

## Case Study 2: Document Retrieval

# Clustering Documents

Machine Learning for Big Data  
CSE547/STAT548, University of Washington

Emily Fox

April 16<sup>th</sup>, 2015

©Emily Fox 2015

1

## Document Retrieval

- **Goal:** Retrieve documents of interest
- **Challenges:**
  - Tons of articles out there
  - How should we measure similarity?



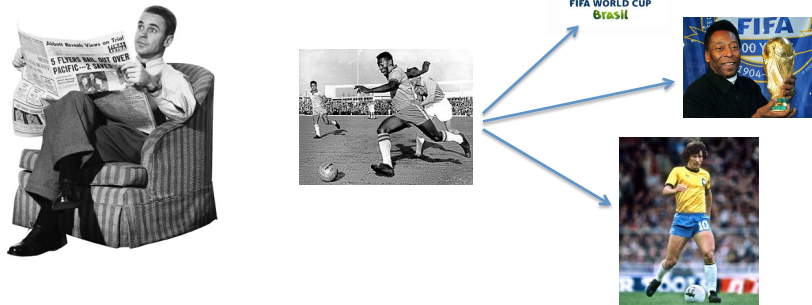
©Emily Fox 2015

2

## Task 1: Find Similar Documents

### ■ So far...

- **Input:** Query article  $x$
- **Output:** Set of  $k$  similar articles  $\{x^{NN_1}, \dots, x^{NN_k}\}$



©Emily Fox 2015

3

## Task 2: Cluster Documents

### ■ Now:

- Cluster documents based on topic



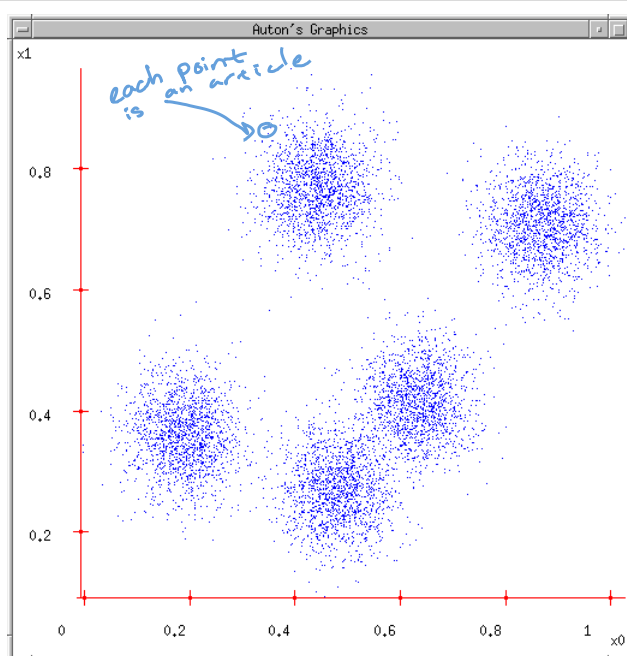
©Emily Fox 2015

4

## Some Data

e.g., doc.  $x$   
 - bag of words  
 - tfidf  
 ...

How to discover  
 clusters?  
 (unsupervised)

