

# Spark

*CS6450: Distributed Systems*

Lecture 18

Ryan Stutsman

Some content adapted from Matei's NSDI talk.

Material taken/derived from Princeton COS-418 materials created by Michael Freedman and Kyle Jamieson at Princeton University.

Licensed for use under a Creative Commons Attribution-NonCommercial-ShareAlike 3.0 Unported License.

Some material taken/derived from MIT 6.824 by Robert Morris, Franz Kaashoek, and Nickolai Zeldovich.

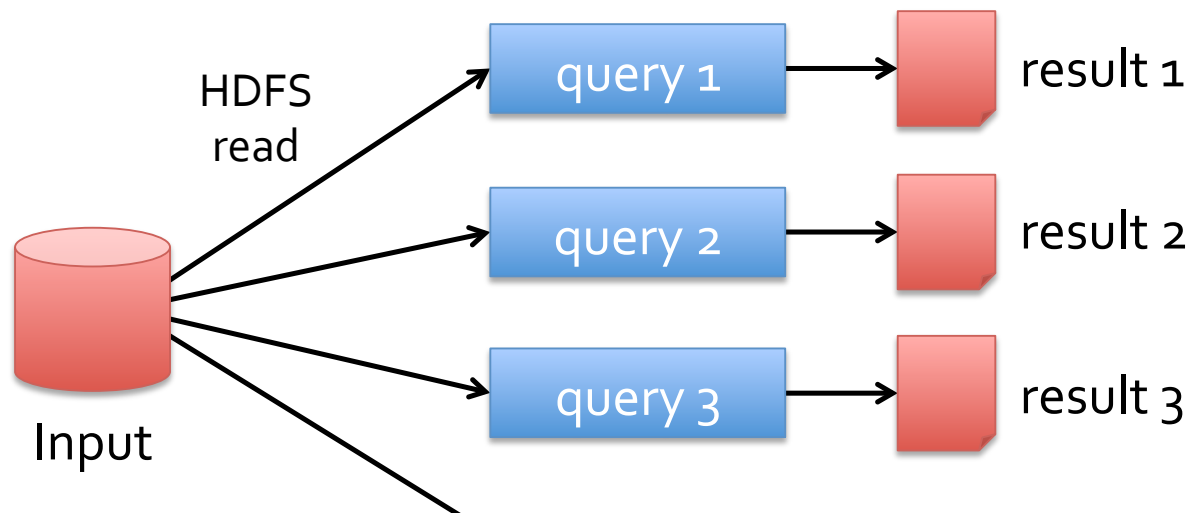
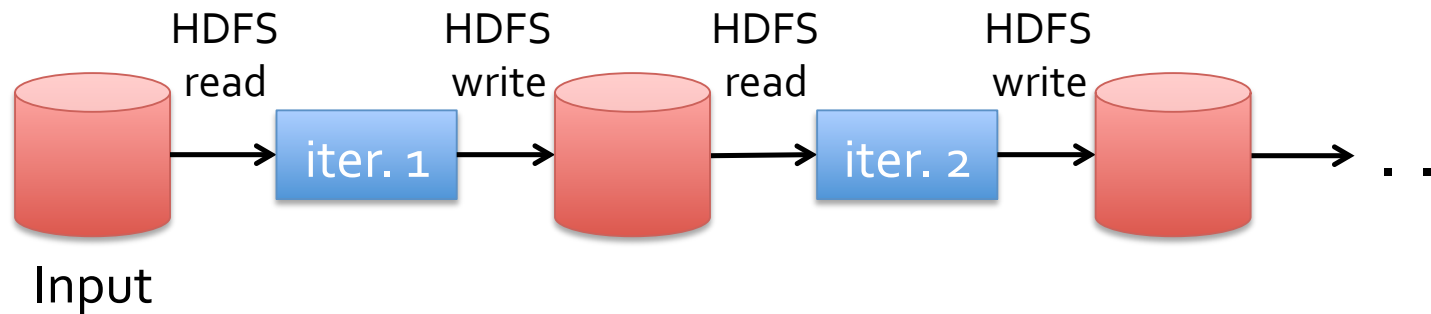
# Problems with Map-Reduce

- Scaled analytics to thousands of machines
- Eliminated fault-tolerance as a concern
- Not very expressive
  - Iterative algorithms  
(PageRank, Logistic Regression, Transitive Closure)
  - Interactive and ad-hoc queries  
(Interactive Log Debugging)
- Lots of specialized frameworks
  - Pregel, GraphLab, PowerGraph, DryadLINQ, HaLoop...

