

# MapReduce Algorithm Design

Based on Jimmy Lin's slides

# Midterm – Take Home Exam

- Take Home Exam, Open Books and Notes
- No collaboration is allowed!
- 1 day (24 hours) to complete and submit using Gradescope
- Based on material before the break
  - Lecture Notes and Slides from Jan 22 to March 4
  - Papers and Book Chapters
  - Topics:
    - B-tree and Linear Hashing
    - Spatial databases and Indexing
    - Temporal databases and Indexing
    - Spatio-temporal Databases
    - Time Series

# Midterm

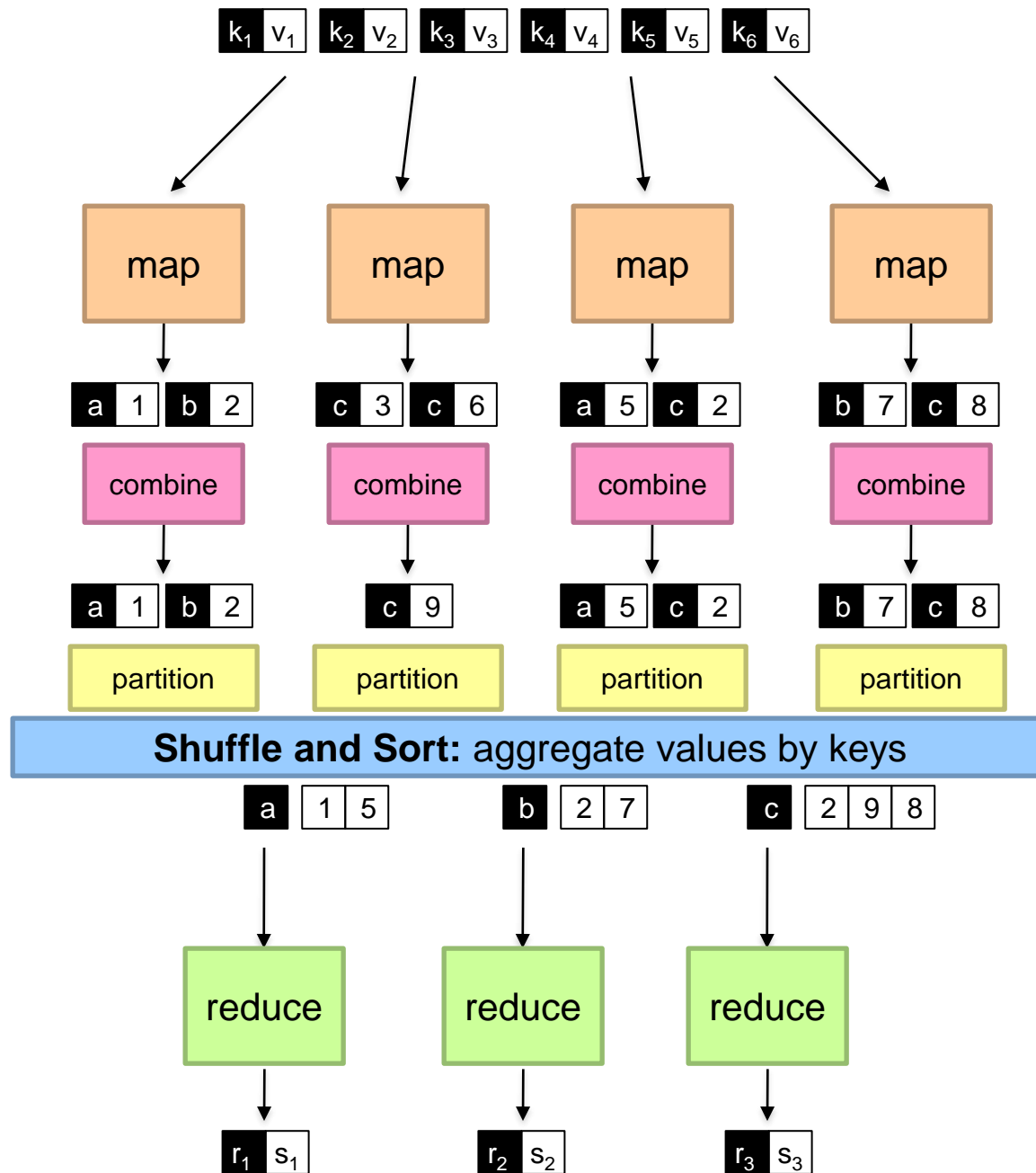
- When?
  - One option April 8, Wednesday
  - Second option, April 13, Monday
  - Other days?
- We decided to do a poll for potential days on Piazza...

# MapReduce: Recap

- Programmers must specify:
  - 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
- Optionally, also:
  - 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
- The execution framework handles everything else...