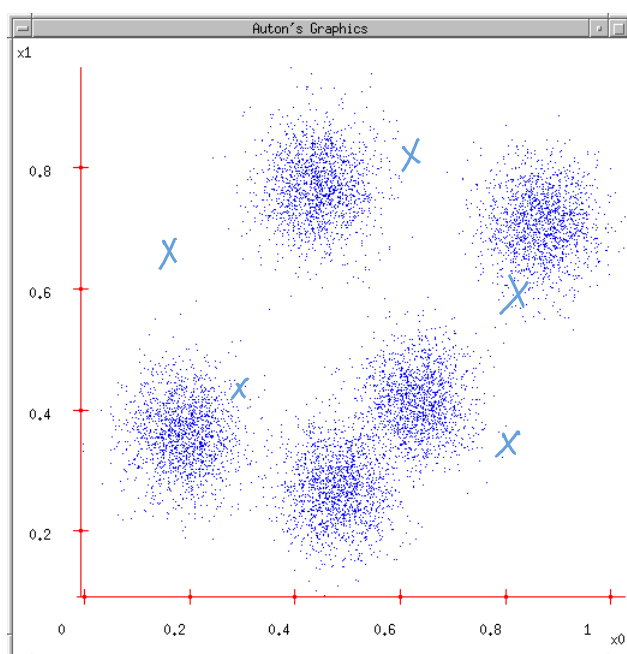


©Emily Fox 2015

5

## K-means

1. Ask user how many clusters they'd like.  
 (e.g.  $k=5$ )

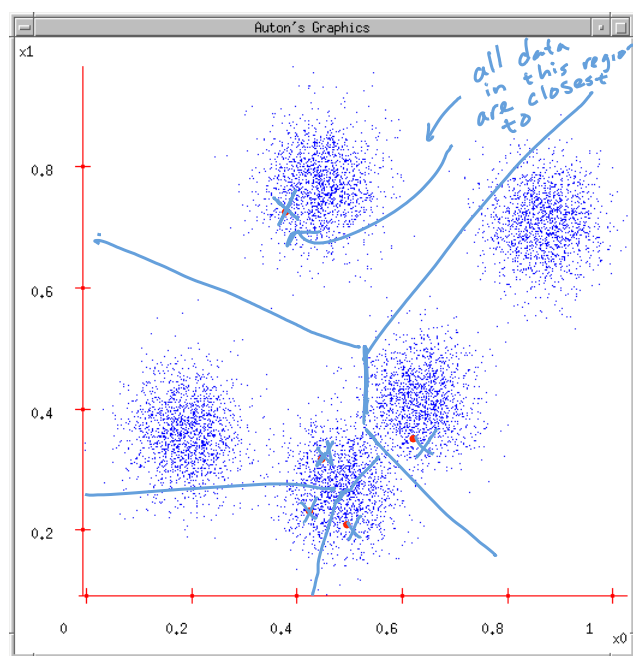


©Emily Fox 2015

6

## K-means

1. Ask user how many clusters they'd like.  
(e.g.  $k=5$ )
2. Randomly guess  $k$  cluster Center locations

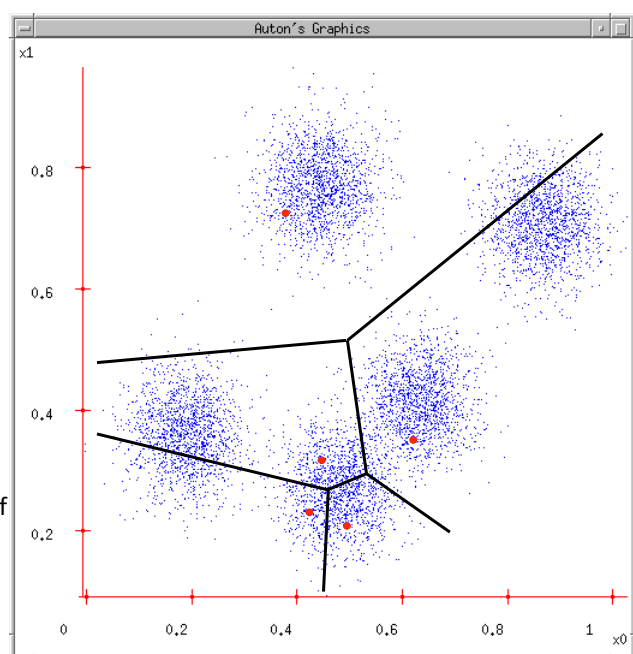


©Emily Fox 2015

7

## K-means

1. Ask user how many clusters they'd like.  
(e.g.  $k=5$ )
2. Randomly guess  $k$  cluster Center locations
3. Each datapoint finds out which Center it's closest to. (Thus each Center "owns" a set of datapoints)

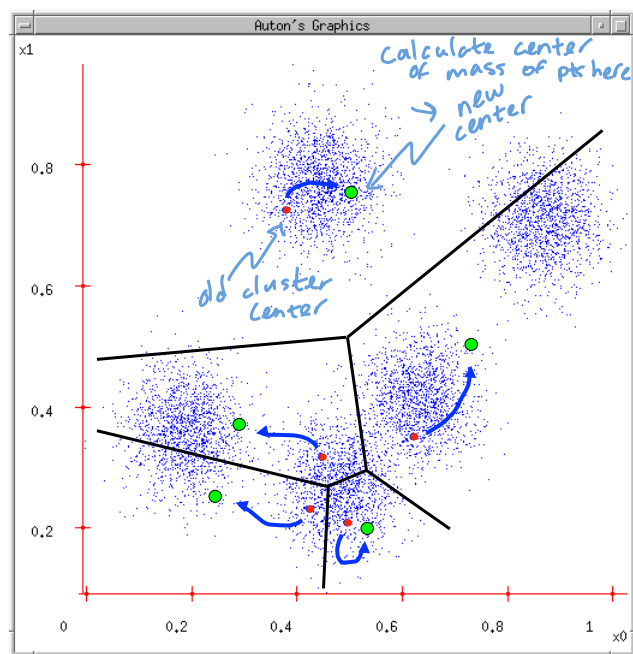


©Emily Fox 2015

8

## K-means

1. Ask user how many clusters they'd like.  
(e.g.  $k=5$ )
2. Randomly guess  $k$  cluster Center locations
3. Each datapoint finds out which Center it's closest to.
4. Each Center finds the centroid of the points it owns

