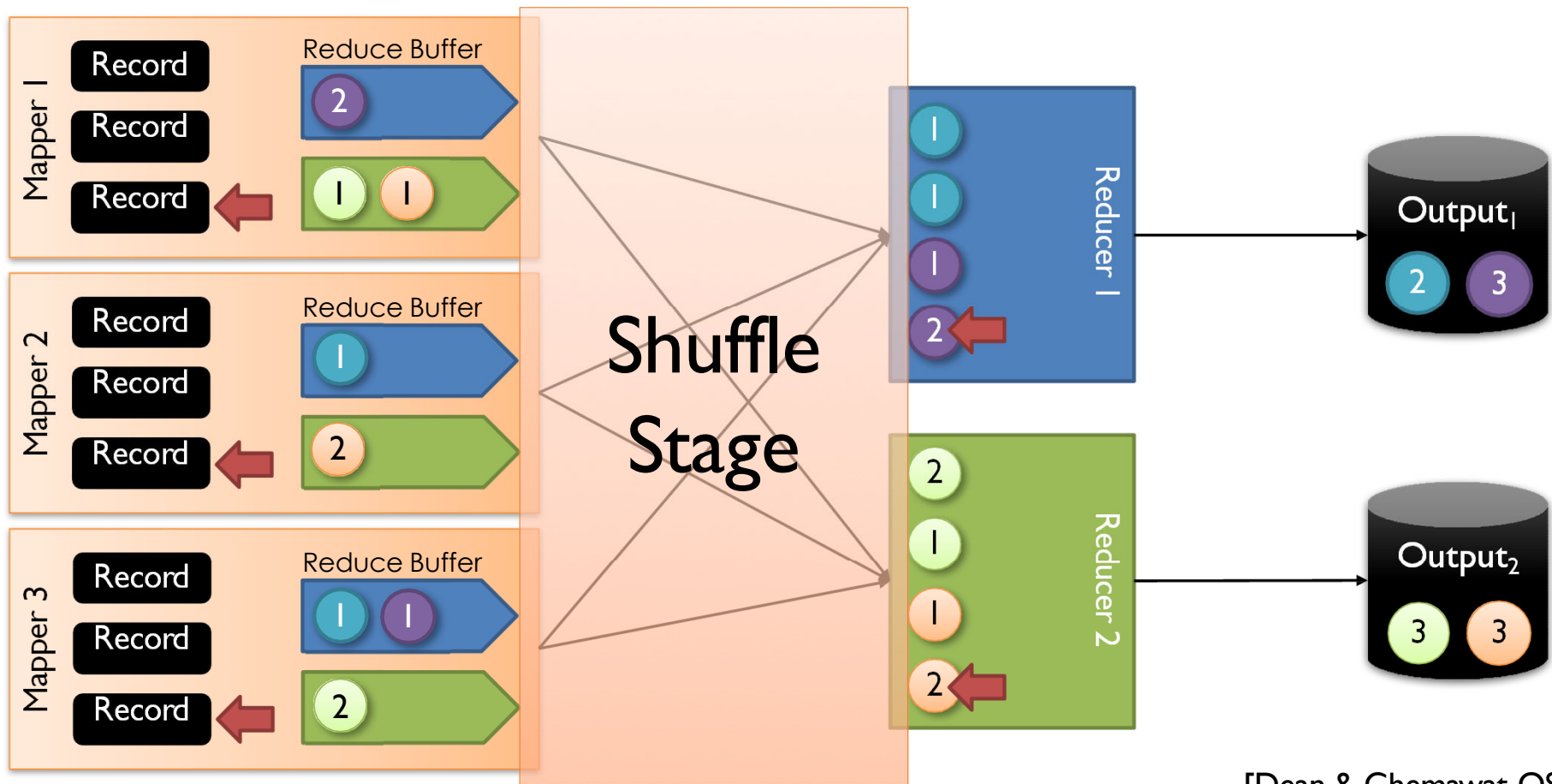
[Dean & Ghemawat, OSDI'04]

The Map Reduce System



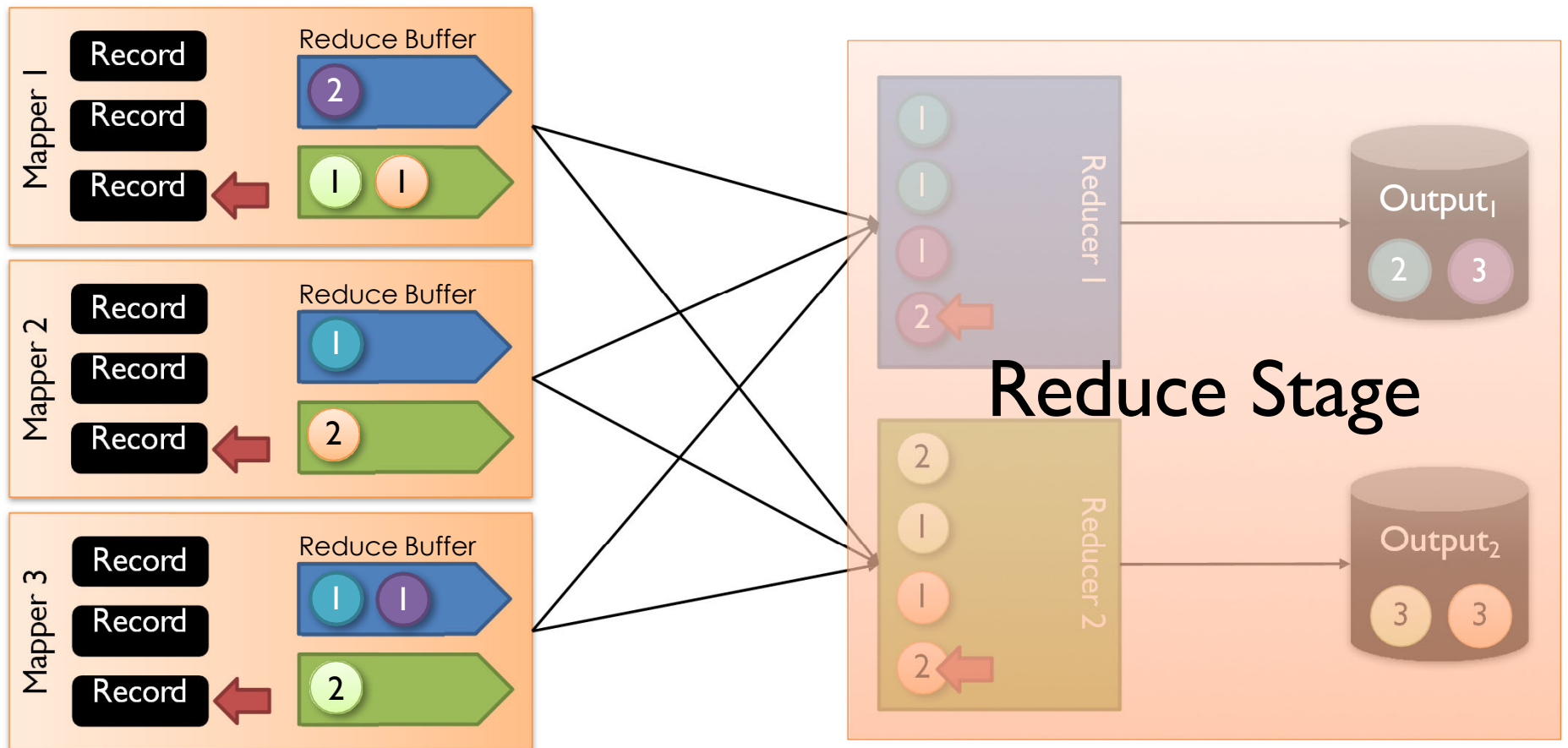
[Dean & Ghemawat, OSDI'04]

The Map Reduce System



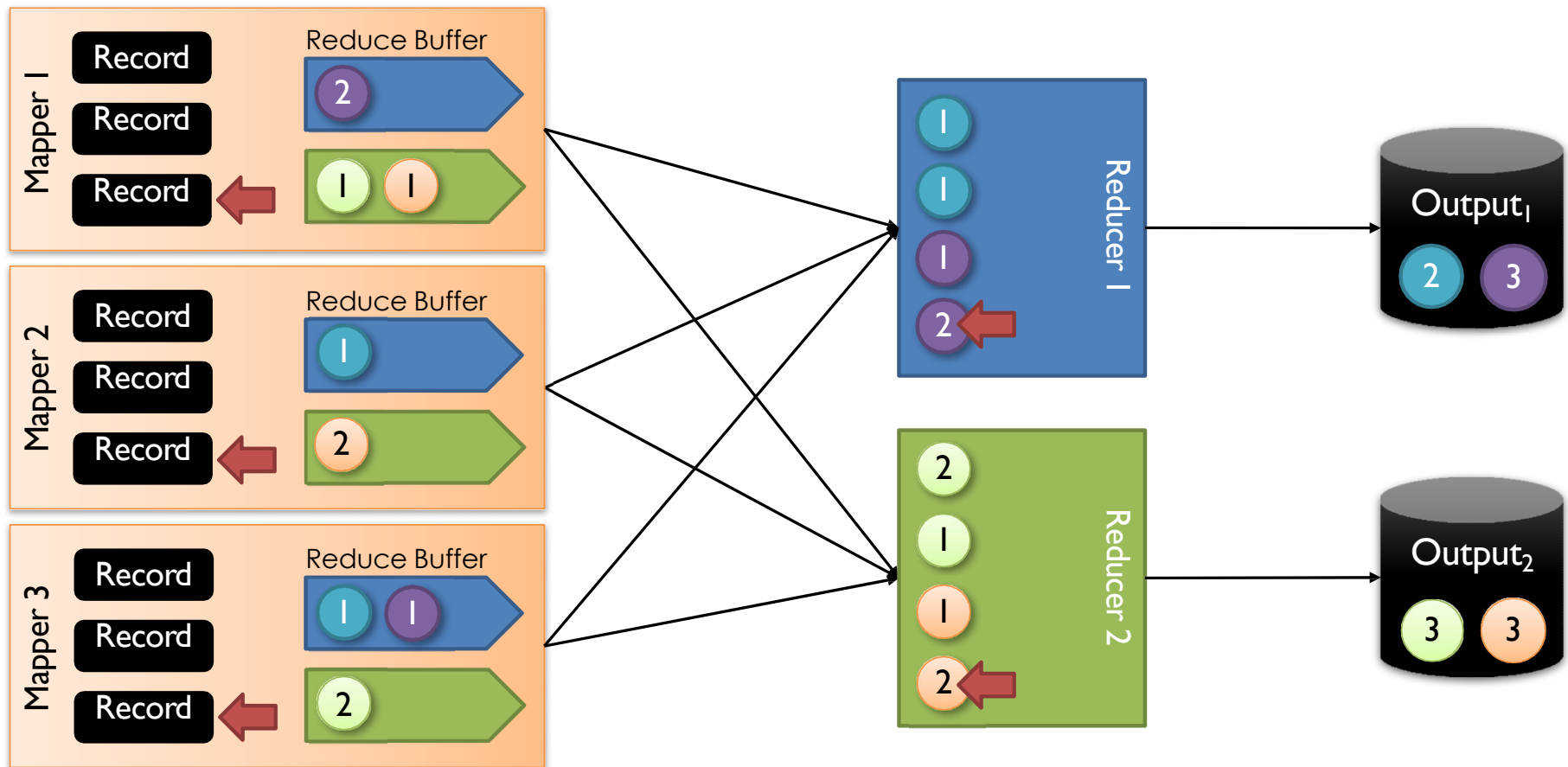
[Dean & Ghemawat, OSDI'04]

The Map Reduce System



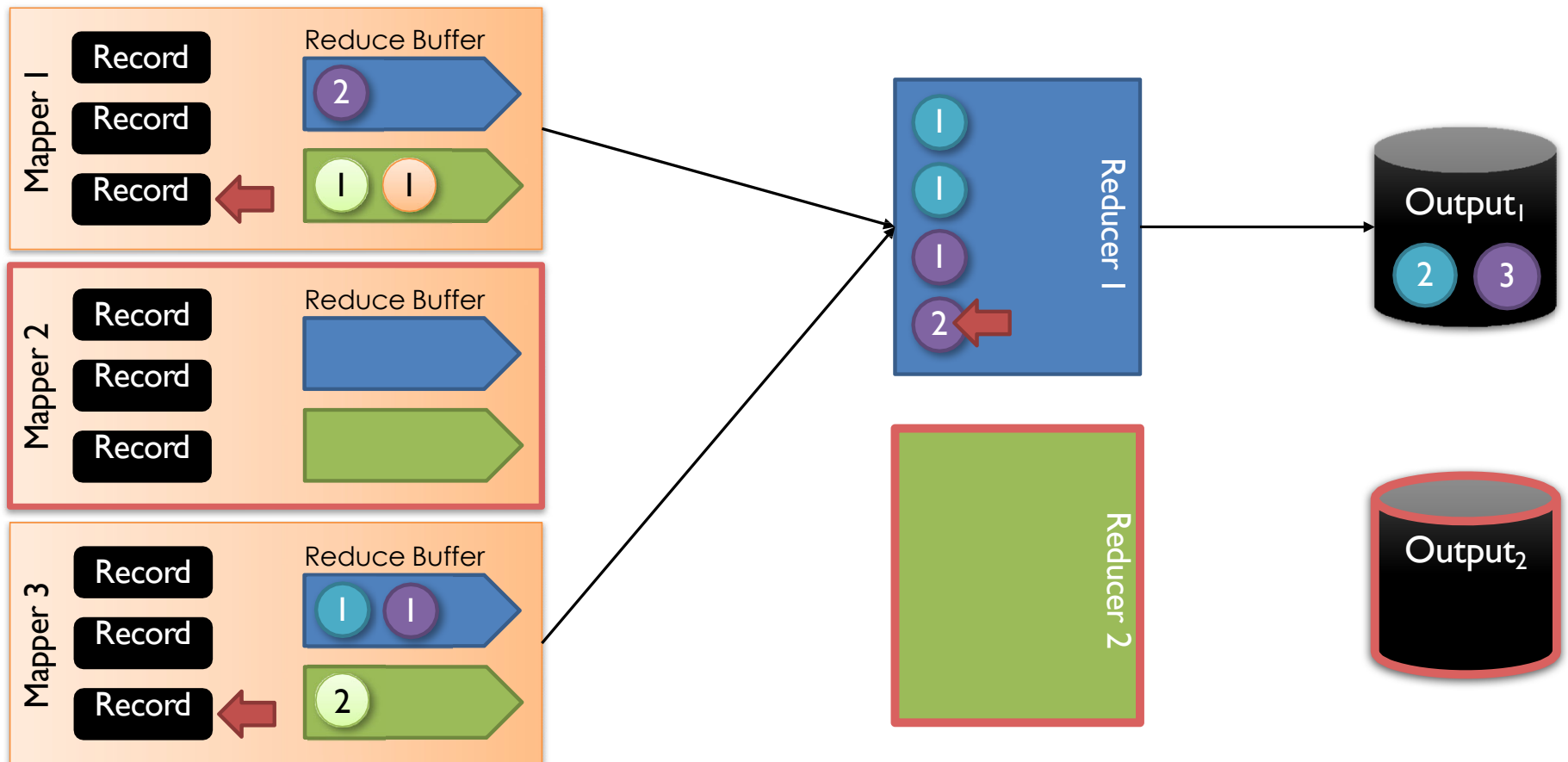
[Dean & Ghemawat, OSDI'04]