# “Everything Else”

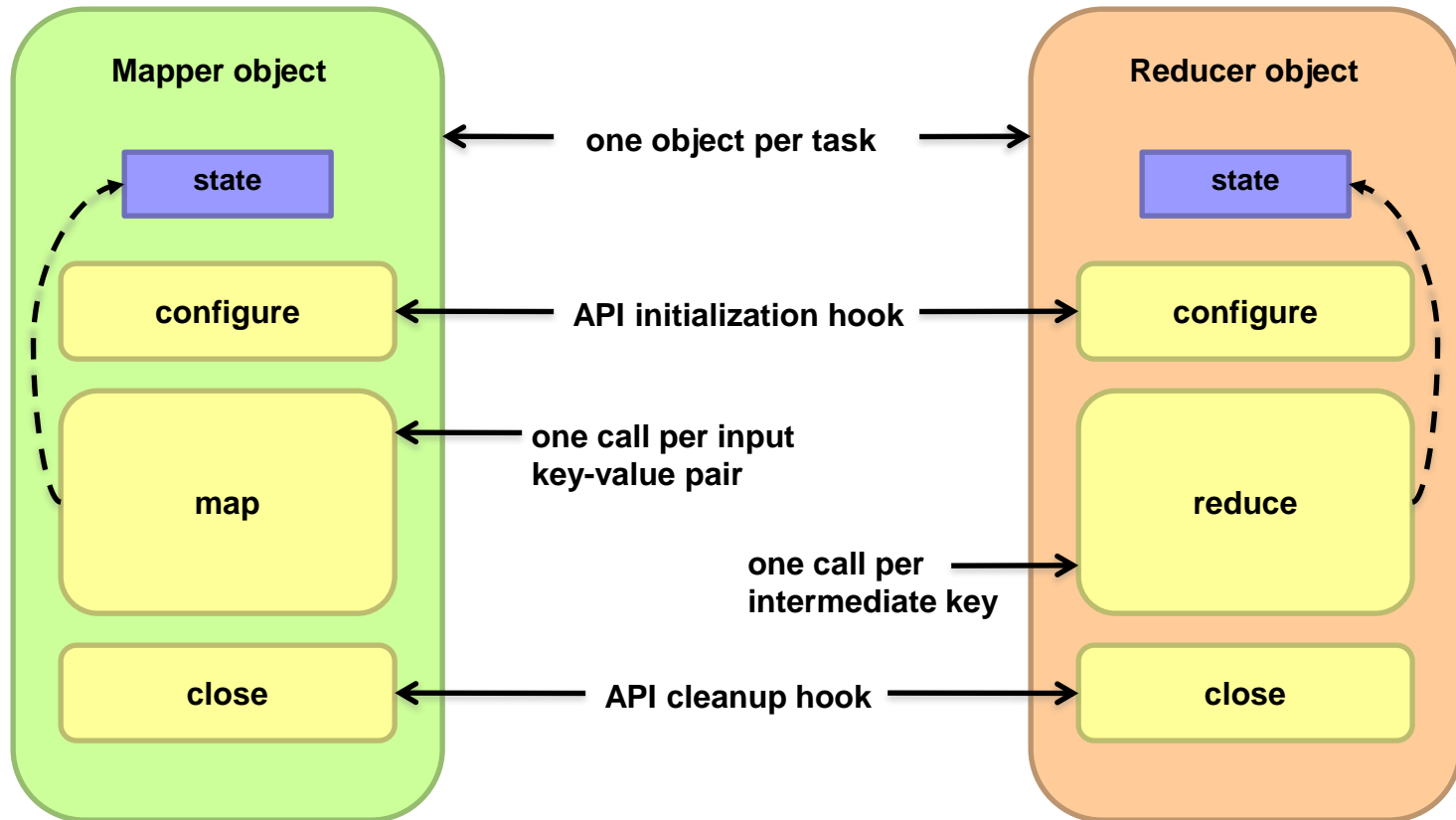
- The execution framework handles everything else...
  - Scheduling: assigns workers to map and reduce tasks
  - “Data distribution”: moves processes to data
  - Synchronization: gathers, sorts, and shuffles intermediate data
  - Errors and faults: detects worker failures and restarts
- Limited control over data and execution flow
  - All algorithms must be expressed in m, r, c, p
- You don’t know:
  - Where mappers and reducers run
  - When a mapper or reducer begins or finishes
  - Which input a particular mapper is processing
  - Which intermediate key a particular reducer is processing



