

Search Engine Architecture

5. Big Data Processing

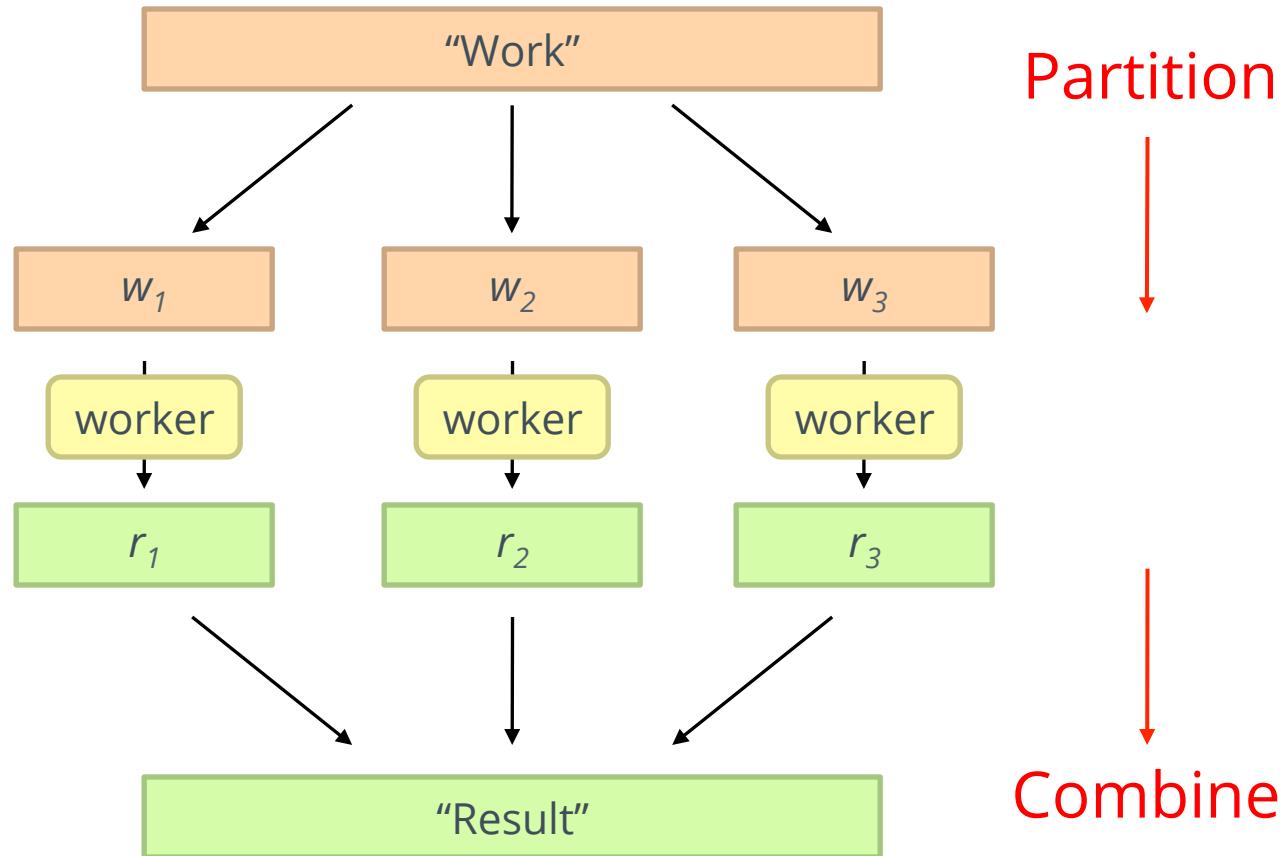


This work is licensed under a Creative Commons Attribution-Noncommercial-Share Alike 3.0 United States
See <http://creativecommons.org/licenses/by-nc-sa/3.0/us/> for details
Noted slides adapted from: Lin et al.'s Big Data Infrastructure, UMD Spring 2015 with cosmetic changes.

Overview

- Parallel computation is hard
 - Because of synchronization
- MapReduce offers a solution
 - (for some applications)
 - Simple, highly-constrained API
 - Sometimes requires creativity to work with
 - No side effects
- Batch index construction can be a good fit

Divide and Conquer



Parallelization Challenges

- How do we assign work units to workers?
- What if we have more work units than workers?
- What if workers need to share partial results?
- How do we aggregate partial results?
- How do we know all the workers have finished?
- What if workers die?

What's the common theme of all of these problems?

Common Theme?

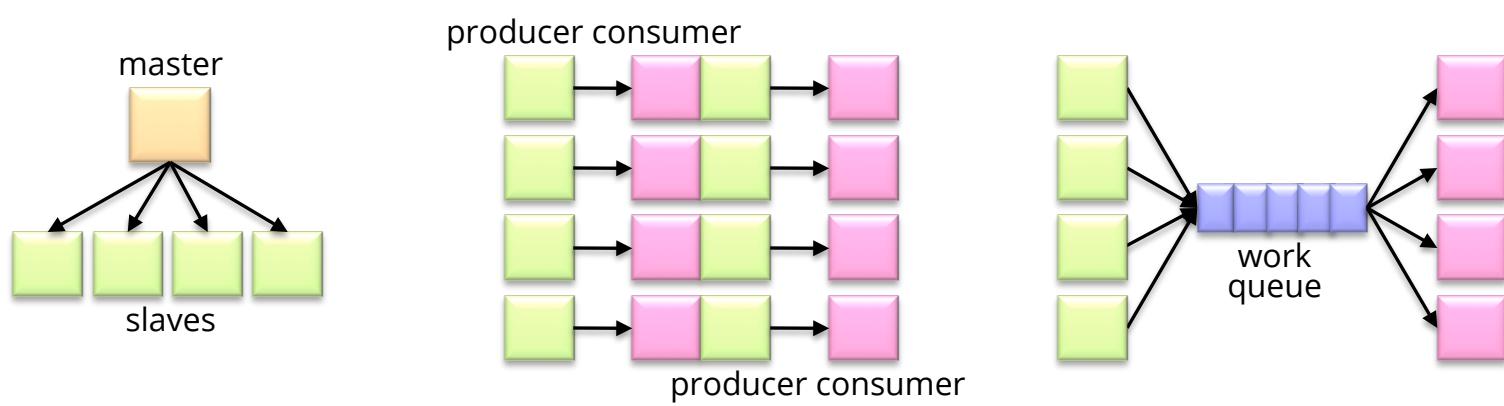
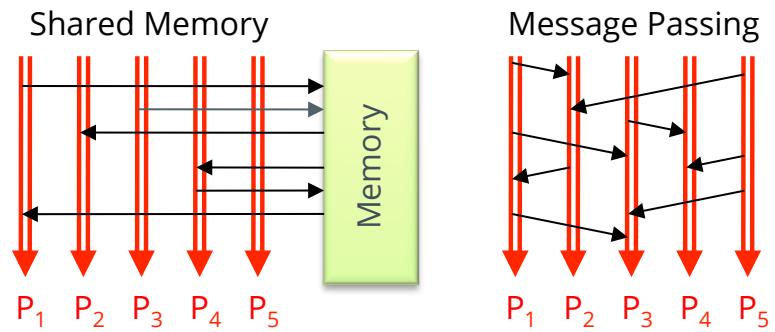
- Parallelization problems arise from:
 - Communication between workers (e.g., to exchange state)
 - Access to shared resources (e.g., data)
- Thus, we need a synchronization mechanism

Managing Multiple Workers

- Difficult because
 - We don't know the order in which workers run
 - We don't know when workers interrupt each other
 - We don't know when workers need to communicate partial results
 - We don't know the order in which workers access shared data
- Thus, we need:
 - Semaphores (lock, unlock)
 - Conditional variables (wait, notify, broadcast)
 - Barriers
- Still, lots of problems:
 - Deadlock, livelock, race conditions...
 - Dining philosophers, sleeping barbers, cigarette smokers...
- Moral of the story: be careful!