# Sharing Data between Iterations/Ops

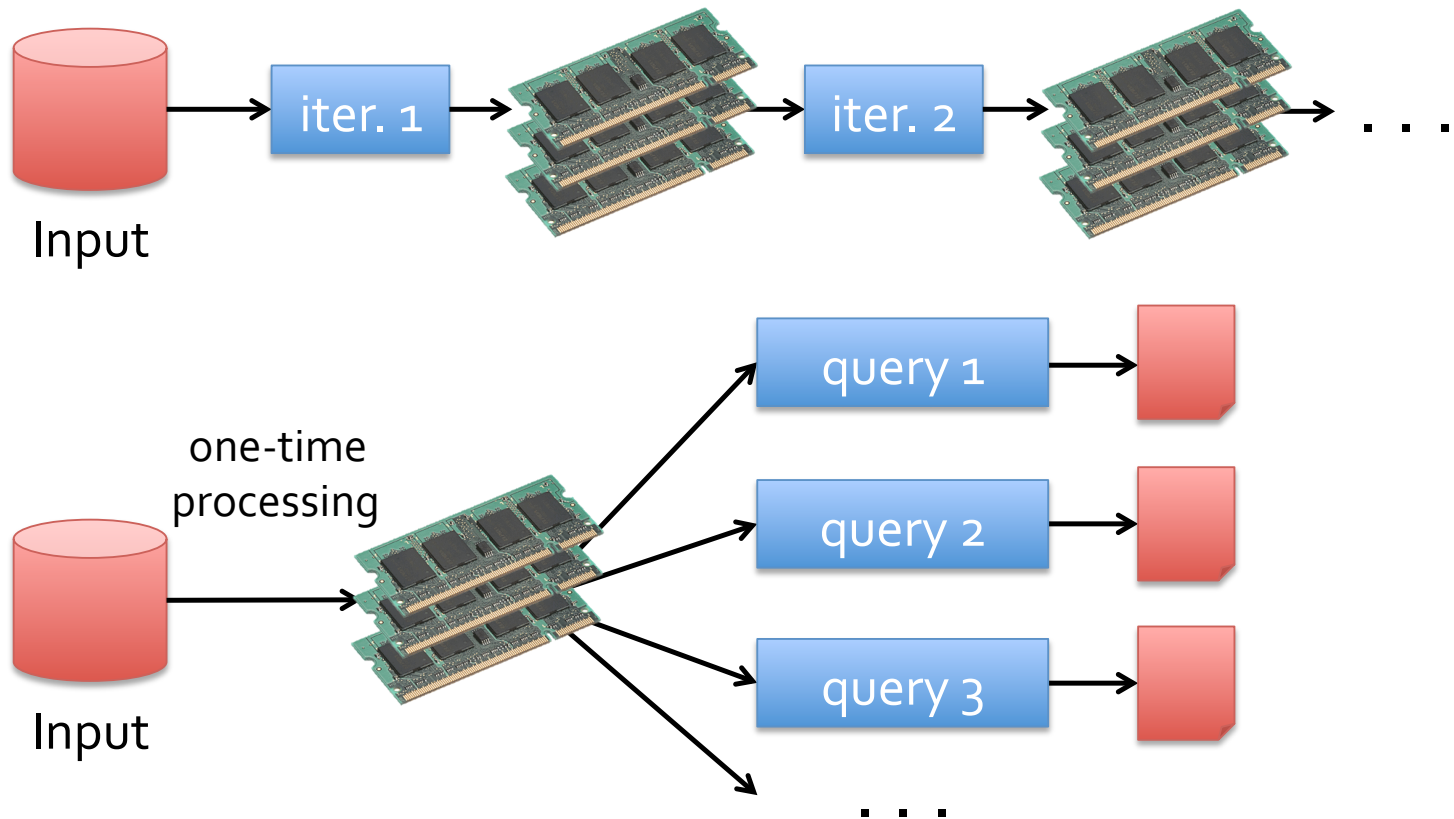
- Only way to share data between iterations/phases is through shared storage
- Allow operations to feed data to one another
  - Ideally, through memory instead of disk-based storage
- Need the "chain" of operations to be exposed to make this work
- Also, does this break the MR fault-tolerance scheme?
  - Retry any Map or Reduce task since idempotent

# Examples



Slow due to replication and disk I/O,  
but necessary for fault tolerance

# Goal: In-memory Data Sharing



10-100x faster than network/disk, but how to get FT?

