

# Batch Processing



COS 518: *Distributed Systems*  
Lecture 10

Andrew Or, Mike Freedman

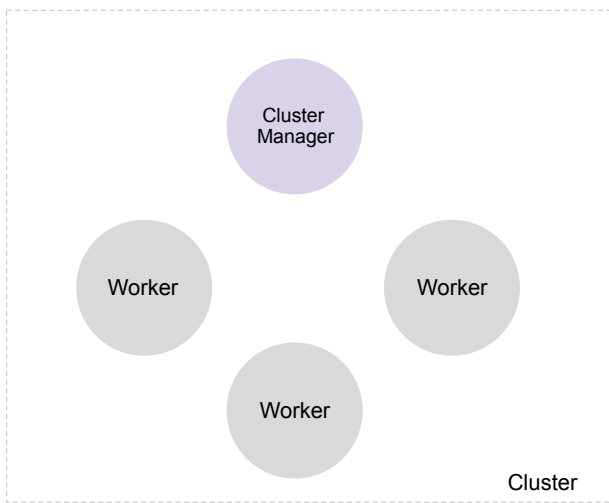
1

# Basic architecture

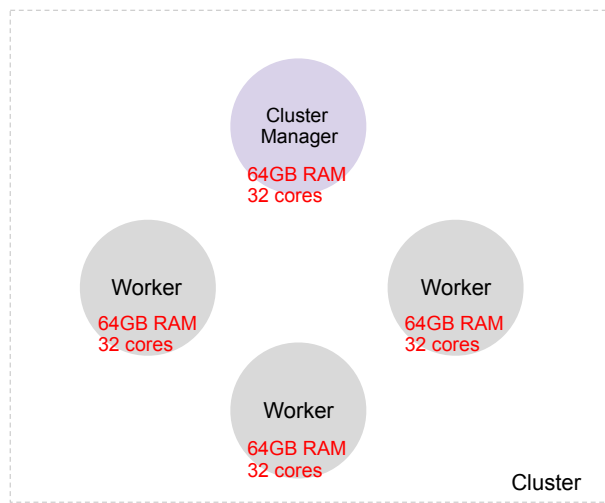
in “big data” systems

2

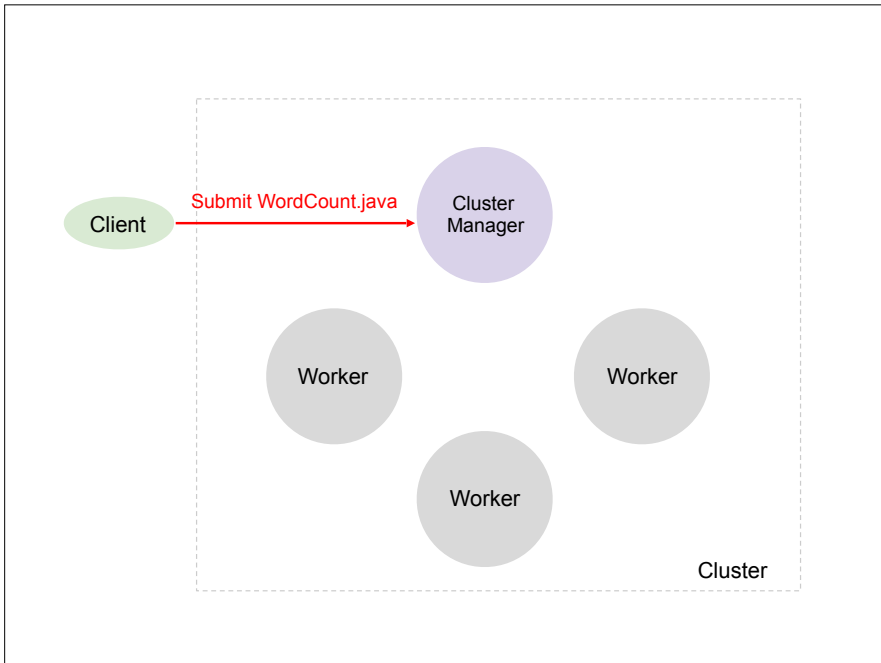
2



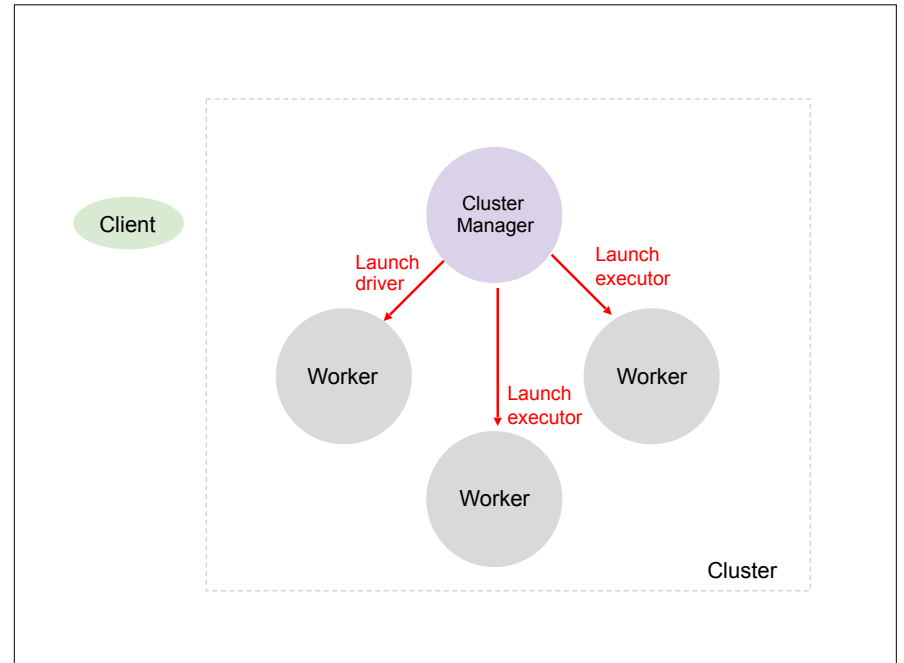
3



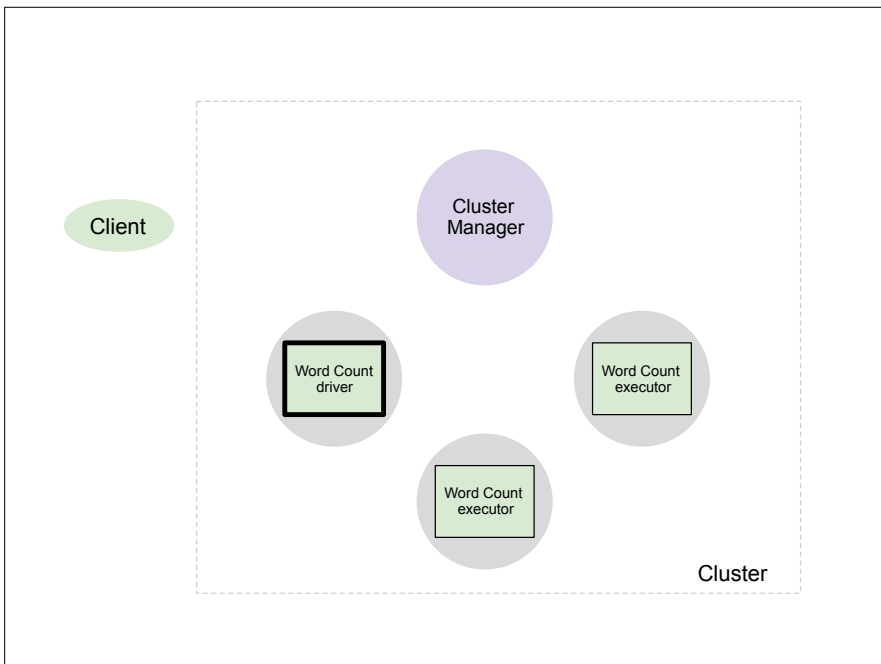
4



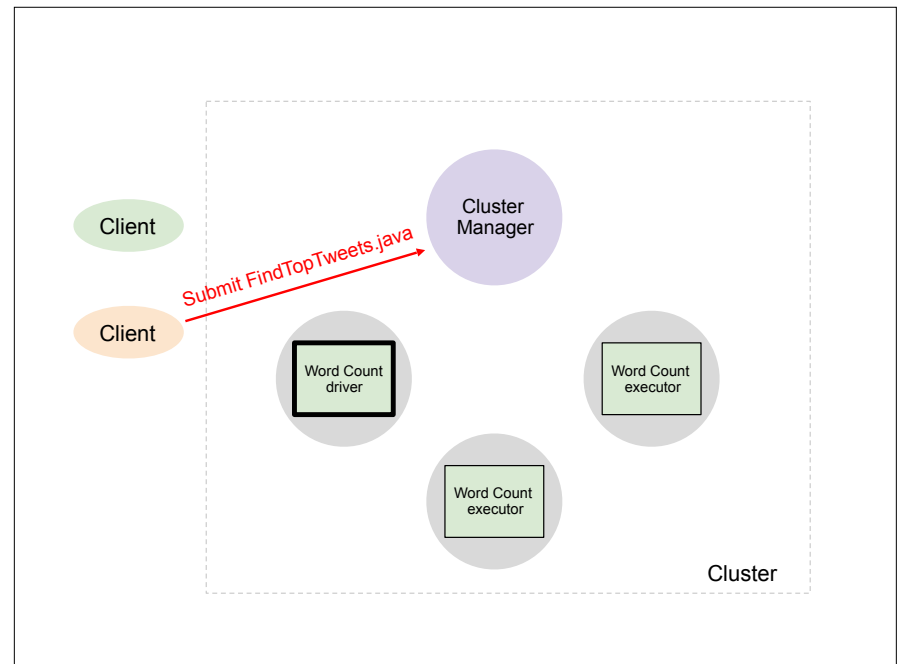
5



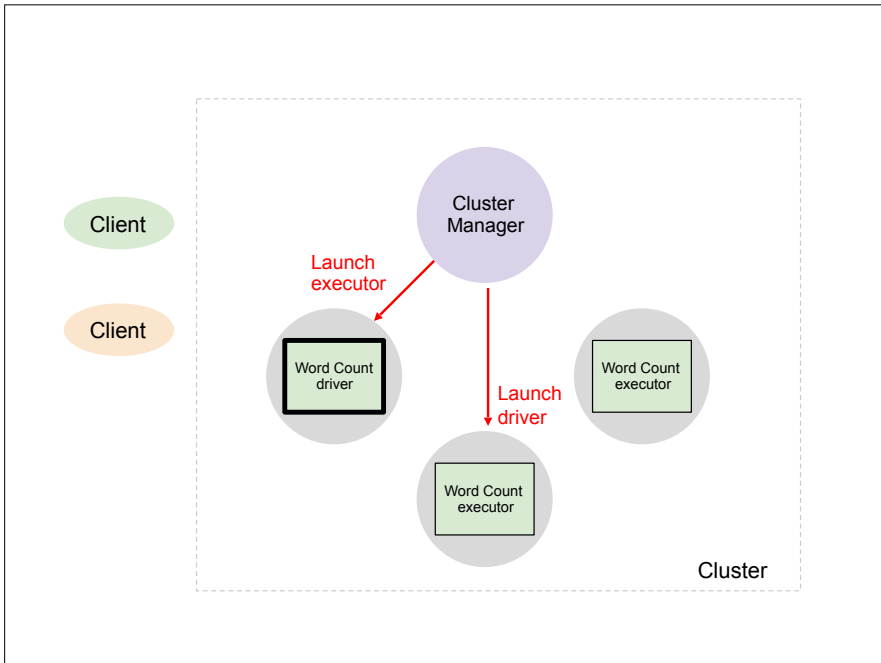
6



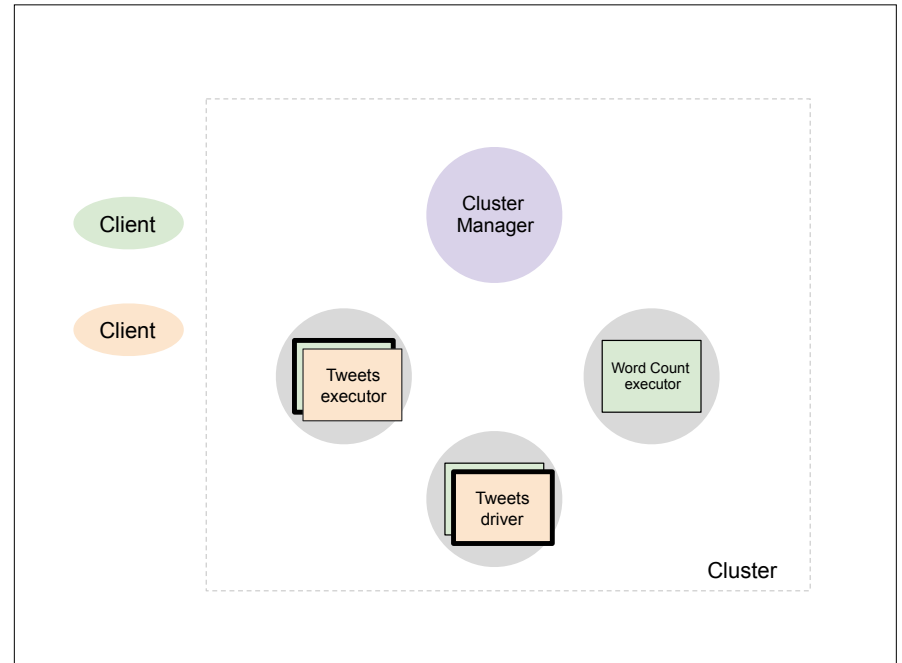
7



