



Big Data Infrastructure

CS 489/698 Big Data Infrastructure (Winter 2016)

Week 3: From MapReduce to Spark (2/2)

January 21, 2016

Jimmy Lin

David R. Cheriton School of Computer Science
University of Waterloo

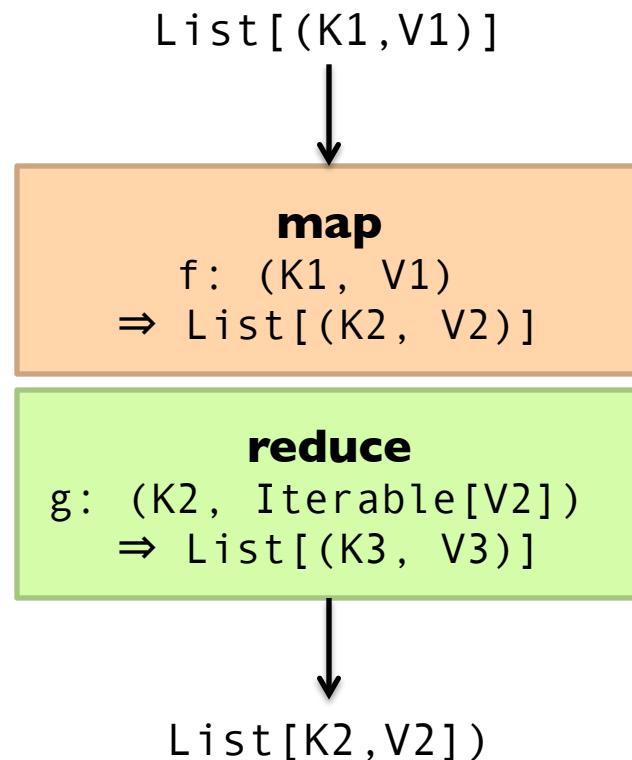
These slides are available at <http://lintool.github.io/bigdata-2016w/>

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

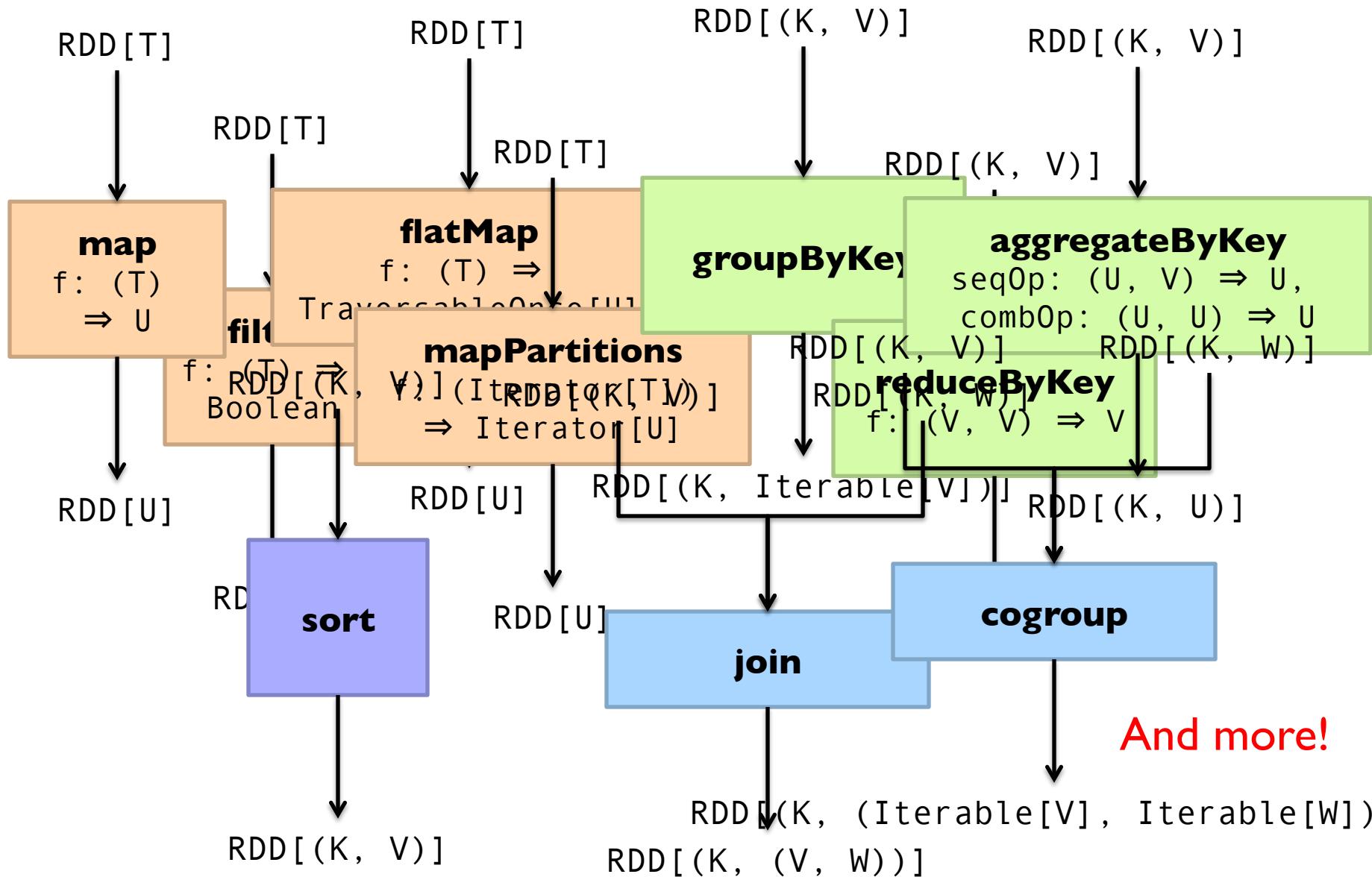
An aerial photograph of a large data center complex during sunset. The sky is a warm orange and yellow. In the foreground, there are several large white industrial buildings, some with flat roofs and others with gabled roofs. A parking lot with many cars is visible in front of one of the buildings. To the right, there is a large field with several white cylindrical storage tanks lined up in rows. In the background, there is a highway with traffic and a small town or cluster of buildings. The overall scene is a mix of industrial infrastructure and rural landscape.

The datacenter *is* the computer!
What's the instruction set?

MapReduce



Spark



What's an RDD?

Resilient Distributed Dataset (RDD)

= immutable = partitioned

Wait, so how do you actually do anything?

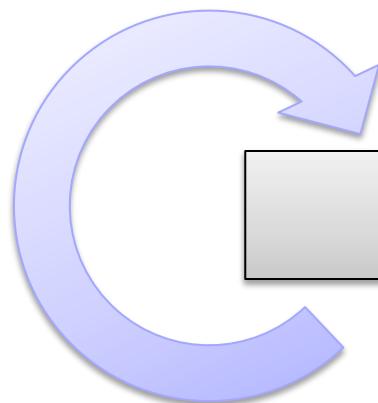
Developers define *transformations* on RDDs
Framework keeps track of lineage

Spark Word Count

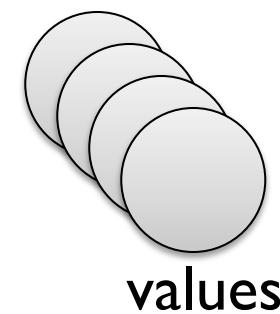
```
val textFile = sc.textFile(args.input())  
  
textFile  
  .flatMap(line => tokenize(line))  
  .map(word => (word, 1))  
  .reduceByKey(_ + _)  
  .saveAsTextFile(args.output())
```

RDD Lifecycle

Transformation



Action



values

Transformations are lazy:
Framework keeps track of lineage

Actions trigger actual execution

Spark Word Count

RDDs

```
→ val textFile = sc.textFile(args.input())
```

```
→ val a = textFile.flatMap(line => line.split(" "))
```

```
→ val b = a.map(word => (word, 1))
```

```
→ val c = b.reduceByKey(_ + _)
```

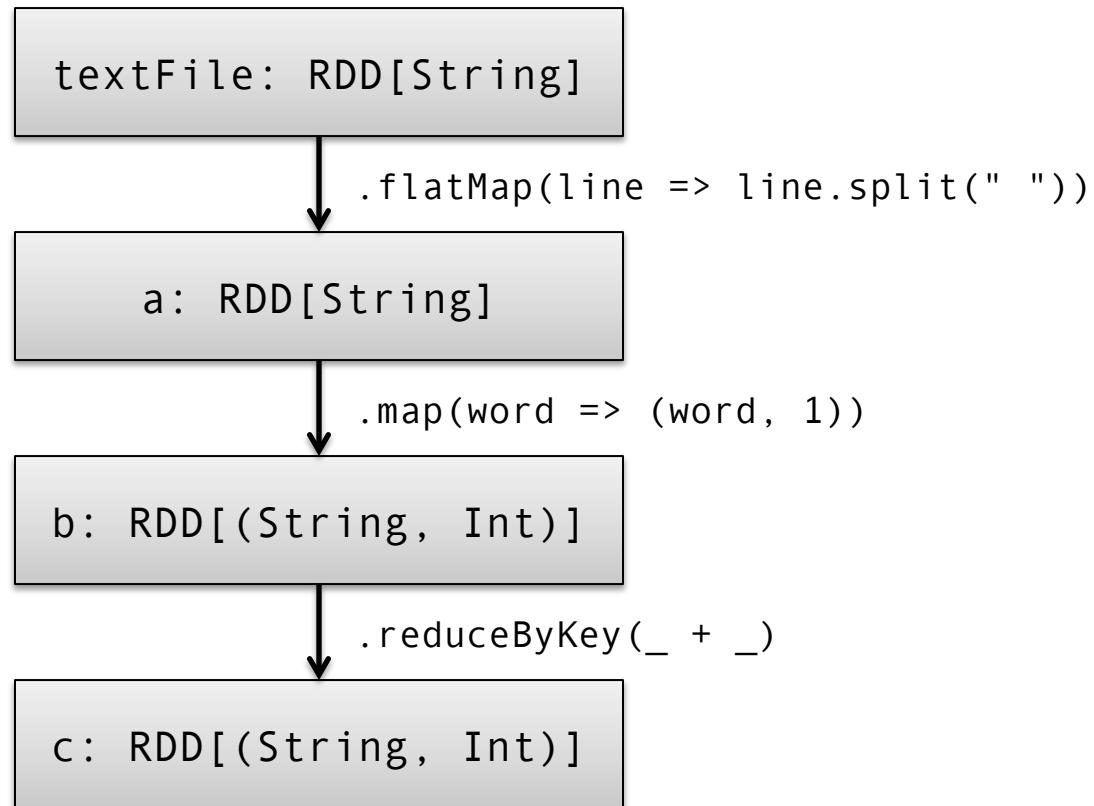
```
c.saveAsTextFile(args.output())
```

Action

Transformations

RDDs and Lineage

On HDFS



Action!