# Challenges

- Want distributed memory abstraction that is both fault-tolerant and efficient
- Existing storage allow fine-grained mutation to state
  - In-memory Key-value stores
  - But, they require costly on-the-fly replication for mutations
- Insight: leverage similar coarse-grained approach that transforms whole data set per op, like MR

# Resilient Distributed Datasets (RDDs)

- Restricted form of distributed shared memory
  - Immutable, partitioned collections of records
  - Can only be built through coarse-grained deterministic *transformations*
  - Map, filter, join...
- Efficient fault recovery using *lineage*
  - Log one operation to apply to many elements
  - Recompute lost partitions on failure
  - No cost if nothing fails

# Spark Programming Interface

- Scala API, exposed within interpreter as well
- RDDs
- Transformations on RDDs (RDD  $\rightarrow$  RDD)
- Actions on RDDs (RDD  $\rightarrow$  output)
- Control over RDD partitioning (how items are split over nodes)
- Control over RDD persistence (in RAM, on disk, or recompute on loss)



# Transformations

<i>map</i> ( $f : T \Rightarrow U$ )	:	$\text{RDD}[T] \Rightarrow \text{RDD}[U]$
<i>filter</i> ( $f : T \Rightarrow \text{Bool}$ )	:	$\text{RDD}[T] \Rightarrow \text{RDD}[T]$
<i>flatMap</i> ( $f : T \Rightarrow \text{Seq}[U]$ )	:	$\text{RDD}[T] \Rightarrow \text{RDD}[U]$
<i>sample</i> ( <i>fraction</i> : Float)	:	$\text{RDD}[T] \Rightarrow \text{RDD}[T]$ (Deterministic sampling)
<i>groupByKey</i> ()	:	$\text{RDD}[(K, V)] \Rightarrow \text{RDD}[(K, \text{Seq}[V])]$
<i>reduceByKey</i> ( $f : (V, V) \Rightarrow V$ )	:	$\text{RDD}[(K, V)] \Rightarrow \text{RDD}[(K, V)]$
<i>union</i> ()	:	$(\text{RDD}[T], \text{RDD}[T]) \Rightarrow \text{RDD}[T]$
<i>join</i> ()	:	$(\text{RDD}[(K, V)], \text{RDD}[(K, W)]) \Rightarrow \text{RDD}[(K, (V, W))]$
<i>cogroup</i> ()	:	$(\text{RDD}[(K, V)], \text{RDD}[(K, W)]) \Rightarrow \text{RDD}[(K, (\text{Seq}[V], \text{Seq}[W]))]$
<i>crossProduct</i> ()	:	$(\text{RDD}[T], \text{RDD}[U]) \Rightarrow \text{RDD}[(T, U)]$
<i>mapValues</i> ( $f : V \Rightarrow W$ )	:	$\text{RDD}[(K, V)] \Rightarrow \text{RDD}[(K, W)]$ (Preserves partitioning)
<i>sort</i> ( $c : \text{Comparator}[K]$ )	:	$\text{RDD}[(K, V)] \Rightarrow \text{RDD}[(K, V)]$
<i>partitionBy</i> ( $p : \text{Partitioner}[K]$ )	:	$\text{RDD}[(K, V)] \Rightarrow \text{RDD}[(K, V)]$

RDDs in terms of Scala types  $\rightarrow$  Scala semantics at workers  
Transformations are lazy "thunks"; cause *no* cluster action

# Actions

*count()* :  $\text{RDD}[T] \Rightarrow \text{Long}$   
*collect()* :  $\text{RDD}[T] \Rightarrow \text{Seq}[T]$   
*reduce*( $f : (T, T) \Rightarrow T$ ) :  $\text{RDD}[T] \Rightarrow T$   
*lookup*( $k : K$ ) :  $\text{RDD}[(K, V)] \Rightarrow \text{Seq}[V]$  (On hash/range partitioned RDDs)  
*save*( $path : \text{String}$ ) : Outputs RDD to a storage system, *e.g.*, HDFS

Consumes an RDD to produce output

either to storage (*save*) or

to interpreter/Scala (*count*, *collect*, *reduce*)

Causes RDD lineage chain to get executed on the cluster to produce the output  
(for any missing piece of the computation)