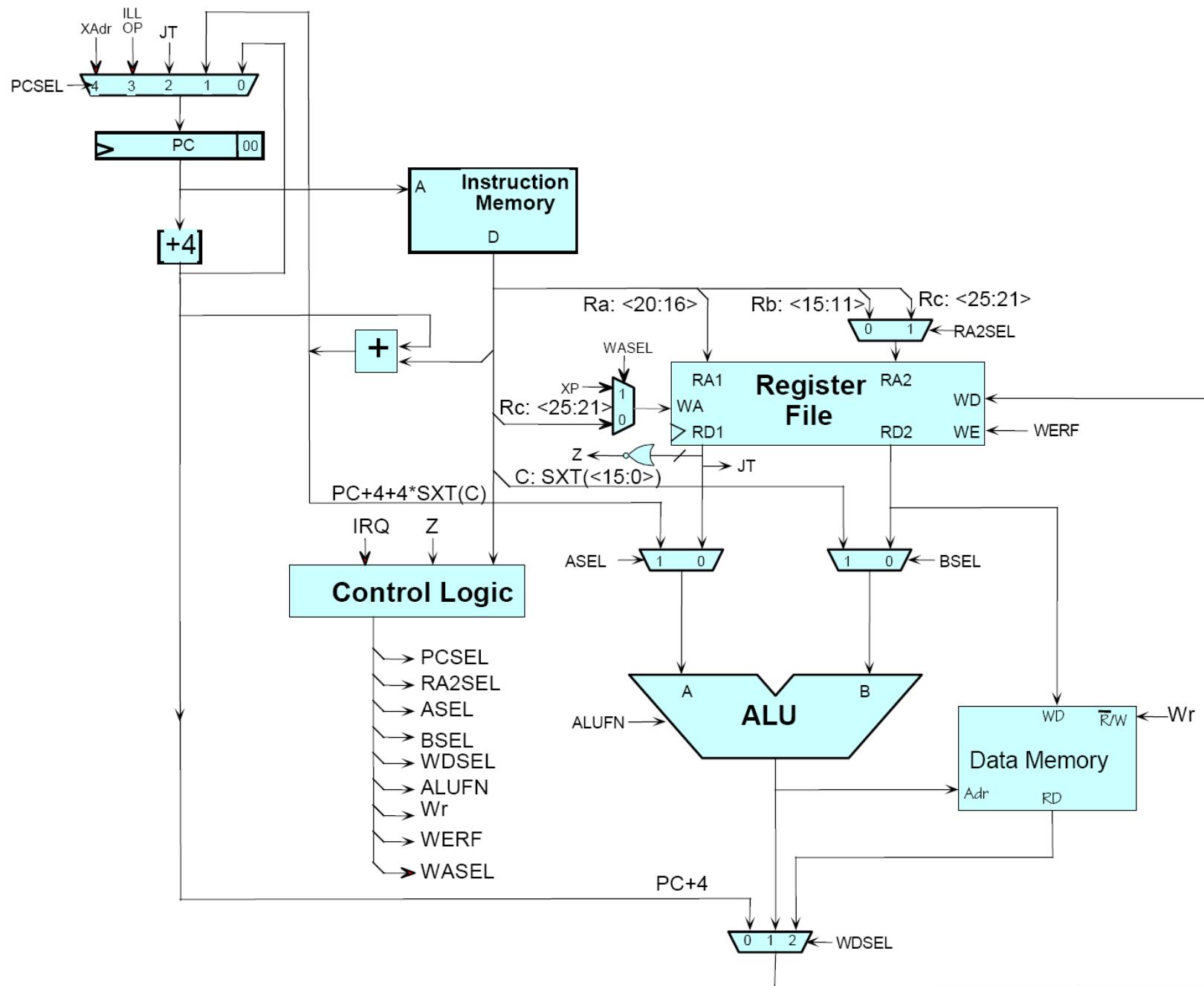
Traditional Tools

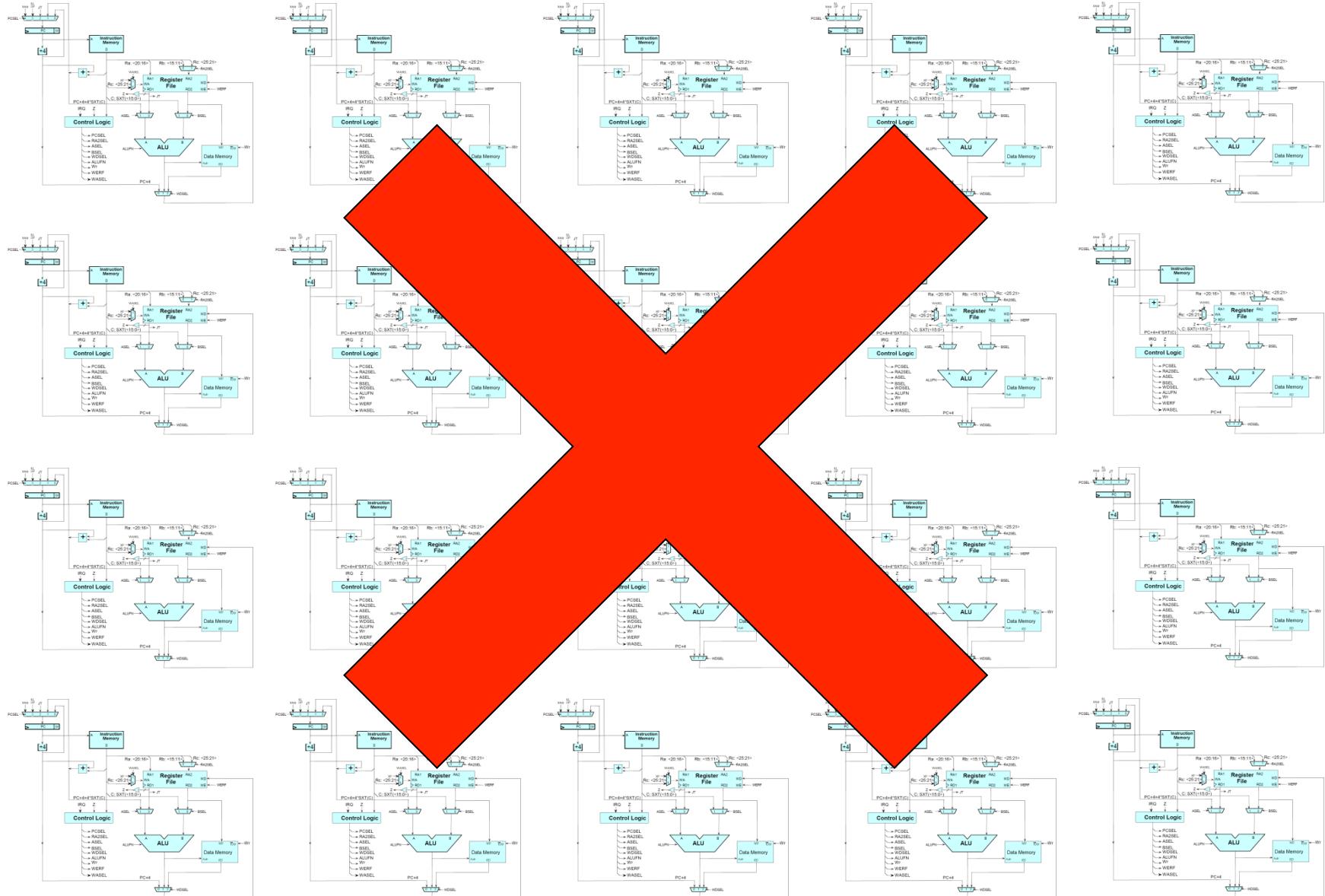
- Programming models
 - Shared memory (pthreads)
 - Message passing (MPI)
- Design Patterns
 - Master-slaves
 - Producer-consumer flows
 - Shared work queues



Where the rubber meets the road

- Concurrency is difficult to reason about
- Concurrency is even more difficult to reason about ...
 - At the scale of datacenters and across datacenters
 - In the presence of failures
 - In terms of multiple interacting services
- Not to mention debugging
- The reality:
 - Lots of one-off solutions, custom code
 - Write your own dedicated library, then program with it
 - Burden on the programmer to explicitly manage everything





Source: Lin et al. Big Data Infrastructure, UMD Spring 2015.



The datacenter *is* the computer!

What's the point?

- It's all about the right level of abstraction
 - Moving beyond the von Neumann architecture
 - We need better programming models
- Hide system-level details from the developers
 - No more race conditions, lock contention, etc.
- Separating the *what* from *how*
 - Developer specifies the computation that needs to be performed
 - Execution framework (“runtime”) handles actual execution

The datacenter *is* the computer!