Remember,
transformations are lazy!

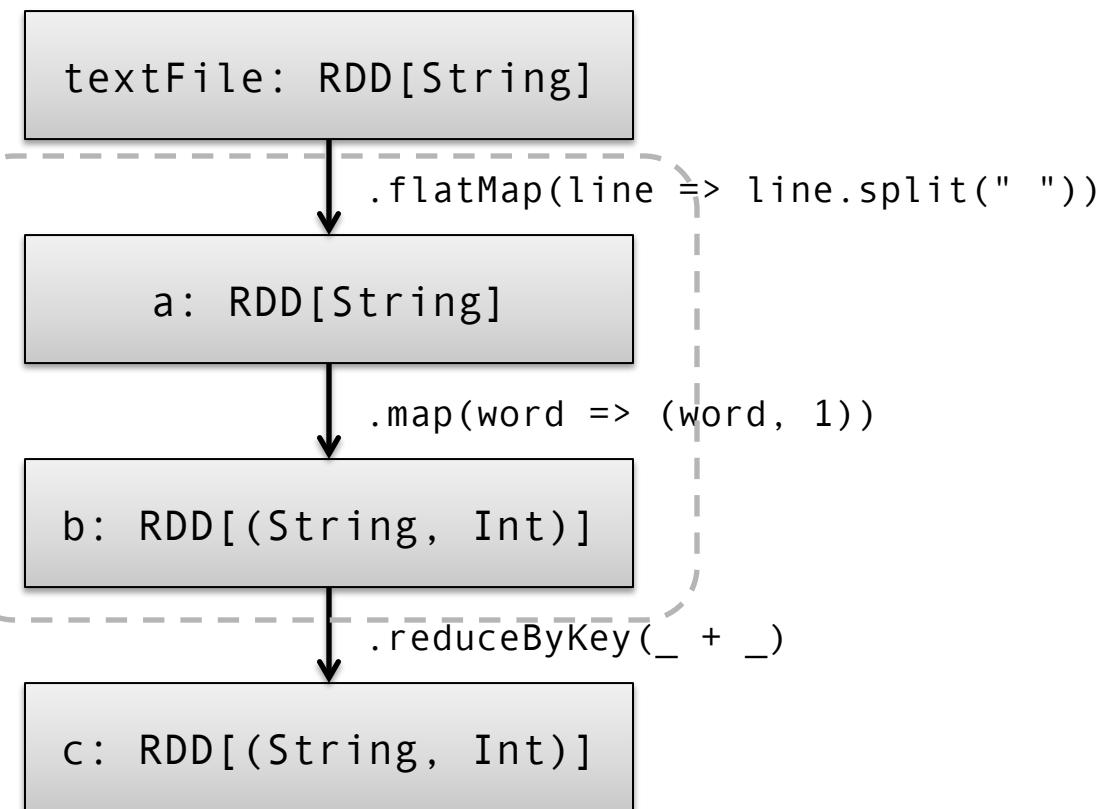
RDDs and Optimizations

Lazy evaluation creates optimization opportunities

On HDFS

RDDs don't need
to be materialized!

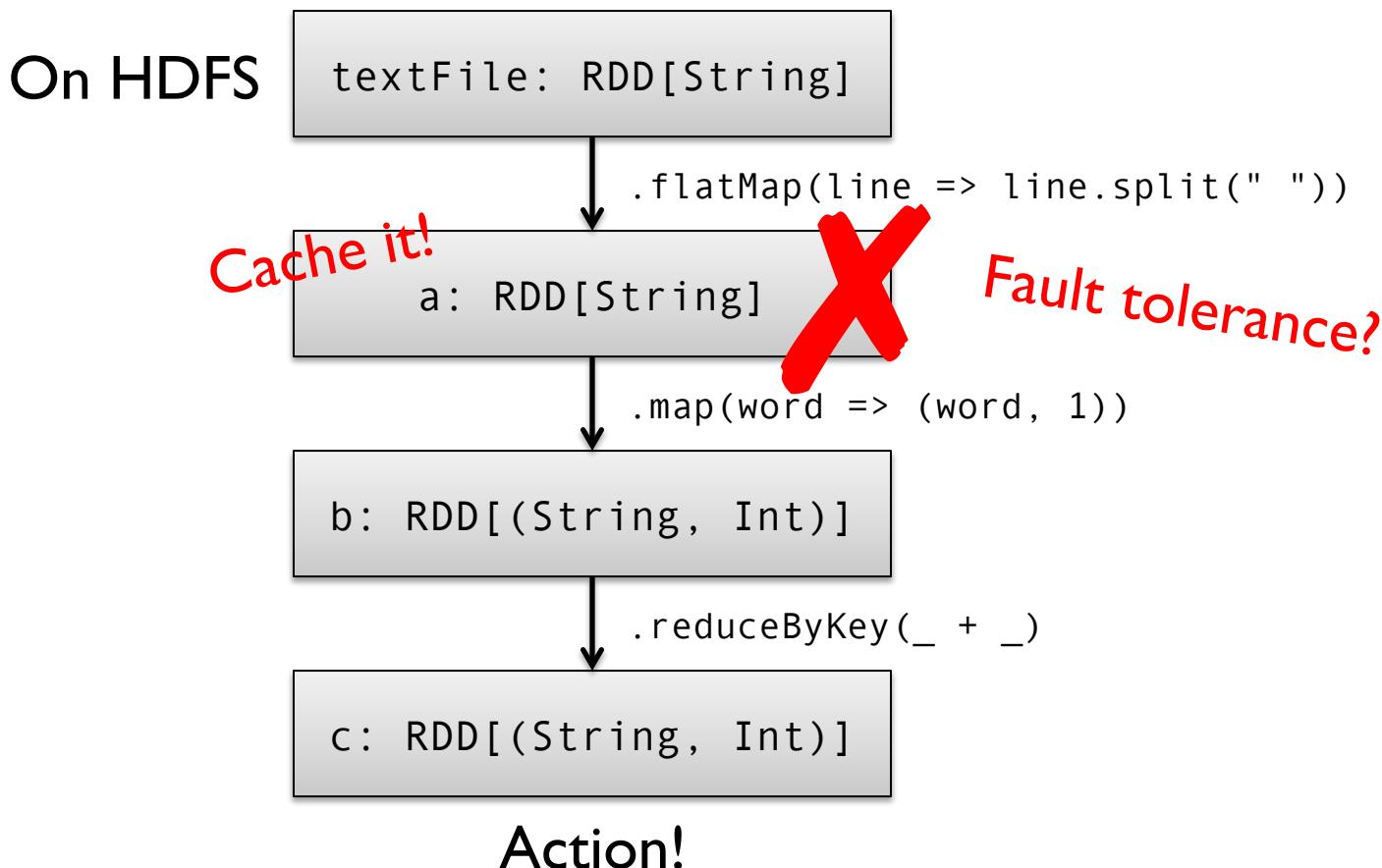
Want MM?



Action!

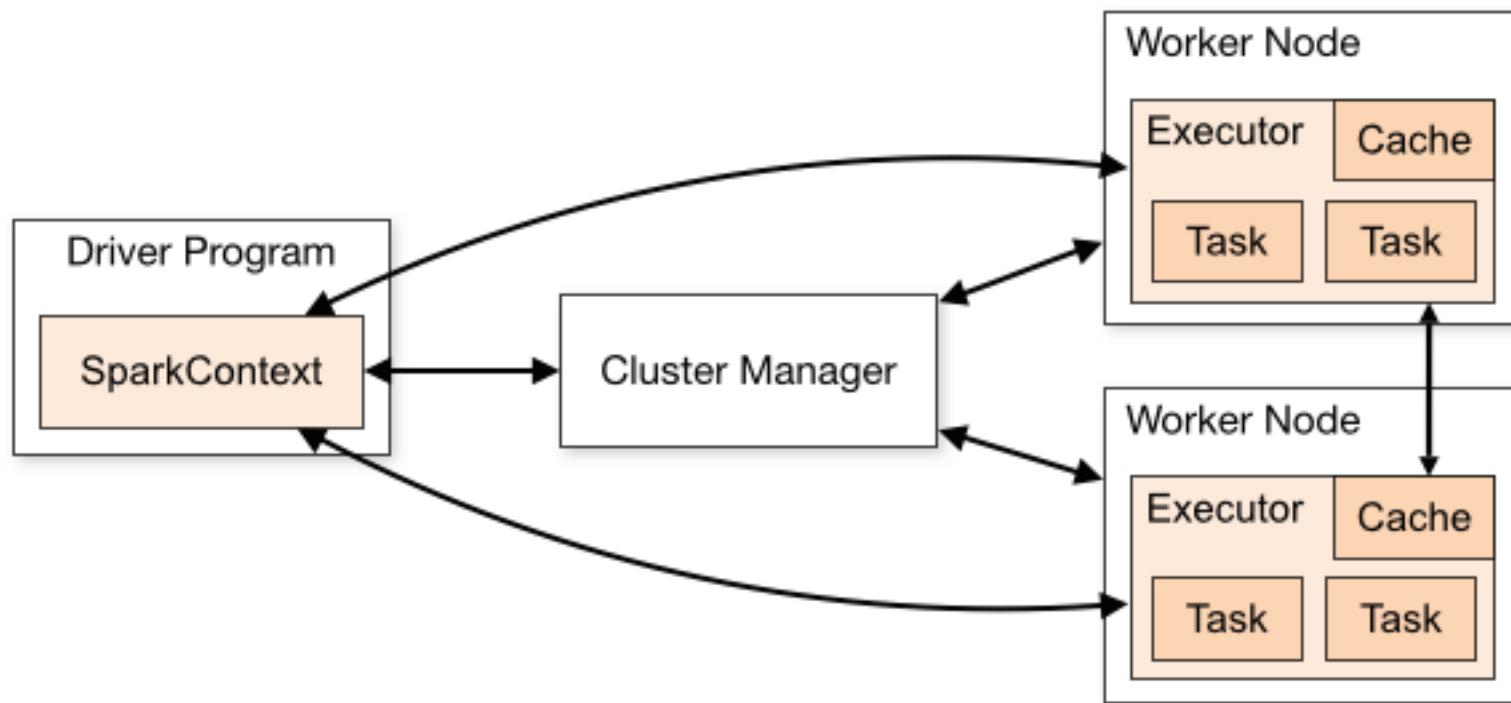
RDDs and Caching

RDDs can be materialized in memory!