# Interactive Debugging

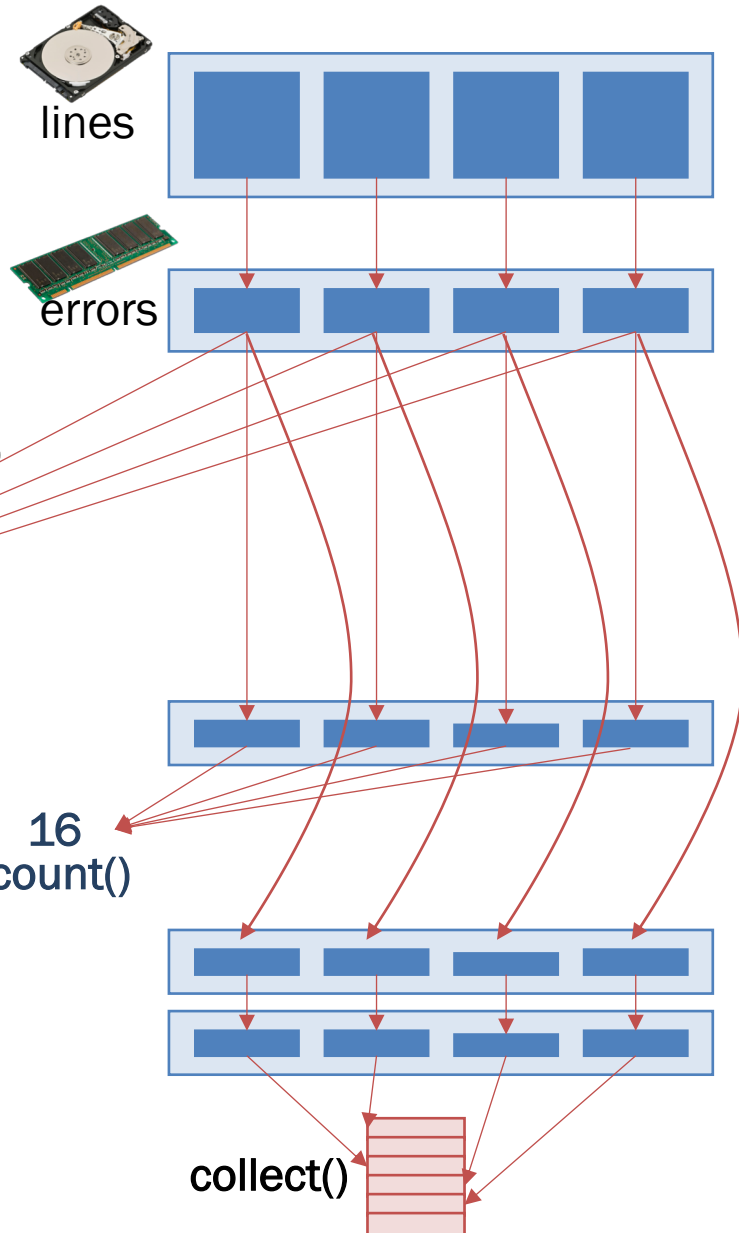
```
lines = textFile("hdfs://foo.log")
errors = lines.filter(
    _.startsWith("ERROR"))
errors.persist()
```

```
errors.count()
errors.filter(
    _.contains("MySQL")).count()
errors.filter(_.contains("HDFS"))
    .map(_.split("\t")(3))
    .collect()
```

32  
count()

16  
count()

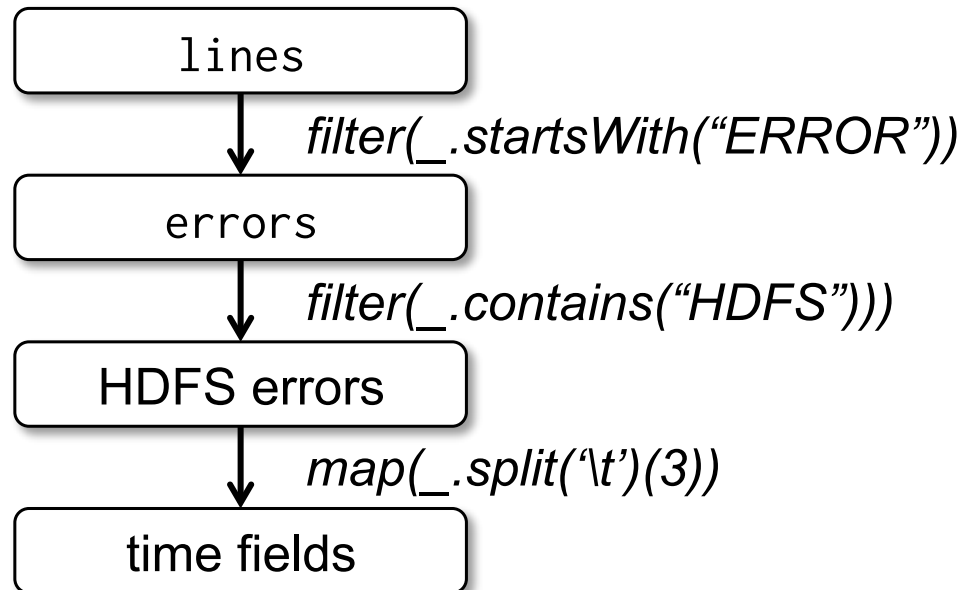
collect()