# 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





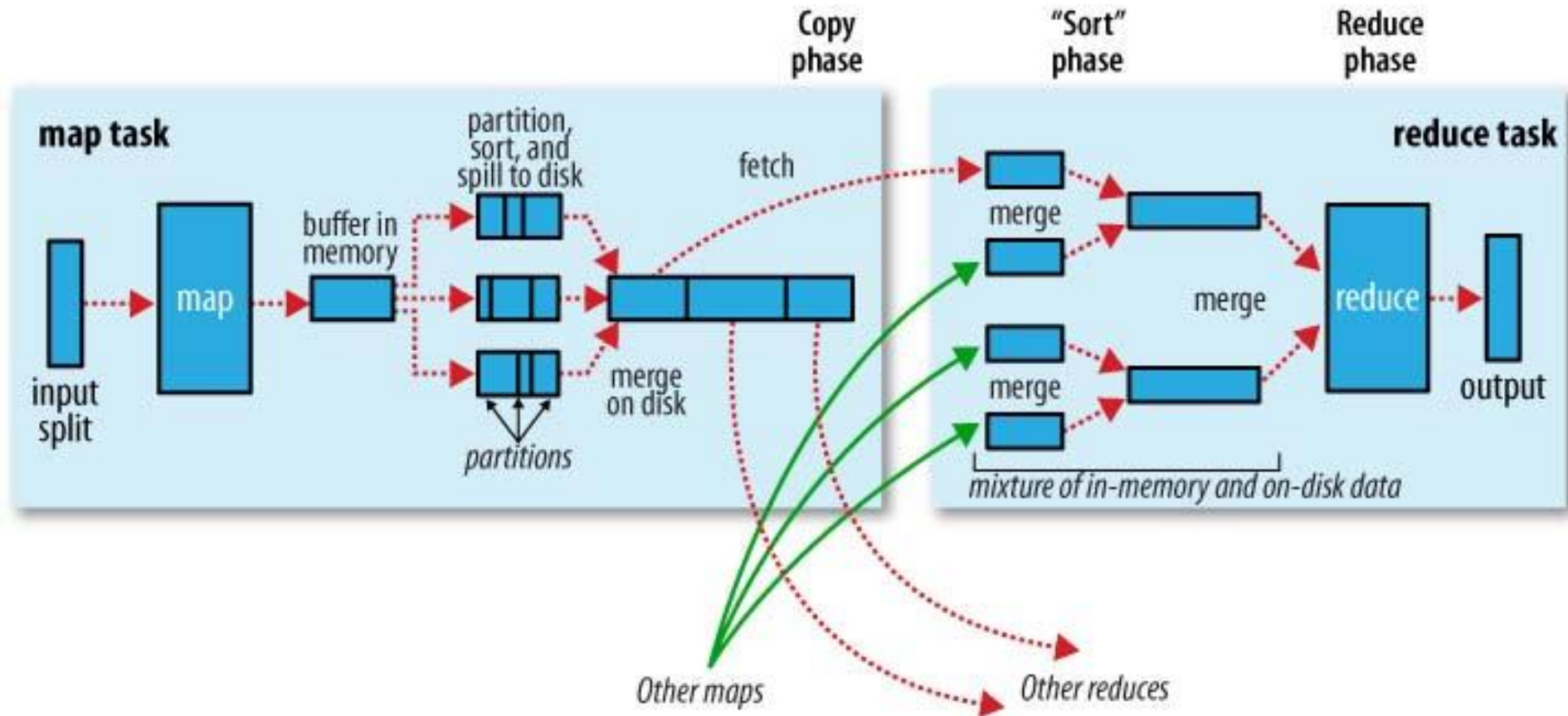
# Scalable Hadoop Algorithms: Themes

- Avoid object creation
  - Inherently costly operation
  - Garbage collection
- Avoid buffering
  - Limited heap size
  - Works for small datasets, but won't scale!

# 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

# Shuffle and Sort



# Word Count: Baseline

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

1: class REDUCER
2:   method REDUCE(term  $t$ , counts  $[c_1, c_2, \dots]$ )
3:      $sum \leftarrow 0$ 
4:     for all count  $c \in \text{counts } [c_1, c_2, \dots]$  do
5:        $sum \leftarrow sum + c$ 
6:     EMIT(term  $t$ , count  $s$ )
```

What's the impact of combiners?

# 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:        $H\{t\} \leftarrow H\{t\} + 1$ 
6:     for all term  $t \in H$  do
7:       EMIT(term  $t$ , count  $H\{t\}$ )
```

▷ Tally counts for entire document

Are combiners still needed?

# Word Count: Version 2

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

Key: preserve state across  
input key-value pairs!

▷ Tally counts *across* documents

Are combiners still needed?

# 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

# Combiner Design

- Combiners and reducers share same method signature
  - Sometimes, reducers can serve as combiners
  - Often, not...
- Remember: combiner are optional optimizations
  - Should not affect algorithm correctness
  - May be run 0, 1, or multiple times
- Example: find average of all integers associated with the same key



# Computing the Mean: Version 1

```
1: class MAPPER
2:     method MAP(string  $t$ , integer  $r$ )
3:         EMIT(string  $t$ , integer  $r$ )

1: class REDUCER
2:     method REDUCE(string  $t$ , integers  $[r_1, r_2, \dots]$ )
3:          $sum \leftarrow 0$ 
4:          $cnt \leftarrow 0$ 
5:         for all integer  $r \in$  integers  $[r_1, r_2, \dots]$  do
6:              $sum \leftarrow sum + r$ 
7:              $cnt \leftarrow cnt + 1$ 
8:              $r_{avg} \leftarrow sum / cnt$ 
9:             EMIT(string  $t$ , integer  $r_{avg}$ )
```

Why can't we use reducer as combiner?

# Computing the Mean: Version 2

```
1: class MAPPER
2:   method MAP(string  $t$ , integer  $r$ )
3:     EMIT(string  $t$ , integer  $r$ )

1: class COMBINER
2:   method COMBINE(string  $t$ , integers  $[r_1, r_2, \dots]$ )
3:      $sum \leftarrow 0$ 
4:      $cnt \leftarrow 0$ 
5:     for all integer  $r \in$  integers  $[r_1, r_2, \dots]$  do
6:        $sum \leftarrow sum + r$ 
7:        $cnt \leftarrow cnt + 1$ 
8:     EMIT(string  $t$ , pair ( $sum, cnt$ ))           ▷ Separate sum and count

1: class REDUCER
2:   method REDUCE(string  $t$ , pairs  $[(s_1, c_1), (s_2, c_2) \dots]$ )
3:      $sum \leftarrow 0$ 
4:      $cnt \leftarrow 0$ 
5:     for all pair  $(s, c) \in$  pairs  $[(s_1, c_1), (s_2, c_2) \dots]$  do
6:        $sum \leftarrow sum + s$ 
7:        $cnt \leftarrow cnt + c$ 
8:      $r_{avg} \leftarrow sum / cnt$ 
9:     EMIT(string  $t$ , integer  $r_{avg}$ )
```

Why doesn't this work?