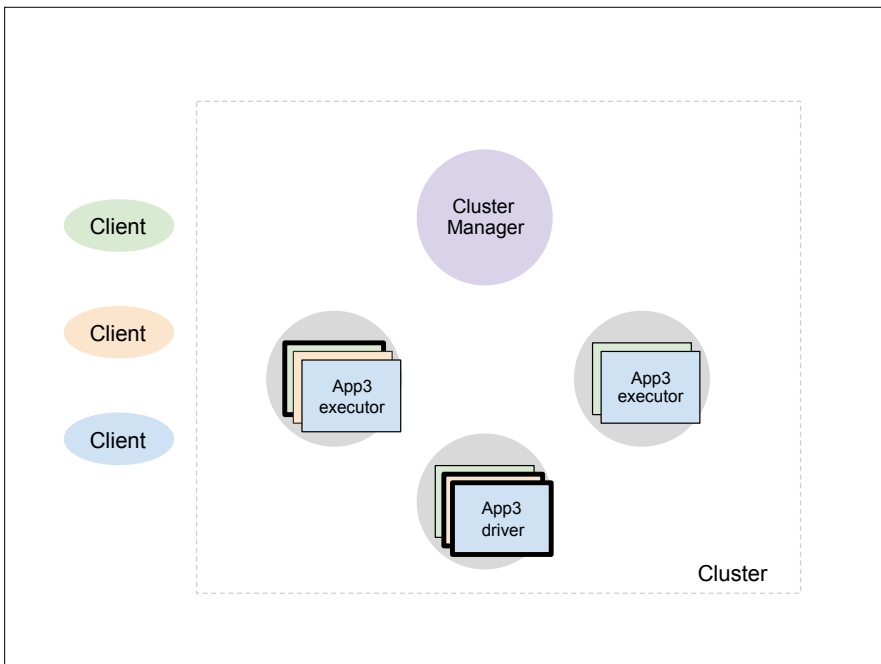
8



9



10



11

## Basic architecture

**Clients** submit applications to the cluster manager

**Cluster manager** assigns cluster resources to applications

Each **Worker** launches containers for each application

**Driver** containers run main method of user program

**Executor** containers run actual computation

*Examples of cluster manager:* YARN, Mesos

*Examples of computing frameworks:* Hadoop MapReduce, Spark

12

## Two levels of scheduling

---

**Cluster-level:** Cluster manager assigns resources to applications

**Application-level:** Driver assigns *tasks* to run on executors

A task is a unit of execution that operates on one *partition*

*Some advantages:*

Applications need not be concerned with resource fairness

Cluster manager need not be concerned with individual tasks

Easy to implement priorities and preemption

13

13

## Case Study: MapReduce

(Data-parallel programming at scale)

14

14

## Application: Word count

---

Hello my love. I love you, my dear. Goodbye.



