

CS109 – Data Science

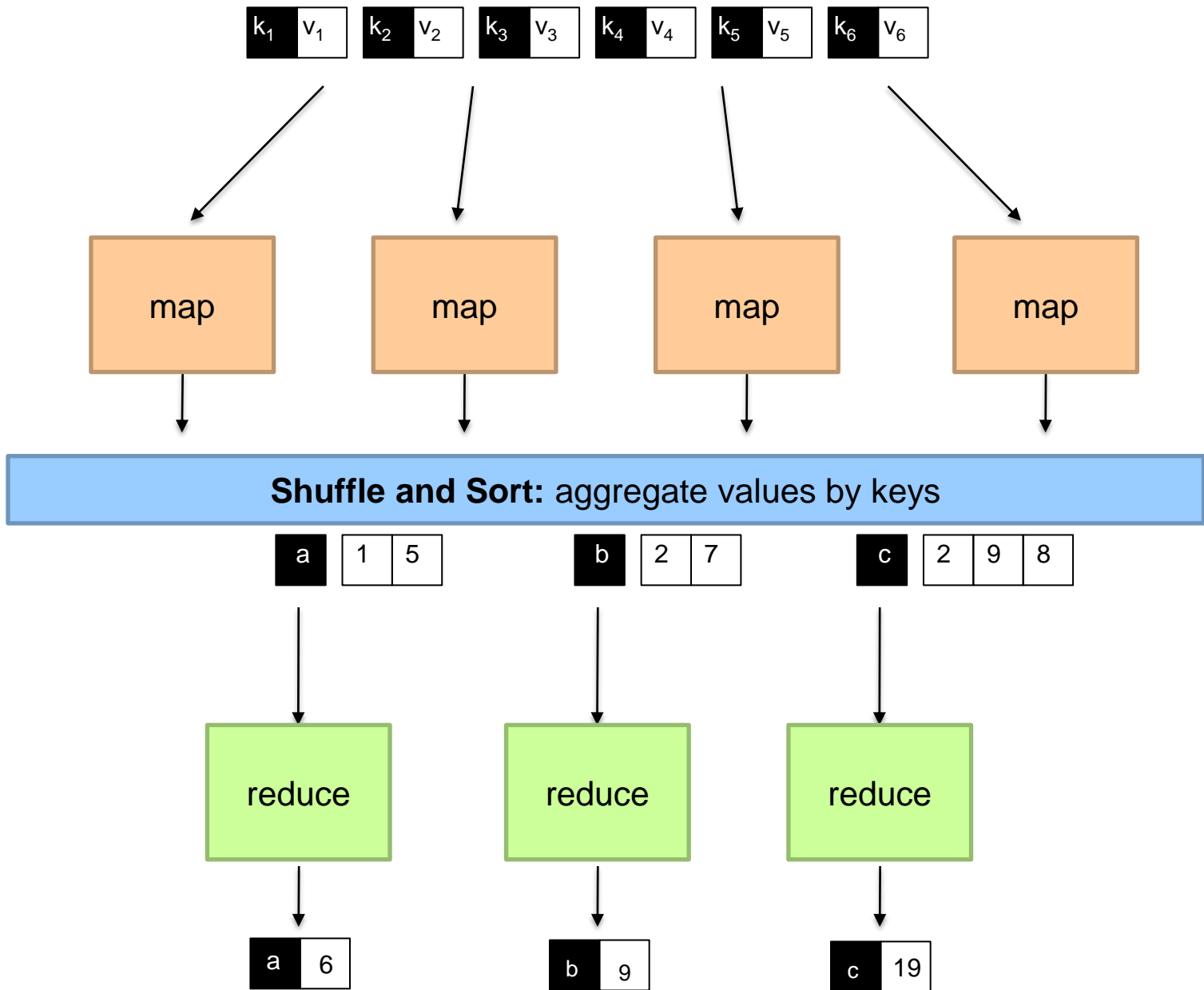
Joe Blitzstein, Hanspeter Pfister, Verena Kaynig-Fittkau

vkaynig@seas.harvard.edu

staff@cs109.org

Announcements

- Homework Collaboration Policy:
 - See Syllabus on CS109.org
 - The work you turn in must be your own
 - This is a data science course. It takes us 20 minutes to get a similarity ranking of all homework submissions.