The Map Reduce: **Fault Tolerance**



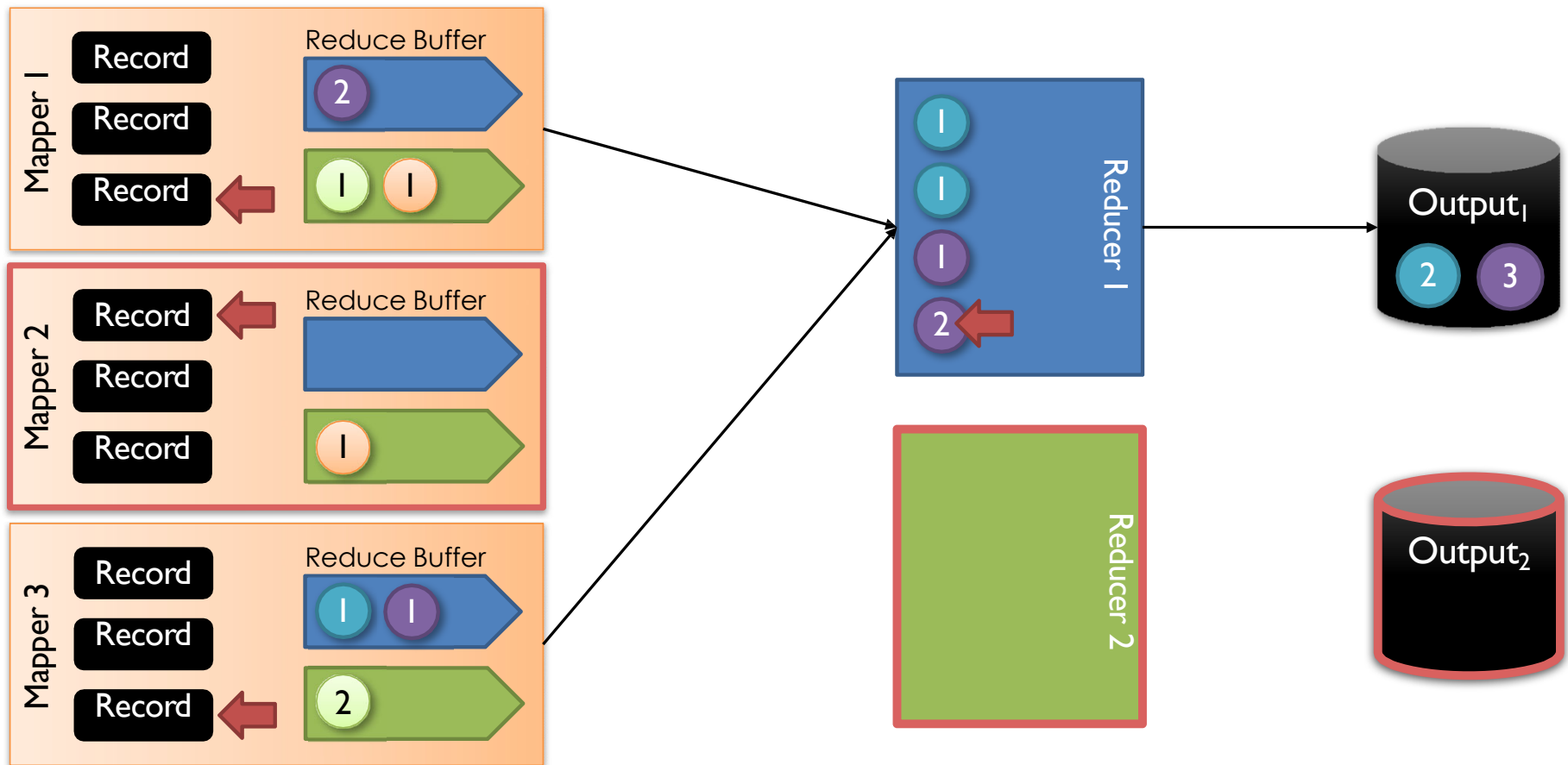
[Dean & Ghemawat, OSDI'04]

The Map Reduce: **Fault Tolerance**



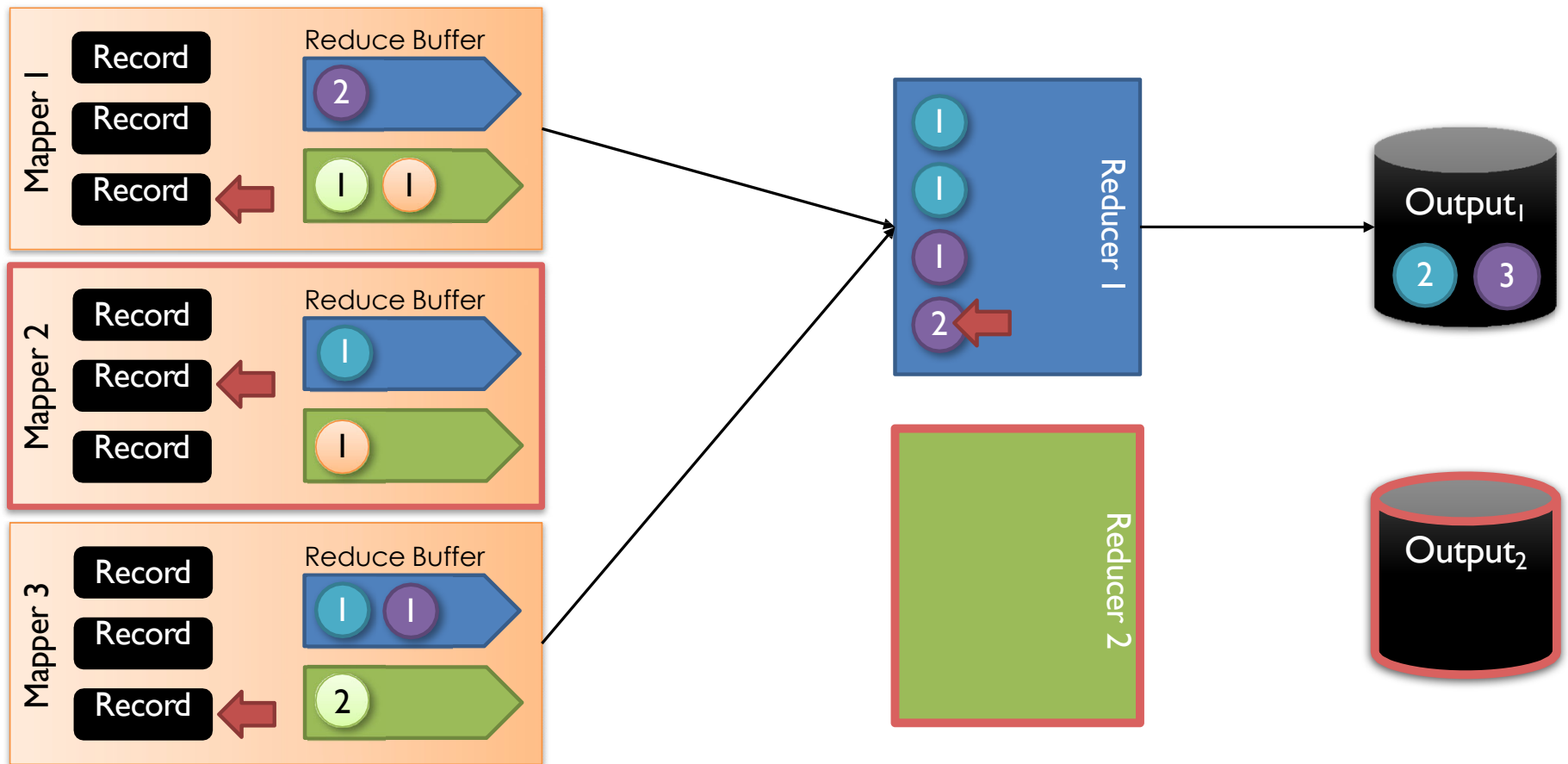
[Dean & Ghemawat, OSDI'04]

The Map Reduce: **Fault Tolerance**



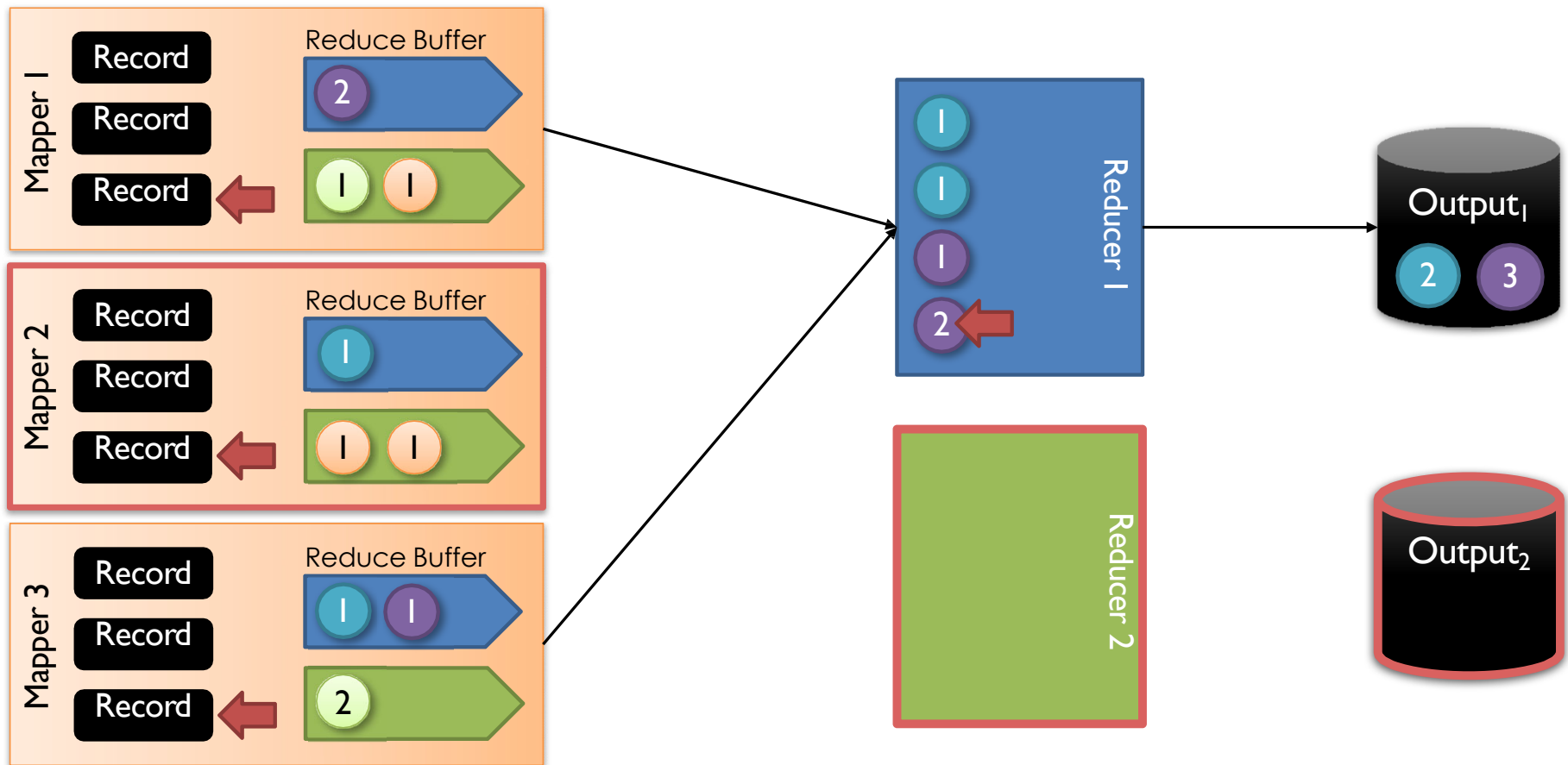
[Dean & Ghemawat, OSDI'04]

The Map Reduce: **Fault Tolerance**



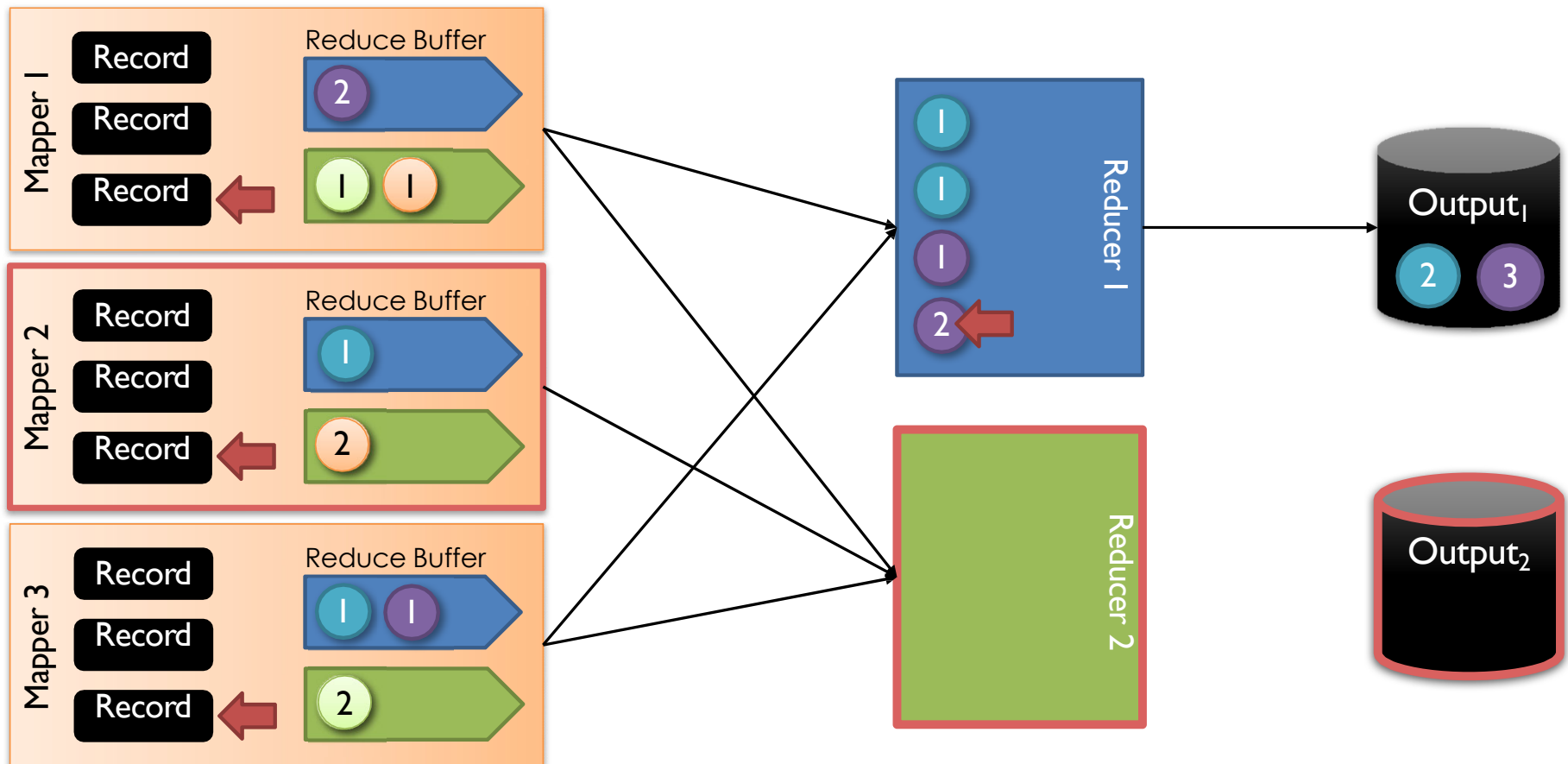
[Dean & Ghemawat, OSDI'04]