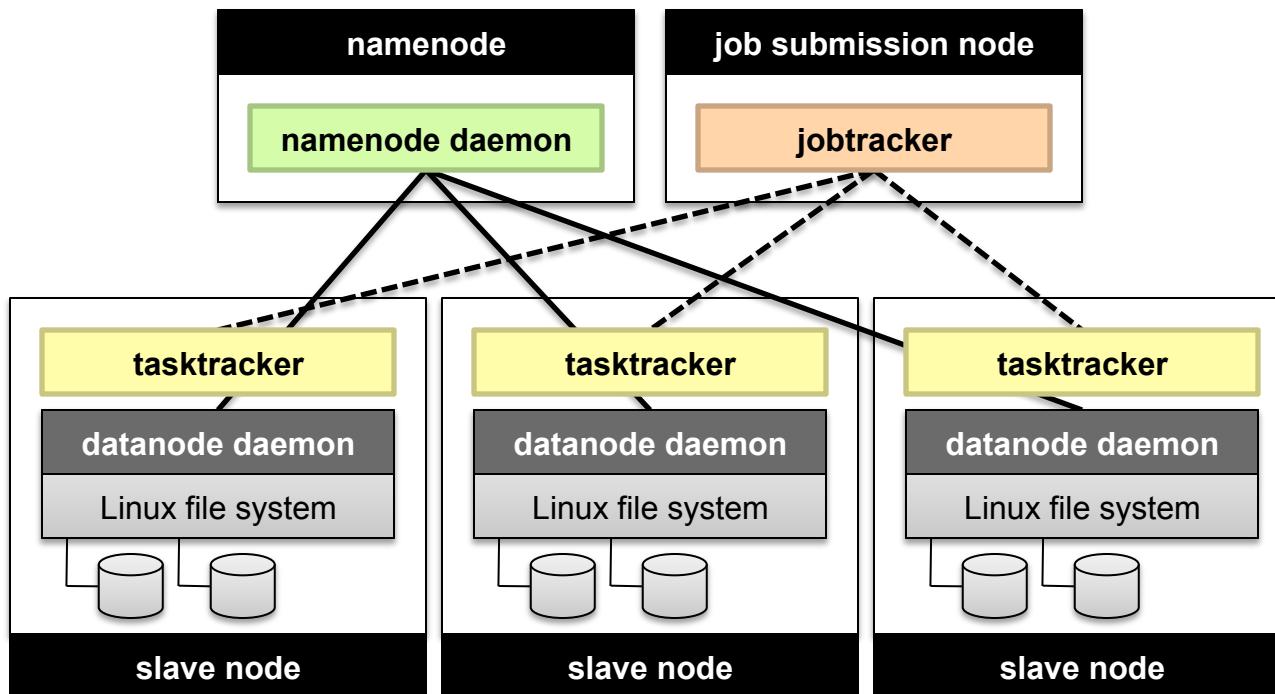


Spark works even if the RDDs are *partially* cached!

Spark Architecture



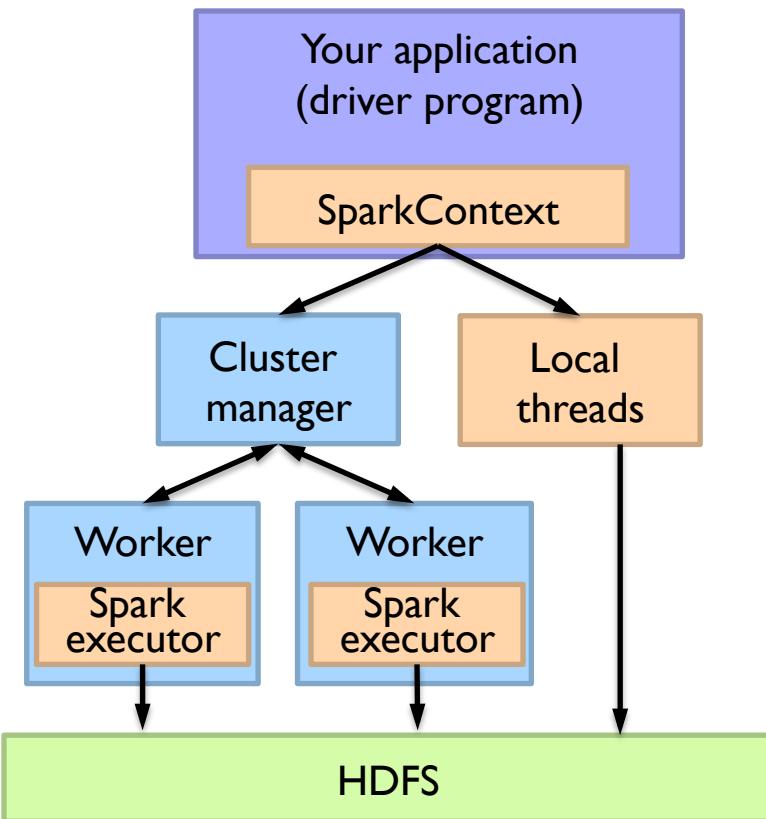
Compare



Spark Programs

Scala, Java, Python, R

spark-shell spark-submit

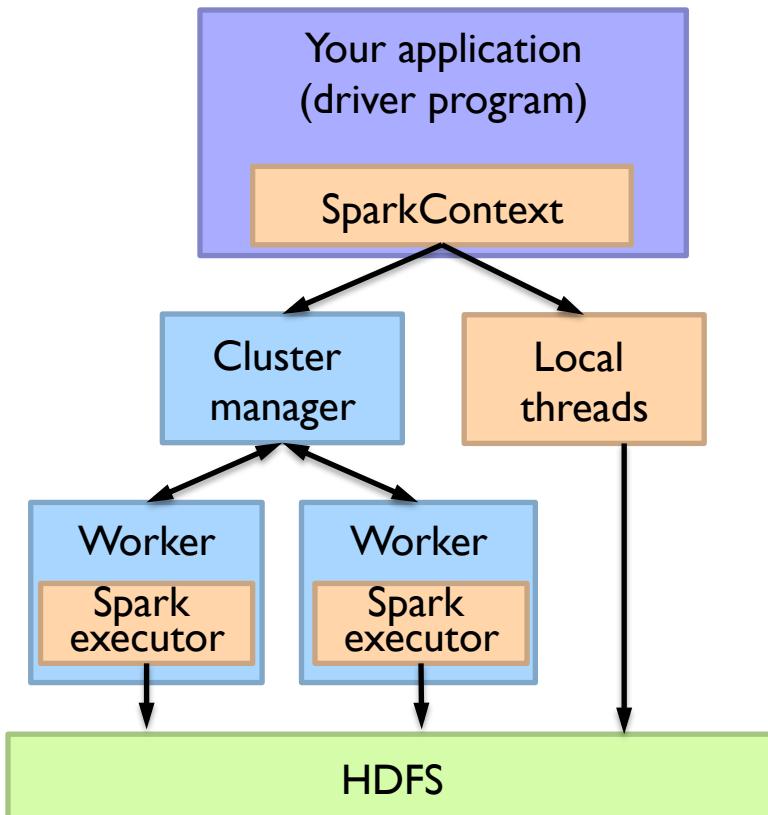


Spark context: tells the framework where to find the cluster

Use the Spark context to create RDDs

Spark Driver

spark-shell spark-submit

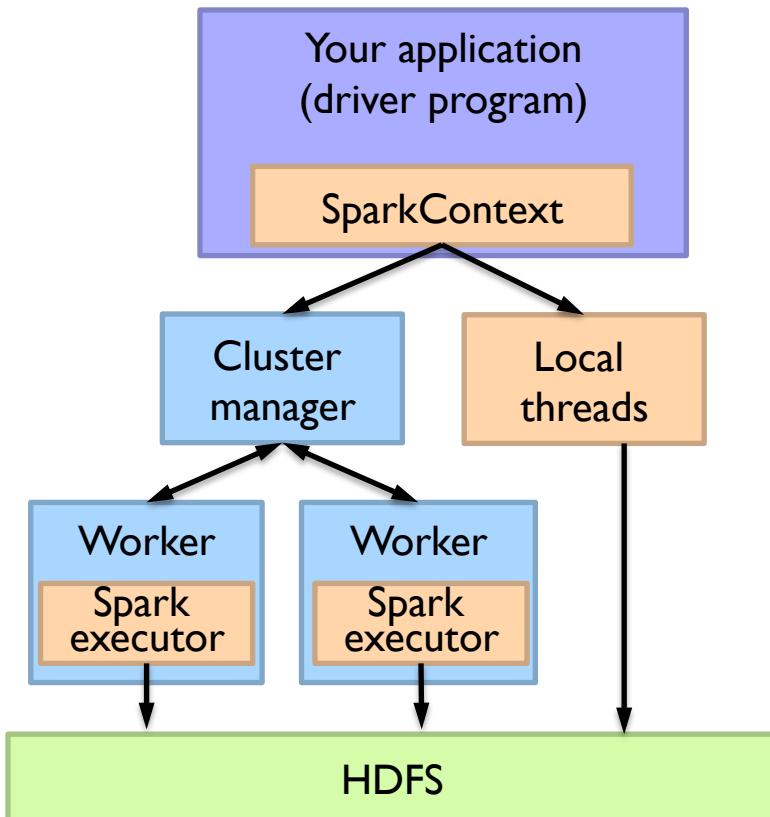


```
val textFile =  
  sc.textFile(args.input())  
  
textFile  
  .flatMap(line => tokenize(line))  
  .map(word => (word, 1))  
  .reduceByKey(_ + _)  
  .saveAsTextFile(args.output())
```

What's happening
to the functions?