# Computing the Mean: Version 3

```
1: class MAPPER
2:   method MAP(string  $t$ , integer  $r$ )
3:     EMIT(string  $t$ , pair ( $r$ , 1))

1: class COMBINER
2:   method COMBINE(string  $t$ , pairs  $[(s_1, c_1), (s_2, c_2) \dots]$ )
3:      $sum \leftarrow 0$ 
4:      $cnt \leftarrow 0$ 
5:     for all pair  $(s, c) \in$  pairs  $[(s_1, c_1), (s_2, c_2) \dots]$  do
6:        $sum \leftarrow sum + s$ 
7:        $cnt \leftarrow cnt + c$ 
8:     EMIT(string  $t$ , pair ( $sum$ ,  $cnt$ ))

1: class REDUCER
2:   method REDUCE(string  $t$ , pairs  $[(s_1, c_1), (s_2, c_2) \dots]$ )
3:      $sum \leftarrow 0$ 
4:      $cnt \leftarrow 0$ 
5:     for all pair  $(s, c) \in$  pairs  $[(s_1, c_1), (s_2, c_2) \dots]$  do
6:        $sum \leftarrow sum + s$ 
7:        $cnt \leftarrow cnt + c$ 
8:      $r_{avg} \leftarrow sum / cnt$ 
9:     EMIT(string  $t$ , pair ( $r_{avg}$ ,  $cnt$ ))
```

Fixed?

# Computing the Mean: Version 4

```
1: class MAPPER
2:   method INITIALIZE
3:      $S \leftarrow \text{new ASSOCIATIVEARRAY}$ 
4:      $C \leftarrow \text{new ASSOCIATIVEARRAY}$ 
5:   method MAP(string  $t$ , integer  $r$ )
6:      $S\{t\} \leftarrow S\{t\} + r$ 
7:      $C\{t\} \leftarrow C\{t\} + 1$ 
8:   method CLOSE
9:     for all term  $t \in S$  do
10:       EMIT(term  $t$ , pair ( $S\{t\}, C\{t\}$ ))
```

Are combiners still needed?

# Algorithm Design: Running 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)  $\rightarrow$  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  $w \in \text{doc } d$  do
4:       for all term  $u \in \text{NEIGHBORS}(w)$  do
5:         EMIT(pair  $(w, u)$ , count 1)      ▷ Emit count for each co-occurrence

1: class REDUCER
2:   method REDUCE(pair  $p$ , counts  $[c_1, c_2, \dots]$ )
3:      $s \leftarrow 0$ 
4:     for all count  $c \in \text{counts } [c_1, c_2, \dots]$  do
5:        $s \leftarrow s + c$                                 ▷ 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: cleverly-constructed data structure  
brings together partial results**