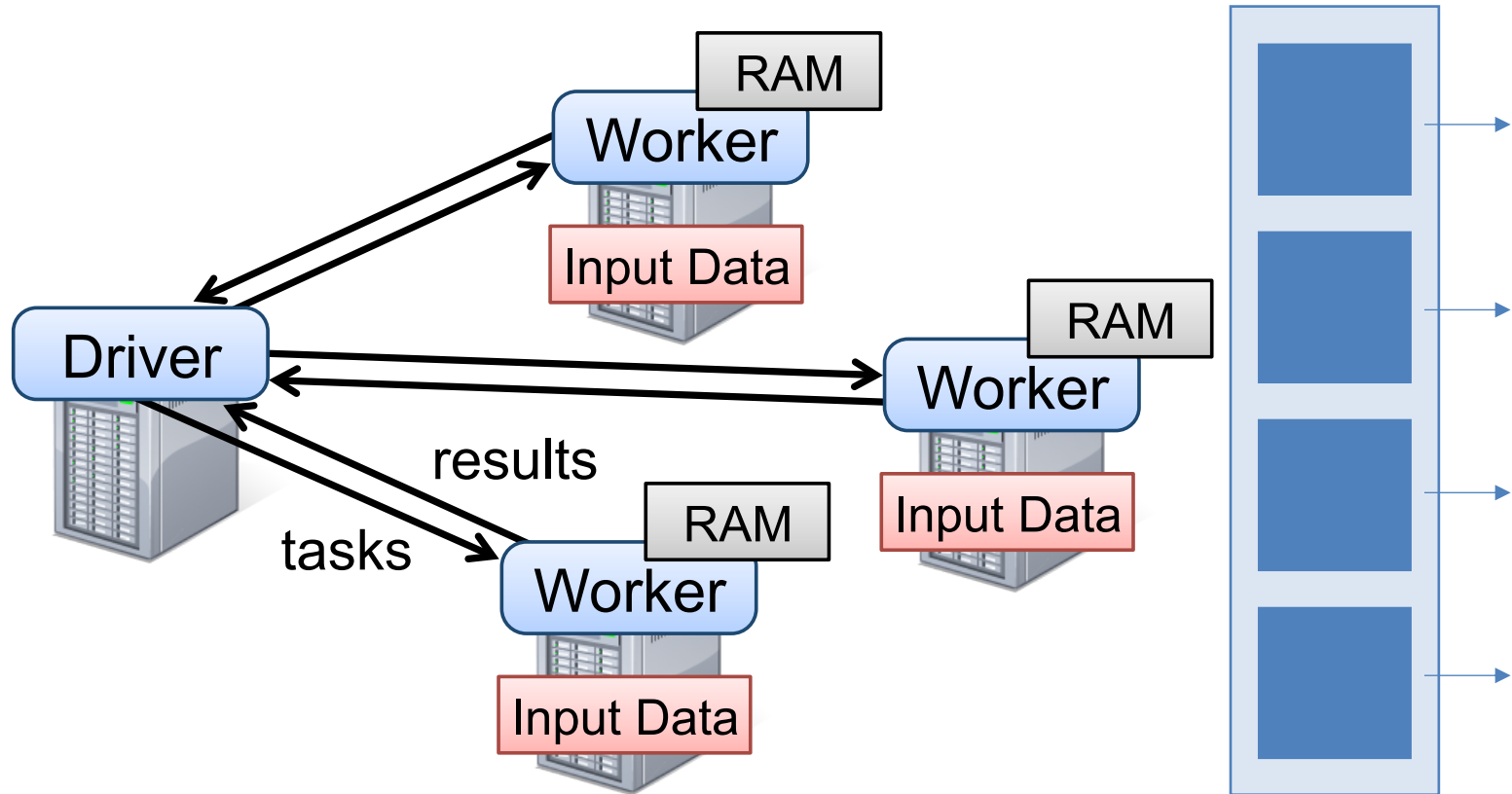
# `persist()`?