Spark Driver

spark-shell spark-submit

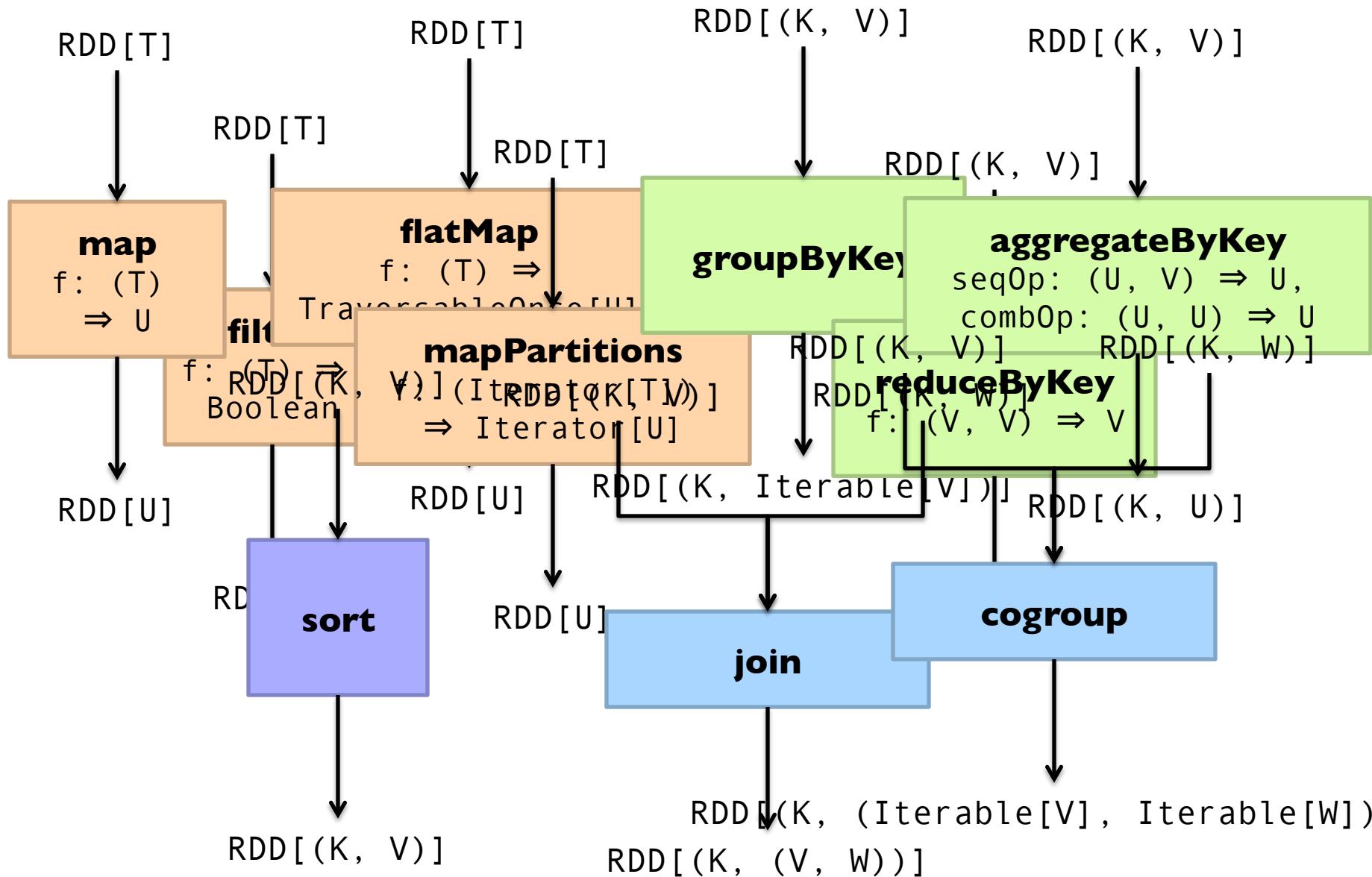


```
val textFile =  
  sc.textFile(args.input())  
  
textFile  
  .flatMap(line => tokenize(line))  
  .map(word => (word, 1))  
  .reduceByKey(_ + _)  
  .saveAsTextFile(args.output())
```

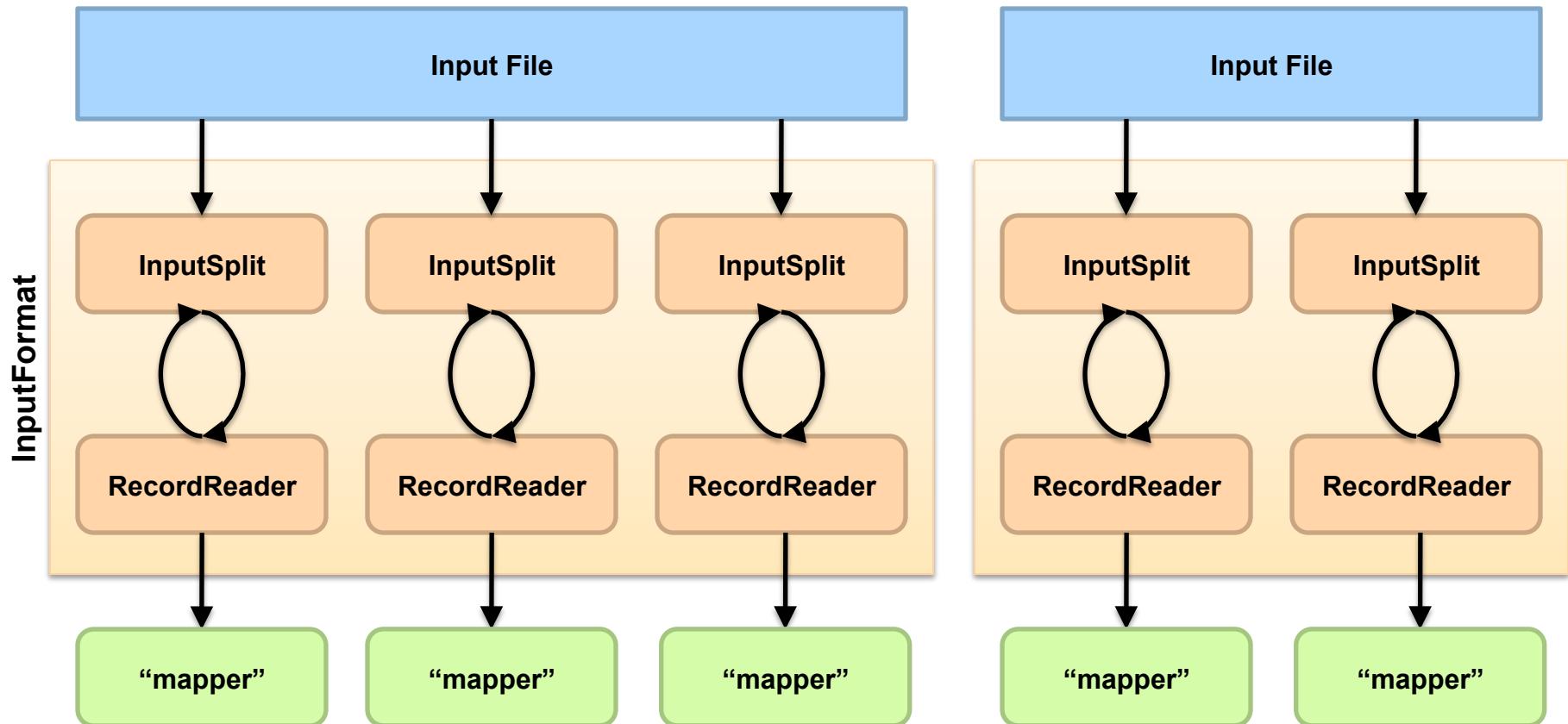
Note: you can run code “locally”,
integrate cluster-computed values!

Beware of the collect action!