The Map Reduce: **Fault Tolerance**



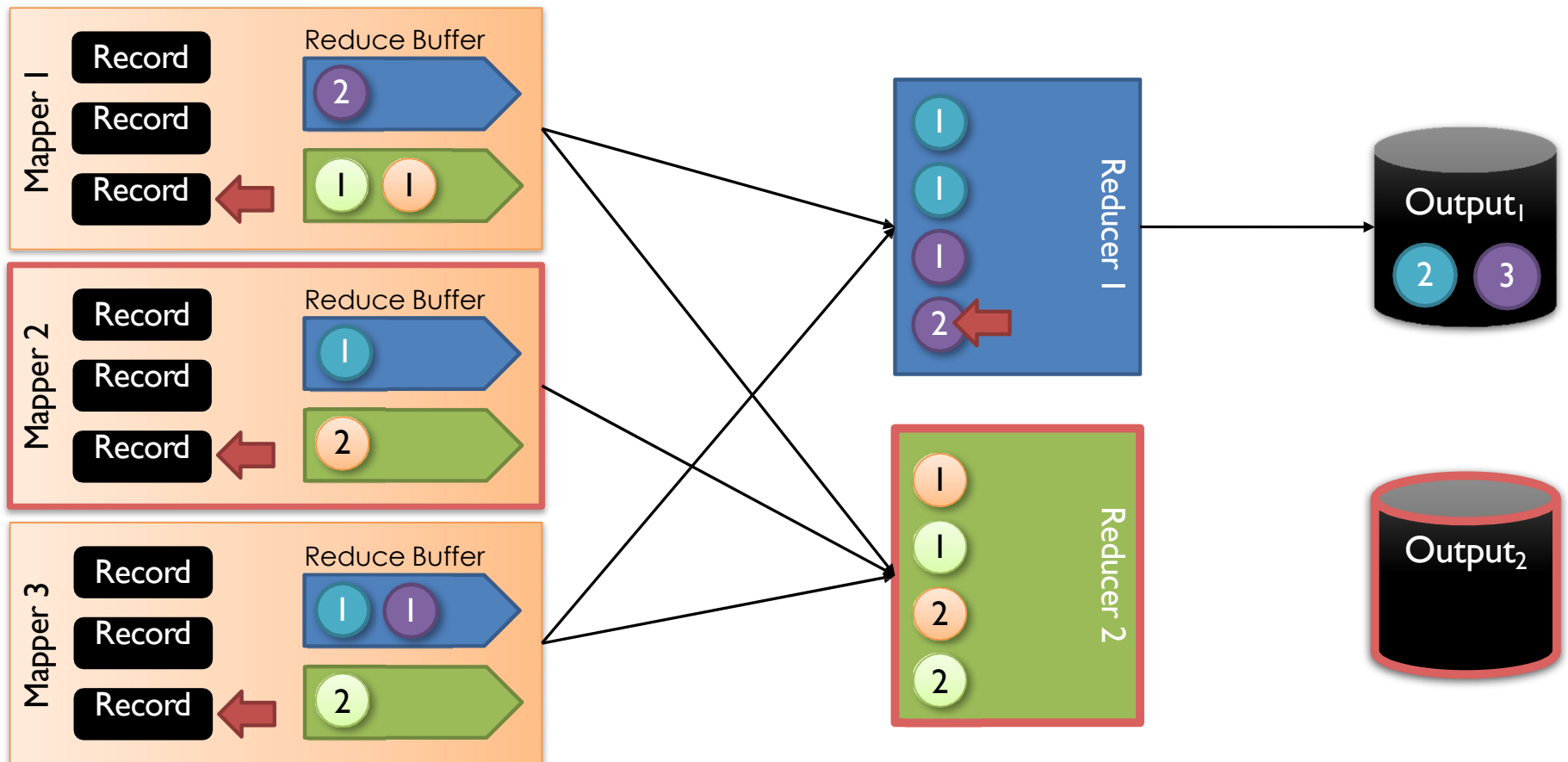
[Dean & Ghemawat, OSDI'04]

The Map Reduce: **Fault Tolerance**



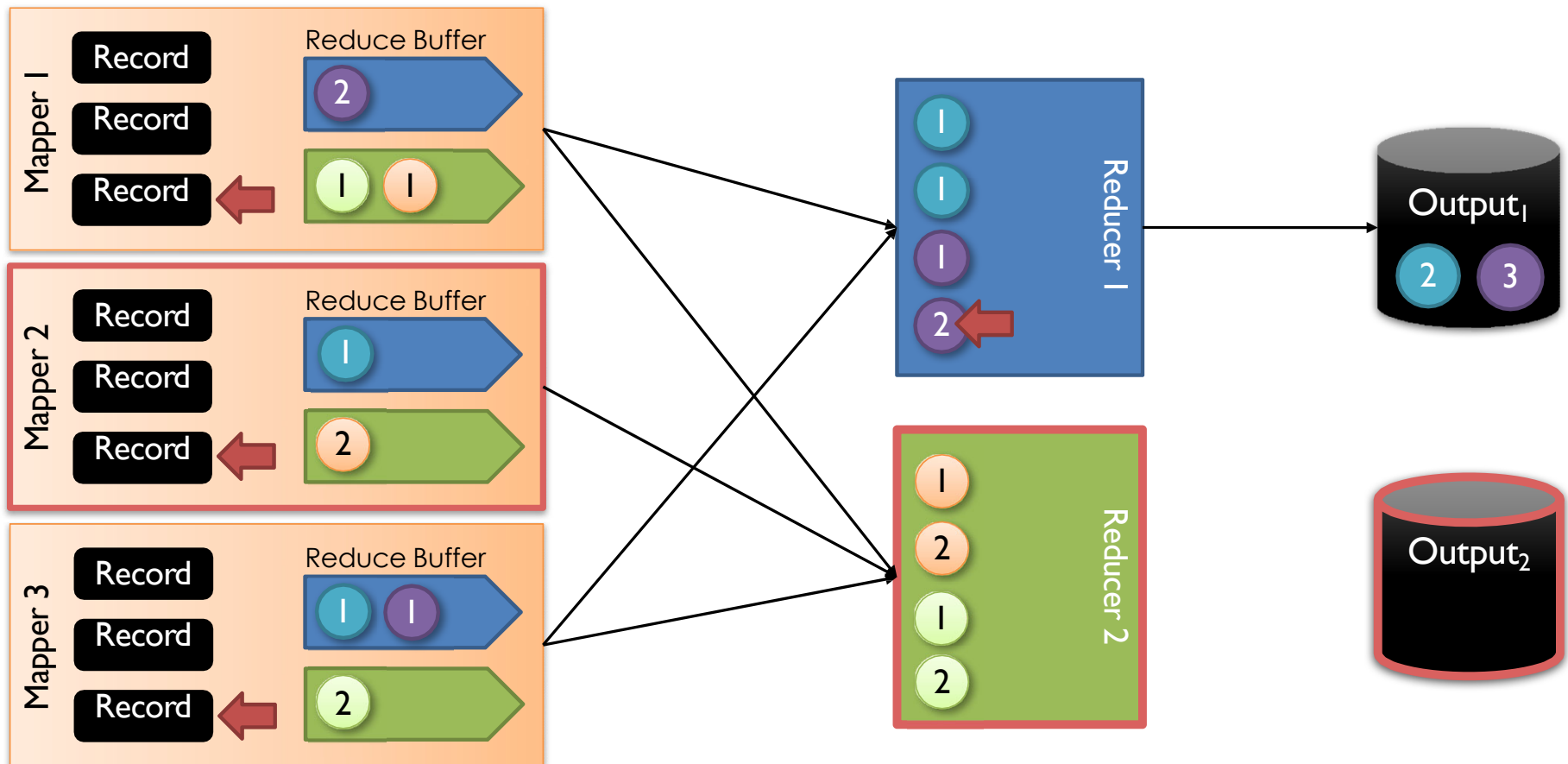
[Dean & Ghemawat, OSDI'04]

The Map Reduce: **Fault Tolerance**



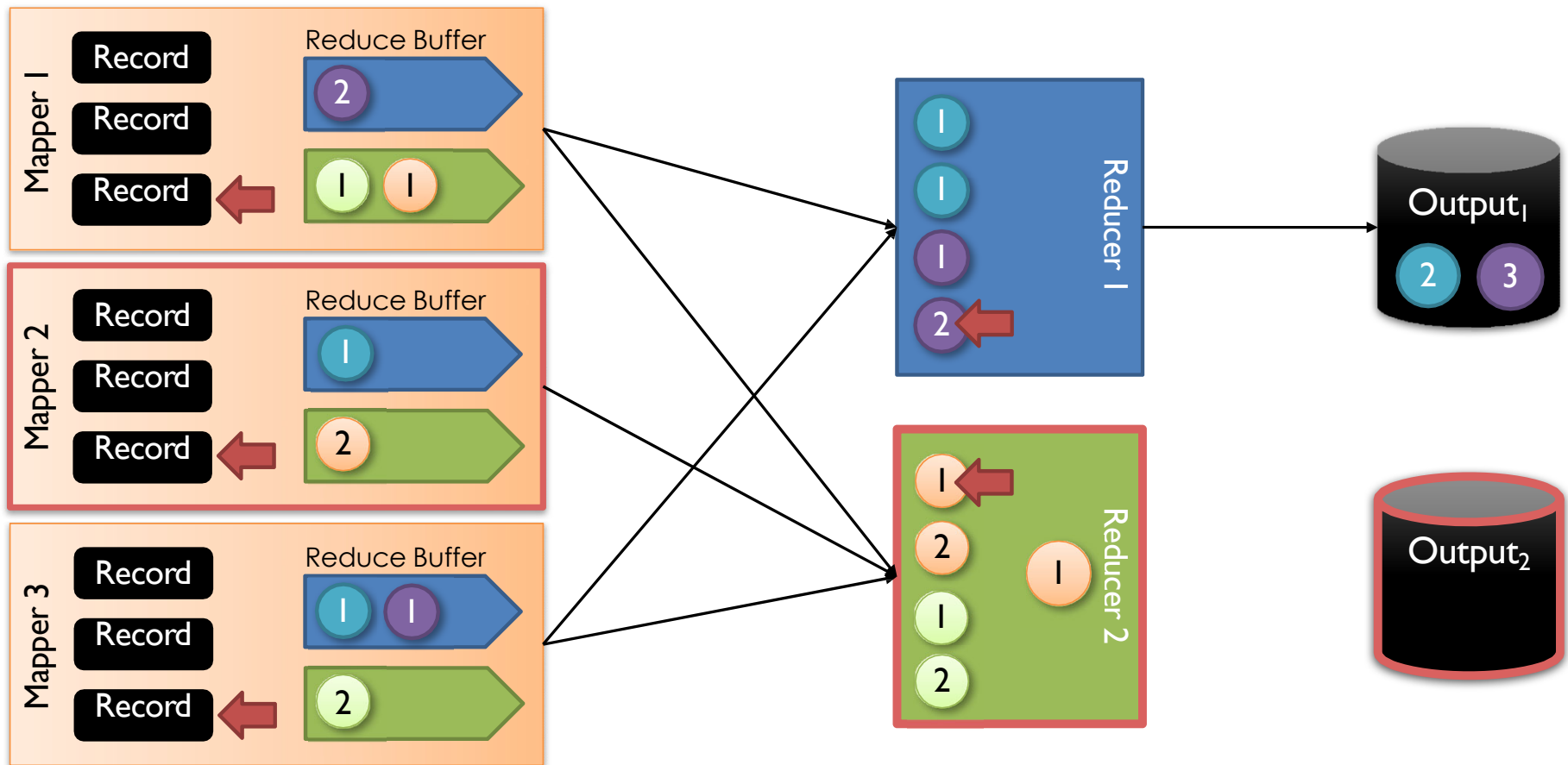
[Dean & Ghemawat, OSDI'04]

The Map Reduce: **Fault Tolerance**



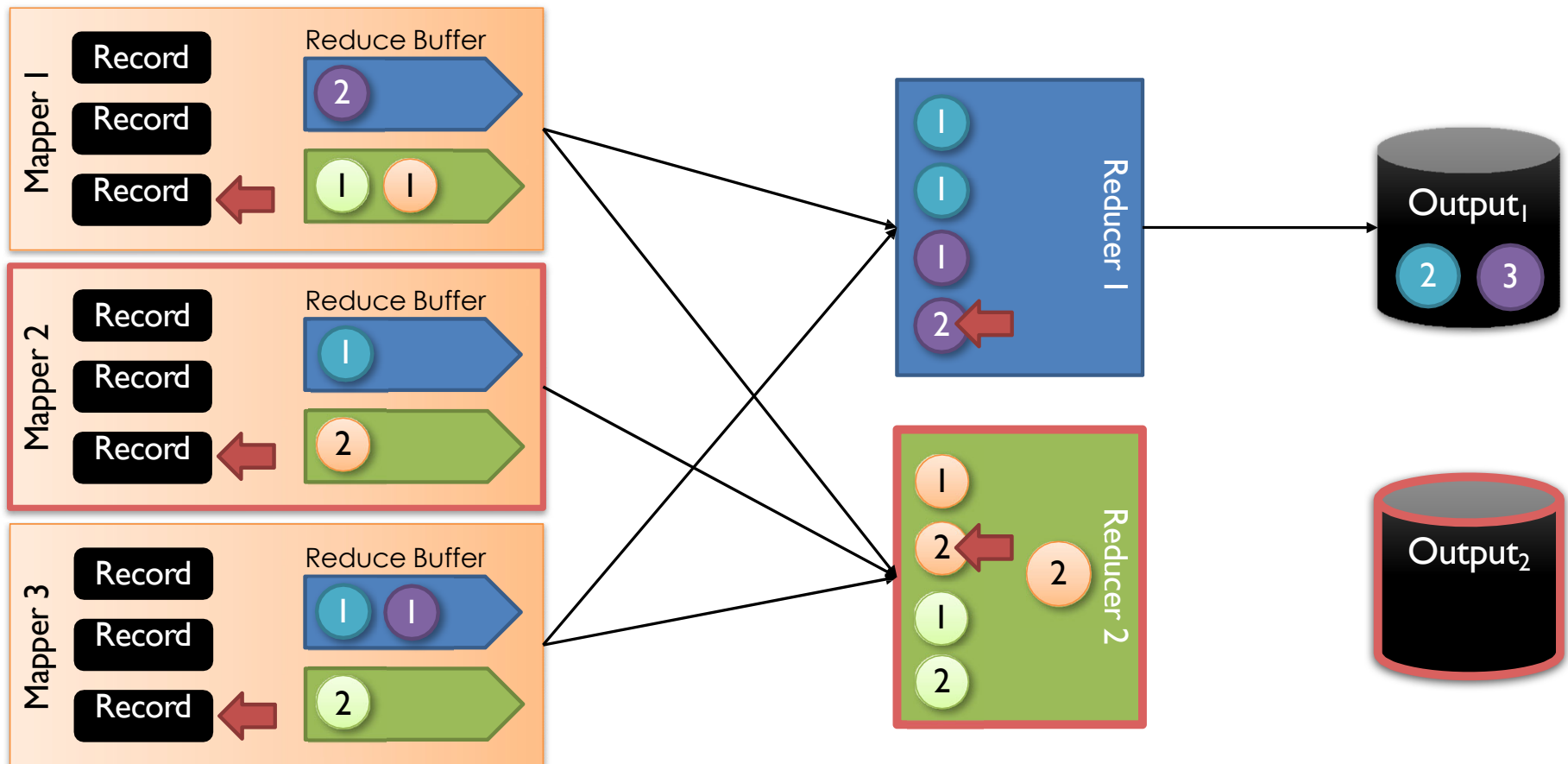
[Dean & Ghemawat, OSDI'04]