# Stripes: Pseudo-Code

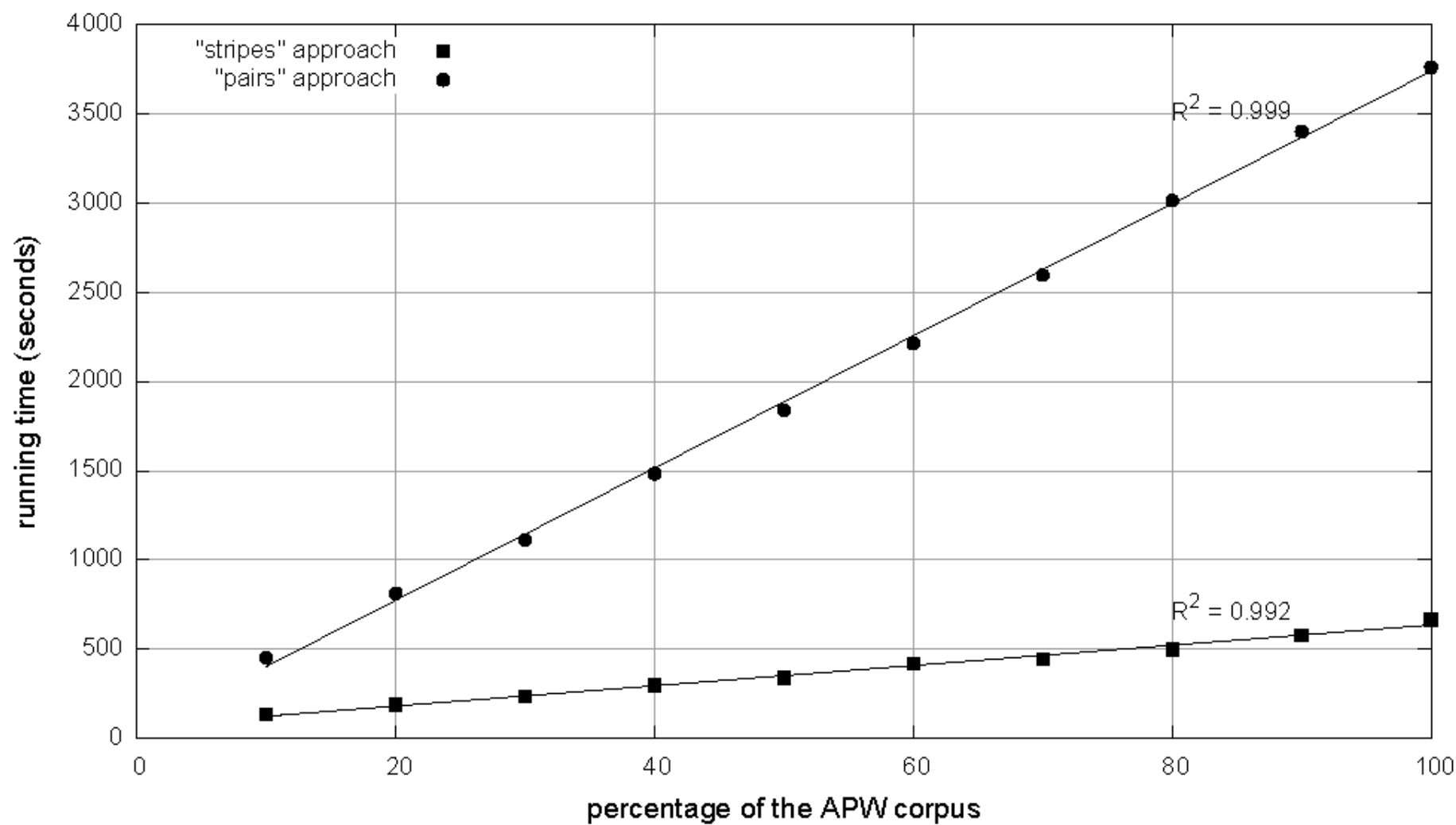
```
1: class MAPPER
2:   method MAP(docid  $a$ , doc  $d$ )
3:     for all term  $w \in \text{doc } d$  do
4:        $H \leftarrow \text{new ASSOCIATIVEARRAY}$ 
5:       for all term  $u \in \text{NEIGHBORS}(w)$  do
6:          $H\{u\} \leftarrow H\{u\} + 1$  ▷ Tally words co-occurring with  $w$ 
7:       EMIT(Term  $w$ , Stripe  $H$ )

1: class REDUCER
2:   method REDUCE(term  $w$ , stripes  $[H_1, H_2, H_3, \dots]$ )
3:      $H_f \leftarrow \text{new ASSOCIATIVEARRAY}$ 
4:     for all stripe  $H \in \text{stripes } [H_1, H_2, H_3, \dots]$  do
5:       SUM( $H_f, H$ ) ▷ Element-wise sum
6:     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)

# Relative Frequencies

- How do we estimate relative frequencies from counts?

$$f(B | A) = \frac{\text{count}(A, B)}{\text{count}(A)} = \frac{\text{count}(A, B)}{\sum_{B'} \text{count}(A, B')}$$

- Why do we want to do this?
- How do we do this with MapReduce?

# $f(B | A)$ : “Stripes”

- Easy!
  - One pass to compute  $(a, *)$
  - Another pass to directly compute  $f(B | A)$

# $f(B|A)$ : “Pairs”

$(a, *) \rightarrow 32$

Reducer holds this value in memory

$(a, b_1) \rightarrow 3$

$(a, b_2) \rightarrow 12$

$(a, b_3) \rightarrow 7$

$(a, b_4) \rightarrow 1$

...



$(a, b_1) \rightarrow 3 / 32$

$(a, b_2) \rightarrow 12 / 32$

$(a, b_3) \rightarrow 7 / 32$

$(a, b_4) \rightarrow 1 / 32$

...

## ○ For this to work:

- Must emit extra  $(a, *)$  for every  $b_n$  in mapper
- Must make sure all  $a$ 's get sent to same reducer (use partitioner)
- Must make sure  $(a, *)$  comes first (define sort order)
- Must hold state in reducer across different key-value pairs



# “Order Inversion”

- Common design pattern
  - Computing relative frequencies requires marginal counts
  - But marginal cannot be computed until you see all counts
  - Buffering is a bad idea!
  - Trick: getting the marginal counts to arrive at the reducer before the joint counts
- Optimizations
  - Apply in-memory combining pattern to accumulate marginal counts

# Synchronization: Pairs vs. Stripes

- Approach 1: turn synchronization into an ordering problem
  - Sort keys into correct order of computation