Spark Transformations

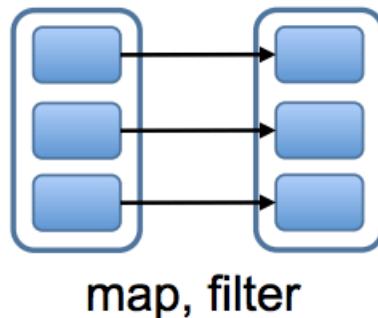


Starting Points

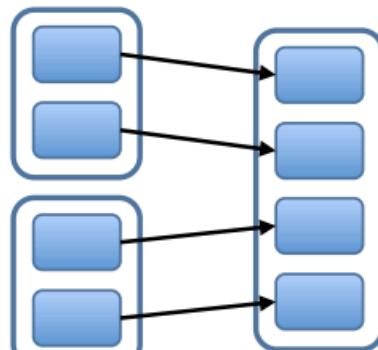


Physical Operators

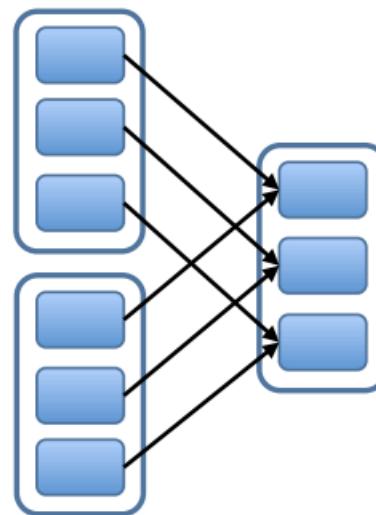
Narrow Dependencies:



map, filter

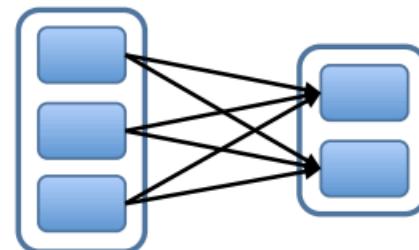


union

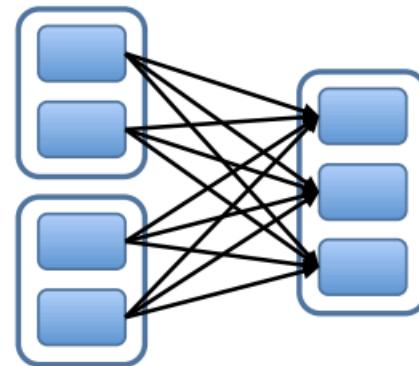


join with inputs
co-partitioned

Wide Dependencies:

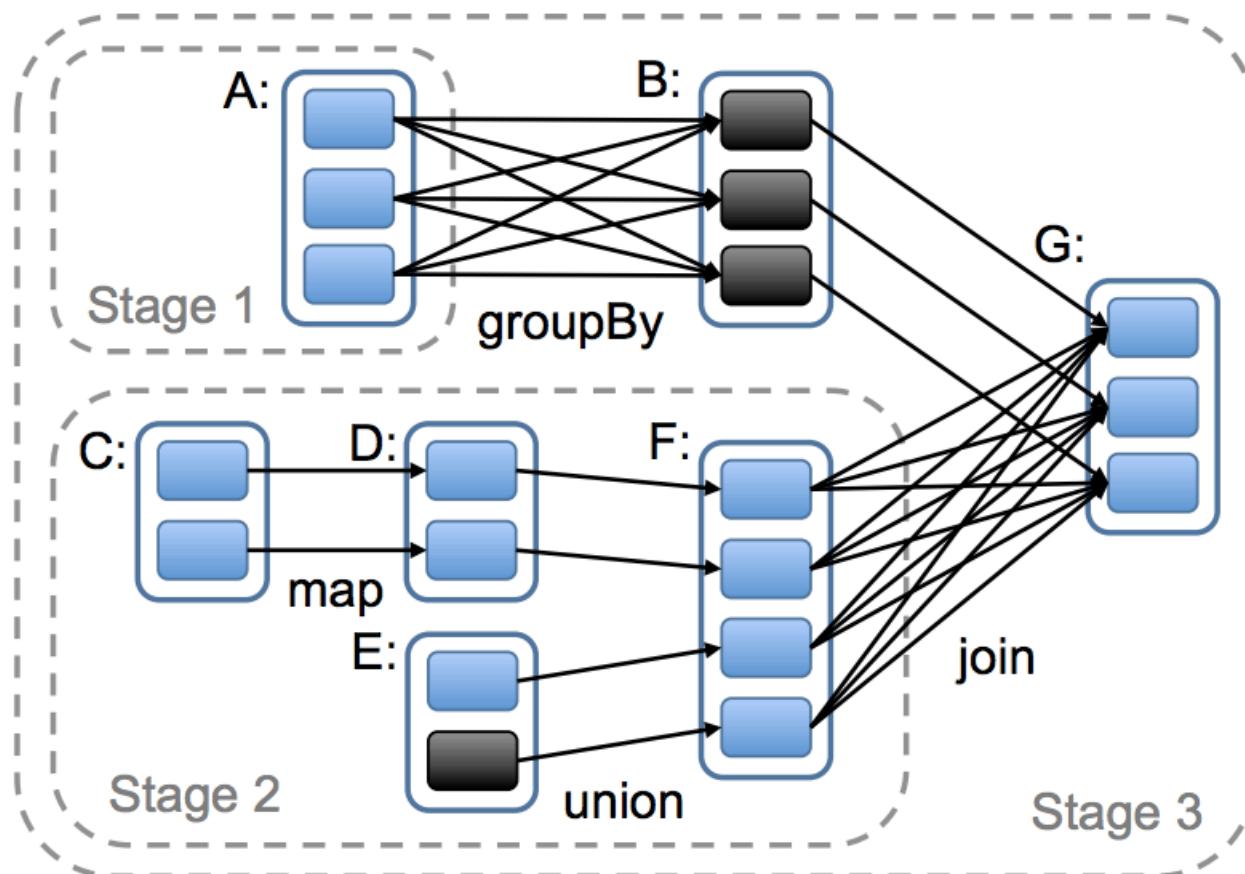


groupByKey

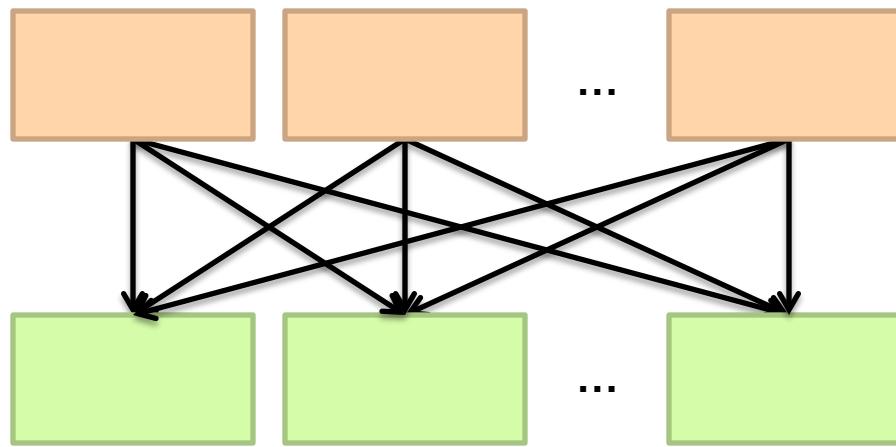


join with inputs not
co-partitioned

Execution Plan

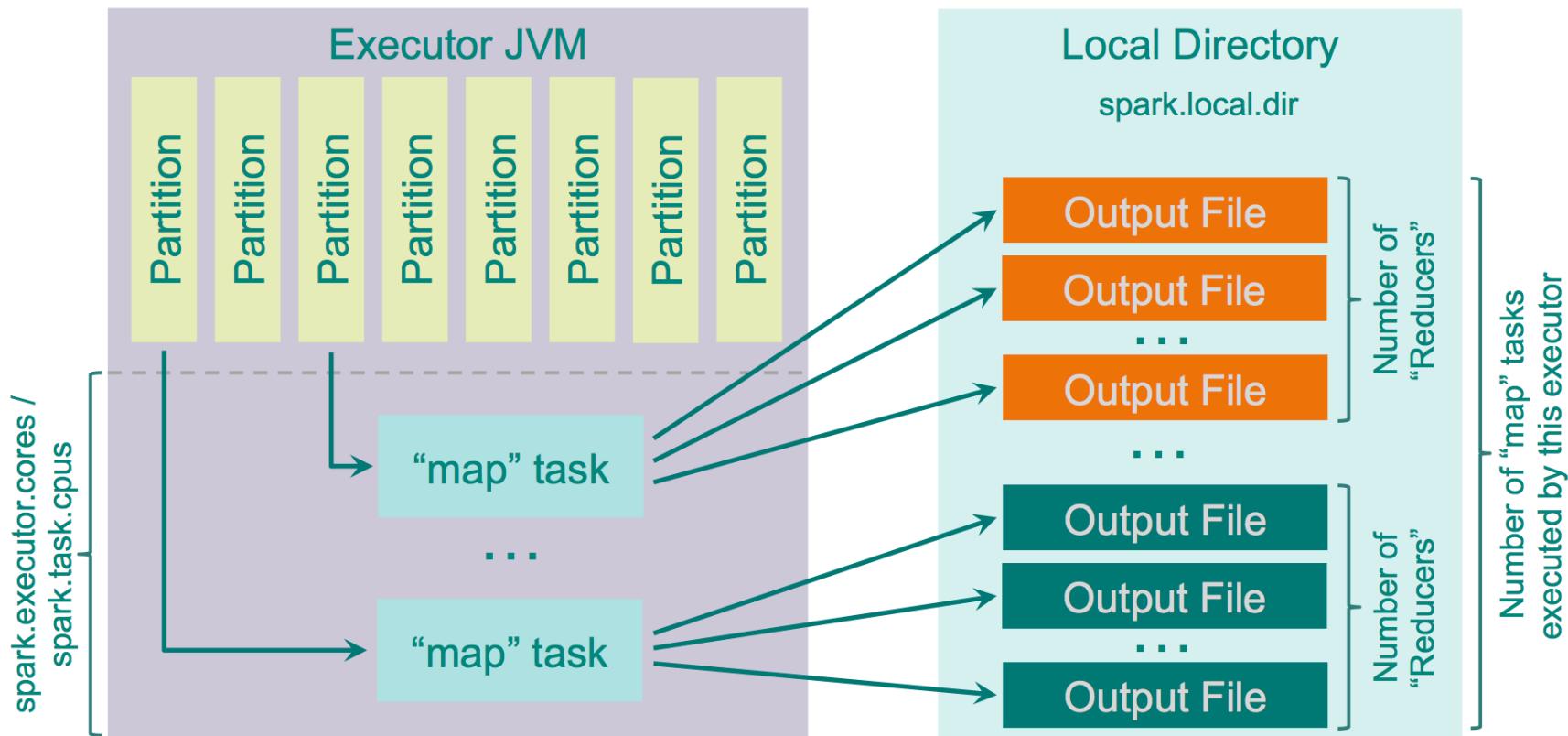


Can't avoid this!