“Big Ideas”

- Scale “out”, not “up”
 - Limits of SMP and large shared-memory machines
- Move processing to the data
 - Cluster has limited bandwidth
- Process data sequentially, avoid random access
 - Seek times are expensive, disk throughput is reasonable
- Seamless scalability
 - From the mythical man-month to the tradable machine-hour

MapReduce



Source: Lin et al. Big Data Infrastructure, UMD Spring 2019

Typical Big Data Problem

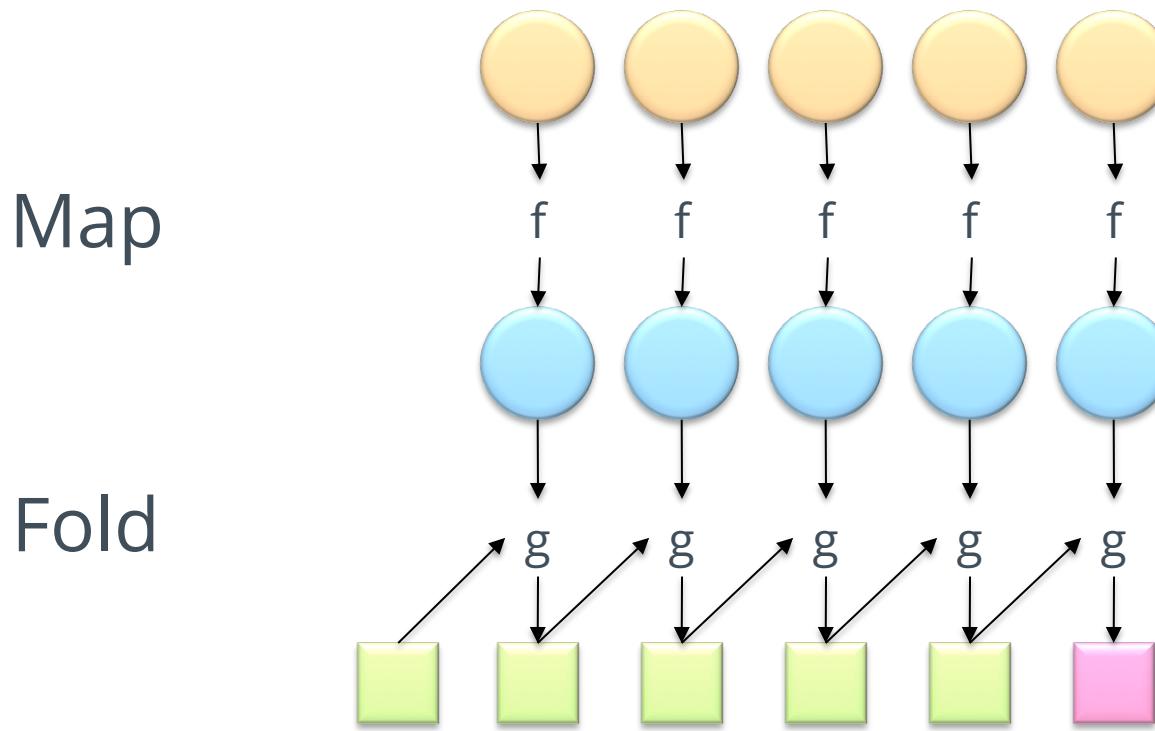
- Iterate over a large number of records
- Extract something of interest from each
- Shuffle and sort intermediate results
- Aggregate intermediate results
- Generate final output

Map

Reduce

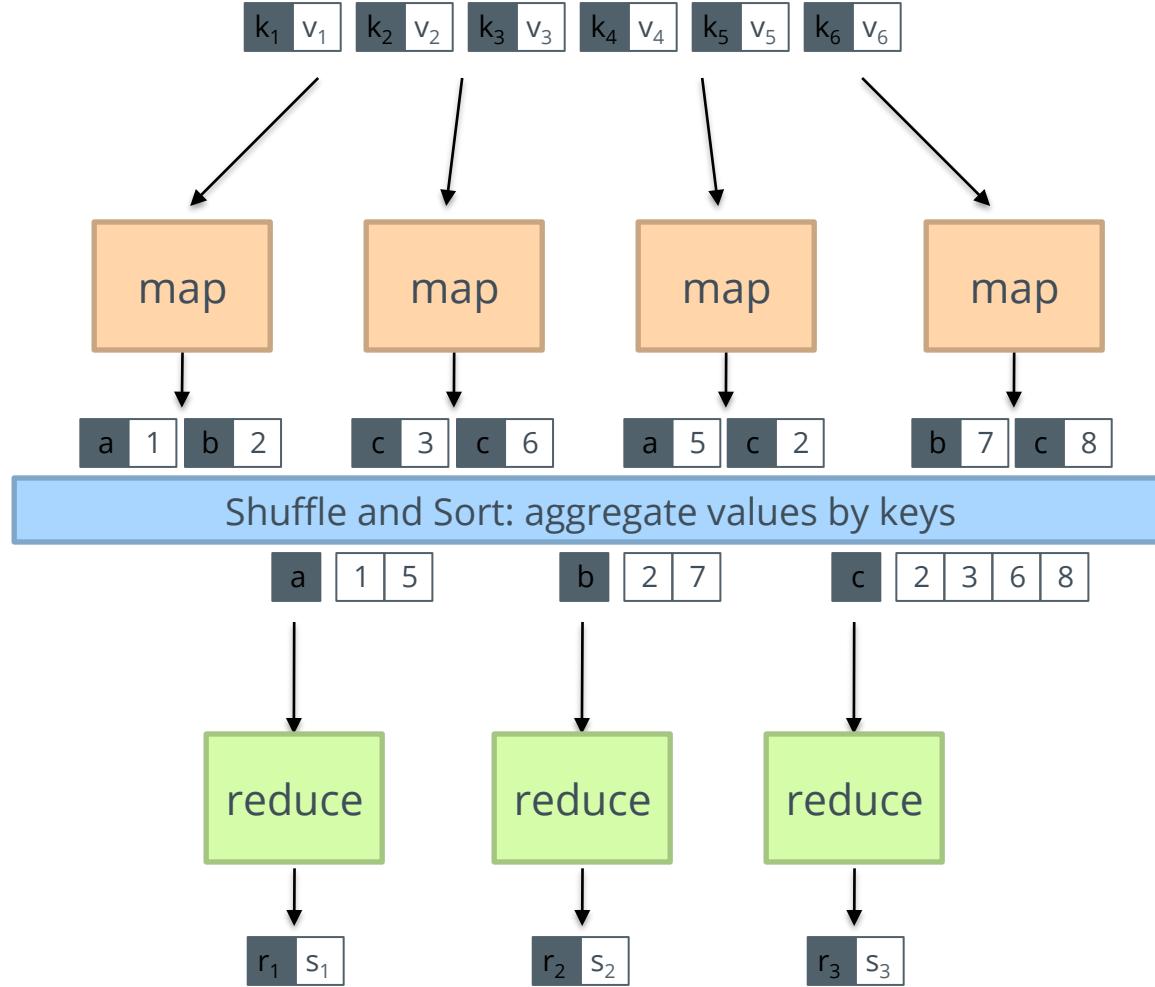
Key idea: provide a functional abstraction for these two operations

Roots in Functional Programming



MapReduce

- Programmers specify two functions:
map (k_1, v_1) \rightarrow [k_2, v_2]
reduce ($k_2, [v_2]$) \rightarrow [k_3, v_3]
 - All values with the same key are sent to the same reducer
 - The execution framework handles everything else