  - Partition key space so that each reducer gets the appropriate set of partial results
  - Hold state in reducer across multiple key-value pairs to perform computation
  - Illustrated by the “pairs” approach
- Approach 2: construct data structures that bring partial results together
  - Each reducer receives all the data it needs to complete the computation
  - Illustrated by the “stripes” approach

# 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

# Recap: Tools for Synchronization

- Cleverly-constructed data structures
  - Bring data 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

# Issues and Tradeoffs

- Number of key-value pairs
  - Object creation overhead
  - Time for sorting and shuffling pairs across the network
- Size of each key-value pair
  - De/serialization overhead
- Local aggregation
  - Opportunities to perform local aggregation varies
  - Combiners make a big difference
  - Combiners vs. in-mapper combining
  - RAM vs. disk vs. network

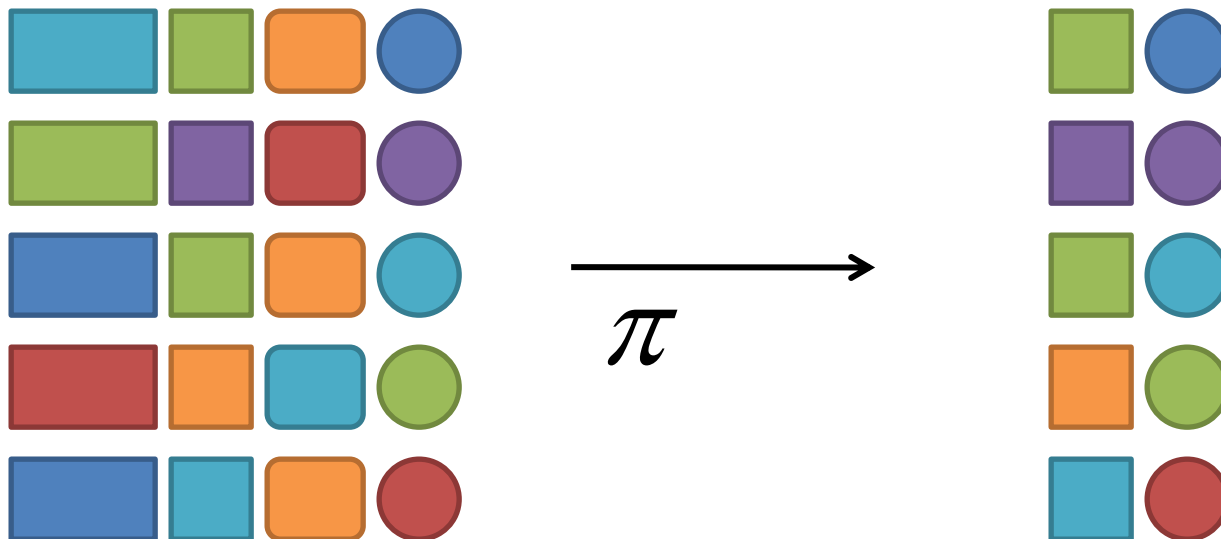
# Mapreduce and Databases

# Relational Algebra

- Primitives
  - Projection ( $\pi$ )
  - Selection ( $\sigma$ )
  - Cartesian product ( $\times$ )
  - Set union ( $\cup$ )
  - Set difference ( $-$ )
  - Rename ( $\rho$ )
- Other operations
  - Join ( $\bowtie$ )
  - Group by... aggregation
  - ...



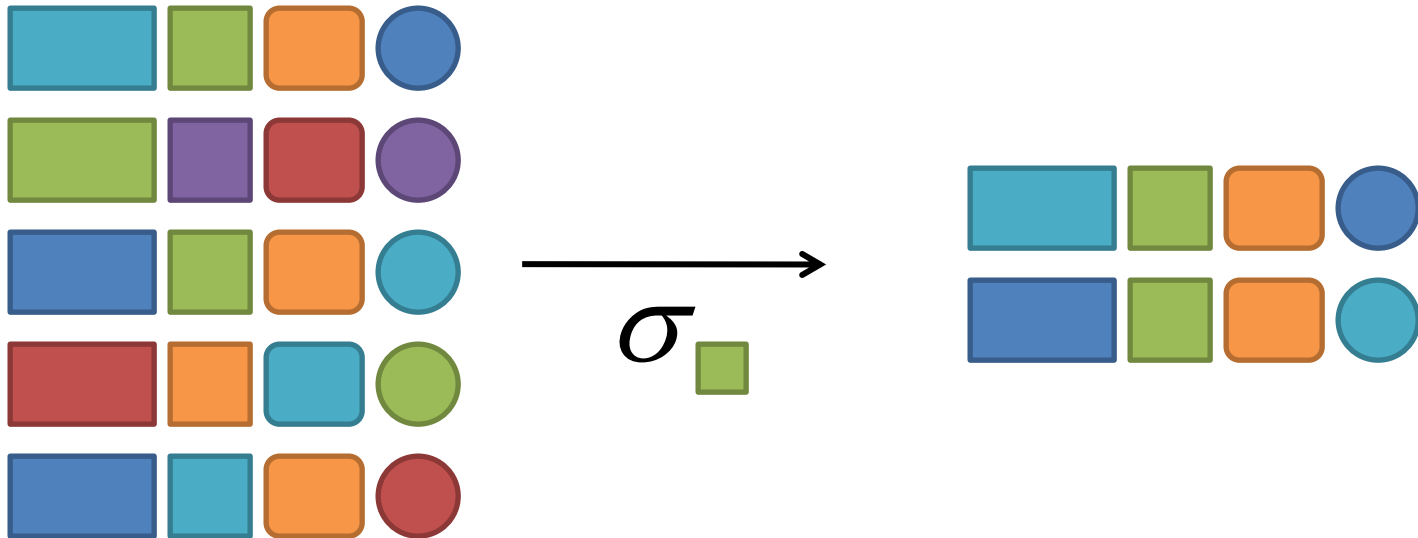
# Projection



# Projection in MapReduce

- Easy!
  - Map over tuples, emit new tuples with appropriate attributes
  - Reduce: take tuples that appear many times and emit only one version (duplicate elimination)
    - Tuple  $t$  in  $R$ :  $\text{Map}(t, t) \rightarrow (t', t')$
    - $\text{Reduce}(t', [t', \dots, t']) \rightarrow [t', t']$
- Basically limited by HDFS streaming speeds
  - Speed of encoding/decoding tuples becomes important
  - Relational databases take advantage of compression

# Selection



# Selection in MapReduce

- Easy!
  - Map over tuples, emit only tuples that meet criteria
  - No reducers, unless for regrouping or resorting tuples (reducers are the identity function)
  - Alternatively: perform in reducer, after some other processing
- But very expensive!!! Has to scan the database
  - Better approaches?