*hello: 1, my: 2, love: 2, i: 1, dear: 1, goodbye: 1*

15

15

## Application: Word count

---

Locally: tokenize and put words in a hash map

**How do you parallelize this?**

Split document by half

Build two hash maps, one for each half

Merge the two hash maps (by key)

16

16

## How do you do this in a distributed environment?



17

When in the Course of human events, it becomes necessary for one people to dissolve the political bands which have connected them with another, and to assume, among the Powers of the earth, the separate and equal station to which the Laws of Nature and of Nature's God entitle them, a decent respect to the opinions of mankind requires that they should declare the causes which impel them to the separation.

**Input document**



18

When in the Course of human events, it becomes necessary for one people to

---

dissolve the political bands which have connected them with another, and to assume, among the Powers of the earth, the separate and equal station to which the Laws of Nature and of Nature's God entitle them, a decent respect to the opinions of mankind

---

requires that they should declare the causes which impel them to the separation.

**Partition**



19

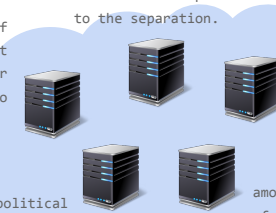
requires that they should declare the causes which impel them to the separation.

When in the Course of human events, it becomes necessary for one people to

Nature and of Nature's God entitle them, a decent respect to the opinions of mankind

dissolve the political bands which have connected them with another, and to assume,

among the Powers of the earth, the separate and equal station to which the Laws of



20

```

requires: 1, that: 1,
they: 1, should: 1,
declare: 1, the: 1,
causes: 1, which: 1 ...

when: 1, in: 1,
the: 1, course: 1,
of: 1, human: 1,
events: 1, it: 1

nature: 2, and: 1, of:
2, god: 1, entitle: 1,
them: 1, decent: 1,
respect: 1, mankind: 1,
opinion: 1 ...

dissolve: 1, the: 2,
political: 1, bands:
1, which: 1, have: 1,
connected: 1, them: 1
...

among: 1, the: 2,
powers: 1, of: 2,
earth: 1, separate: 1,
equal: 1, and: 1 ...

```

Compute word counts locally

21

```

requires: 1, that: 1,
they: 1, should: 1,
declare: 1, the: 1,
causes: 1, which: 1 ...

when: 1, in: 1,
the: 1, course: 1,
of: 1, human: 1,
events: 1, it: 1

nature: 2, and: 1, of:
2, god: 1, entitle: 1,
them: 1, decent: 1,
respect: 1, mankind: 1,
opinion: 1 ...

dissolve: 1, the: 2,
political: 1, bands:
1, which: 1, have: 1,
connected: 1, them: 1
...

among: 1, the: 2,
powers: 1, of: 2,
earth: 1, separate: 1,
equal: 1, and: 1 ...

```

Now what...  
How to merge results?

Compute word counts locally

22

## Merging results computed locally

Several options

Don't merge — requires additional computation for correct results

Send everything to one node — what if data is too big? Too slow...

Partition key space among nodes in cluster (e.g. [a-e], [f-j], [k-p] ...)

1. Assign a key space to each node
2. Partition local results by the key spaces
3. Fetch and merge results that correspond to the node's key space