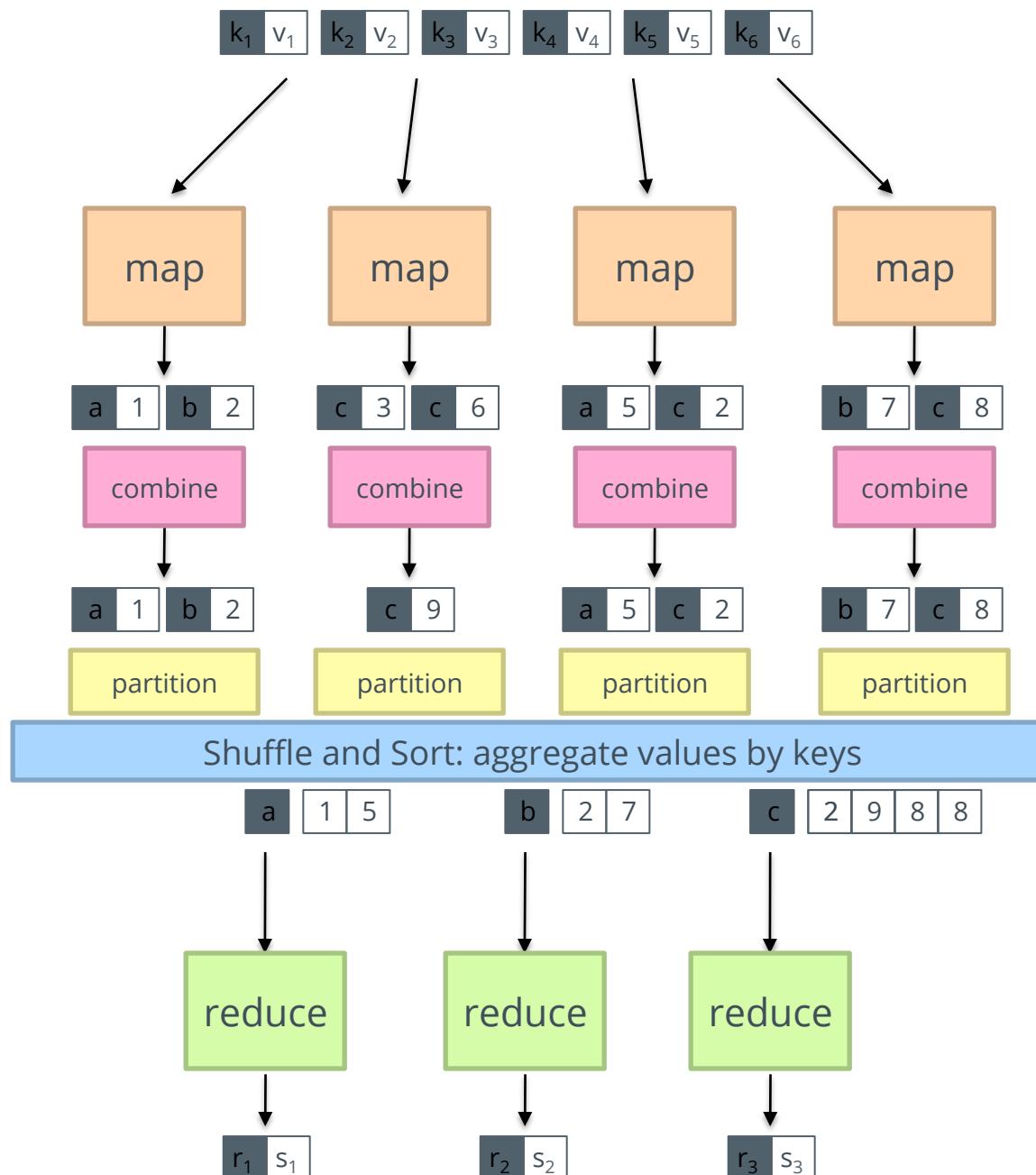


MapReduce “Runtime”

- Handles scheduling
 - Assigns workers to map and reduce tasks
- Handles “data distribution”
 - Moves processes to data
- Handles synchronization
 - Gathers, sorts, and shuffles intermediate data
- Handles errors and faults
 - Detects worker failures and restarts
- Everything happens on top of distributed FS

MapReduce

- Programmers specify two functions:
map $(k, v) \rightarrow \langle k', v' \rangle^*$
reduce $(k', v') \rightarrow \langle k', v' \rangle^*$
 - All values with the same key are reduced together
- The execution framework handles everything else...
- (Not quite.) Usually, programmers also specify:
partition $(k', \text{number of partitions}) \rightarrow \text{partition for } k'$
 - Often a simple hash of the key, e.g., $\text{hash}(k') \bmod n$
 - Divides up key space for parallel reduce operations
combine $(k', v') \rightarrow \langle k', v' \rangle^*$
 - Mini-reducers that run in memory after the map phase
 - Used as an optimization to reduce network traffic



Source: Lin et al. Big Data Infrastructure, UMD Spring 2015.

Two more details...

- Barrier between map and reduce phases
 - But we can begin copying intermediate data earlier
- Keys arrive at each reducer in sorted order
 - No enforced ordering *across* reducers

“Hello World”: Word Count

```
1: class MAPPER
2:   method MAP(docid a, doc d)
3:     for all term t  $\in$  doc d do
4:       EMIT(term t, count 1)

1: class REDUCER
2:   method REDUCE(term t, counts [c1, c2, ...])
3:     sum  $\leftarrow$  0
4:     for all count c  $\in$  counts [c1, c2, ...] do
5:       sum  $\leftarrow$  sum + c
6:     EMIT(term t, count s)
```

MapReduce can refer to...

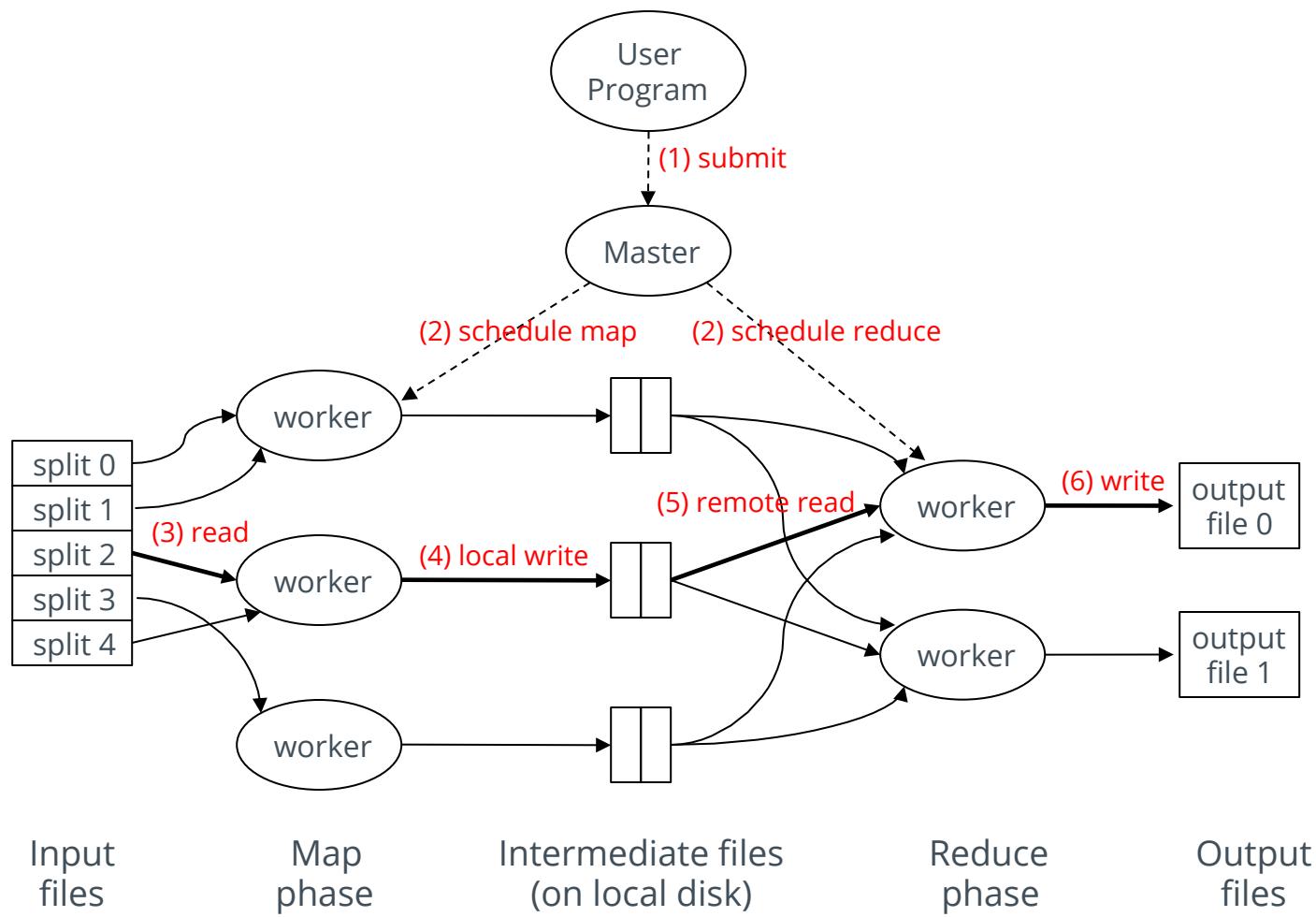
- The programming model
- The execution framework (aka “runtime”)
- The specific implementation

Usage is usually clear from context!