Example Input File

I am Sam

I am Sam
Sam I am

That Sam I am
That Sam I am
I do not like
that Sam I am

Do you like
green eggs and ham

I do not like them
Sam I am
I do not like
green eggs and ham

```
from mrjob.job import MRJob
```

```
class mrWordCount(MRJob):
```

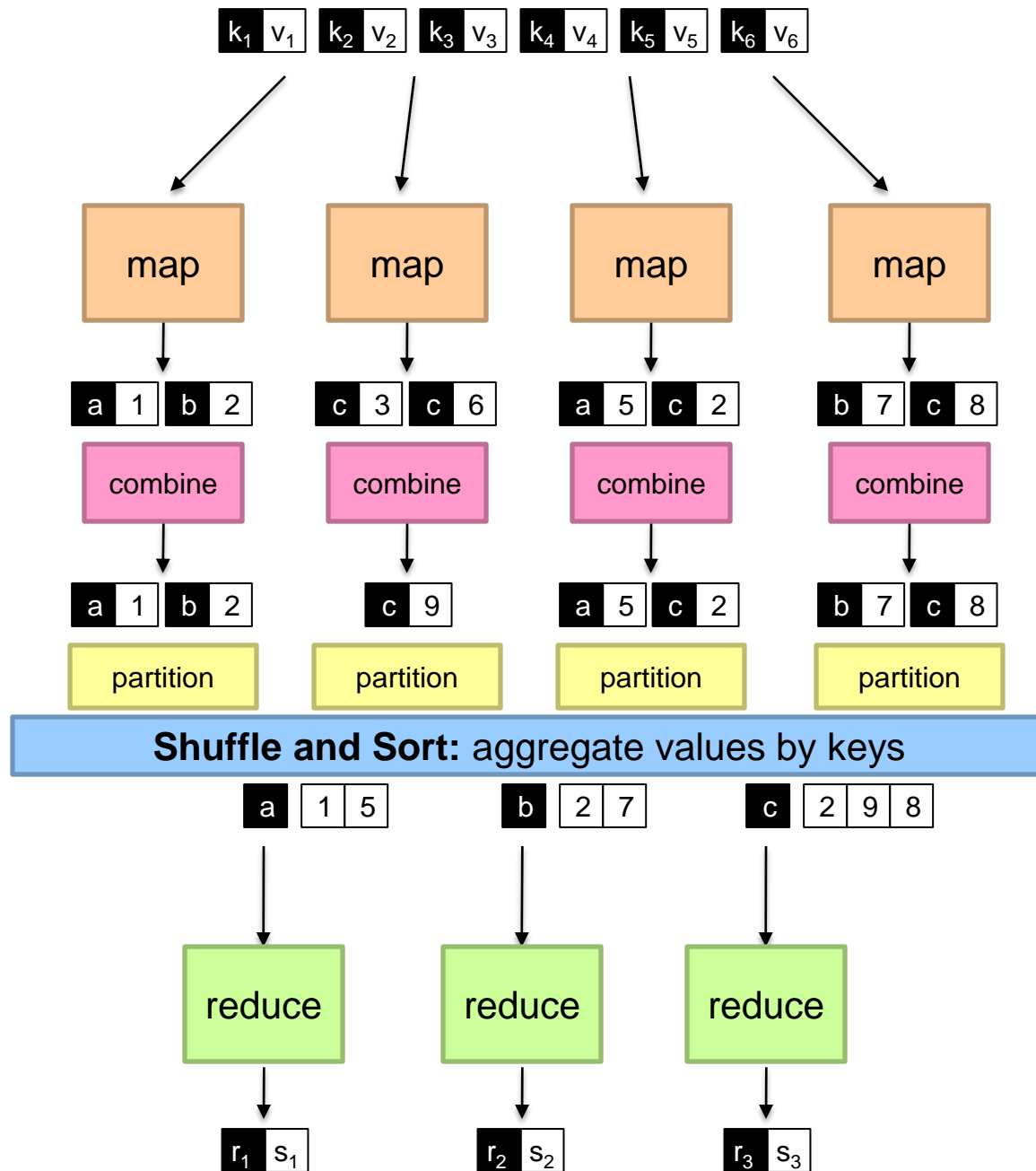
```
    def mapper(self, key, line):  
        for word in line.split(' '):  
            yield word.lower(), 1
```

```
    def reducer(self, word, occurrences):  
        yield word, sum(occurrences)
```

```
if __name__ == '__main__':  
    mrWordCount.run()
```