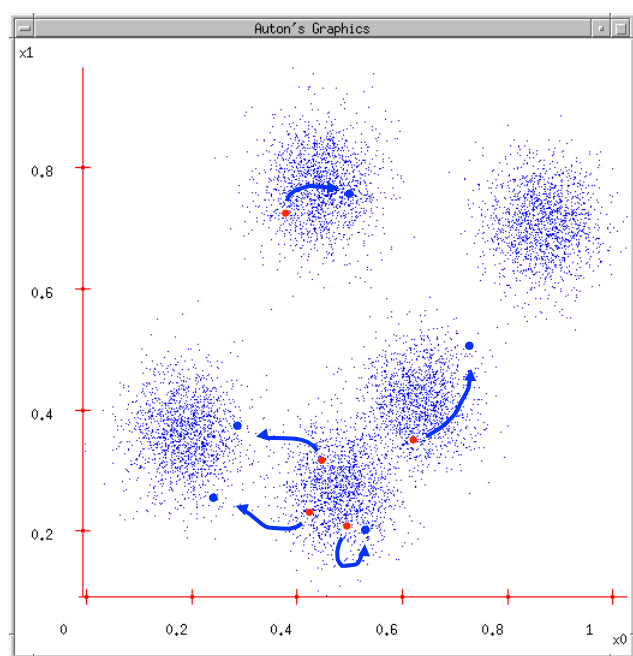


©Emily Fox 2015

9

## K-means

1. Ask user how many clusters they'd like.  
(e.g.  $k=5$ )
2. Randomly guess  $k$  cluster Center locations
3. Each datapoint finds out which Center it's closest to.
4. Each Center finds the centroid of the points it owns...
5. ...and jumps there
6. ...Repeat until terminated!



©Emily Fox 2015

10

## K-means

- Randomly initialize  $k$  centers

$$\mu^{(0)} = \mu_1^{(0)}, \dots, \mu_k^{(0)}$$

- Classify:** Assign each point  $j \in \{1, \dots, N\}$  to nearest center:

$$z^j \leftarrow \arg \min_i \|\mu_i - \mathbf{x}^j\|_2^2$$

- Recenter:**  $\mu_i$  becomes centroid of its point:

$$\mu_i^{(t+1)} \leftarrow \arg \min_{\mu} \sum_{j: z^j = i} \|\mu - \mathbf{x}^j\|_2^2$$

– Equivalent to  $\mu_i \leftarrow$  average of its points!

©Emily Fox 2015

$z^j \in \{1, \dots, k\}$  is indicator of nearest center for obs.  $j$

all pts assigned to  $i$ th cluster  

$$\mu_i = \frac{\sum_{j: z^j = i} \mathbf{x}^j}{|\{j: z^j = i\}|}$$

## Case Study 2: Document Retrieval

### Parallel Programming Map-Reduce

Machine Learning for Big Data  
CSE547/STAT548, University of Washington

Emily Fox

April 16<sup>th</sup>, 2015

©Emily Fox 2015

12

## Needless to Say, We Need Machine Learning for Big Data



6 Billion  
Flickr Photos



28 Million  
Wikipedia Pages



1 Billion  
Facebook Users



72 Hours a Minute  
YouTube

The New York Times  
**Sunday Review**

WORLD U.S. N.Y. / REGION BUSINESS TEC

NEWS ANALYSIS  
**The Age of Big Data**

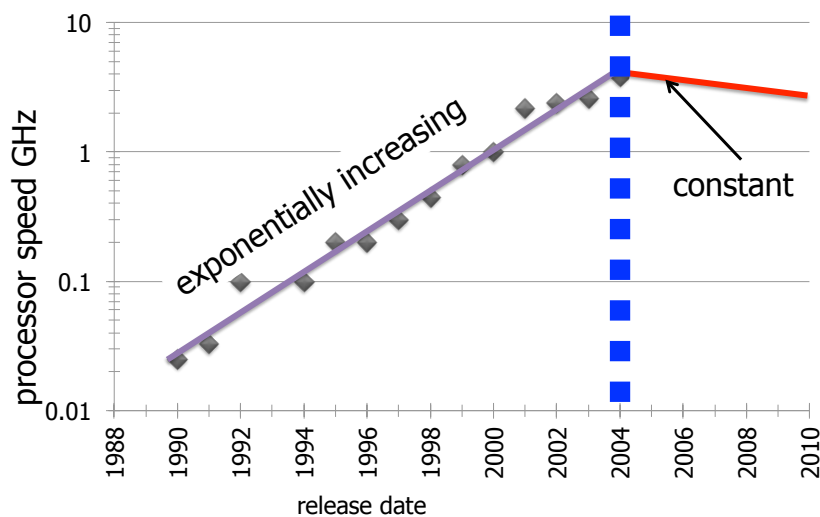
By STEVE LOHR  
Published: February 11, 2012

"... data a new class of economic asset,  
like currency or gold."

©Emily Fox 2015

13

## CPUs Stopped Getting Faster...



©Emily Fox 2015

14

## ML in the Context of Parallel Architectures

*use more processors*



- But scalable ML in these systems is hard, especially in terms of:
  1. Programmability
  2. Data distribution
  3. Failures

*we'll go through these ideas...*

©Emily Fox 2015

15

## Programmability Challenge 1: Designing Parallel Programs

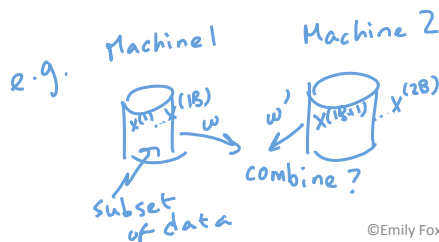
- SGD for LR:
  - For each data point  $\mathbf{x}^{(t)}$ :

$$w_i^{(t+1)} \leftarrow w_i^{(t)} + \eta_t \left\{ -\lambda w_i^{(t)} + \phi_i(\mathbf{x}^{(t)}) [y^{(t)} - P(Y=1|\phi(\mathbf{x}^{(t)}), \mathbf{w}^{(t)})] \right\}$$

$$w^{(0)} \xrightarrow{x^{(1)}} w^{(1)} \xrightarrow{x^{(2)}} w^{(2)} \xrightarrow{x^{(3)}} w^{(3)} \rightarrow \dots$$