MapReduce Implementations

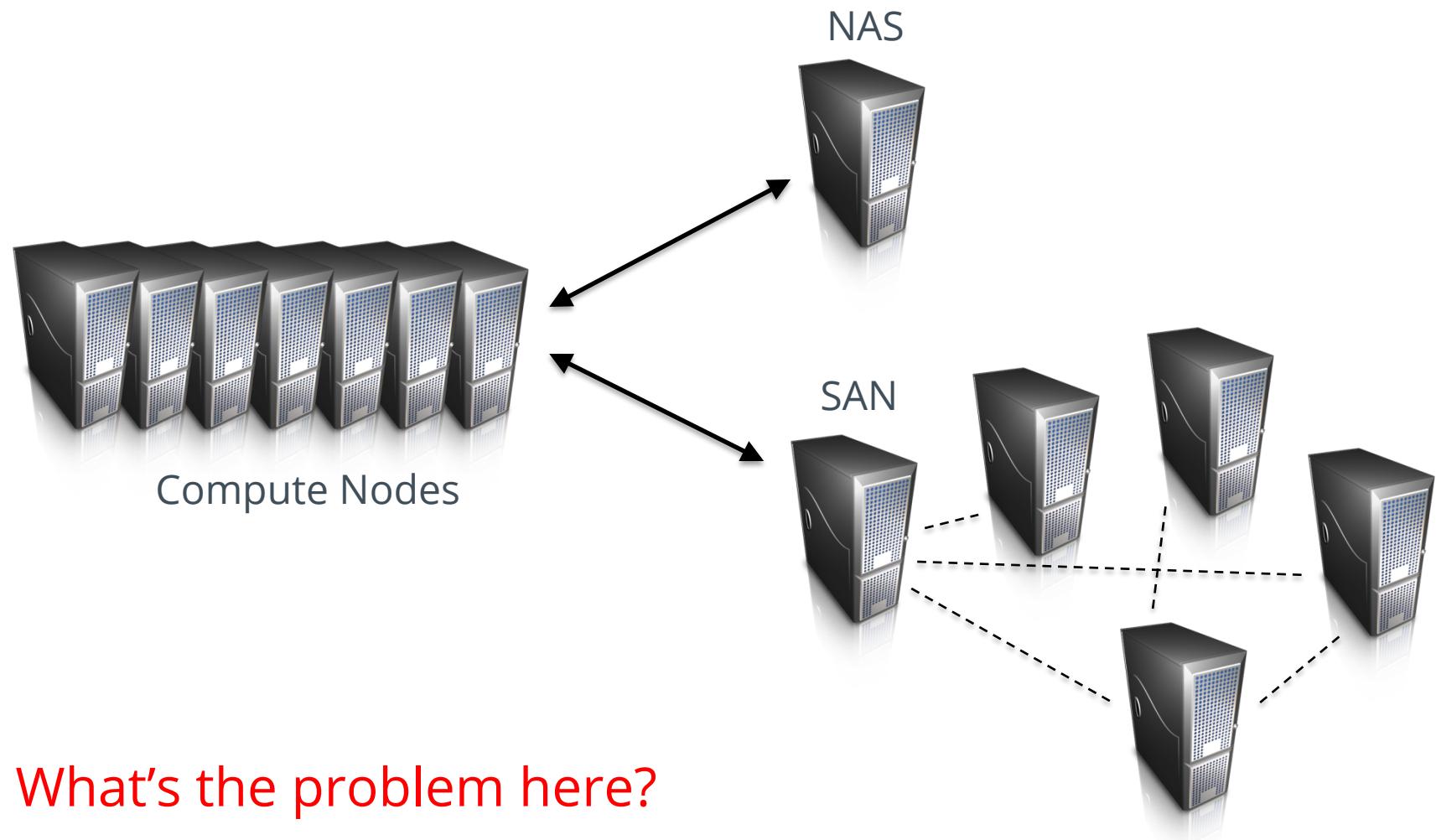
- Google has a proprietary implementation in C++
 - Bindings in Java, Python
- Hadoop is an open-source implementation in Java
 - Development led by Yahoo, now an Apache project
 - Used in production at Yahoo, Facebook, Twitter, LinkedIn, Netflix, ...
 - The *de facto* big data processing platform
 - Rapidly expanding software ecosystem
- Lots of custom research implementations
 - For GPUs, cell processors, etc.





Adapted from (Dean and Ghemawat, OSDI 2004) via Lin et al. Big Data Infrastructure, UMD Spring 2015.

How do we get data to the workers?



What's the problem here?

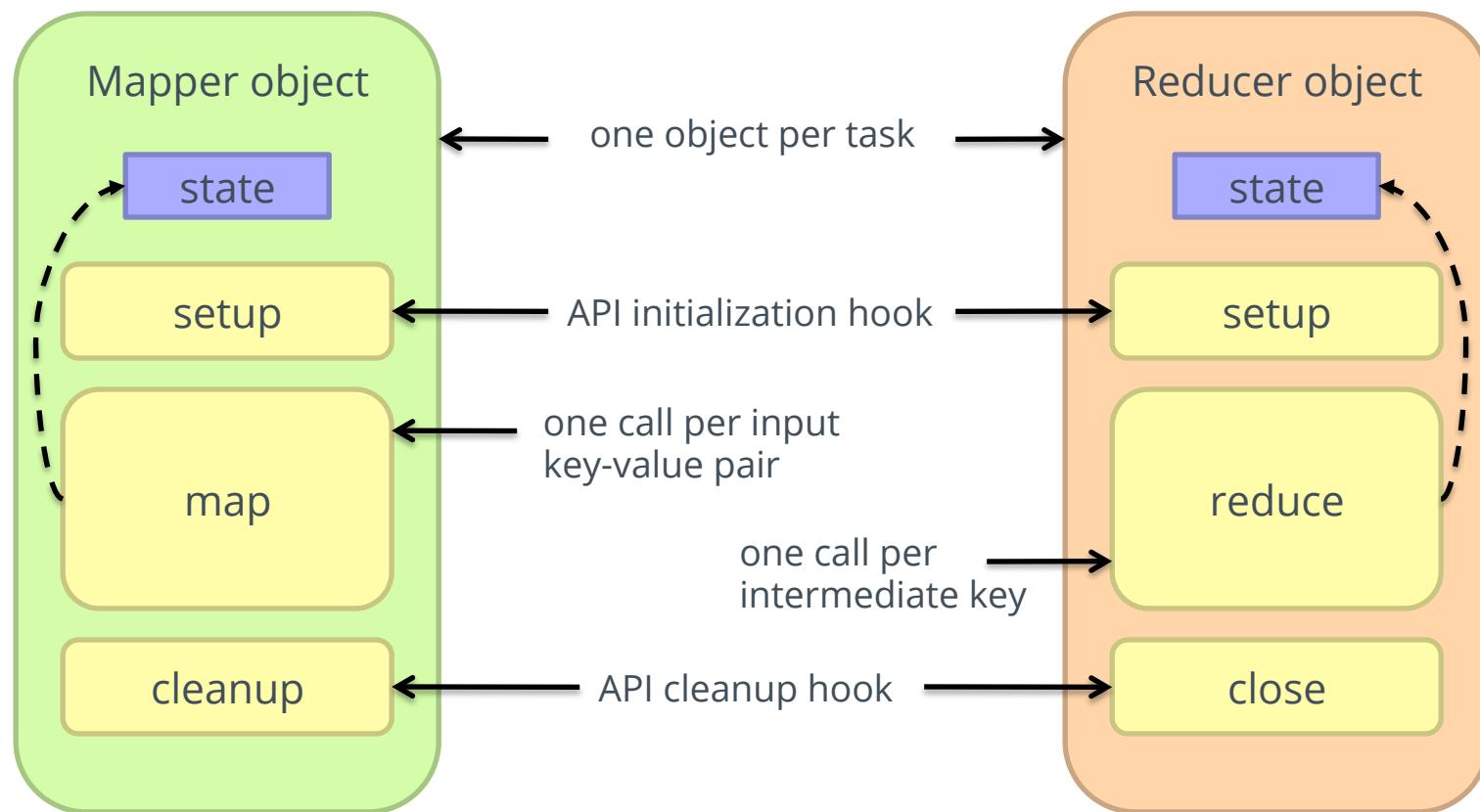
Distributed File System

- Don't move data to workers... move workers to the data!
 - Store data on the local disks of nodes in the cluster
 - Start up the workers on the node that has the data local
- Why?
 - Not enough RAM to hold all the data in memory
 - Disk access is slow, but disk throughput is reasonable
- A distributed file system is the answer
 - GFS (Google File System) for Google's MapReduce
 - HDFS (Hadoop Distributed File System) for Hadoop

Tools for Synchronization

- Cleverly-constructed data structures
 - Bring partial results together
- Sort order of intermediate keys
 - Control order in which reducers process keys
- Partitioner
 - Control which reducer processes which keys
- Preserving state in mappers and reducers
 - Capture dependencies across multiple keys and values

Preserving State



Importance of Local Aggregation

- Ideal scaling characteristics:
 - Twice the data, twice the running time
 - Twice the resources, half the running time
- Why can't we achieve this?
 - Synchronization requires communication
 - Communication kills performance
- Thus... avoid communication!
 - Reduce intermediate data via local aggregation
 - Combiners can help

Word Count: Baseline

```
1: class MAPPER
2:   method MAP(docid a, doc d)
3:     for all term t  $\in$  doc d do
4:       EMIT(term t, count 1)

1: class REDUCER
2:   method REDUCE(term t, counts [c1, c2, ...])
3:     sum  $\leftarrow$  0
4:     for all count c  $\in$  counts [c1, c2, ...] do
5:       sum  $\leftarrow$  sum + c
6:     EMIT(term t, count s)
```

Word Count: Version 1

```
1: class MAPPER
2:   method MAP(docid a, doc d)
3:     H  $\leftarrow$  new ASSOCIATIVEARRAY
4:     for all term t  $\in$  doc d do
5:       H{t}  $\leftarrow$  H{t} + 1                                 $\triangleright$  Tally counts for entire document
6:     for all term t  $\in$  H do
7:       EMIT(term t, count H{t})
```

Word Count: Version 2

```
1: class MAPPER
2:   method INITIALIZE
3:      $H \leftarrow$  new ASSOCIATIVEARRAY
4:   method MAP(docid  $a$ , doc  $d$ )
5:     for all term  $t \in$  doc  $d$  do
6:        $H\{t\} \leftarrow H\{t\} + 1$ 
7:   method CLOSE
8:     for all term  $t \in H$  do
9:       EMIT(term  $t$ , count  $H\{t\}$ )
```

Key idea: preserve state across
input key-value pairs!