The Map Reduce: **Fault Tolerance**



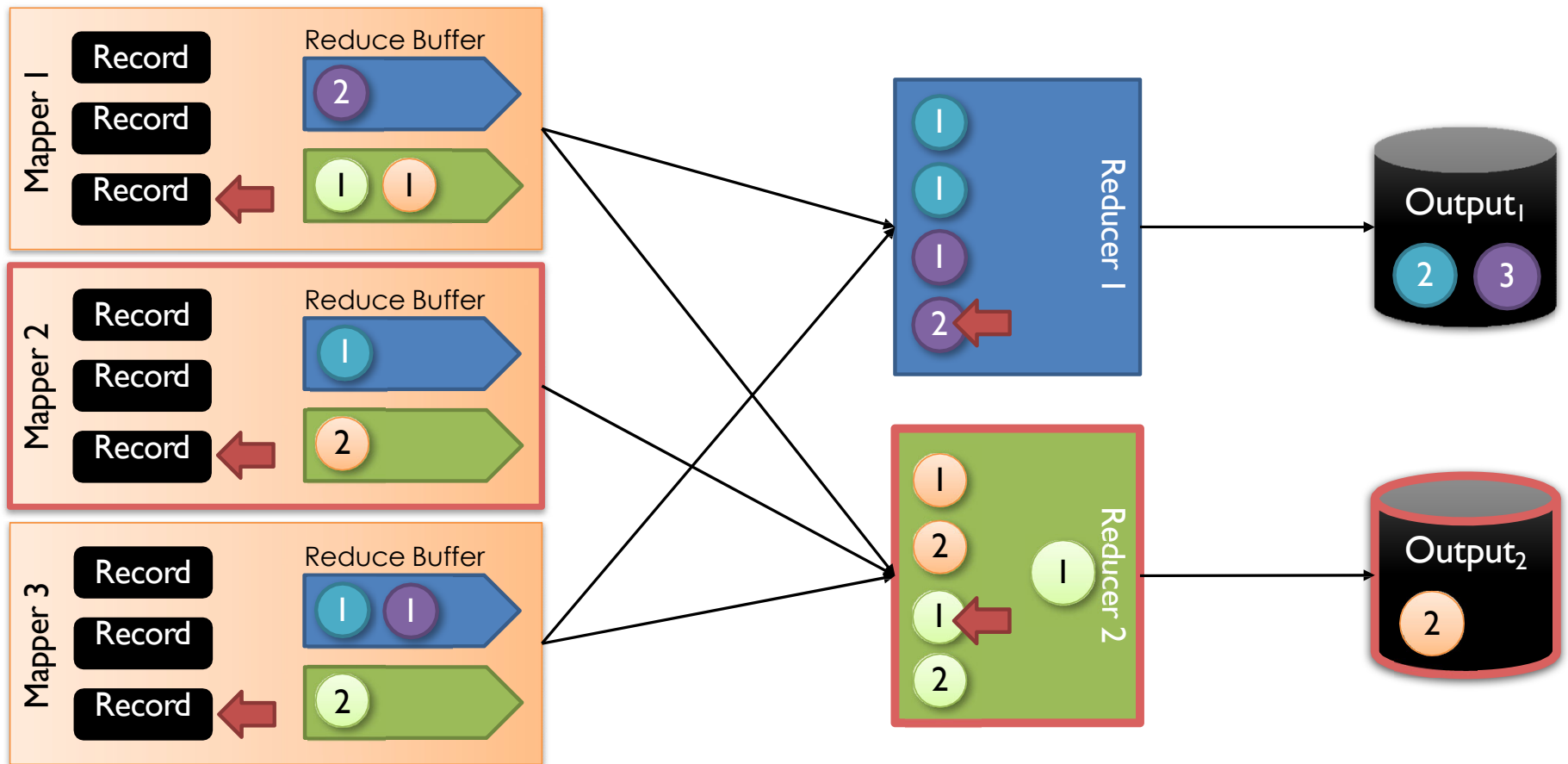
[Dean & Ghemawat, OSDI'04]

The Map Reduce: **Fault Tolerance**



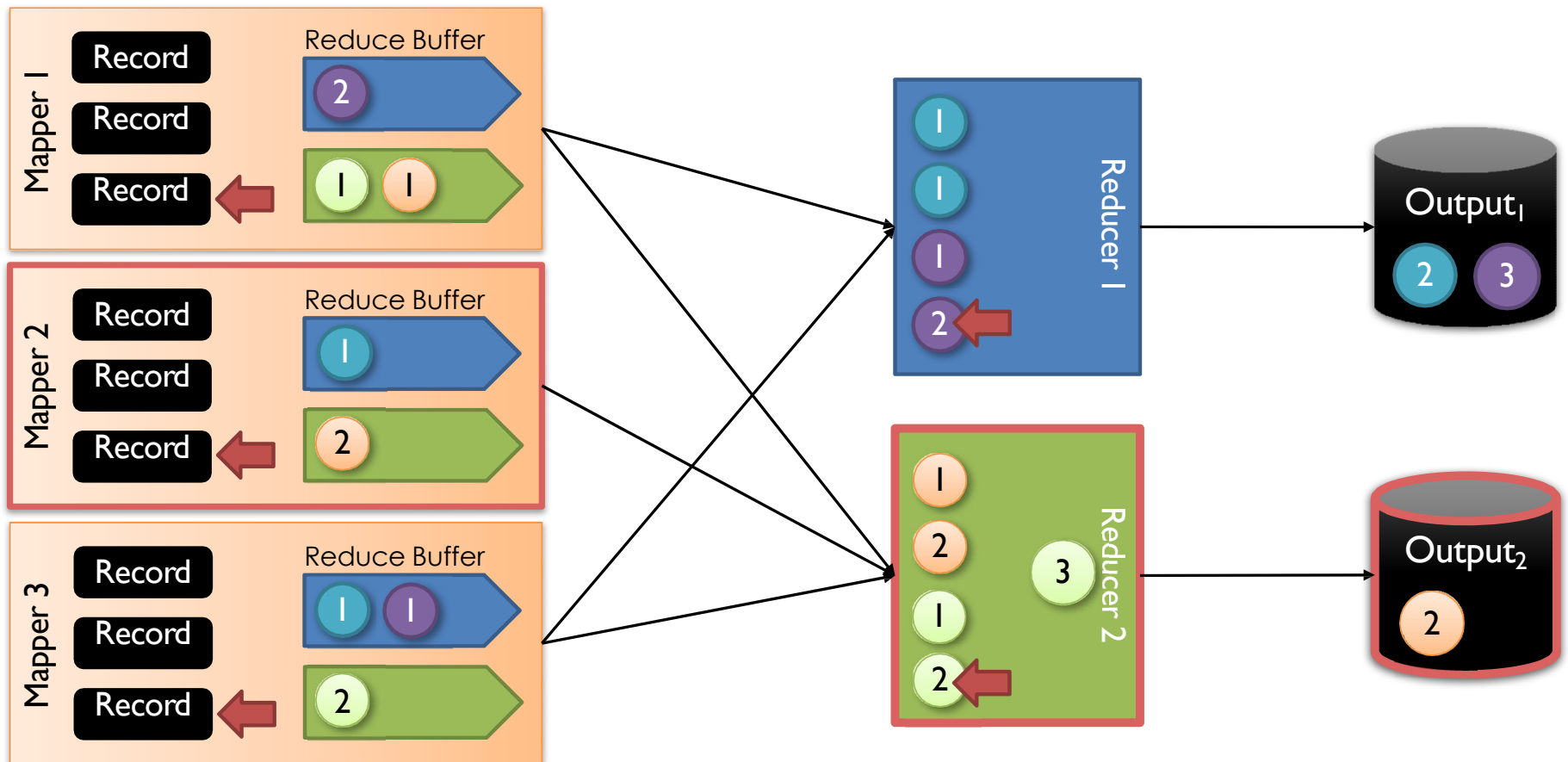
[Dean & Ghemawat, OSDI'04]

The Map Reduce: **Fault Tolerance**



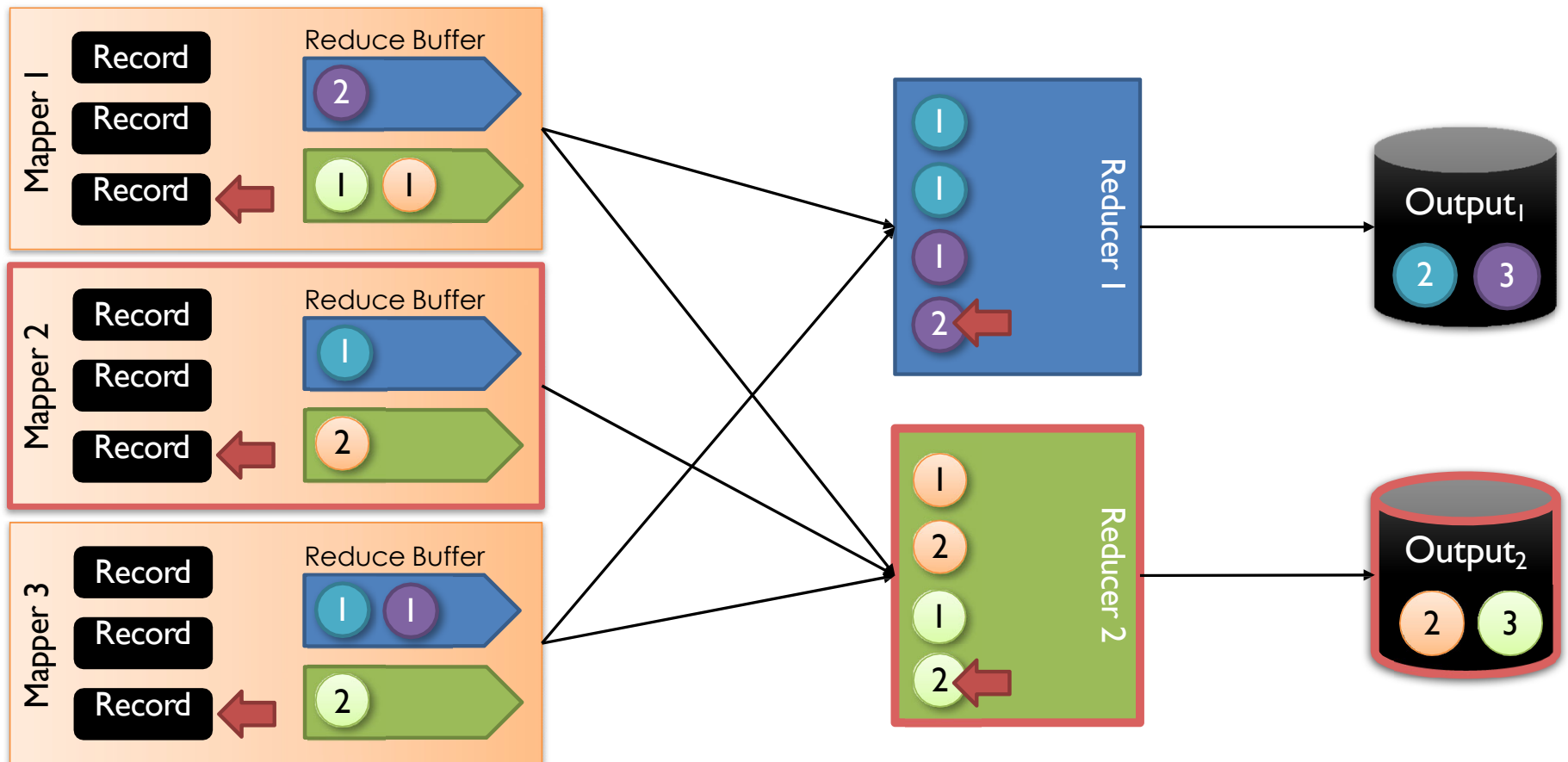
[Dean & Ghemawat, OSDI'04]

The Map Reduce: **Fault Tolerance**



[Dean & Ghemawat, OSDI'04]

The Map Reduce: **Fault Tolerance**

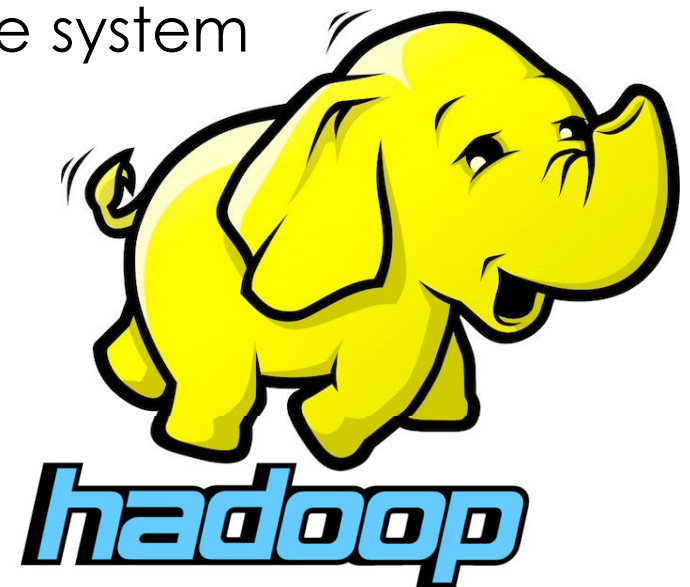


[Dean & Ghemawat, OSDI'04]