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
 - Two possibilities:
 - Combiners
 - In-mapper combining



Combiner

- “mini-reducers”
- Takes mapper output before shuffle and sort
- Can significantly reduce network traffic
- No access to other mappers
- Not guaranteed to get all values for a key
- Not guaranteed to run at all!
- Key and value output must match mapper

Why does the key and value output have to match the mapper output?

Word Count with Combiner

```
from mrjob.job import MRJob

class mrWordCount(MRJob):

    def mapper(self, key, line):
        for word in line.split(' '):
            yield word.lower(), 1

    def combiner(self, word, occurrences):
        yield word, sum(occurrences)

    def reducer(self, word, occurrences):
        yield word, sum(occurrences)

if __name__ == '__main__':
    mrWordCount.run()
```


Combiner Design

- Combiners and reducers share same method signature
 - Sometimes, reducers can serve as combiners
 - Often, not...
- Remember: combiners 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? What if combiner does not run?

In-Mapper Combining

- “Fold the functionality of the combiner into the mapper by preserving state across multiple map calls

```
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 ( $S\{t\}$ ,  $C\{t\}$ ))
```

In-Mapper Combining

- Advantages
 - Speed
 - Why is this faster than actual combiners?
- Disadvantages
 - Explicit memory management required
 - Potential for order-dependent bugs

Word Count with In-Mapper-Comb.

```
from collections import defaultdict
from mrjob.job import MRJob

class mrWordCount(MRJob):
    def __init__(self, *args, **kwargs):
        super(mrWordCount, self).__init__(*args, **kwargs)
        self.localWordCount = defaultdict(int)

    def mapper(self, key, line):
        if False:
            yield
        for word in line.split(' '):
            self.localWordCount[word.lower()] += 1

    def mapper_final(self):
        for (word, count) in self.localWordCount.iteritems():
            yield word, count

    def reducer(self, word, occurrences):
        yield word, sum(occurrences)

if __name__ == '__main__':
    mrWordCount.run()
```

Which is better?

- For large dictionaries?
 - Combiner has no memory problems
- For skewed word distributions (“the”)?
 - In-mapper reduces load on reducer

Pairs and Stripes:

- 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
 - Context can be a sentence, sequence of m words, etc.
 - In this case co-occurrence matrix is symmetric

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

First Try: “Pairs”

- Each mapper takes a sentence:
 - Generate all co-occurring term pairs
 - For all pairs, emit (a, b) → 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
 - Not many opportunities for combiners to work

Another Try: “Stripes”