# Union, Set Intersection and Set Difference

- Similar ideas: each map outputs the tuple pair  $(t,t)$ . For union, we output it once, for intersection only when in the reduce we have  $(t, [t,t])$
- For Set difference?

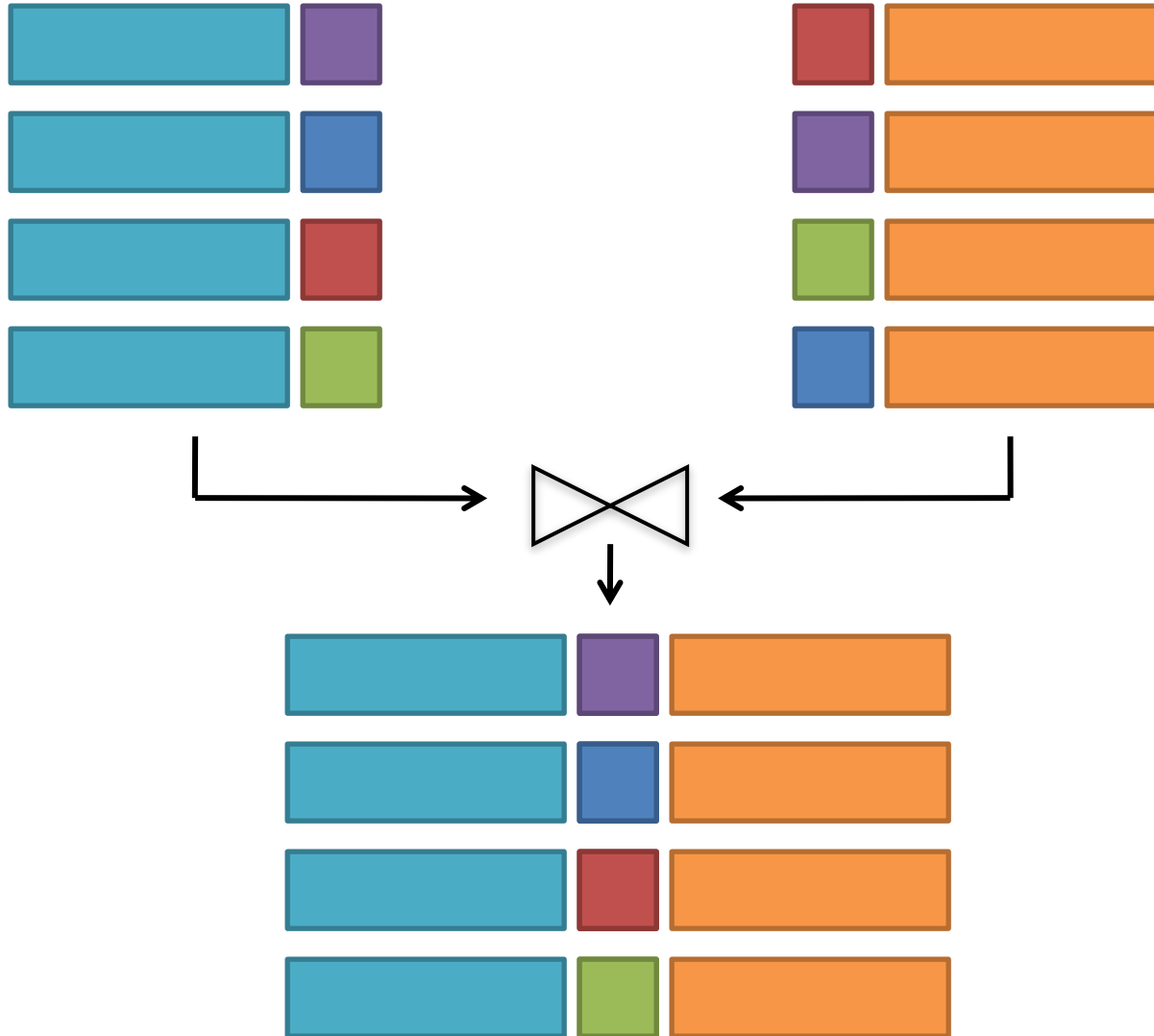
# Set Difference

- Map Function: For a tuple  $t$  in  $R$ , produce key-value pair  $(t, R)$ , and for a tuple  $t$  in  $S$ , produce key-value pair  $(t, S)$ .
- Reduce Function: For each key  $t$ , do the following.
  1. If the associated value list is  $[R]$ , then produce  $(t, t)$ .
  2. If the associated value list is anything else, which could only be  $[R, S]$ ,  $[S, R]$ , or  $[S]$ , produce  $(t, \text{NULL})$ .

# Group by... Aggregation

- Example: What is the average time spent per URL?
- In SQL:
  - `SELECT url, AVG(time) FROM visits GROUP BY url`
- In MapReduce:
  - Map over tuples, emit time, keyed by url
  - Framework automatically groups values by keys
  - Compute average in reducer
  - Optimize with combiners

# Relational Joins





# Join Algorithms in MapReduce

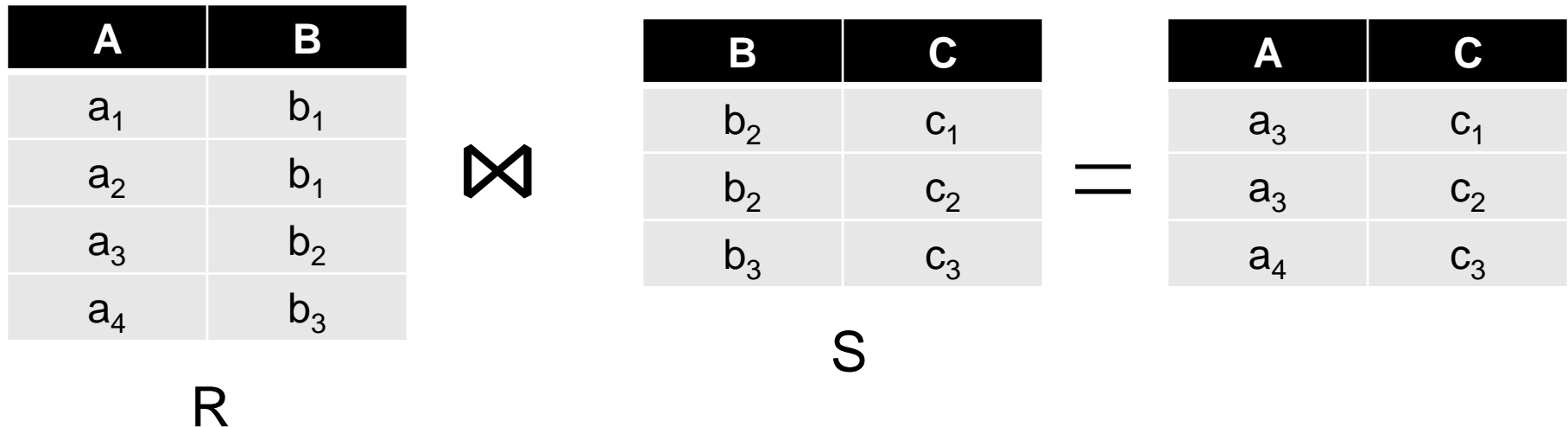
- Reduce-side join
- Map-side join
- In-memory join
  - Striped variant
  - Memcached variant

# Reduce-side Join

- Basic idea: group by join key
  - Map over both sets of tuples
  - Emit tuple as value with join key as the intermediate key
  - Execution framework brings together tuples sharing the same key
  - Perform actual join in reducer
  - Similar to a “sort-merge join” in database terminology

# Example: Join By Map-Reduce

- Compute the natural join  $R(A,B) \bowtie S(B,C)$
- $R$  and  $S$  are each stored in files
- Tuples are pairs  $(a,b)$  or  $(b,c)$



# Map-Reduce Join

- Use a hash function  $h$  from B-values to  $1\dots k$
- **A Map process turns:**
  - Each input tuple  $R(a,b)$  into key-value pair  $(b,(a,R))$
  - Each input tuple  $S(b,c)$  into  $(b,(c,S))$
- **Map processes** send each key-value pair with key  $b$  to Reduce process  $h(b)$ 
  - Hadoop does this automatically; just tell it what  $k$  is.
- Each **Reduce process** matches all the pairs  $(b,(a,R))$  with all  $(b,(c,S))$  and outputs  $(a,b,c)$ .