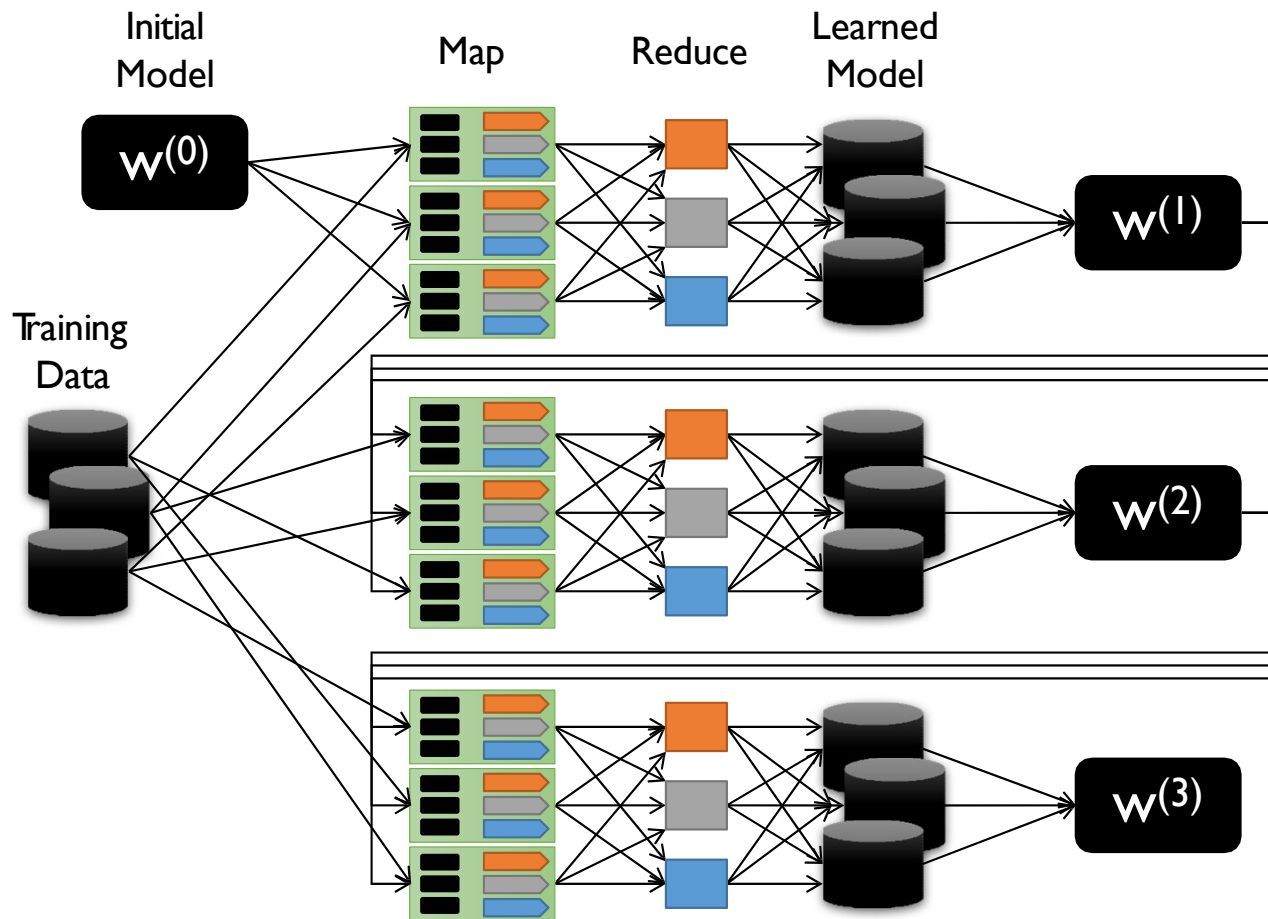
Map Reduce Technologies

Hadoop

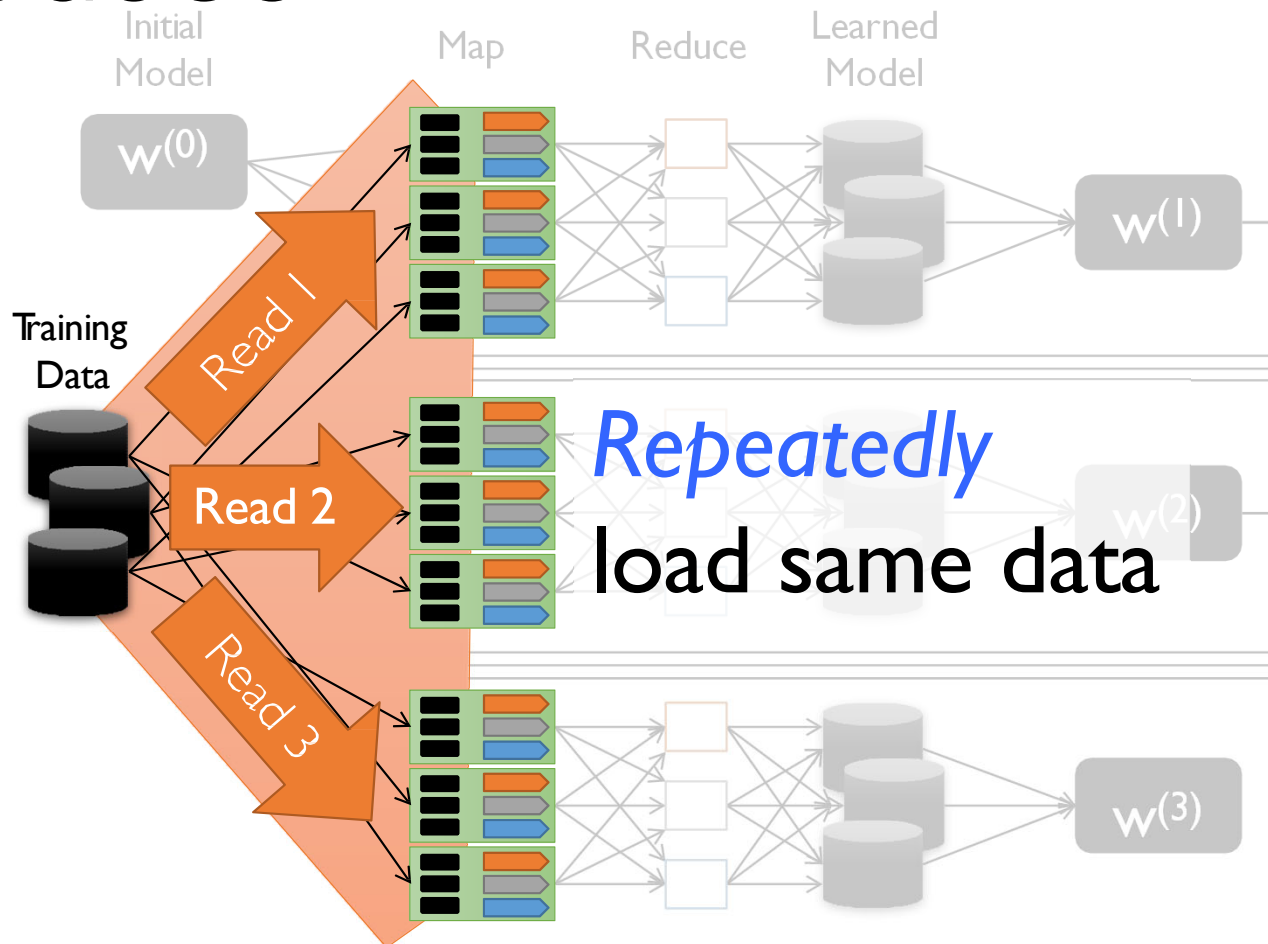
- First open-source map-reduce software system
 - Managed by Apache foundation
- Based on Google's
 - Map-Reduce
 - Google File System
- Several key technologies
 - **HDFS:** Hadoop File System
 - **Yarn:** Yet another resource negotiator
 - **MapReduce:** map-reduce compute framework