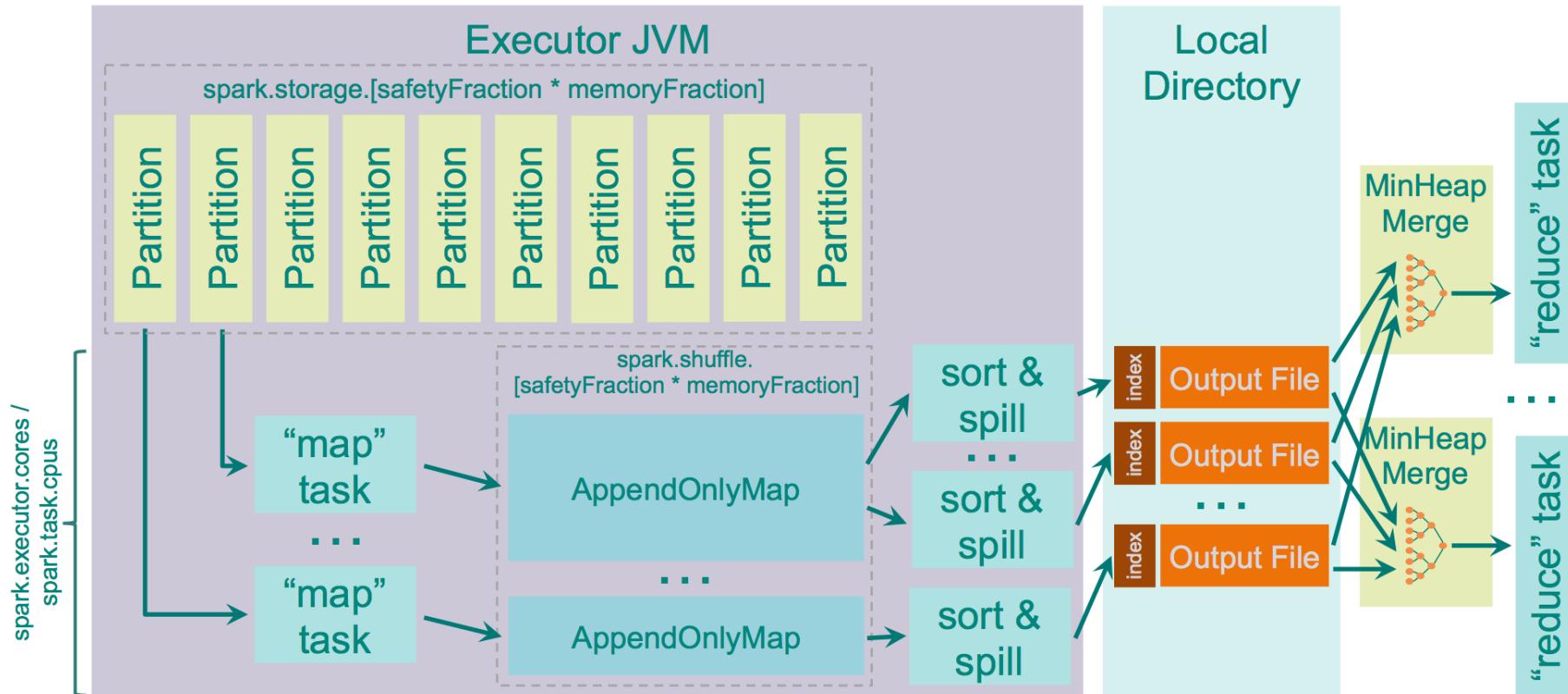
Spark Shuffle Implementations

Hash shuffle

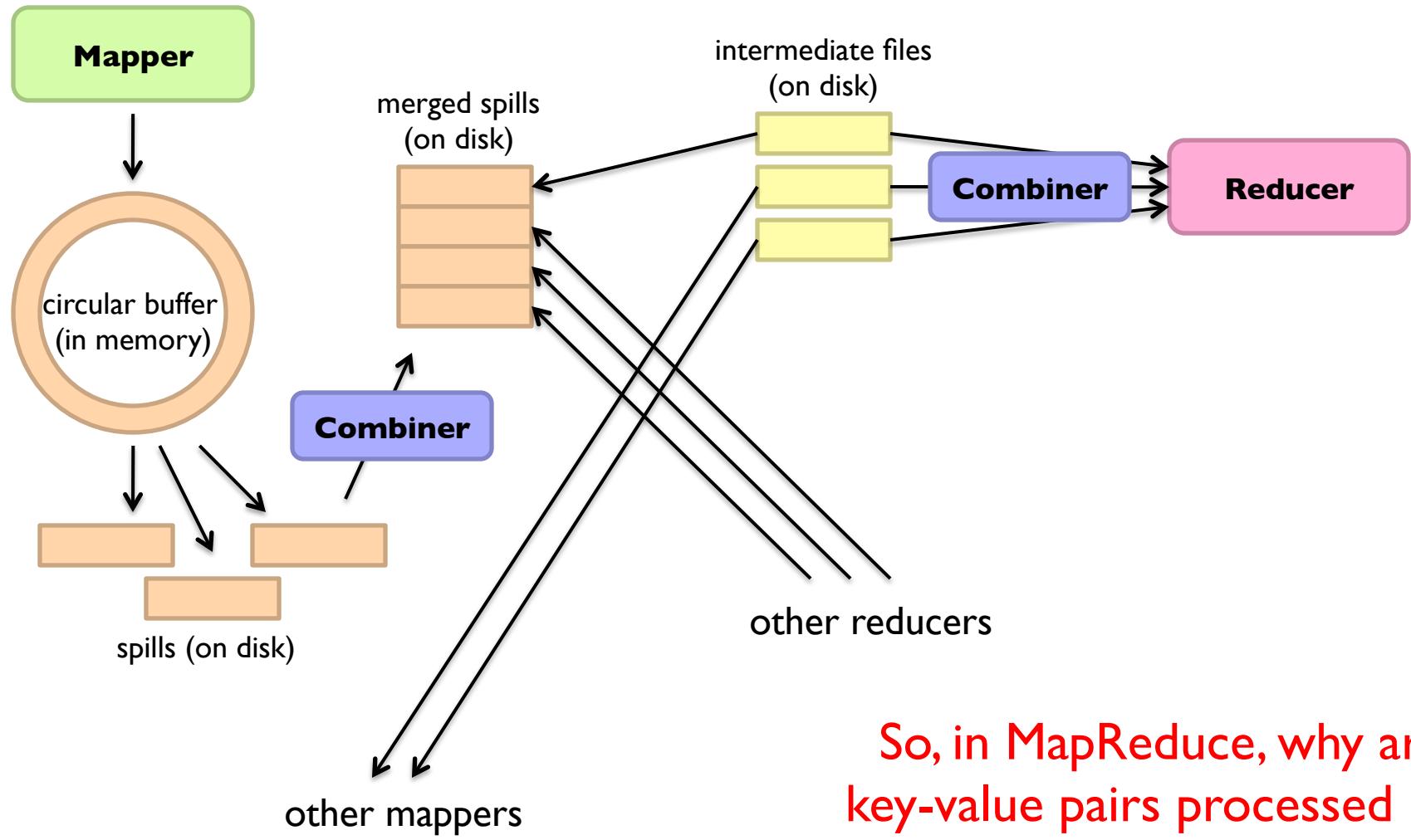


Spark Shuffle Implementations

Sort shuffle

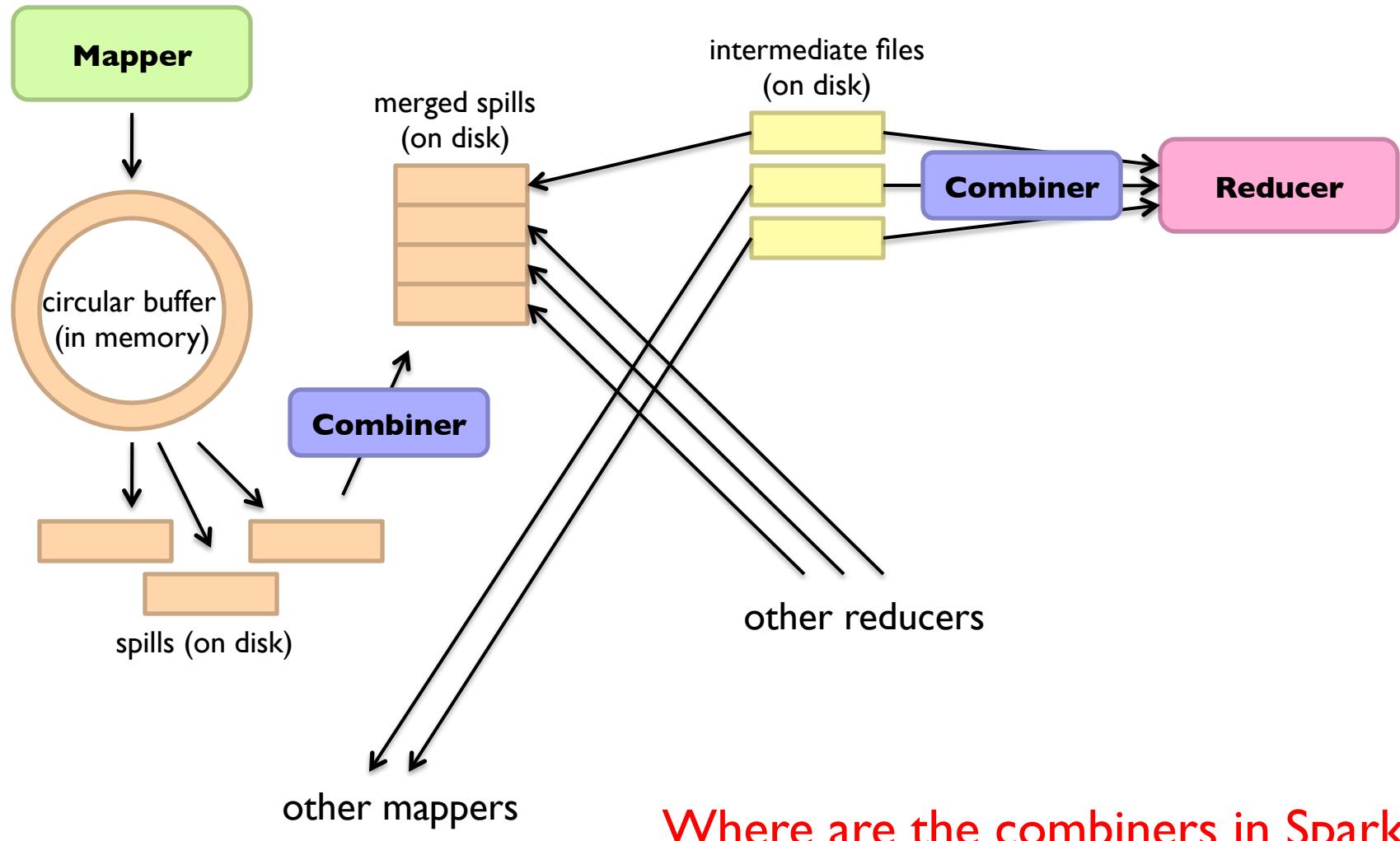


Remember this?