# Map-side Join: Parallel Scans

- If datasets are sorted by join key, join can be accomplished by a scan over both datasets
- How can we accomplish this in parallel?
  - Partition and sort both datasets in the same manner
- In MapReduce:
  - Map over one dataset, read from other corresponding partition
  - No reducers necessary (unless to repartition or resort)
- Consistently partitioned datasets: realistic to expect?

# In-Memory Join

- Basic idea: load one dataset into memory, stream over other dataset
  - Works if  $R \ll S$  and  $R$  fits into memory
  - Similar to “hash Join” algorithm
- MapReduce implementation
  - Distribute  $R$  to all nodes
  - Map over  $S$ , each mapper loads  $R$  in memory, hashed by join key
  - For every tuple in  $S$ , look up join key in  $R$
  - No reducers, unless for regrouping or resorting tuples

# In-Memory Join: Variants

- Striped variant:
  - R too big to fit into memory?
  - Divide R into  $R_1, R_2, R_3, \dots$  s.t. each  $R_n$  fits into memory
  - Perform in-memory join:  $\forall n, R_n \bowtie S$
  - Take the union of all join results
- Memcached join:
  - Load R into memcached
  - Replace in-memory hash lookup with memcached lookup

# Memcached Join

- Memcached join:
  - Load R into memcached
  - Replace in-memory hash lookup with memcached lookup
- Capacity and scalability?
  - Memcached capacity  $\gg$  RAM of individual node
  - Memcached scales out with cluster
- Latency?
  - Memcached is fast (basically, speed of network)
  - Batch requests to amortize latency costs



# Which join to use?

- In-memory join > map-side join > reduce-side join
  - Why?
- Limitations of each?
  - In-memory join: memory
  - Map-side join: sort order and partitioning
  - Reduce-side join: general purpose

# Processing Relational Data: Summary

- MapReduce algorithms for processing relational data:
  - Group by, sorting, partitioning are handled automatically by shuffle/sort in MapReduce
  - Selection, projection, and other computations (e.g., aggregation), are performed either in mapper or reducer
  - Multiple strategies for relational joins
- Complex operations require multiple MapReduce jobs
  - Example: top ten URLs in terms of average time spent
  - Opportunities for automatic optimization