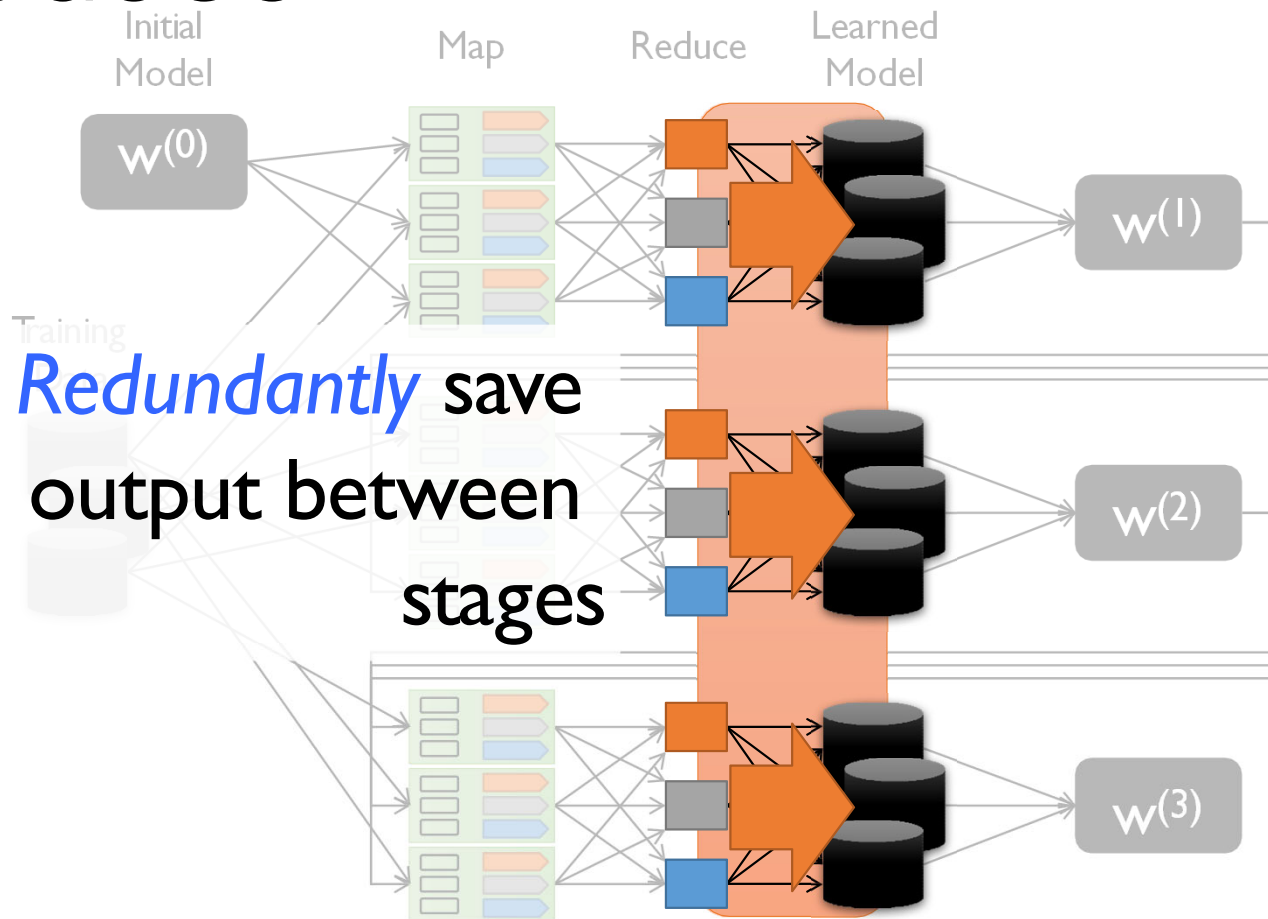
Iteration in Map-Reduce



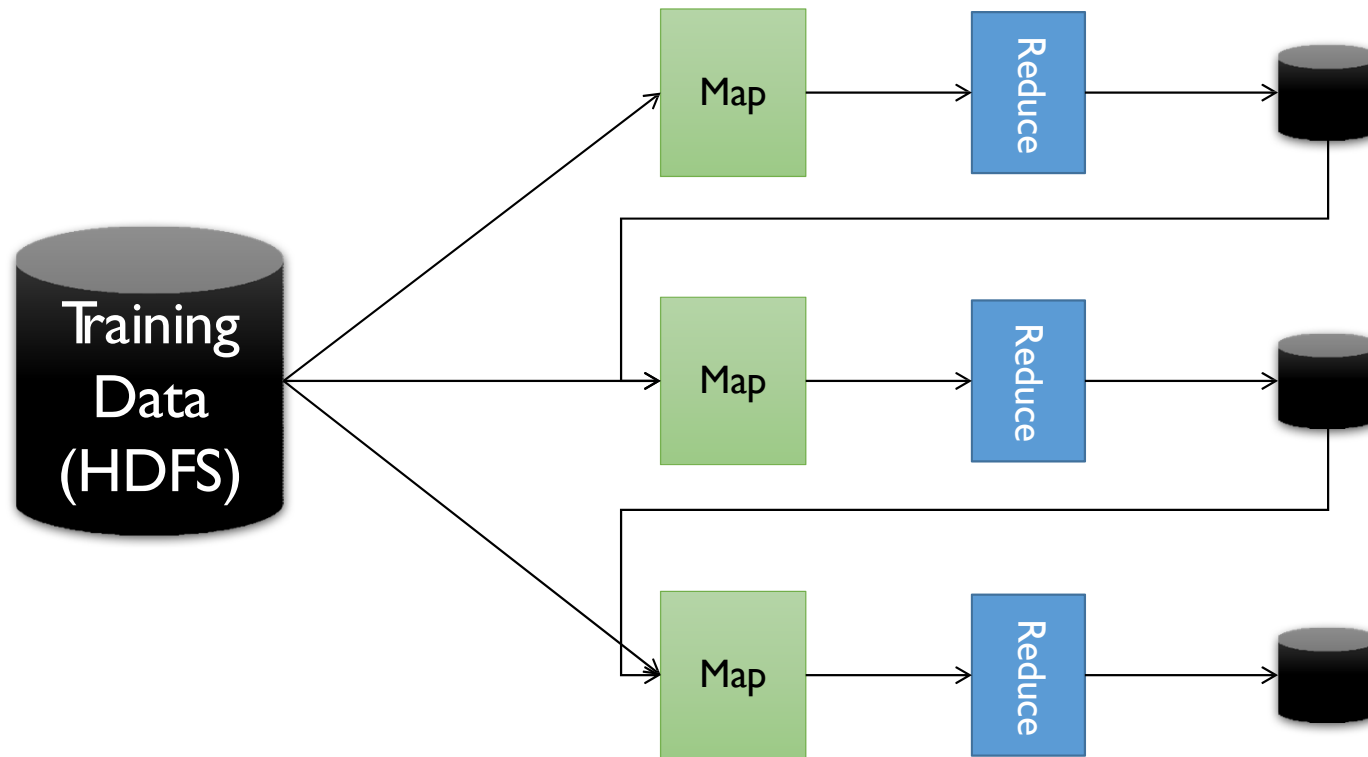
Cost of Iteration in Map-Reduce



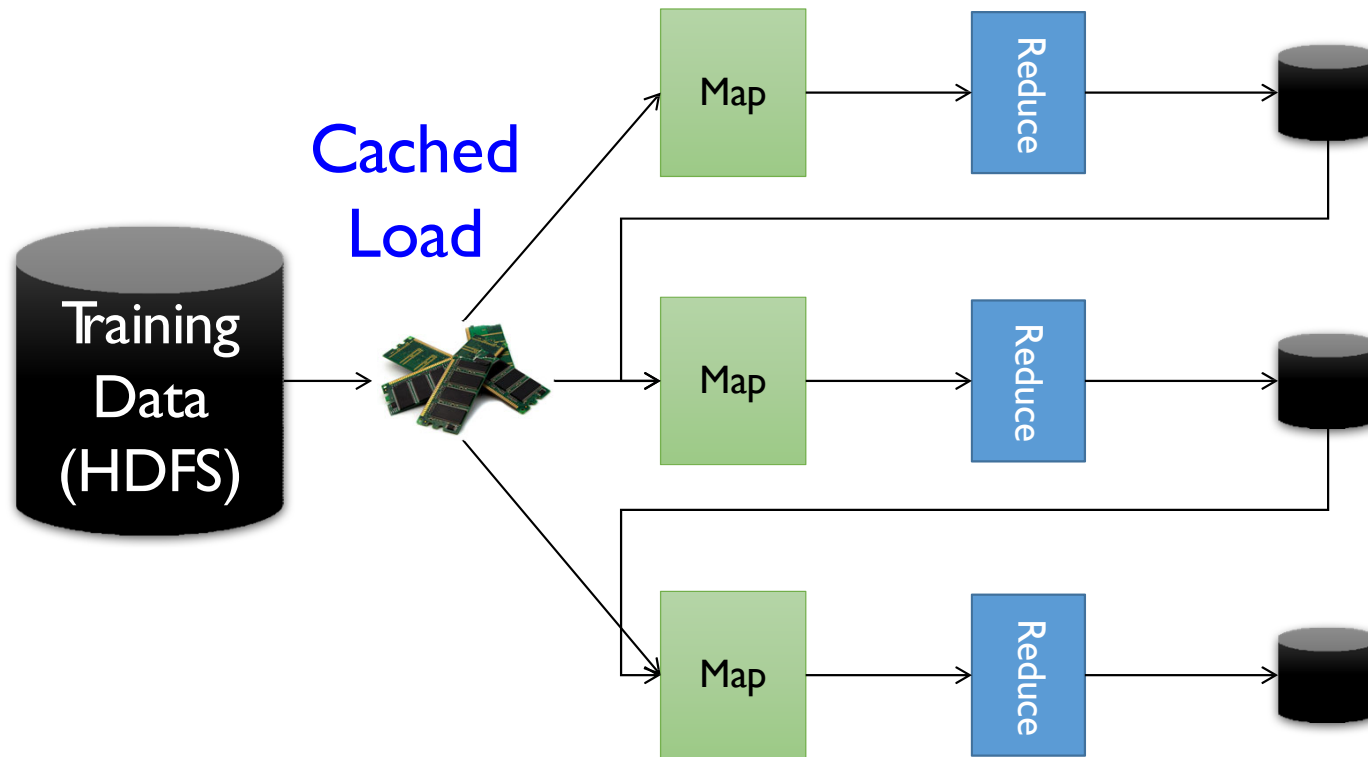
Cost of Iteration in Map-Reduce



Dataflow View

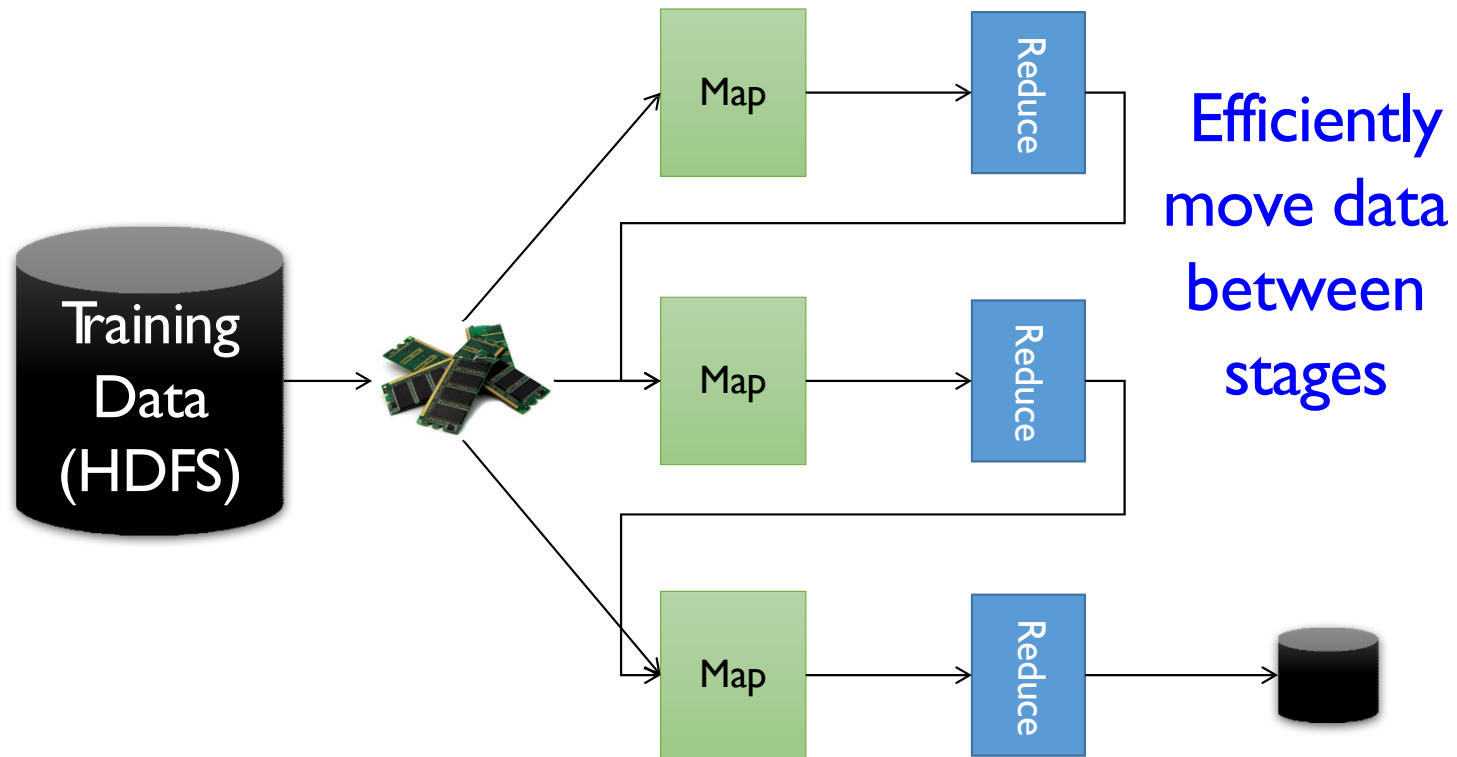


Memory Opt. Dataflow



10-100× faster than network and disk

Memory Opt. Dataflow View





In-Memory Dataflow System

Developed at the UC Berkeley AMP Lab

M. Zaharia, M. Chowdhury, M. J. Franklin, S. Shenker, and I. Stoica. *Spark: cluster computing with working sets*. HotCloud'10

M. Zaharia, M. Chowdhury, T. Das, A. Dave, J. Ma, M. McCauley, M. J. Franklin, S. Shenker, I. Stoica. *Resilient Distributed Datasets: A Fault-Tolerant Abstraction for In-Memory Cluster Computing*, NSDI 2012

What Is *Spark*

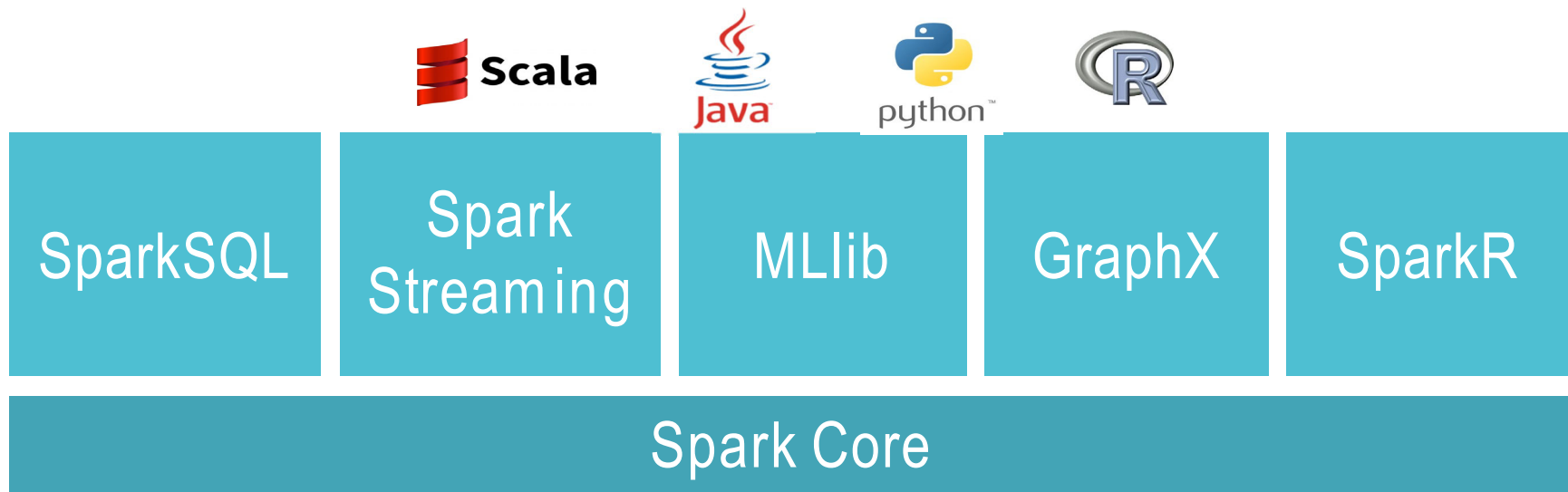
- Parallel execution engine for big data processing
- **General**: efficient support for multiple workloads
- **Easy** to use: 2-5x less code than Hadoop MR
 - High level API's in Python, Java, and Scala
- **Fast**: up to 100x faster than Hadoop MR
 - Can exploit in-memory when available
 - Low overhead scheduling, optimized engine

Spark Programming Abstraction

- *Write programs in terms of transformations on distributed datasets*
- Resilient Distributed Datasets (RDDs)
 - Distributed collections of objects that can be stored in memory or on disk
 - Built via parallel transformations (map, filter, ...)
 - Automatically rebuilt on failure

General

- Unifies **batch, interactive, streaming** workloads
- Easy to build sophisticated applications
 - Support iterative, graph-parallel algorithms
 - Powerful APIs in Scala, Python, Java, R



Easy to Write Code

```
1 public class WordCount {
2     public static class TokenizerMapper
3         extends Mapper<Object, Text, Text, IntWritable> {
4
5         private final static IntWritable one = new IntWritable(1);
6         private Text word = new Text();
7
8         public void map(Object key, Text value, Context context
9             ) throws IOException, InterruptedException {
10             StringTokenizer itr = new StringTokenizer(value.toString());
11             while (itr.hasMoreTokens()) {
12                 word.set(itr.nextToken());
13                 context.write(word, one);
14             }
15         }
16     }
17
18     public static class IntSumReducer
19         extends Reducer<Text, IntWritable, Text, IntWritable> {
20         private IntWritable result = new IntWritable();
21
22         public void reduce(Text key, Iterable<IntWritable> values,
23             Context context
24             ) throws IOException, InterruptedException {
25             int sum = 0;
26             for (IntWritable val : values) {
27                 sum += val.get();
28             }
29             result.set(sum);
30             context.write(key, result);
31         }
32     }
33
34     public static void main(String[] args) throws Exception {
35         Configuration conf = new Configuration();
36         String[] otherArgs = new GenericOptionsParser(conf, args).getRemainingArgs();
37         if (otherArgs.length < 2) {
38             System.err.println("Usage: wordcount <input> <output>");
39             System.exit(2);
40         }
41         Job job = new Job(conf, "word count");
42         job.setJarByClass(WordCount.class);
43         job.setMapperClass(TokenizerMapper.class);
44         job.setCombinerClass(IntSumReducer.class);
45         job.setReducerClass(IntSumReducer.class);
46         job.setOutputKeyClass(Text.class);
47         job.setOutputValueClass(IntWritable.class);
48         for (int i = 0; i < otherArgs.length - 1; i++) {
49             FileInputFormat.addInputPath(job, new Path(otherArgs[i]));
50         }
51         FileOutputFormat.setOutputPath(job,
52             new Path(otherArgs[otherArgs.length - 1]));
53         System.exit(job.waitForCompletion(true) ? 0 : 1);
54     }
55 }
```

```
1 val f = sc.textFile(inputPath)
2 val w = f.flatMap(l => l.split(" ")).map(word => (word, 1)).cache()
3 w.reduceByKey(_ + _).saveAsText(outputPath)
```

WordCount in 3 lines of Spark

WordCount in 50+ lines of Java MR

Fast: Time to sort 100TB

2013 Record:
Hadoop

2100 machines



72 minutes



2014 Record:
Spark

207 machines



23 minutes



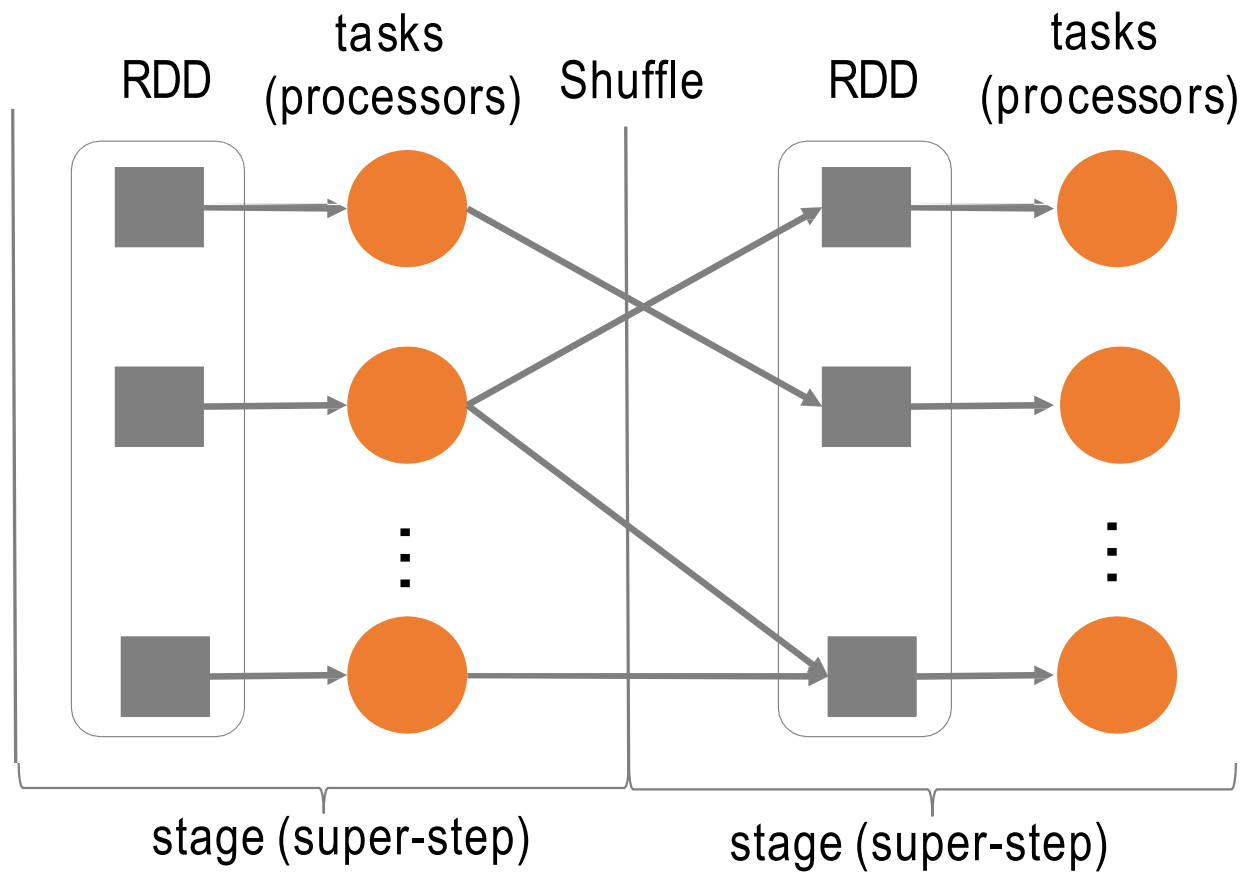
Also sorted 1PB in 4 hours (bottlenecked by network)