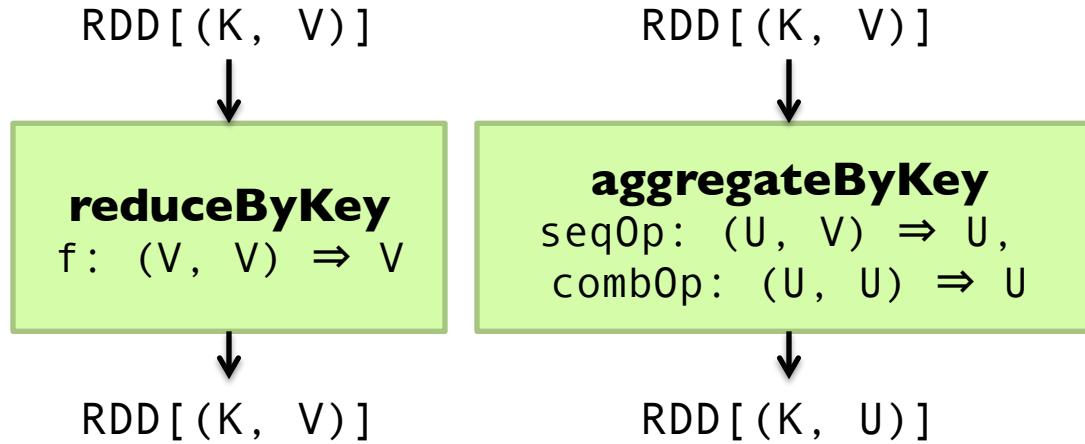


So, in MapReduce, why are key-value pairs processed in sorted order in the reducer?

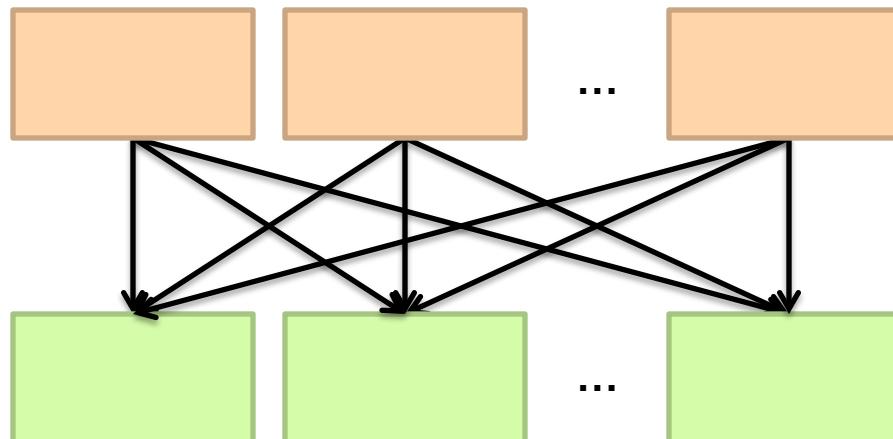
Remember this?



Reduce-like Operations



How can we optimize?
What happened to combiners?



Spark #wins

Richer operators

RDD abstraction supports
optimizations (pipelining, caching, etc.)

Scala, Java, Python, R, bindings



Algorithm design, redux



Two superpowers:

Associativity
Commutativity
(sorting)

What follows... very basic category theory...

The Power of Associativity

You can put parenthesis where ever you want!

$$(v_1 \oplus v_2 \oplus v_3) \oplus (v_4 \oplus v_5 \oplus v_6 \oplus v_7) \oplus (v_8 \oplus v_9)$$

$$(v_1 \oplus v_2) \oplus (v_3 \oplus v_4 \oplus v_5) \oplus (v_6 \oplus v_7 \oplus v_8 \oplus v_9)$$

$$(v_1 \oplus v_2 \oplus (v_3 \oplus v_4 \oplus v_5)) \oplus (v_6 \oplus v_7 \oplus v_8 \oplus v_9)$$

The Power of Commutativity

You can swap order of operands however you want!

$$(v_1 \oplus v_2 \oplus v_3) \oplus (v_4 \oplus v_5 \oplus v_6 \oplus v_7) \oplus (v_8 \oplus v_9)$$

$$(v_4 \oplus v_5 \oplus v_6 \oplus v_7) \oplus (v_1 \oplus v_2 \oplus v_3) \oplus (v_8 \oplus v_9)$$

$$(v_8 \oplus v_9) \oplus (v_4 \oplus v_5 \oplus v_6 \oplus v_7) \oplus (v_1 \oplus v_2 \oplus v_3)$$

Implications for distributed processing?

You don't know when the tasks begin

You don't know when the tasks end

You don't know when the tasks interrupt each other

You don't know when intermediate data arrive

...

It's okay!

Fancy Labels for Simple Concepts...

Semigroup = (M , \oplus)

$\oplus : M \times M \rightarrow M$, s.t., $\forall m_1, m_2, m_3 \in M$

$$(m_1 \oplus m_2) \oplus m_3 = m_1 \oplus (m_2 \oplus m_3)$$

Monoid = Semigroup + identity

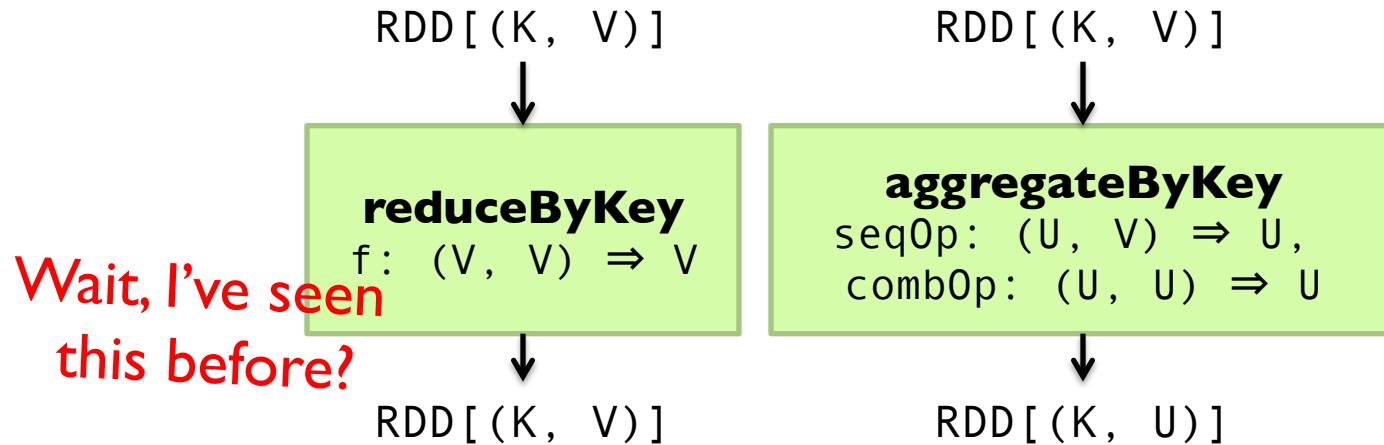
ε s.t., $\varepsilon \oplus m = m \oplus \varepsilon = m$, $\forall m \in M$

Commutative Monoid = Monoid + commutativity

$\forall m_1, m_2 \in M, m_1 \oplus m_2 = m_2 \oplus m_1$

A few examples?

Back to these...