- Idea: group together pairs into an associative array

(a, b) → 1

(a, c) → 2

(a, d) → 5

(a, e) → 3

(a, f) → 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

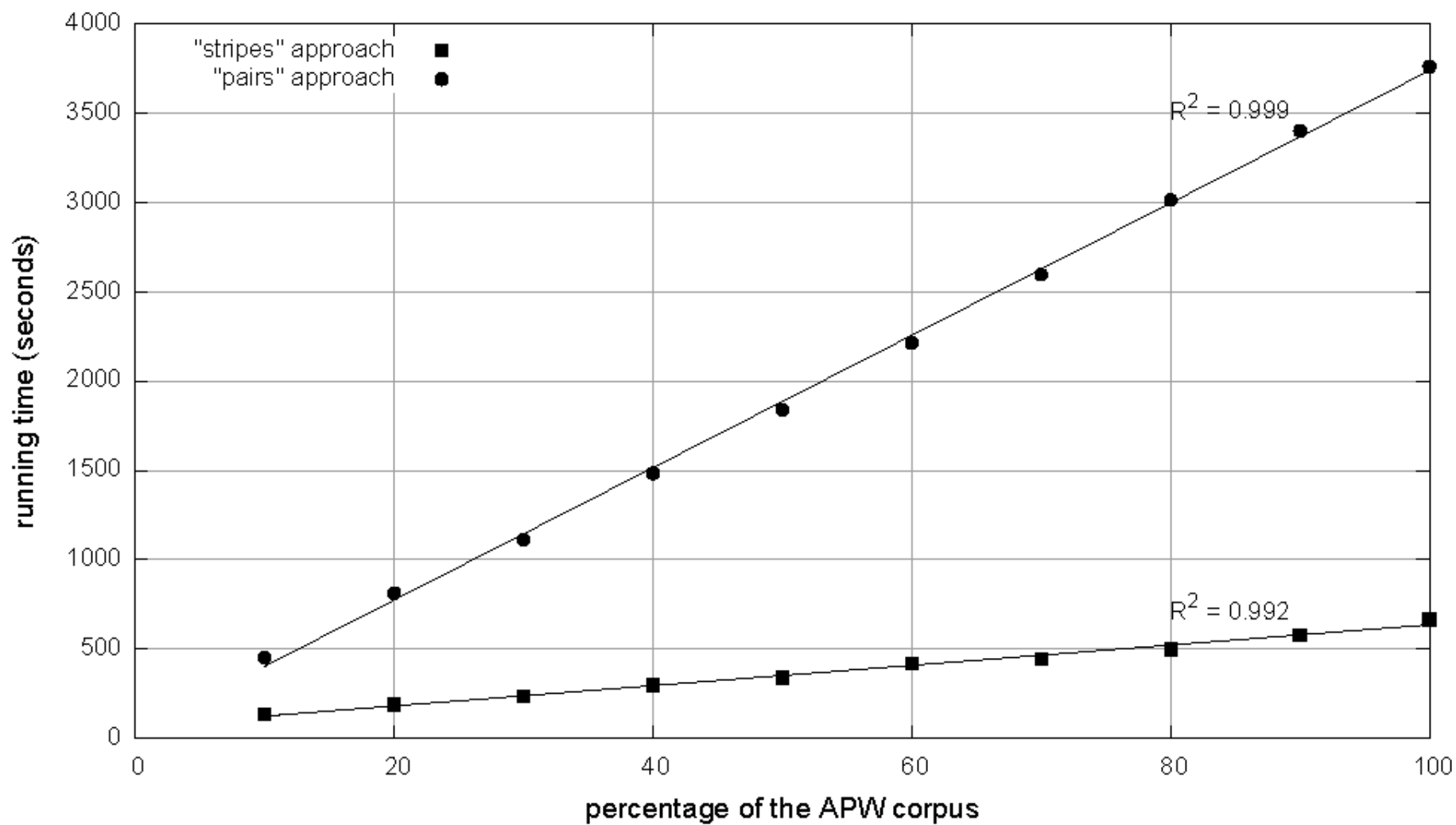
```
1: class MAPPER
2:   method MAP(docid  $a$ , doc  $d$ )
3:     for all term  $w \in \text{doc } d$  do
4:        $H \leftarrow \text{new ASSOCIATIVEARRAY}$ 
5:       for all term  $u \in \text{NEIGHBORS}(w)$  do
6:          $H\{u\} \leftarrow H\{u\} + 1$  ▷ Tally words co-occurring with  $w$ 
7: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

Map Reduce for Machine Learning

- Random Forest?
- SVM?