23

```

requires: 1, that: 1,
they: 1, should: 1,
declare: 1, the: 1,
causes: 1, which: 1 ...

when: 1, in: 1,
the: 1, course: 1,
of: 1, human: 1,
events: 1, it: 1

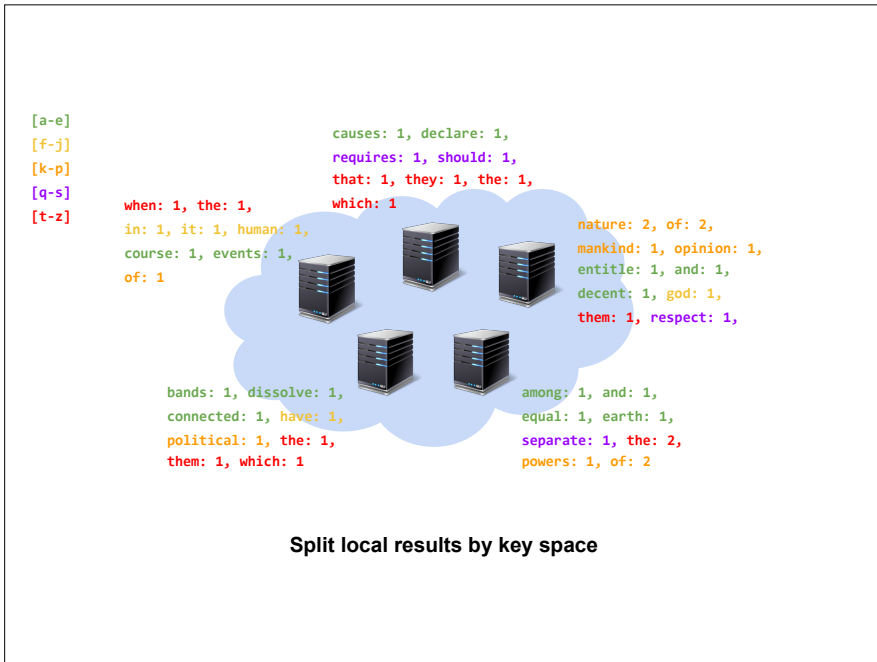
nature: 2, and: 1, of:
2, god: 1, entitle: 1,
them: 1, decent: 1,
respect: 1, mankind: 1,
opinion: 1 ...

dissolve: 1, the: 2,
political: 1, bands:
1, which: 1, have: 1,
connected: 1, them: 1
...

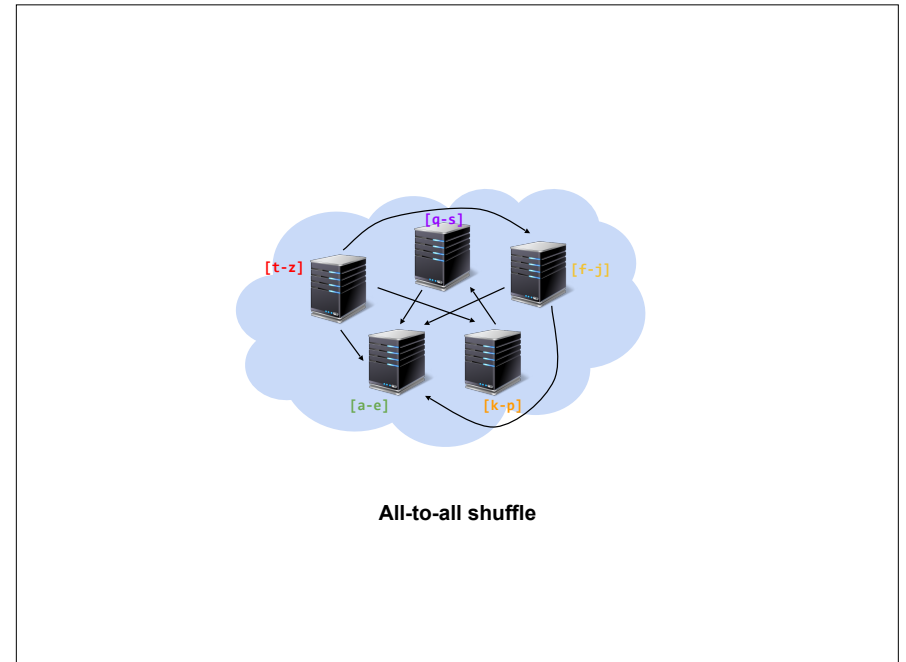
among: 1, the: 2,
powers: 1, of: 2,
earth: 1, separate: 1,
equal: 1, and: 1 ...