*How do we parallelize?*

*weights  $w^{(t)}$  updated sequentially*






©Emily Fox 2015

16



## Programmability Challenge 2: Race Conditions

- We are used to sequential programs:
    - Read data, think, write data, read data, think, write data, read data, think, write data, read data, think, write data, read data, think, write data...
  - But, in parallel, you can have non-deterministic effects:
    - One machine reading data while other is writing
- ①  → [I hate you] → ... → [I love you]
- ②  → [I love you]
- 
- Called a race-condition:
    - Very annoying
    - One of the hardest problems to debug in practice:
      - because of non-determinism, bugs are hard to reproduce

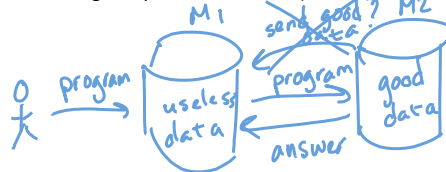
©Emily Fox 2015

17

## Data Distribution Challenge

- Accessing data:
  - Main memory reference: 100ns ( $10^{-7}$ s)
  - Round trip time within data center: 500,000ns ( $5 * 10^{-4}$ s) ← net access
  - Disk seek: 10,000,000ns ( $10^{-2}$ s)
- Reading 1MB sequentially:
  - Local memory: 250,000ns ( $2.5 * 10^{-4}$ s)
  - Network: 10,000,000ns ( $10^{-2}$ s)
  - Disk: 30,000,000ns ( $3 * 10^{-2}$ s)

2 orders of magnitude diff. of memory vs. network / disk
- Conclusion: Reading data from local memory is **much faster** → Must have data locality:
  - Good data partitioning strategy fundamental!
  - “Bring computation to data” (rather than moving data around)



©Emily Fox 2015

18




## Robustness to Failures Challenge

- From Google's Jeff Dean, about their clusters of 1800 servers, in first year of operation:
  - 1,000 individual machine failures
  - thousands of hard drive failures
  - one power distribution unit will fail, bringing down 500 to 1,000 machines for about 6 hours
  - 20 racks will fail, each time causing 40 to 80 machines to vanish from the network
  - 5 racks will "go wonky," with half their network packets missing in action
  - the cluster will have to be rewired once, affecting 5 percent of the machines at any given moment over a 2-day span
  - 50% chance cluster will overheat, taking down most of the servers in less than 5 minutes and taking 1 to 2 days to recover
- How do we design distributed algorithms and systems robust to failures?
  - It's not enough to say: run, if there is a failure, do it again... because you may never finish

©Emily Fox 2015

19

## Move Towards Higher-Level Abstraction

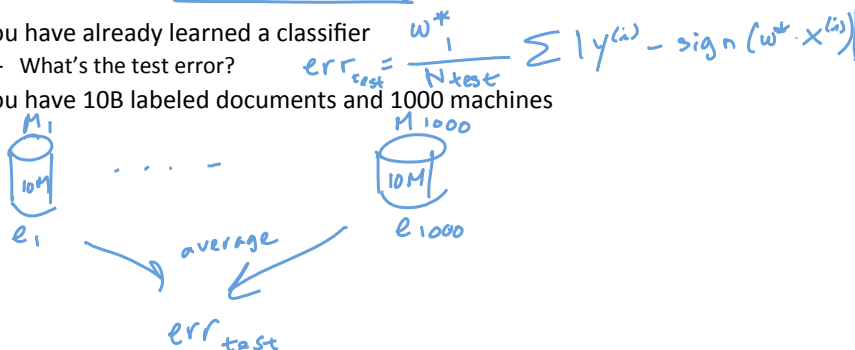
- Distributed computing challenges are hard and annoying!
  1. Programmability
  2. Data distribution
  3. Failures
- High-level abstractions try to simplify distributed programming by hiding challenges:
  - Provide different levels of robustness to failures, optimizing data movement and communication, protect against race conditions...
  - ✱ Generally, you are still on your own WRT designing parallel algorithms
- Some common parallel abstractions:
  - Lower-level:
    - Pthreads: abstraction for distributed threads on single machine
    - MPI: abstraction for distributed communication in a cluster of computers
  - Higher-level:
    - Map-Reduce (Hadoop: open-source version): mostly data-parallel problems
    - GraphLab: for graph-structured distributed problems


©Emily Fox 2015

20

## Simplest Type of Parallelism: Data Parallel Problems

- You have already learned a classifier
  - What's the test error?
- You have 10B labeled documents and 1000 machines

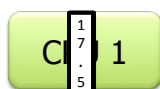


- Problems that can be broken into independent subproblems are called data-parallel (or embarrassingly parallel)
- Map-Reduce is a great tool for this...
  - Focus of today's lecture
  - but first a simple example

©Emily Fox 2015

21

## Data Parallelism (MapReduce)



*Solve a huge number of **independent** subproblems,  
e.g., extract features in images*