Source: Daytona GraySort benchmark, sortbenchmark.org

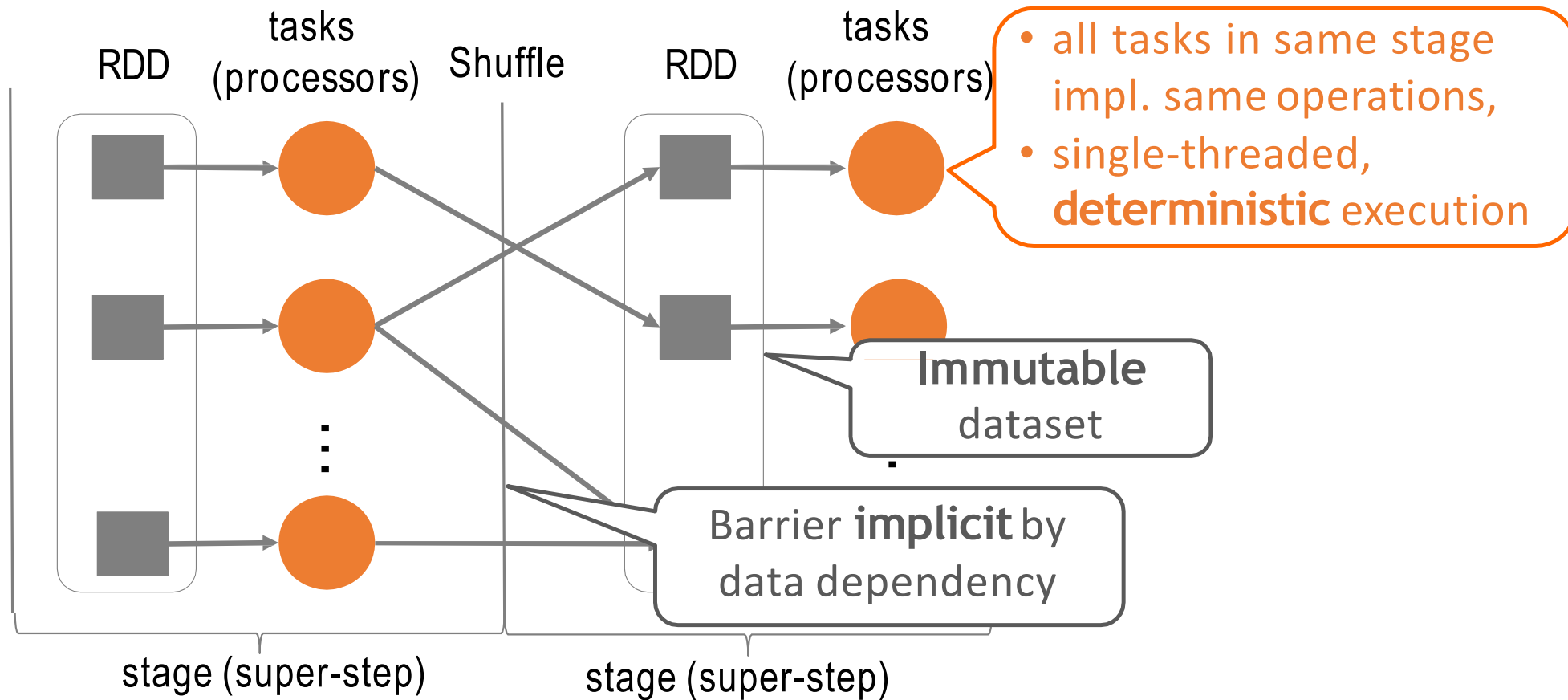
RDD: Resilient Distributed Datasets

- Collections of objects partitioned & distributed across a cluster
 - Stored in RAM or on Disk
 - Resilient to failures
- Operations
 - Transformations
 - Actions

Spark



Spark



Operations on RDDs

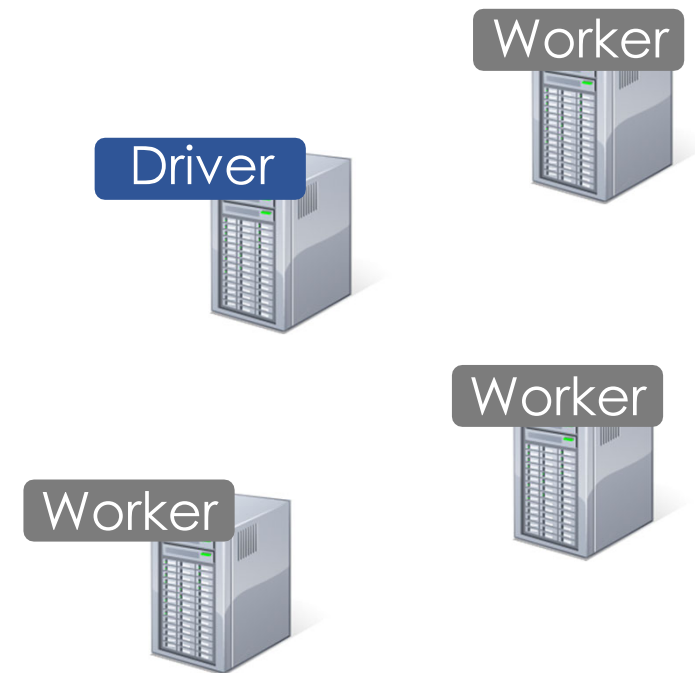
- Transformations $f(\text{RDD}) \Rightarrow \text{RDD}$
 - Lazy (not computed immediately)
 - E.g., “map”, “filter”, “groupBy”
- Actions:
 - Triggers computation
 - E.g. “count”, “collect”, “saveAsTextFile”

Example: Log Mining

Load error messages from a log into memory,
then interactively search for various patterns

Example: Log Mining

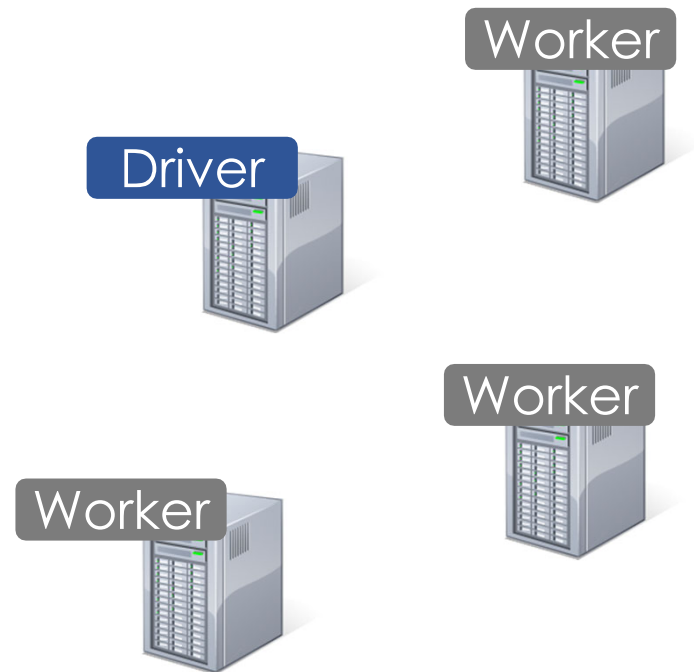
Load error messages from a log into memory,
then interactively search for various patterns



Example: Log Mining

Load error messages from a log into memory,
then interactively search for various patterns

```
lines = spark.textFile("hdfs://...")
```

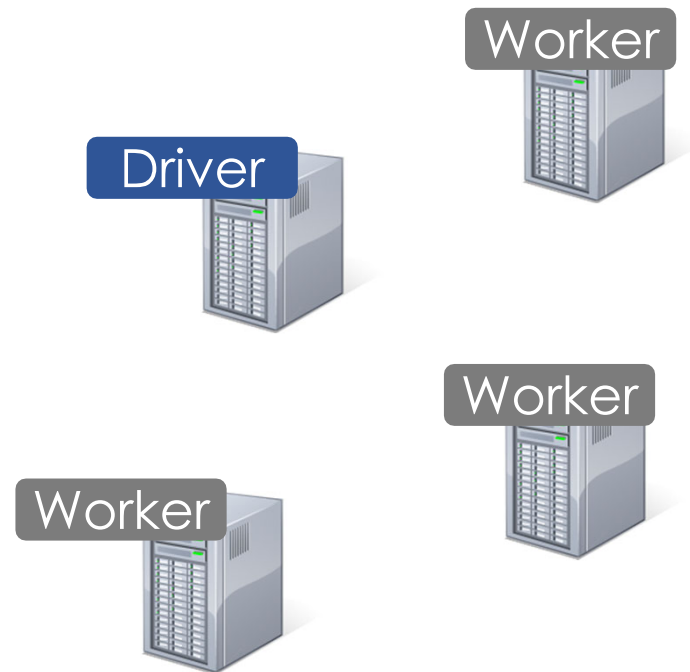


Example: Log Mining

Load error messages from a log into memory,
then interactively search for various patterns

Base RDD

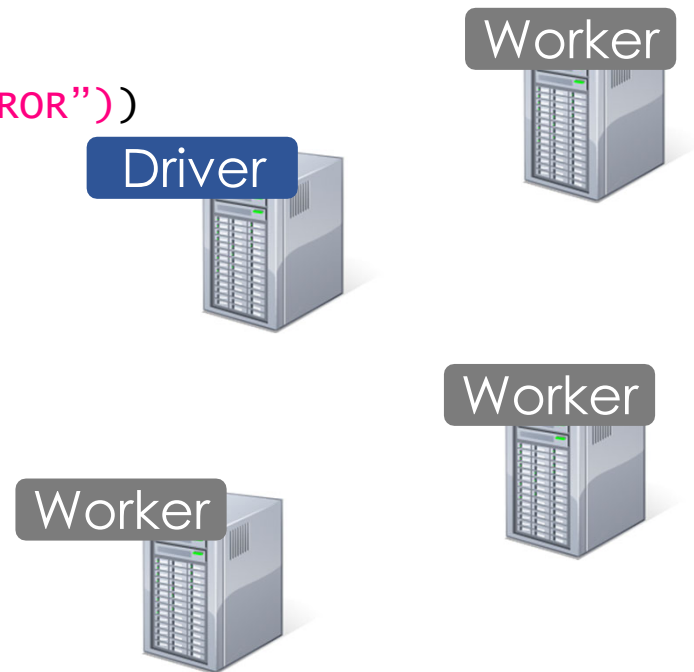
```
lines = spark.textFile("hdfs://...")
```



Example: Log Mining

Load error messages from a log into memory,
then interactively search for various patterns

```
lines = spark.textFile("hdfs://...")  
errors = lines.filter(lambda s: s.startswith("ERROR"))
```



Example: Log Mining

Load error messages from a log into memory,
then interactively search for various patterns

Transformed RDD

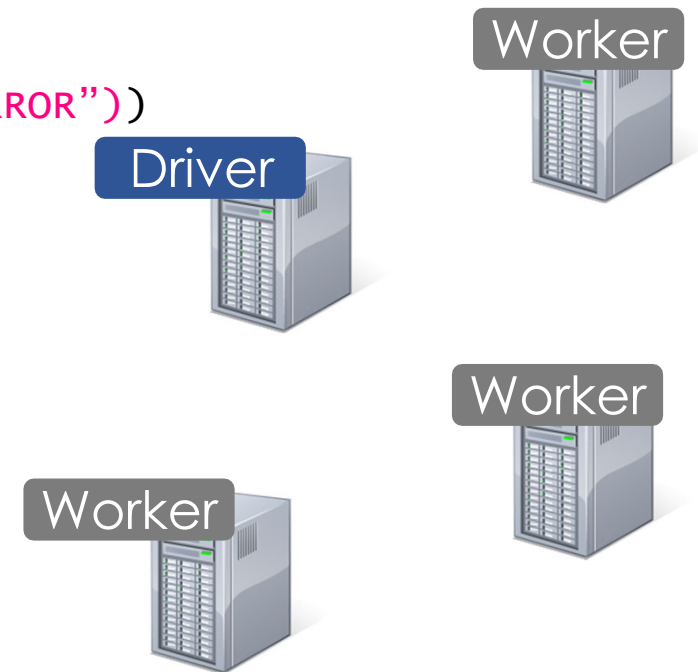
```
lines = spark.textFile("hdfs://...")  
errors = lines.filter(lambda s: s.startswith("ERROR"))
```

Driver

Worker

Worker

Worker



Example: Log Mining

Load error messages from a log into memory, then interactively search for various patterns

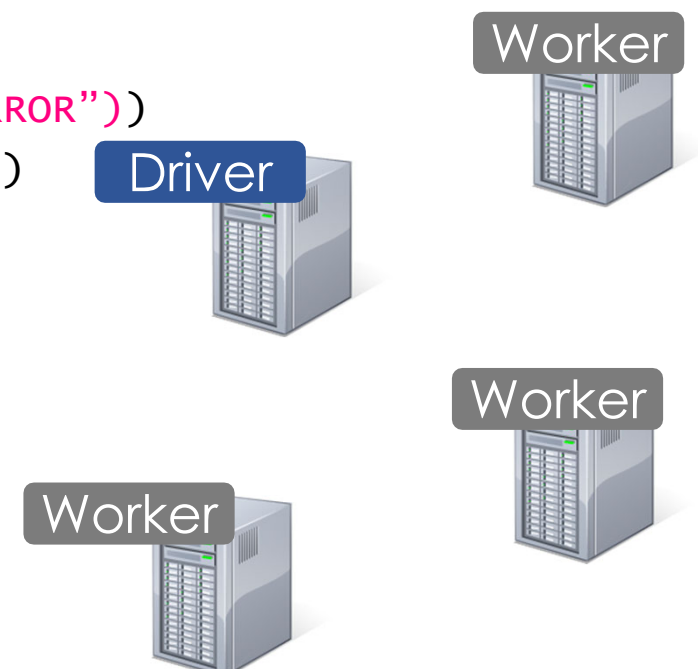
```
lines = spark.textFile("hdfs://...")  
errors = lines.filter(lambda s: s.startswith("ERROR"))  
messages = errors.map(lambda s: s.split("\t")[2])  
messages.cache()  
  
messages.filter(lambda s: "mysql" in s).count()
```

Driver

Worker

Worker

Worker



Example: Log Mining

Load error messages from a log into memory, then interactively search for various patterns

```
lines = spark.textFile("hdfs://...")
errors = lines.filter(lambda s: s.startswith("ERROR"))
messages = errors.map(lambda s: s.split("\t")[2])
messages.cache()

messages.filter(lambda s: "mysql" in s).count()
```

Driver

Action

Worker

Worker

Worker



Example: Log Mining

Load error messages from a log into memory, then interactively search for various patterns

```
lines = spark.textFile("hdfs://...")
errors = lines.filter(lambda s: s.startswith("ERROR"))
messages = errors.map(lambda s: s.split("\t")[2])
messages.cache()
```

```
messages.filter(lambda s: "mysql" in s).count()
```

Driver

Worker

Partition 1

Worker

Partition 2

Worker

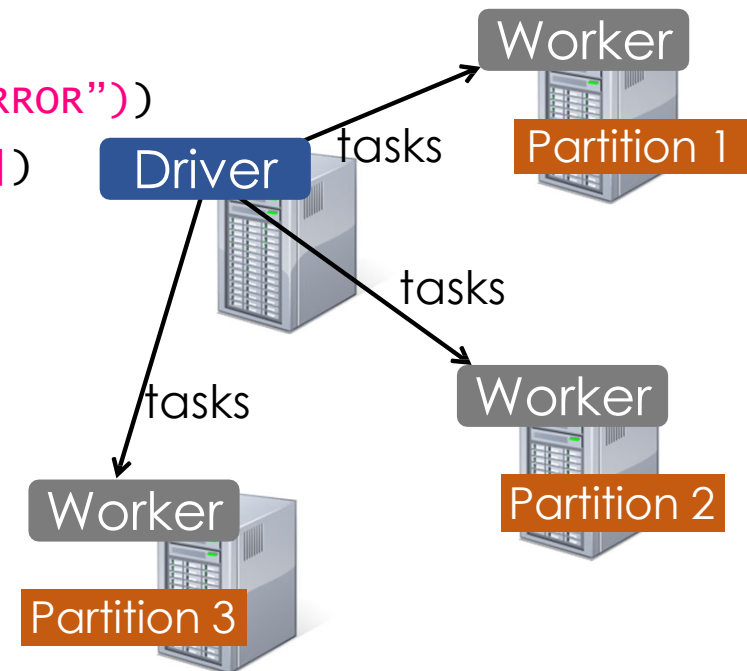
Partition 3

Example: Log Mining

Load error messages from a log into memory, then interactively search for various patterns

```
lines = spark.textFile("hdfs://...")
errors = lines.filter(lambda s: s.startswith("ERROR"))
messages = errors.map(lambda s: s.split("\t")[2])
messages.cache()

messages.filter(lambda s: "mysql" in s).count()
```