```

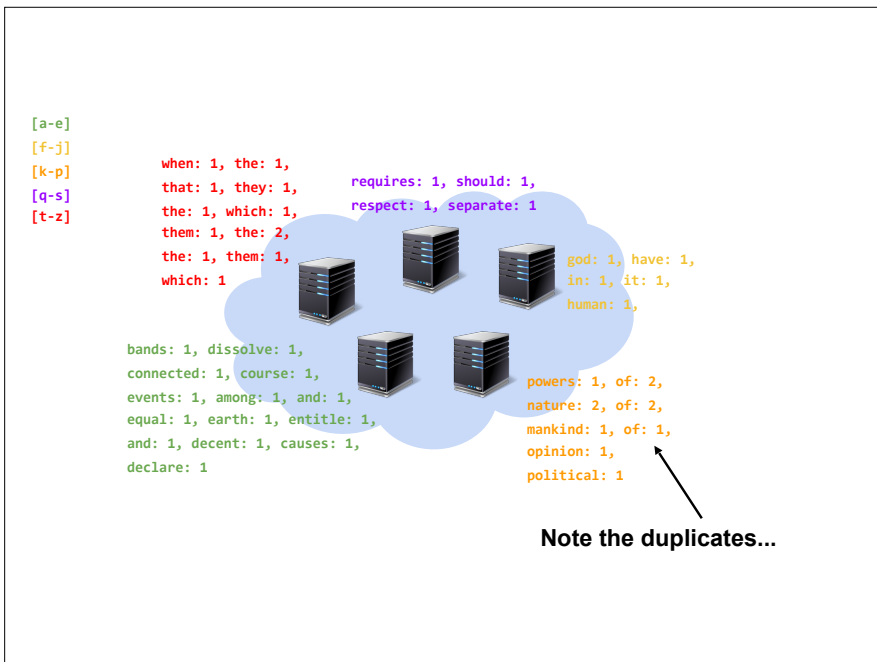
24



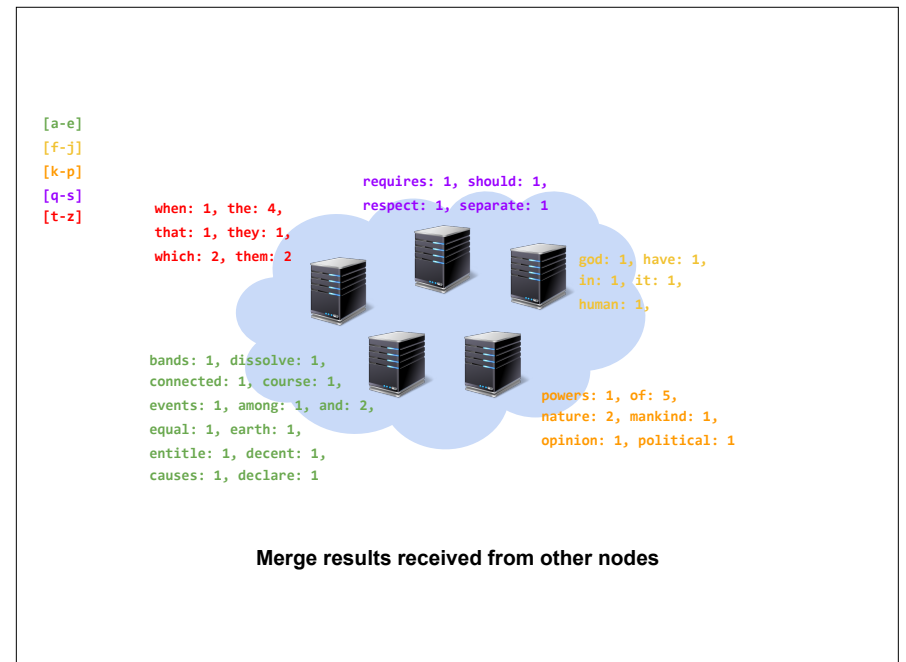
25



26



27



28

## MapReduce

Partition dataset into many chunks

**Map stage:** Each node processes one or more chunks locally

**Reduce stage:** Each node fetches and merges partial results from all other nodes

29

29

## MapReduce Interface

**map(key, value) -> list(<k', v'>)**

Apply function to (key, value) pair

Outputs set of intermediate pairs

**reduce(key, list<value>) -> <k', v'>**

Applies aggregation function to values

Outputs result

30

30

## MapReduce: Word count

**map(key, value):**

```
// key = document name
// value = document contents
for each word w in value:
    emit (w, 1)
```

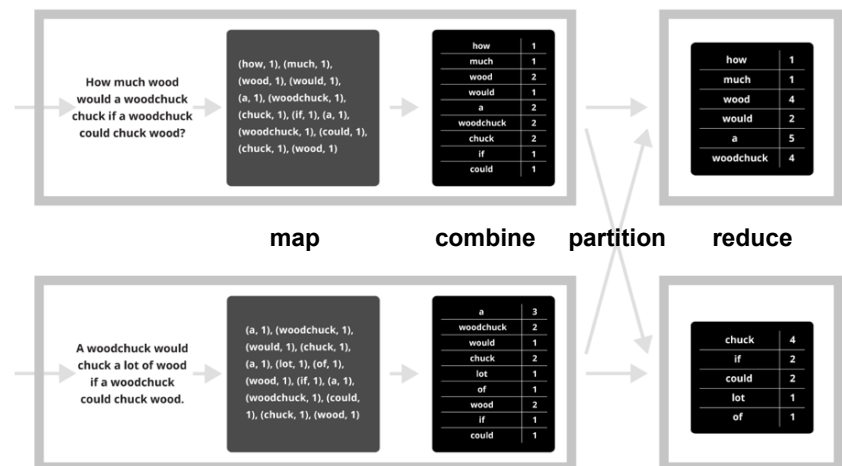
**reduce(key, values):**

```
// key = the word
// values = number of occurrences of that word
count = sum(values)
emit (key, count)
```

31

31

## MapReduce: Word count

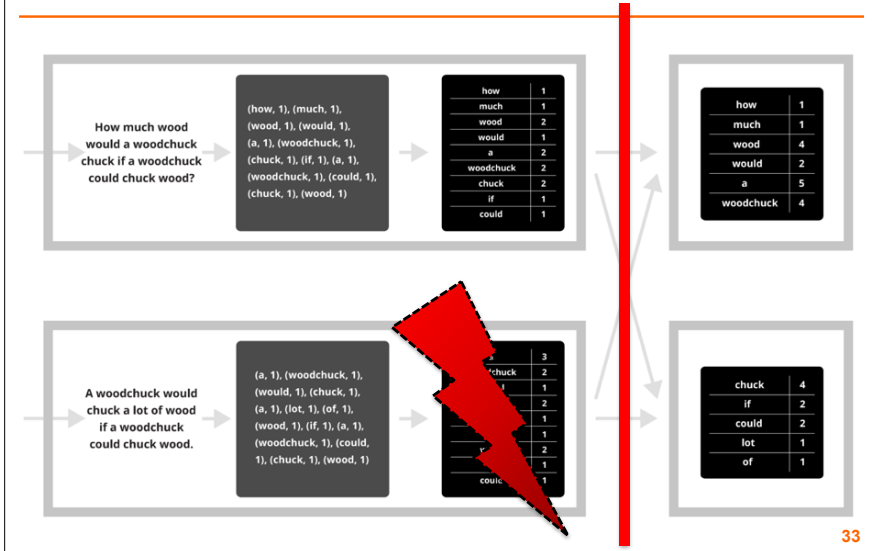


32

32

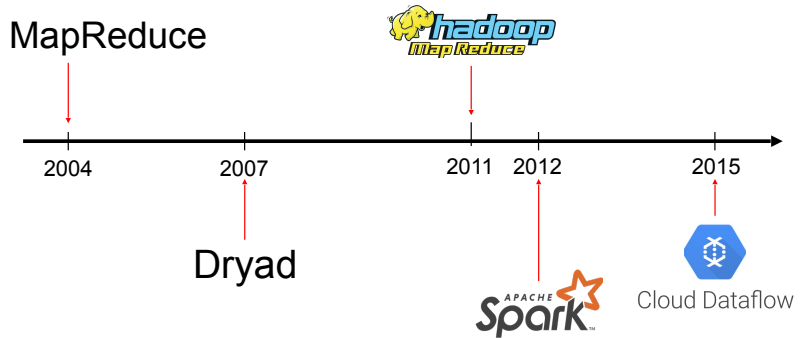


## Synchronization Barrier



33

## MapReduce



34

## Brainstorm: Top K

Top K is the problem of finding the largest K values from a set of numbers

How would you express this as a distributed application?