▷ Tally counts *across* documents

Design Pattern for Local Aggregation

- “In-mapper combining”
 - Fold the functionality of the combiner into the mapper by preserving state across multiple map calls
- Advantages
 - Speed
 - Why is this faster than actual combiners?
- Disadvantages
 - Explicit memory management required
 - Potential for order-dependent bugs

Algorithm Design: Example

- Term co-occurrence matrix for a text collection
 - $M = N \times N$ matrix ($N =$ vocabulary size)
 - M_{ij} : number of times i and j co-occur in some context
(for concreteness, let's say context = sentence)
- Why?
 - Distributional profiles as a way of measuring semantic distance
 - Semantic distance useful for many language processing tasks

MapReduce: Large Counting Problems

- Term co-occurrence matrix for a text collection
= specific instance of a large counting problem
 - A large event space (number of terms)
 - A large number of observations (the collection itself)
 - Goal: keep track of interesting statistics about the events
- Basic approach
 - Mappers generate partial counts
 - Reducers aggregate partial counts

How do we aggregate partial counts efficiently?

First Try: “Pairs”

- Each mapper takes a sentence:
 - Generate all co-occurring term pairs
 - For all pairs, emit $(a, b) \rightarrow \text{count}$
- Reducers sum up counts associated with these pairs
- Use combiners!

Pairs: Pseudo-Code

```
1: class MAPPER
2:   method MAP(docid a, doc d)
3:     for all term w  $\in$  doc d do
4:       for all term u  $\in$  NEIGHBORS(w) do
5:         EMIT(pair (w, u), count 1)       $\triangleright$  Emit count for each co-occurrence

1: class REDUCER
2:   method REDUCE(pair p, counts [c1, c2, ...])
3:     s  $\leftarrow$  0
4:     for all count c  $\in$  counts [c1, c2, ...] do
5:       s  $\leftarrow$  s + c                       $\triangleright$  Sum co-occurrence counts
6:     EMIT(pair p, count s)
```

“Pairs” Analysis

- Advantages
 - Easy to implement, easy to understand
- Disadvantages
 - Lots of pairs to sort and shuffle around (upper bound?)
 - Not many opportunities for combiners to work

Another Try: “Stripes”

- Idea: group together pairs into an associative array

$(a, b) \rightarrow 1$

$(a, c) \rightarrow 2$

$(a, d) \rightarrow 5$

$(a, e) \rightarrow 3$

$(a, f) \rightarrow 2$

$a \rightarrow \{ b: 1, c: 2, d: 5, e: 3, f: 2 \}$

- Each mapper takes a sentence:
 - Generate all co-occurring term pairs
 - For each term, emit $a \rightarrow \{ b: \text{count}_b, c: \text{count}_c, d: \text{count}_d \dots \}$
- Reducers perform element-wise sum of associative arrays

$$\begin{array}{r} a \rightarrow \{ b: 1, \quad d: 5, e: 3 \} \\ + \quad a \rightarrow \{ b: 1, c: 2, d: 2, \quad f: 2 \} \\ \hline a \rightarrow \{ b: 2, c: 2, d: 7, e: 3, f: 2 \} \end{array}$$

Key idea: cleverly-constructed data structure
brings together partial results

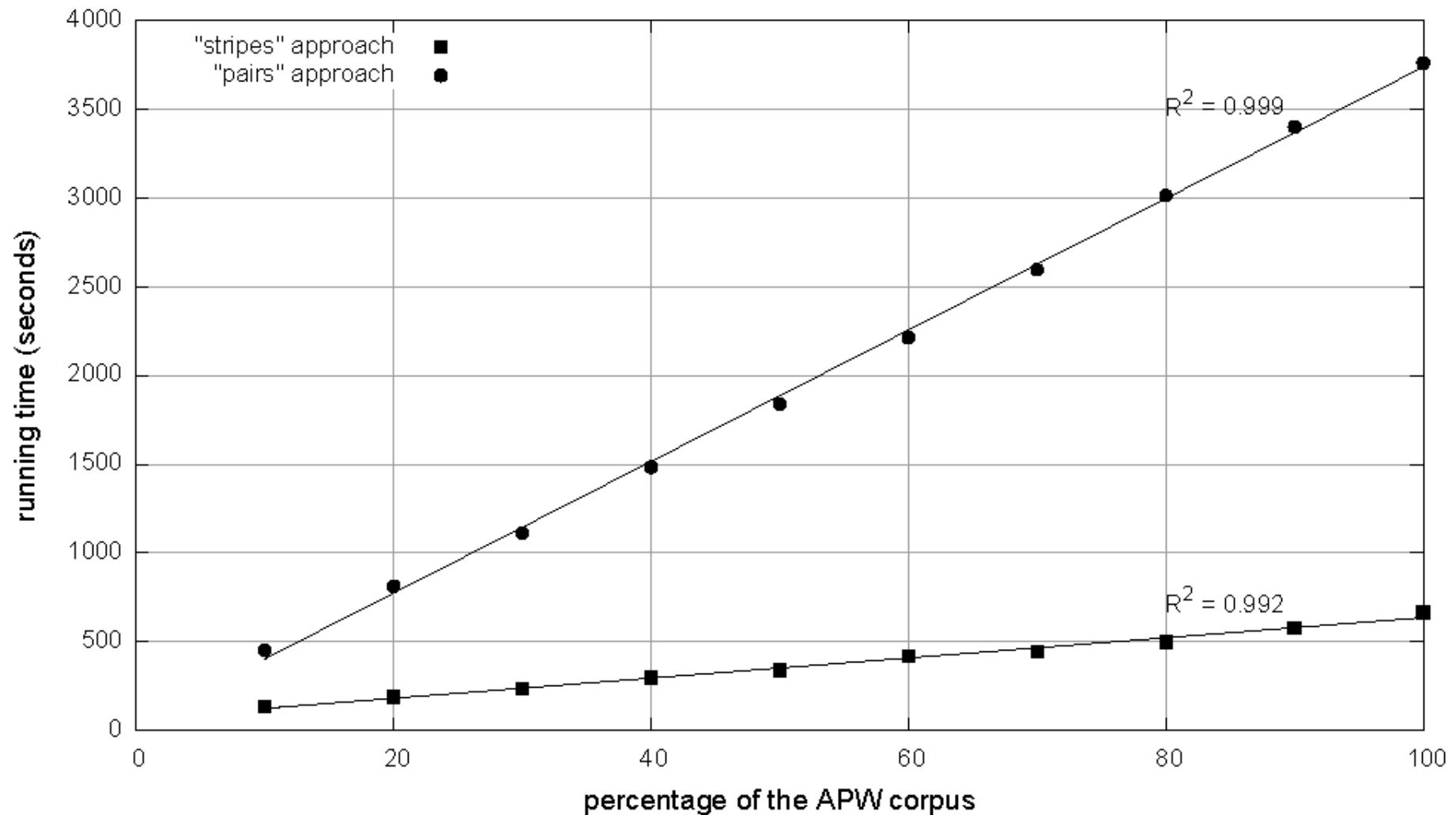
Stripes: Pseudo-Code

```
1: class MAPPER
2:   method MAP(docid  $a$ , doc  $d$ )
3:     for all term  $w \in$  doc  $d$  do
4:        $H \leftarrow$  new ASSOCIATIVEARRAY
5:       for all term  $u \in$  NEIGHBORS( $w$ ) do
6:          $H\{u\} \leftarrow H\{u\} + 1$             $\triangleright$  Tally words co-occurring with  $w$ 
7:       EMIT(Term  $w$ , Stripe  $H$ )
8:
9: class REDUCER
10:   method REDUCE(term  $w$ , stripes [ $H_1, H_2, H_3, \dots$ ])
11:      $H_f \leftarrow$  new ASSOCIATIVEARRAY
12:     for all stripe  $H \in$  stripes [ $H_1, H_2, H_3, \dots$ ] do
13:       SUM( $H_f, H$ )                    $\triangleright$  Element-wise sum
14:     EMIT(term  $w$ , stripe  $H_f$ )
```