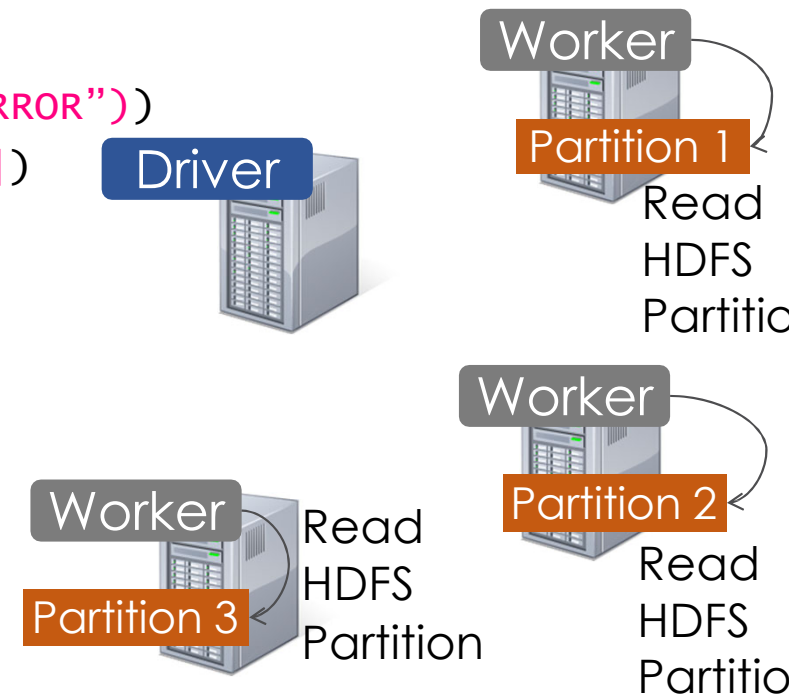


Example: Log Mining

Load error messages from a log into memory, then interactively search for various patterns

```
lines = spark.textFile("hdfs://...")
errors = lines.filter(lambda s: s.startswith("ERROR"))
messages = errors.map(lambda s: s.split("\t")[2])
messages.cache()

messages.filter(lambda s: "mysql" in s).count()
```

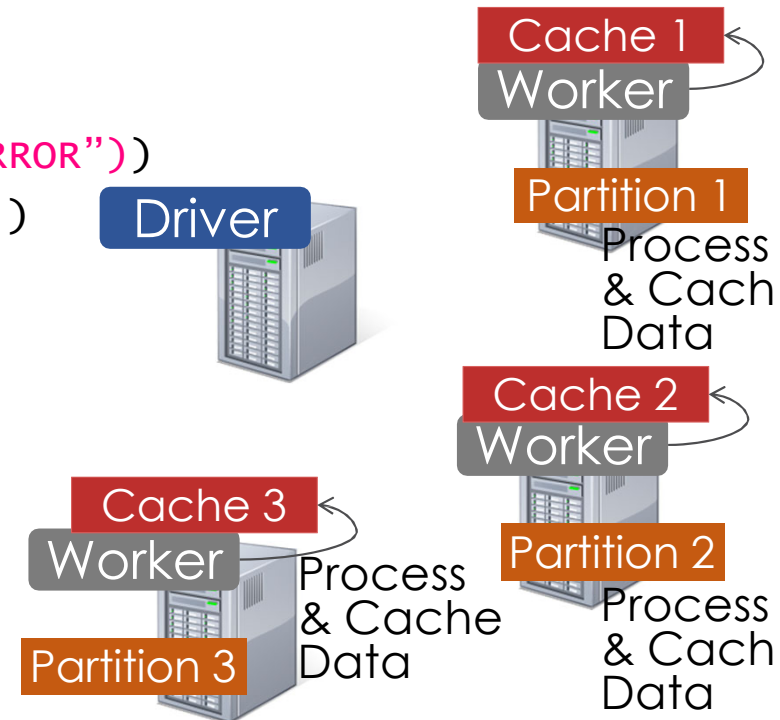


Example: Log Mining

Load error messages from a log into memory, then interactively search for various patterns

```
lines = spark.textFile("hdfs://...")
errors = lines.filter(lambda s: s.startswith("ERROR"))
messages = errors.map(lambda s: s.split("\t")[2])
messages.cache()

messages.filter(lambda s: "mysql" in s).count()
```

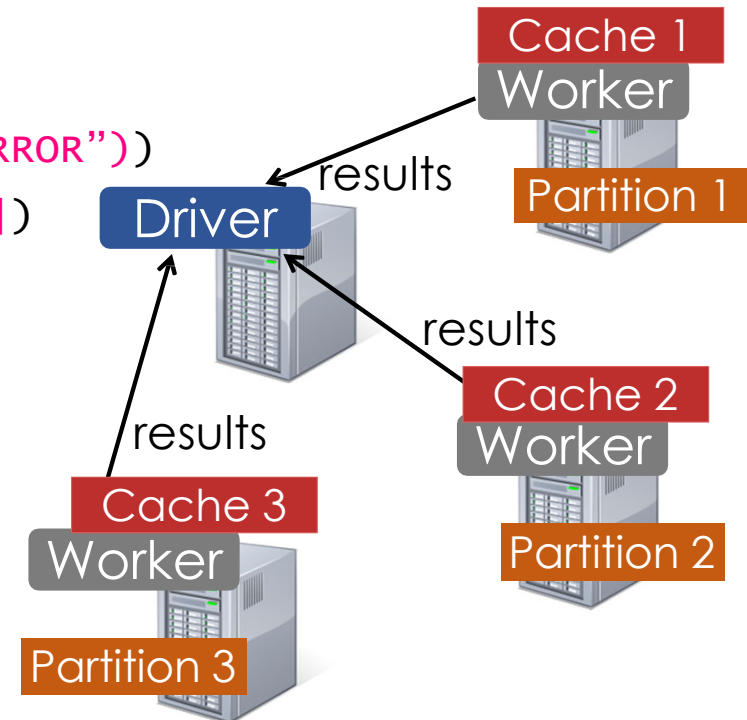


Example: Log Mining

Load error messages from a log into memory, then interactively search for various patterns

```
lines = spark.textFile("hdfs://...")
errors = lines.filter(lambda s: s.startswith("ERROR"))
messages = errors.map(lambda s: s.split("\t")[2])
messages.cache()

messages.filter(lambda s: "mysql" in s).count()
```

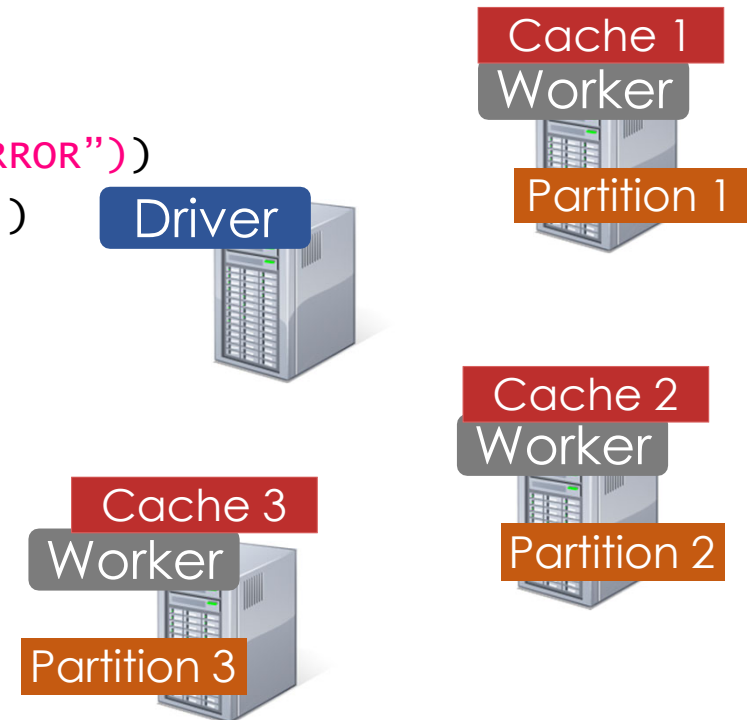


Example: Log Mining

Load error messages from a log into memory,
then interactively search for various patterns

```
lines = spark.textFile("hdfs://...")  
errors = lines.filter(lambda s: s.startswith("ERROR"))  
messages = errors.map(lambda s: s.split("\t")[2])  
messages.cache()
```

```
messages.filter(lambda s: "mysql" in s).count()  
messages.filter(lambda s: "php" in s).count()
```

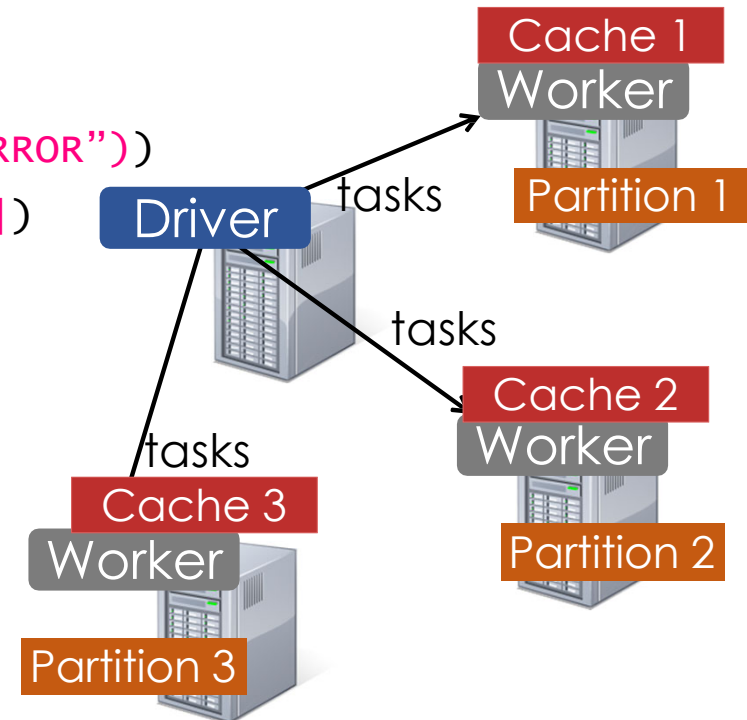


Example: Log Mining

Load error messages from a log into memory, then interactively search for various patterns

```
lines = spark.textFile("hdfs://...")
errors = lines.filter(lambda s: s.startswith("ERROR"))
messages = errors.map(lambda s: s.split("\t")[2])
messages.cache()

messages.filter(lambda s: "mysql" in s).count()
messages.filter(lambda s: "php" in s).count()
```

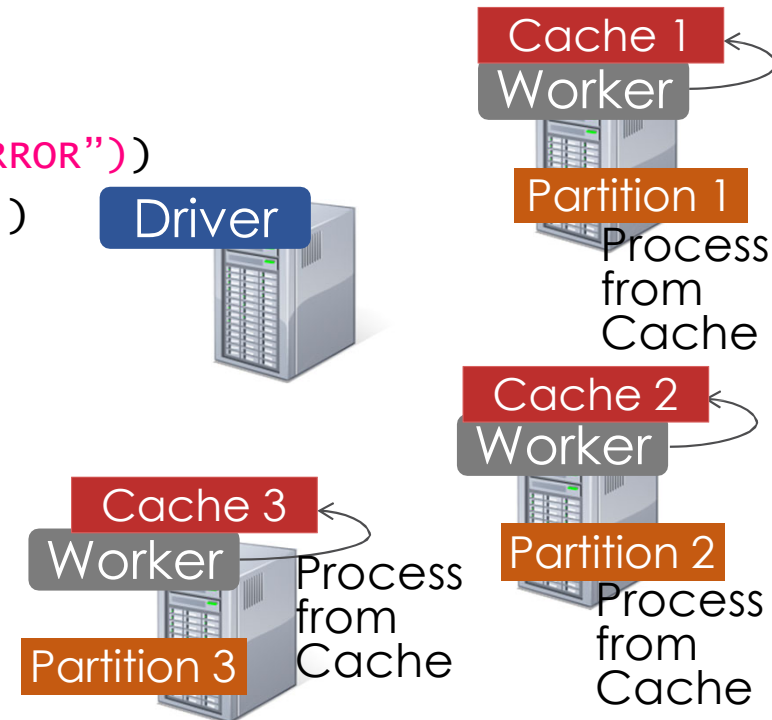


Example: Log Mining

Load error messages from a log into memory, then interactively search for various patterns

```
lines = spark.textFile("hdfs://...")
errors = lines.filter(lambda s: s.startswith("ERROR"))
messages = errors.map(lambda s: s.split("\t")[2])
messages.cache()
```

```
messages.filter(lambda s: "mysql" in s).count()
messages.filter(lambda s: "php" in s).count()
```

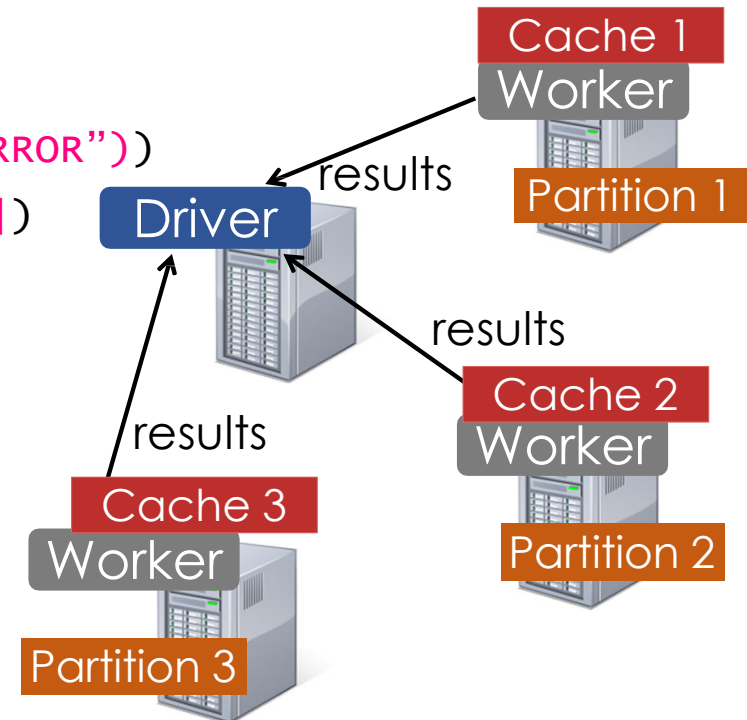


Example: Log Mining

Load error messages from a log into memory, then interactively search for various patterns

```
lines = spark.textFile("hdfs://...")
errors = lines.filter(lambda s: s.startswith("ERROR"))
messages = errors.map(lambda s: s.split("\t")[2])
messages.cache()
```

```
messages.filter(lambda s: "mysql" in s).count()
messages.filter(lambda s: "php" in s).count()
```



Example: Log Mining

Load error messages from a log into memory, then interactively search for various patterns

```
lines = spark.textFile("hdfs://...")
errors = lines.filter(lambda s: s.startswith("ERROR"))
messages = errors.map(lambda s: s.split("\t")[2])
messages.cache()
```

```
messages.filter(lambda s: "mysql" in s).count()
messages.filter(lambda s: "php" in s).count()
```

Cache your data → Faster Results

Full-text search of Wikipedia

- 60GB on 20 EC2 machines
- 0.5 sec from mem vs. 20s for on-disk

Driver

Cache 1
Worker

Partition 1

Cache 2
Worker

Partition 2

Cache 3
Worker

Partition 3