In particular, what would the map and reduce phases look like?

*Hint: use a heap...*

35

35

## Brainstorm: Top K

*Assuming that a set of K integers fit in memory...*

Key idea...

Map phase: everyone maintains a heap of K elements

Reduce phase: merge the heaps until you're left with one

36

36

## Brainstorm: Top K

---

Problem: What are the keys and values here?

No notion of key here, just assign the same key to all the values (e.g. key = 1)

Map task 1: [10, 5, 3, 700, 18, 4] → (1, heap(700, 18, 10))

Map task 2: [16, 4, 523, 100, 88] → (1, heap(523, 100, 88))

Map task 3: [3, 3, 3, 3, 300, 3] → (1, heap(300, 3, 3))

Map task 4: [8, 15, 20015, 89] → (1, heap(20015, 89, 15))

Then all the heaps will go to a single reducer responsible for the key 1

This works, but clearly not scalable...

37

37

## Brainstorm: Top K

---

Idea: Use X different keys to balance load (e.g. X = 2 here)

Map task 1: [10, 5, 3, 700, 18, 4] → (1, heap(700, 18, 10))

Map task 2: [16, 4, 523, 100, 88] → (1, heap(523, 100, 88))

Map task 3: [3, 3, 3, 3, 300, 3] → (2, heap(300, 3, 3))

Map task 4: [8, 15, 20015, 89] → (2, heap(20015, 89, 15))

Then all the heaps will (hopefully) go to X different reducers

Rinse and repeat (*what's the runtime complexity?*)

38

38

## Case Study: Spark

(Data-parallel programming at scale)

39

39

## What is Spark ?

---

General distributed data execution engine

Key optimizations

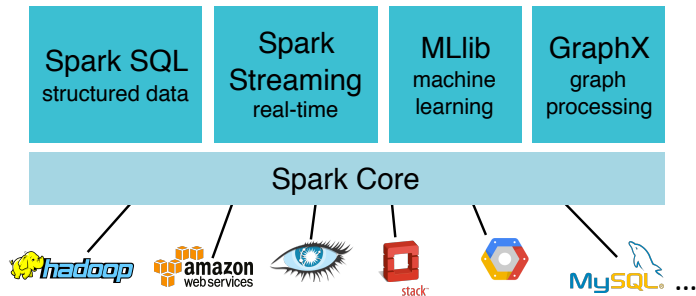
- General computation graphs (pipelining, lazy execution)
- In-memory data sharing (caching)
- Fine-grained fault tolerance

#1 most active big data project

40

40

# What is Spark?



41

41

# Spark computational model

Most computation can be expressed in terms of two phases

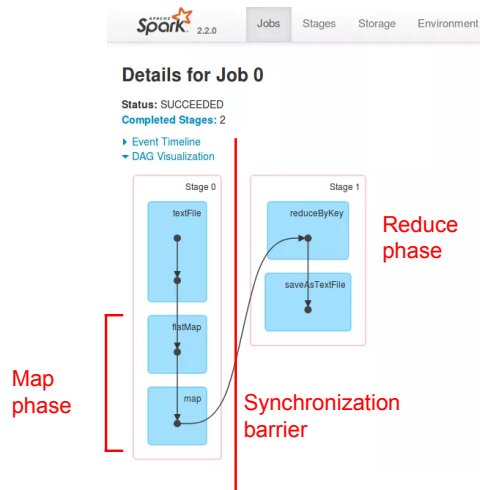
**Map phase** defines how each machine processes its individual partition

**Reduce phase** defines how to merge map outputs from previous phase

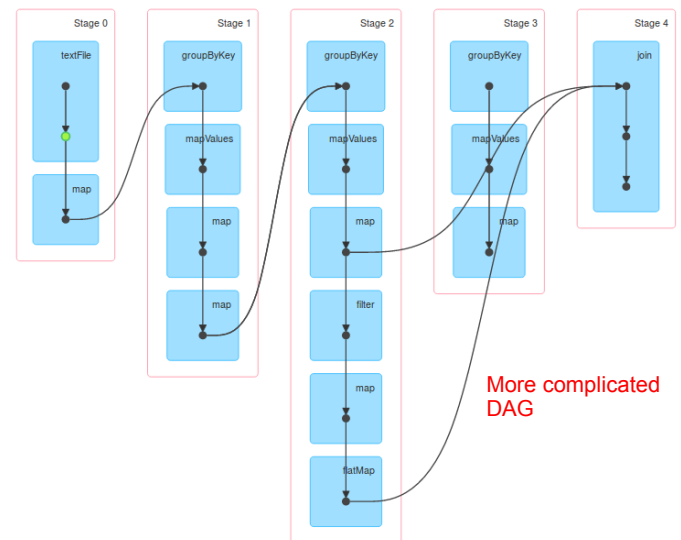
*Spark expresses computation as a DAG of maps and reduces*

42

42



43



44

## Spark: Word count

---

**Transformations** express how to process a dataset

**Actions** express how to turn a transformed dataset into results

```
sc.textFile("declaration-of-independence.txt")
  .flatMap { line => line.split(" ") }
  .map { word => (word, 1) }
  .reduceByKey { case (counts1, counts2) => counts1 + counts2 }
  .collect()
```

45

45

## Pipelining + Lazy execution

---

Transformations can be pipelined until we hit

A synchronization barrier (e.g. reduce), or

An action

*Example:*

```
data.map { ... }.filter { ... }.flatMap { ... }.groupByKey().count()
```

These **three operations** can all be run in the same task

This allows **lazy execution**; we don't need to eagerly execute map

46

46

## In-memory caching

---

Store intermediate results in **memory** to bypass disk access

Important for iterative workloads (e.g. machine learning!)