## Counting Words on a Single Processor

- (This is the “Hello World!” of Map-Reduce)
- Suppose you have 10B documents and 1 machine
- You want to count the number of appearances of each word in this corpus
  - Similar ideas useful for, e.g., building Naïve Bayes classifiers and computing TF-IDF
- Code:

```

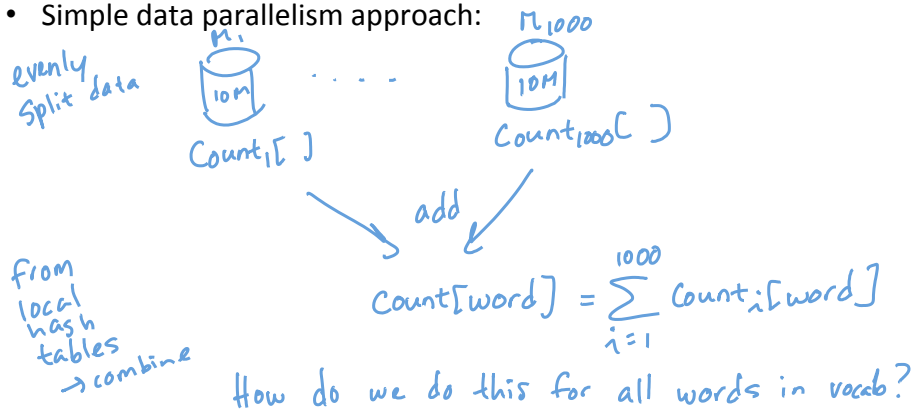
count[ ] ← init a hash table
for d in documents
  for word in d
    count[word] += 1
  
```

©Emily Fox 2015

23

## Naïve Parallel Word Counting

- Simple data parallelism approach:



- Merging hash tables: annoying, potentially not parallel → no gain from parallelism??? have to merge hash tables sequentially

©Emily Fox 2015

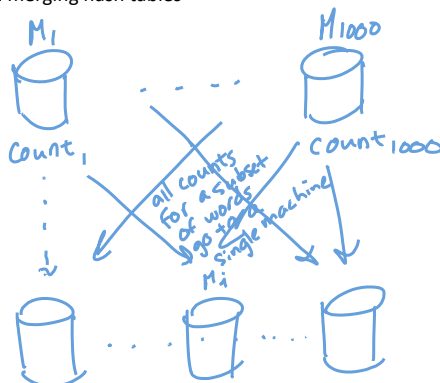
24

## Counting Words in Parallel & Merging Hash Tables in Parallel

- Generate pairs (word,count) ('uw', 17)
- ★ Merge counts for each word in parallel
  - Thus parallel merging hash tables

Phase 1:  
dist. counting

Phase 2:  
parallel sum over words



which words go to machine  $i$ ?  
 $h: \text{vocab} \rightarrow [1, \dots, \# \text{machines}]$   
 send counts of word 'uw' to machine  $h['uw']$

©Emily Fox 2015

25

## Map-Reduce Abstraction

- Map: Transform a data element
  - Data-parallel over elements, e.g., documents
  - Generate (key,value) pairs
    - "value" can be any data type

In our example: ('uw', 1), ('Mary', 1), ('uw', 1)

- Reduce: Take all values associated with a key and aggregate
  - Aggregate values for each key
  - Must be commutative-associate operation
  - Data-parallel over keys
  - Generate (key,value) pairs

reduce ('uw', [1, 17, 0, 0, 12])  
 emit ('uw', 30)

Example:

word count  
 map (document)  
 for word in doc  
 emit (word, 1)

Reduce (word, count: list(int))  
 $c = 0$   
 for  $i$  in count  
 $c += \text{count}[i]$   
 emit (word, c)

- Map-Reduce has long history in functional programming
  - But popularized by Google, and subsequently by open-source Hadoop implementation from Yahoo!

©Emily Fox 2015

26

## Map Code (Hadoop): Word Count

```

public static class Map extends Mapper<LongWritable, Text, Text, IntWritable> {
    private final static IntWritable one = new IntWritable(1);
    private Text word = new Text();

    public void map(LongWritable key, Text value, Context context) throws <stuff>
    {
        String line = value.toString();
        StringTokenizer tokenizer = new StringTokenizer(line);
        while (tokenizer.hasMoreTokens()) {
            word.set(tokenizer.nextToken());
            context.write(word, one);
        }
    }
}

```

Handwritten annotations for Map Code:

- name* (points to `Map`)
- class* (points to `Mapper`)
- For each word in string* (points to the `while` loop)
- word object* (points to `word.set(tokenizer.nextToken());`)
- emit(word, 1)* (points to `context.write(word, one);`)