“Stripes” Analysis

- Advantages
 - Far less sorting and shuffling of key-value pairs
 - Can make better use of combiners
- Disadvantages
 - More difficult to implement
 - Underlying object more heavyweight
 - Fundamental limitation in terms of size of event space

Comparison of "pairs" vs. "stripes" for computing word co-occurrence matrices

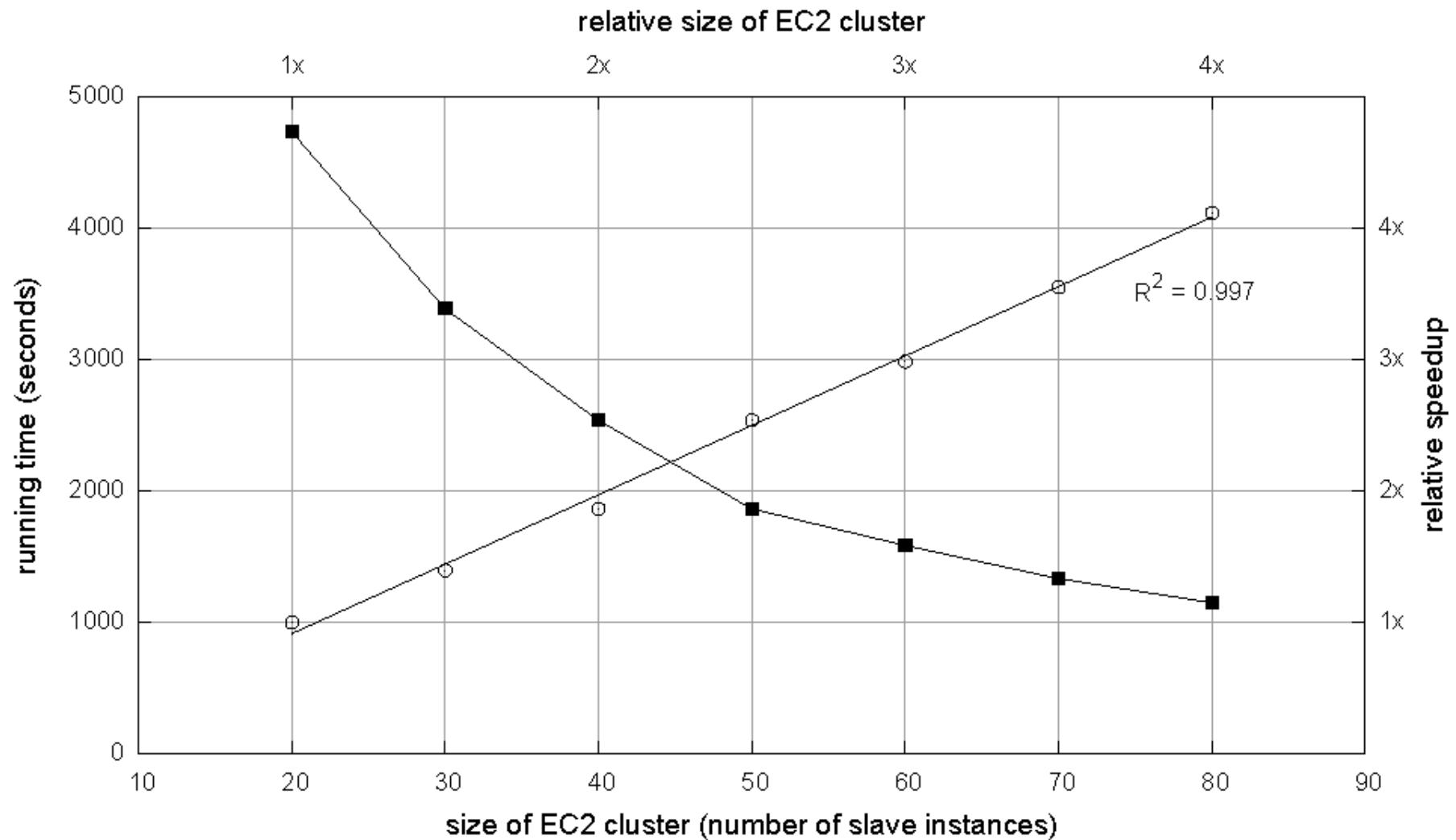


Cluster size: 38 cores

Data Source: Associated Press Worldstream (APW) of the English Gigaword Corpus (v3), which contains 2.27 million documents (1.8 GB compressed, 5.7 GB uncompressed)

Source: Lin et al. Big Data Infrastructure, UMD Spring 2015.

Effect of cluster size on "stripes" algorithm



Debugging at Scale

- Works on small datasets, won't scale... why?
 - Memory management issues (buffering and object creation)
 - Too much intermediate data
 - Mangled input records
- Real-world data is messy!
 - There's no such thing as "consistent data"
 - Watch out for corner cases
 - Isolate unexpected behavior, bring local