*Example:*

```
val cached = data.map { ... }.filter { ... }.cache()
(1 to 100).foreach { i =>
  cached.reduceByKey { ... }.saveAsTextFile(...)
}
```

47

47

## Reusing map outputs

---

Reusing map outputs allows Spark to avoid redoing map stages

Along with caching, this makes iterative workloads much faster

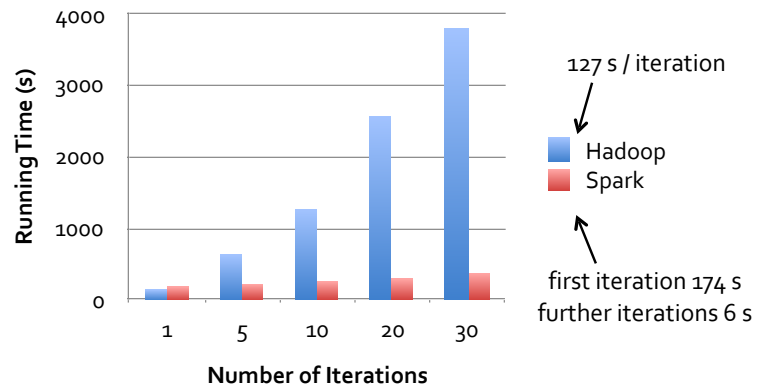
*Example:*

```
val transformed = data.map { ... }.reduceByKey { ... }
transformed.collect()
transformed.collect() // does not run map phase again
```

48

48

## Logistic Regression Performance



[Borrowed from Matei Zaharia's talk, "Spark: Fast, Interactive, Language-Integrated Cluster Computing"]

49

# Monday

## Stream processing

50

50

## Application: Word Count

```
SELECT count(word) FROM data  
GROUP BY word
```

```
cat data.txt  
| tr -s '[:punct:][:space:]' '\n'  
| sort | uniq -c
```

51

51

## Using partial aggregation

1. Compute word counts from individual files
2. Then merge intermediate output
3. Compute word count on merged outputs

52

52

## Using partial aggregation

---

1. In parallel, send to worker:
  - Compute word counts from individual files
  - Collect result, wait until all finished
2. Then merge intermediate output
3. Compute word count on merged intermediates

53

53

## MapReduce: Programming Interface

---

```
map(key, value) -> list(<k', v'>)
```

- Apply function to (key, value) pair and produces set of intermediate pairs

```
reduce(key, list<value>) -> <k', v'>
```

- Applies aggregation function to values
- Outputs result

54

54

## MapReduce: Programming Interface

---

```
map(key, value):  
  for each word w in value:  
    EmitIntermediate(w, "1");
```

```
reduce(key, list(values):  
  int result = 0;  
  for each v in values:  
    result += ParseInt(v);  
  Emit(AsString(result));
```

55

55

## MapReduce: Optimizations

---

```
combine(list<key, value>) -> list<k, v>
```

- Perform partial aggregation on mapper node:  
 $\langle \text{the}, 1 \rangle, \langle \text{the}, 1 \rangle, \langle \text{the}, 1 \rangle \rightarrow \langle \text{the}, 3 \rangle$
- reduce() should be commutative and associative

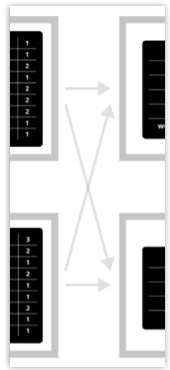
```
partition(key, int) -> int
```

- Need to aggregate intermediate vals with same key
- Given n partitions, map key to partition  $0 \leq i < n$
- Typically via  $\text{hash}(\text{key}) \bmod n$

56

56

## Fault Tolerance in MapReduce



- Map worker writes intermediate output to local disk, separated by partitioning. Once completed, tells master node.
- Reduce worker told of location of map task outputs, pulls their partition's data from each mapper, execute function across data
- Note:
  - “All-to-all” shuffle b/w mappers and reducers
  - Written to disk (“materialized”) b/w *each* stage

57

57

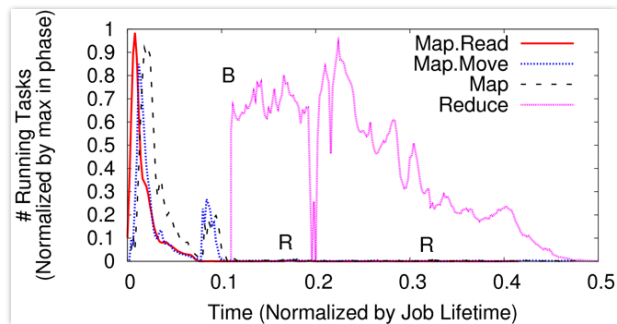
## Fault Tolerance in MapReduce

- Master node monitors state of system
  - If master failures, job aborts and client notified
- Map worker failure
  - Both in-progress/completed tasks marked as idle
  - Reduce workers notified when map task is re-executed on another map worker
- Reducer worker failure
  - In-progress tasks are reset to idle (and re-executed)
  - Completed tasks had been written to global file system

58

58

## Straggler Mitigation in MapReduce



- Tail latency means some workers finish late
- For slow map tasks, execute in parallel on second map worker as “backup”, race to complete task

59

59