Computing the mean via (again)

```
1: class MAPPER
2:   method MAP(string t, integer r)
3:     EMIT(string t, integer r)
4:
5: class REDUCER
6:   method REDUCE(string t, integers [r1, r2, ...])
7:     sum  $\leftarrow$  0
8:     cnt  $\leftarrow$  0
9:     for all integer r  $\in$  integers [r1, r2, ...] do
10:       sum  $\leftarrow$  sum + r
11:       cnt  $\leftarrow$  cnt + 1
12:     ravg  $\leftarrow$  sum/cnt
13:     EMIT(string t, integer ravg)
```

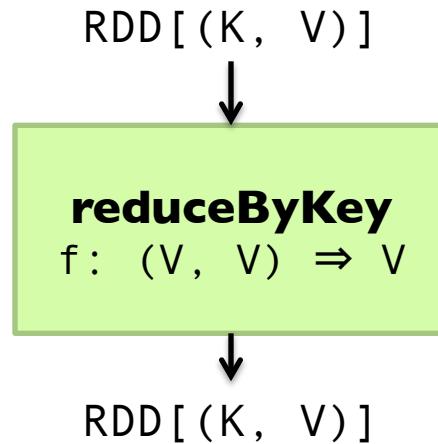
Computing the mean v2 (again)

```
1: class MAPPER
2:   method MAP(string t, integer r)
3:     EMIT(string t, integer r)
4:
5: class COMBINER
6:   method COMBINE(string t, integers [r1, r2, ...])
7:     sum  $\leftarrow$  0
8:     cnt  $\leftarrow$  0
9:     for all integer r  $\in$  integers [r1, r2, ...] do
10:       sum  $\leftarrow$  sum + r
11:       cnt  $\leftarrow$  cnt + 1
12:     EMIT(string t, pair (sum, cnt))            $\triangleright$  Separate sum and count
13:
14: class REDUCER
15:   method REDUCE(string t, pairs [(s1, c1), (s2, c2) ...])
16:     sum  $\leftarrow$  0
17:     cnt  $\leftarrow$  0
18:     for all pair (s, c)  $\in$  pairs [(s1, c1), (s2, c2) ...] do
19:       sum  $\leftarrow$  sum + s
20:       cnt  $\leftarrow$  cnt + c
21:     ravg  $\leftarrow$  sum/cnt
22:     EMIT(string t, integer ravg)
```

Computing the mean v3 (again)

```
1: class MAPPER
2:     method MAP(string t, integer r)
3:         EMIT(string t, pair (r, 1))
4:
5: class COMBINER
6:     method COMBINE(string t, pairs [(s1, c1), (s2, c2) ...])
7:         sum ← 0
8:         cnt ← 0
9:         for all pair (s, c) ∈ pairs [(s1, c1), (s2, c2) ...] do
10:             sum ← sum + s
11:             cnt ← cnt + c
12:         EMIT(string t, pair (sum, cnt))
13:
14: class REDUCER
15:     method REDUCE(string t, pairs [(s1, c1), (s2, c2) ...])
16:         sum ← 0
17:         cnt ← 0
18:         for all pair (s, c) ∈ pairs [(s1, c1), (s2, c2) ...] do
19:             sum ← sum + s
20:             cnt ← cnt + c
21:             ravg ← sum / cnt
22:         EMIT(string t, pair (ravg, cnt))
```

Wait, I've seen
this before?

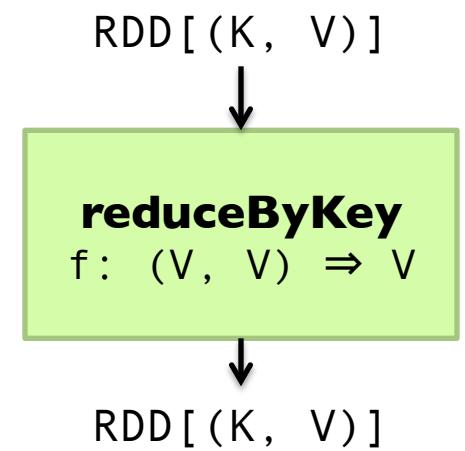


Co-occurrence Matrix: Stripes

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

Wait, I've seen this before?

```
1: class REDUCER
2:   method REDUCE(term w, stripes [H1, H2, H3, ...])
3:     Hf ← new ASSOCIATIVEARRAY
4:     for all stripe H ∈ stripes [H1, H2, H3, ...] do
5:       SUM(Hf, H)
6:     EMIT(term w, stripe Hf)
```



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

What about this?

Commutative monoids!

$f(B|A)$: “Pairs”

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

$(a, *) \rightarrow 32$

Reducer holds this value in memory

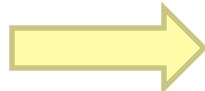
$(a, b_1) \rightarrow 3$

$(a, b_2) \rightarrow 12$

$(a, b_3) \rightarrow 7$

$(a, b_4) \rightarrow 1$

...



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

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

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

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

...

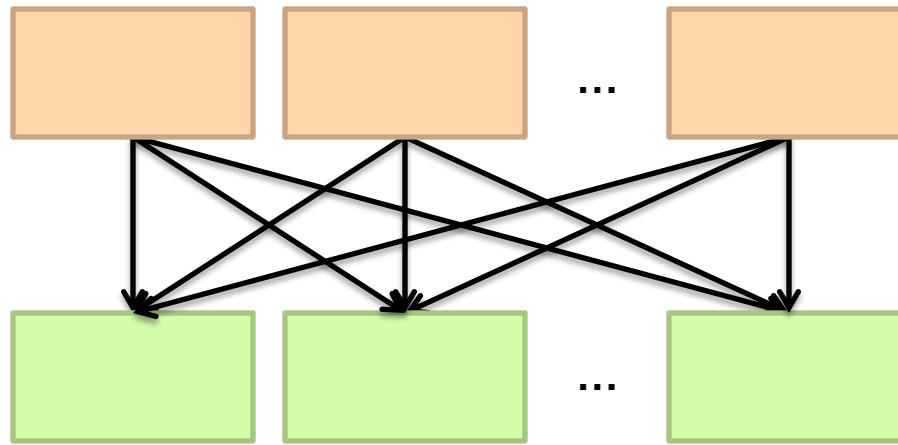
- For this to work:

- 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



Two superpowers:
Associativity
Commutativity
(sorting)

Because you can't avoid this...



And sort-based shuffling is pretty efficient!

An aerial photograph of a large data center complex during sunset. The sky is a warm orange and yellow. In the foreground, there are several large white industrial buildings, some with flat roofs and others with gabled roofs. A parking lot with many cars is visible in front of one of the buildings. To the right, there is a large field with several white cylindrical storage tanks lined up in rows. In the background, there is a highway with traffic and a small town or cluster of buildings. The overall scene is a mix of industrial infrastructure and rural landscape.

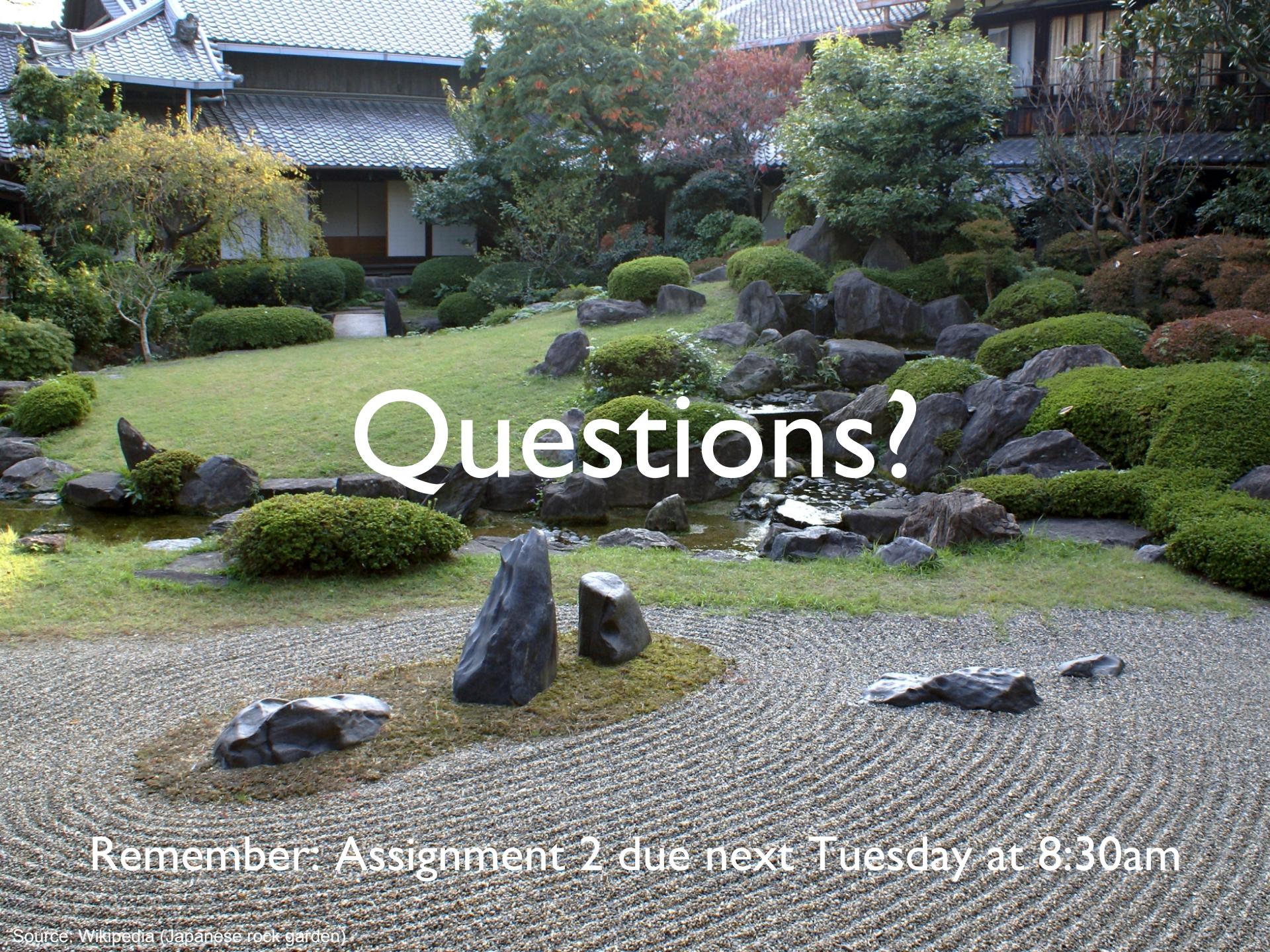
The datacenter *is* the computer!
What's the instruction set?

Algorithm design in a nutshell...



Exploit associativity and commutativity
via commutative monoids (if you can)

Exploit framework-based sorting to
sequence computations (if you can't)

A photograph of a traditional Japanese rock garden. The foreground features a gravel path with raked patterns. Several large, dark, irregular stones are scattered across the garden. In the middle ground, there is a small pond with rocks and a waterfall. The background consists of a building with a tiled roof and various trees and shrubs, some with autumn-colored leaves.

Questions?

Remember: Assignment 2 due next Tuesday at 8:30am