Demo

Index Construction with MapReduce

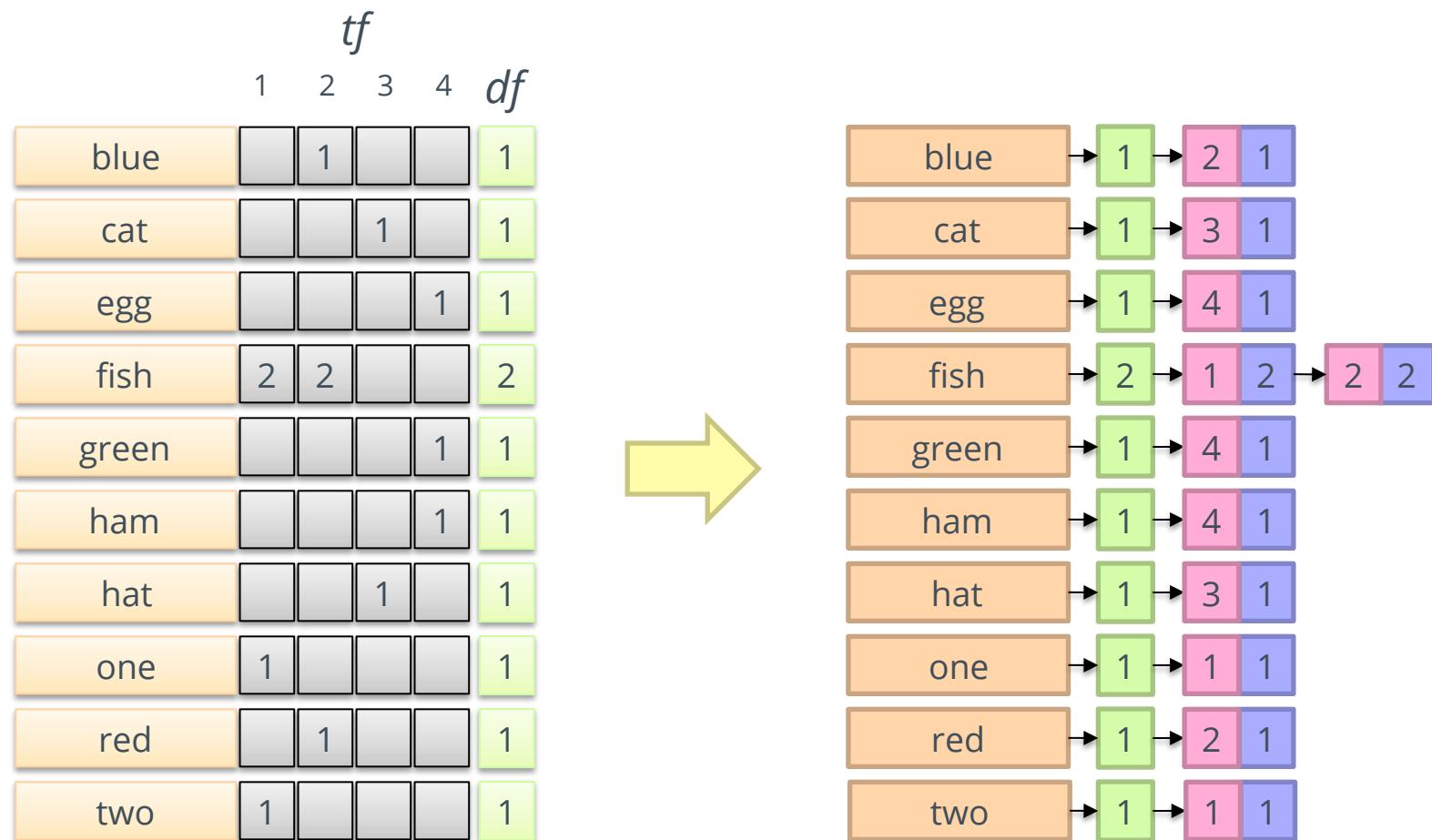
Inverted Index: TF.IDF

Doc 1
one fish, two fish

Doc 2
red fish, blue fish

Doc 3
cat in the hat

Doc 4
green eggs and ham



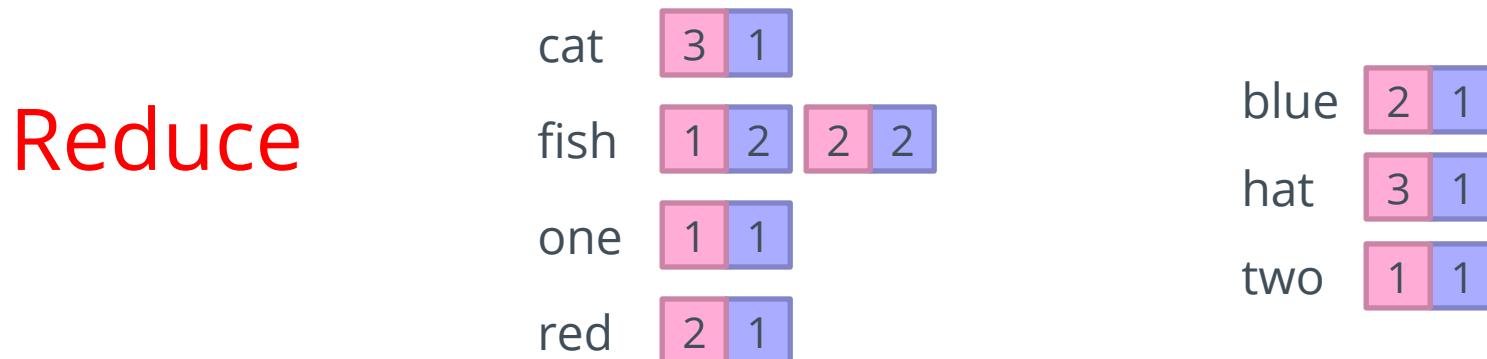
MapReduce: Index Construction

- Map over all documents
 - Emit *term* as key, $(docno, tf)$ as value
 - Emit other information as necessary (e.g., term position)
- Sort/shuffle: group postings by term
- Reduce
 - Gather and sort the postings (e.g., by *docno* or *tf*)
 - Write postings to disk
- MapReduce does all the heavy lifting!

Inverted Indexing with MapReduce



Shuffle and Sort: aggregate values by keys



Inverted Indexing: Pseudo-Code

```
1: class MAPPER
2:   procedure MAP(docid  $n$ , doc  $d$ )
3:      $H \leftarrow$  new ASSOCIATIVEARRAY
4:     for all term  $t \in$  doc  $d$  do
5:        $H\{t\} \leftarrow H\{t\} + 1$ 
6:     for all term  $t \in H$  do
7:       EMIT(term  $t$ , posting  $\langle n, H\{t\} \rangle$ )
```



```
1: class REDUCER
2:   procedure REDUCE(term  $t$ , postings [ $\langle a_1, f_1 \rangle, \langle a_2, f_2 \rangle \dots$ ])
3:      $P \leftarrow$  new LIST
4:     for all posting  $\langle a, f \rangle \in$  postings [ $\langle a_1, f_1 \rangle, \langle a_2, f_2 \rangle \dots$ ] do
5:       APPEND( $P, \langle a, f \rangle$ )
6:     SORT( $P$ )
7:     EMIT(term  $t$ , postings  $P$ )
```

Inverted Indexing: Pseudo-Code

```
1: class MAPPER
2:   procedure MAP(docid  $n$ , doc  $d$ )
3:      $H \leftarrow$  new ASSOCIATIVEARRAY
4:     for all term  $t \in$  doc  $d$  do
5:        $H\{t\} \leftarrow H\{t\} + 1$ 
6:     for all term  $t \in H$  do
7:       EMIT(term  $t$ , posting  $\langle n, H\{t\} \rangle$ )
```



```
1: class REDUCER
2:   procedure REDUCE(term  $t$ , postings [ $\langle a_1, f_1 \rangle, \langle a_2, f_2 \rangle \dots$ ])
3:      $P \leftarrow$  new LIST
4:     for all posting  $\langle a, f \rangle \in$  postings [ $\langle a_1, f_1 \rangle, \langle a_2, f_2 \rangle \dots$ ] do
5:       APPEND( $P, \langle a, f \rangle$ )
6:     SORT( $P$ )
7:     EMIT(term  $t$ , postings  $P$ )
```

Inverted Indexing: Pseudo-Code

```
1: class MAPPER
2:   procedure MAP(docid  $n$ , doc  $d$ )
3:      $H \leftarrow$  new ASSOCIATIVEARRAY
4:     for all term  $t \in$  doc  $d$  do
5:        $H\{t\} \leftarrow H\{t\} + 1$ 
6:     for all term  $t \in H$  do
7:       EMIT(term  $t$ , posting  $\langle n, H\{t\} \rangle$ )
```



```
1: class REDUCER
2:   procedure REDUCE(term  $t$ , postings [ $\langle a_1, f_1 \rangle, \langle a_2, f_2 \rangle \dots$ ])
3:      $P \leftarrow$  new LIST
4:     for all posting  $\langle a, f \rangle \in$  postings [ $\langle a_1, f_1 \rangle, \langle a_2, f_2 \rangle \dots$ ] do
5:       APPEND( $P, \langle a, f \rangle$ )
6:     SORT( $P$ )
7:     EMIT(term  $t$ , postings  $P$ )
```

What's the problem?

Scalability Bottleneck

- Initial implementation: terms as keys, postings as values
 - Reducers must buffer all postings associated with key (to sort)
 - What if we run out of memory to buffer postings

Secondary Sorting

- MapReduce sorts input to reducers by key
 - Values may be arbitrarily ordered
- What if want to sort value also?
 - E.g., $k \rightarrow (v_1, r), (v_3, r), (v_4, r), (v_8, r) \dots$

Secondary Sorting: Solutions

- Solution 1:
 - Buffer values in memory, then sort
 - Why is this a bad idea?
- Solution 2:
 - “Value-to-key conversion” design pattern: form composite intermediate key, (k, v_1)
 - Let execution framework do the sorting
 - Preserve state across multiple key-value pairs to handle processing
 - Anything else we need to do?

Another Try...

```
1: class MAPPER
2:   method MAP(docid  $n$ , doc  $d$ )
3:      $H \leftarrow$  new ASSOCIATIVEARRAY
4:     for all term  $t \in$  doc  $d$  do
5:        $H\{t\} \leftarrow H\{t\} + 1$ 
6:     for all term  $t \in H$  do
7:       EMIT(tuple  $\langle t, n \rangle$ , tf  $H\{t\}$ )
```



```
1: class REDUCER
2:   method INITIALIZE
3:      $t_{prev} \leftarrow \emptyset$ 
4:      $P \leftarrow$  new POSTINGSLIST
5:   method REDUCE(tuple  $\langle t, n \rangle$ , tf  $[f]$ )
6:     if  $t \neq t_{prev} \wedge t_{prev} \neq \emptyset$  then
7:       EMIT(term  $t_{prev}$ , postings  $P$ )
8:        $P.\text{RESET}()$ 
9:        $P.\text{ADD}(\langle n, f \rangle)$ 
10:       $t_{prev} \leftarrow t$ 
11:    method CLOSE
12:      EMIT(term  $t$ , postings  $P$ )
```

Another Try...

(key)	(values)		
fish	<table><tr><td>1</td><td>2</td></tr></table>	1	2
1	2		
	<table><tr><td>34</td><td>1</td></tr></table>	34	1
34	1		
	<table><tr><td>21</td><td>3</td></tr></table>	21	3
21	3		
	<table><tr><td>35</td><td>2</td></tr></table>	35	2
35	2		
	<table><tr><td>80</td><td>3</td></tr></table>	80	3
80	3		
	<table><tr><td>9</td><td>1</td></tr></table>	9	1
9	1		



(keys)	(values)		
fish	<table><tr><td>1</td><td>2</td></tr></table>	1	2
1	2		
fish	<table><tr><td>9</td><td>1</td></tr></table>	9	1
9	1		
fish	<table><tr><td>21</td><td>3</td></tr></table>	21	3
21	3		
fish	<table><tr><td>34</td><td>1</td></tr></table>	34	1
34	1		
fish	<table><tr><td>35</td><td>2</td></tr></table>	35	2
35	2		
fish	<table><tr><td>80</td><td>3</td></tr></table>	80	3
80	3		

How is this different?

- Let the framework do the sorting
- Directly write postings to disk!

Questions?