©Emily Fox 2015

27

## Reduce Code (Hadoop): Word Count

```

public static class Reduce extends Reducer<Text, IntWritable,
Text, IntWritable> {
    public void reduce(Text key, Iterable<IntWritable> values,
        Context context)
        throws IOException, InterruptedException {
        int sum = 0;
        for (IntWritable val : values) {
            sum += val.get();
        }
        context.write(key, new IntWritable(sum));
    }
}

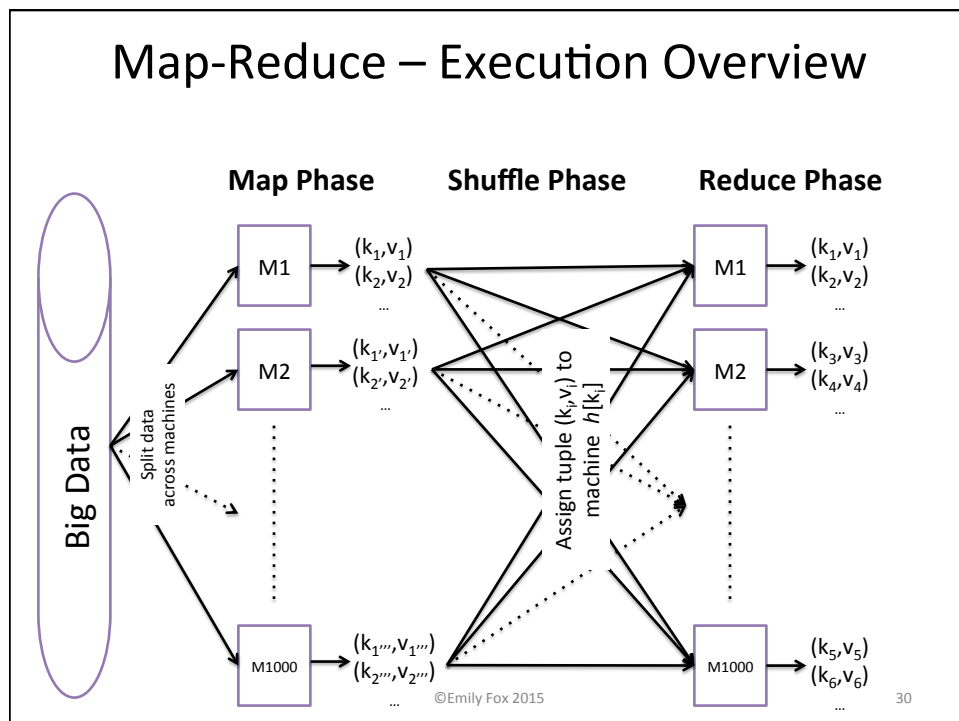
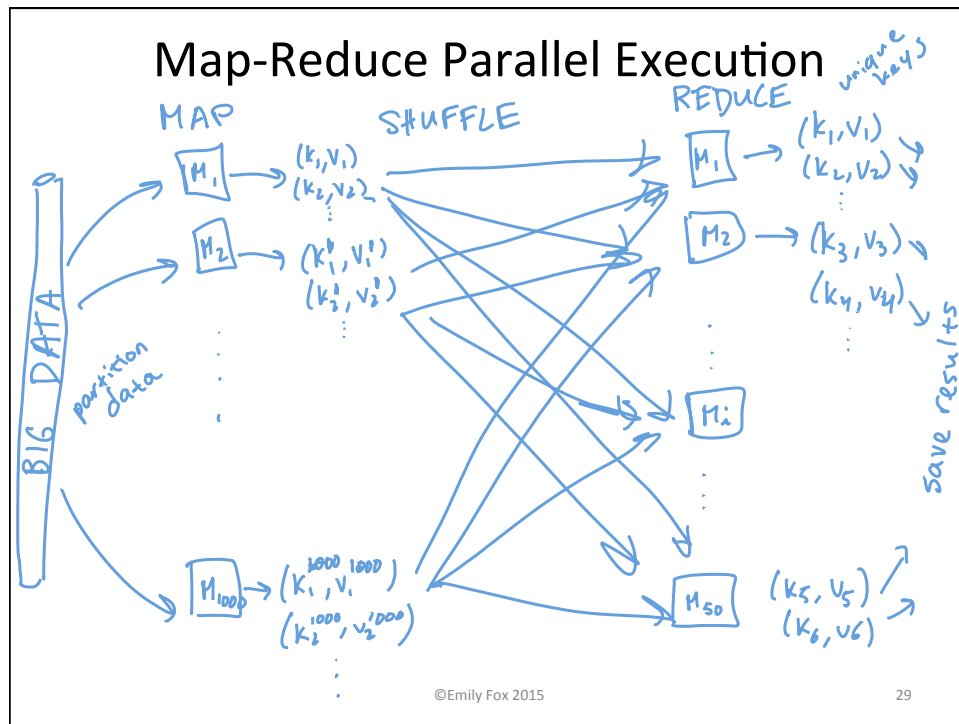
```

Handwritten annotations for Reduce Code:

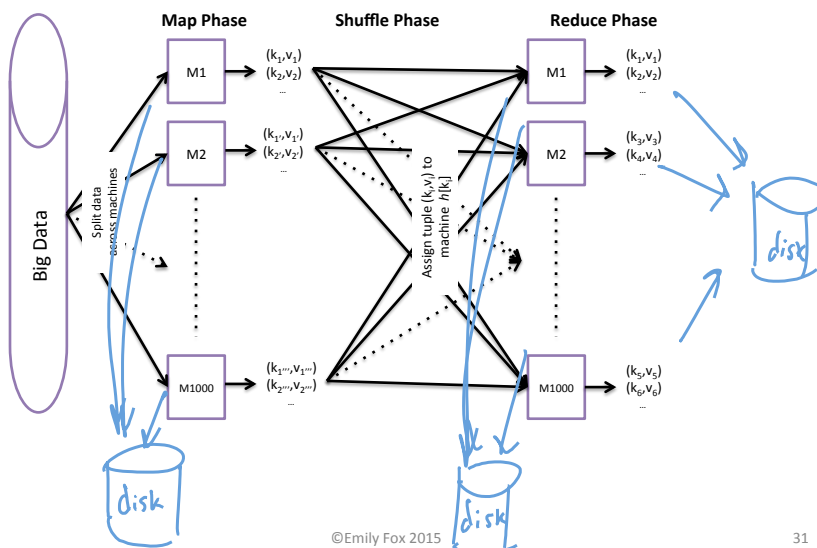
- word: 'uv'* (points to `Text key`)
- int: count* (points to `int sum`)
- values: {5, 17, 0, 0, 12}* (points to `Iterable<IntWritable> values`)
- emit(word, 30)* (points to `context.write(key, new IntWritable(sum));`)

©Emily Fox 2015

28



## Map-Reduce – Robustness to Failures 1: Protecting Data: **Save To Disk Constantly**



## Distributed File Systems

- Saving to disk locally is not enough  $\rightarrow$  If disk or machine fails, all data is lost
- Replicate data among multiple machines!

### Distributed File System (DFS)

- Write a file from anywhere  $\rightarrow$  automatically replicated
- Can read a file from anywhere  $\rightarrow$  read from closest copy

- If failure, try next closest copy

Handwritten notes and diagrams illustrate replication and failure handling:

- usually 3* (referring to the number of replicas)
- write('foo.txt')  $\rightarrow h_i[foo.txt]$  ie  $\{1, 2, 3\}$*  (showing data being written to multiple machines)
- get('foo.txt') :  $h_i('foo.txt') \rightarrow$  where it is next one if failure* (showing a sequence of machines where one is crossed out as 'failed' and the next is used)

### Common implementations:

- Google File System (GFS)
- Hadoop File System (HDFS)

### Important practical considerations:

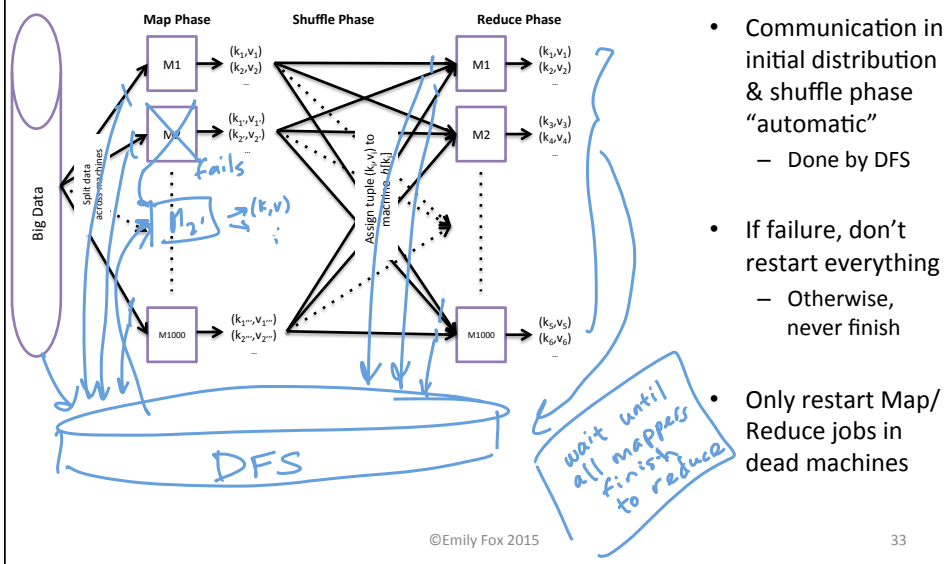
- Write large files
  - Many small files  $\rightarrow$  becomes way too slow
- Typically, files can't be "modified", just "replaced"  $\rightarrow$  makes robustness much simpler

©Emily Fox 2015

32

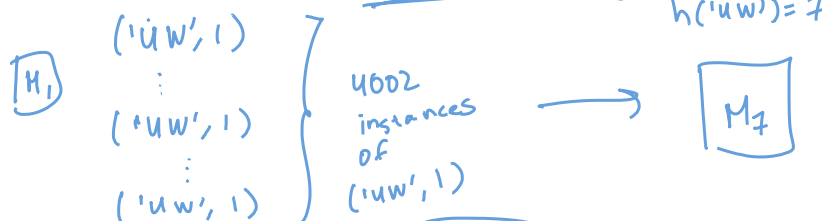


## Map-Reduce – Robustness to Failures 2: Recovering From Failures: **Read from DFS**



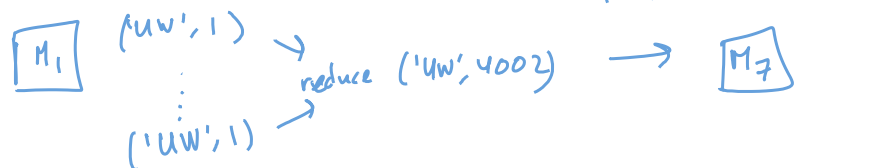
## Improving Performance: Combiners

- Naïve implementation of M-R very wasteful in communication during shuffle:



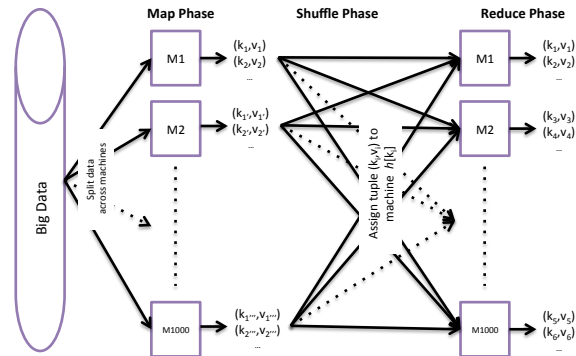
- Combiner:** Simple solution, perform reduce locally before communicating for global reduce

- Works because reduce is commutative-associative



## (A few of the) Limitations of Map-Reduce

- “Too much synchrony”
  - E.g., reducers don't start until all mappers are done
- “Too much” robustness
  - Writing to disk all the time
- Not all problems fit in Map-Reduce
  - E.g., you can't communicate between mappers
- Oblivious to structure in data
  - E.g., if data is a graph, can be much more efficient
    - For example, no need to shuffle nearly as much
- Nonetheless, extremely useful; industry standard for Big Data
  - Though many many companies are moving away from Map-Reduce (Hadoop)



©Emily Fox 2015

35

## What you need to know about Map-Reduce

- Distributed computing challenges are **hard** and **annoying!**
  1. Programmability ⚡
  2. Data distribution
  3. Failures
- High-level abstractions help a lot!
- Data-parallel problems & Map-Reduce
- Map:
  - Data-parallel transformation of data
    - Parallel over data points
- Reduce:
  - Data-parallel aggregation of data
    - Parallel over keys
- Combiner helps reduce communication
- Distributed execution of Map-Reduce:
  - Map, shuffle, reduce
  - Robustness to failure by writing to disk
  - Distributed File Systems

©Emily Fox 2015

36

## Case Study 2: Document Retrieval

# Parallel K-Means on Map-Reduce

Machine Learning for Big Data  
CSE547/STAT548, University of Washington

Emily Fox

April 16<sup>th</sup>, 2015

©Emily Fox 2015

37

## Map-Reducing One Iteration of K-Means

- **Classify:** Assign each point  $j \in \{1, \dots, N\}$  to nearest center:

$$z^j \leftarrow \arg \min_i \|\mu_i - \mathbf{x}^j\|_2^2$$

- **Recenter:**  $\mu_i$  becomes centroid of its point:

$$\mu_i^{(t+1)} \leftarrow \arg \min_{\mu} \sum_{j: z^j = i} \|\mu - \mathbf{x}^j\|_2^2$$

- Equivalent to  $\mu_i \leftarrow$  average of its points!

- **Map:** data parallel : classify phase  
for each data point, given  $(\{\mu_i\}, \mathbf{x}^j)$ , emit  $(z^j, \mathbf{x}^j)$

- **Reduce:** Recenter phase  
average over all pts in cluster  $i$

©Emily Fox 2015

$z^j = i$

38

## Classification Step as Map

- **Classify:** Assign each point  $j \in \{1, \dots, N\}$  to nearest center:

$$z^j \leftarrow \arg \min_i \|\mu_i - \mathbf{x}^j\|_2^2$$

- **Map:**  $\text{map}([\mu_1, \dots, \mu_k], \mathbf{x}^j)$   

$$z^j \leftarrow \arg \min_i \|\mu_i - \mathbf{x}^j\|_2^2$$
  

$$\text{emit}(z^j, \mathbf{x}^j)$$
  
 e.g.  $\text{emit}(2, [17, 0, 0, 1, 7, 2])$   

$$z^j = 2 \text{ (assigned to cluster 2)}$$
  

$$\underbrace{[17, 0, 0, 1, 7, 2]}_{\mathbf{x}^j \text{ vector}}$$

©Emily Fox 2015

39

## Recenter Step as Reduce

- **Recenter:**  $\mu_i$  becomes centroid of its point:

$$\mu_i^{(t+1)} \leftarrow \arg \min_{\mu} \sum_{j: z^j = i} \|\mu - \mathbf{x}^j\|_2^2$$

- Equivalent to  $\mu_i \leftarrow$  average of its points!

- **Reduce:**  $\text{Reduce}(i, \text{list\_x}: [\mathbf{x}^1, \mathbf{x}^3, \dots])$   

$$\text{count} = 0$$
  

$$\text{sum} = 0$$
  
 For  $x$  in list\_x  

$$\text{sum} += x$$
  

$$\text{count} += 1$$
  

$$\text{emit}(i, \text{sum} / \text{count})$$
  

$$\text{avg. of pts in cluster } i$$

©Emily Fox 2015

40

## Some Practical Considerations

- K-Means needs an iterative version of Map-Reduce
  - Not standard formulation
- Mapper needs to get data point and all centers
  - A lot of data!
  - Better implementation: mapper gets many data points

©Emily Fox 2015

41

## What you need to know about Parallel K-Means on Map-Reduce

- Map: **classification step**;  
data parallel over data points
- Reduce: **recompute means**;  
data parallel over centers

©Emily Fox 2015

42