- Not an action and not a transformation
- A scheduler hint
- Tells which RDDs the Spark schedule should *materialize* and whether in memory or storage
- Gives the user control over reuse/recompute/recovery tradeoffs
- Q: If `persist()` asks for the materialization of an RDD why isn't it an action?

# Lineage Graph of RDDs



# Physical Execution of Tasks over RDDs



# Example: PageRank

1. Start each page with a rank of 1
2. On each iteration, update each page's rank to

$$\sum_{i \in \text{neighbors}} \text{rank}_i / |\text{neighbors}_i|$$

```
RDD[(URL, Seq[URL])]
links = // RDD of (url, neighbors) pairs
ranks = // RDD of (url, rank) pairs ← RDD[(URL, Rank)]

for (i <- 1 to ITERATIONS) {
  ranks = links.join(ranks).flatMap {
    (url, (links, rank)) =>
      links.map(dest => (dest, rank/links.size))
  }.reduceByKey(_ + _)
}
Reduce to RDD[(URL, Rank)]
```

For each neighbor in links, emits (URL, RankContrib)

# Join ( $\bowtie$ )

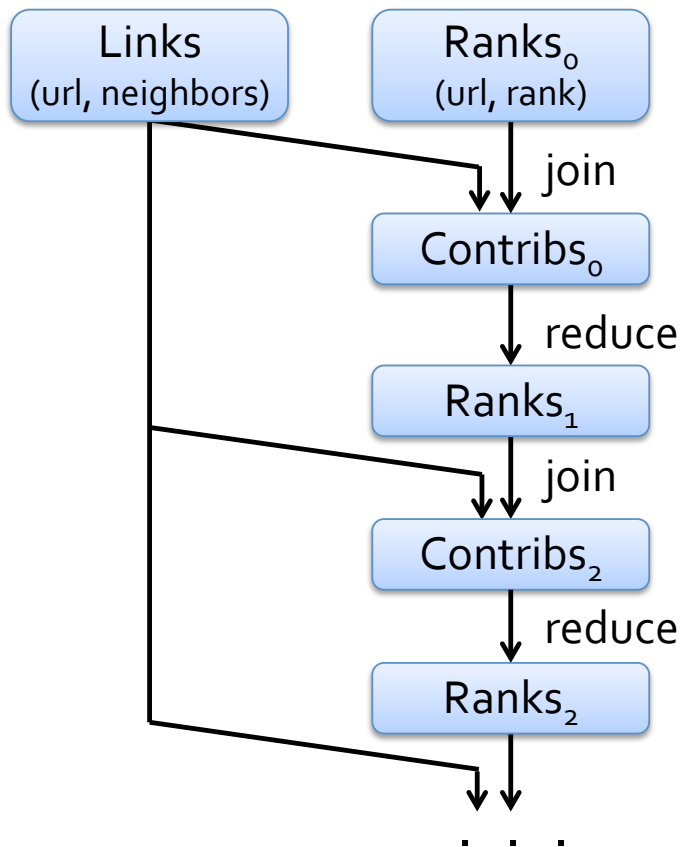
Ryan	5		Ryan	M		Ryan	5	M
Claire	6		Claire	F		Claire	6	F
Elliott	3		Elliott	M		Elliott	3	M

A	5		C	M
A	3		B	M
A	5		A	M
B	4		B	F
B	7		A	F
C	1		C	F
C	9		B	F

If partitioning  
doesn't match,  
then need to reshuffle  
to match pairs.  
Same problem in reduce()  
for Map-Reduce.



# Optimizing Placement



Links & ranks repeatedly joined

Can *co-partition* them (e.g. hash both on URL) to avoid shuffles

Can also use app knowledge, e.g., hash on DNS name

```
links = links.partitionBy(  
    new URLPartitioner())
```

Q: Where might we have placed `persist()`?

# Co-partitioning Example

From join of “top” partitions

foo.com	5
foo.com	4
widget.com	3
bar.com	2
foo.com	2

bar.com	2
bad.com	0
foo.com	11
widget.com	3

From join of “bottom” partitions

bar.com	foo.com
bad.com	foo.com
foo.com	widget.com
widget.com	bar.com foo.com

bar.com	5
bad.com	4
foo.com	3
widget.com	4

- Can avoid shuffle on join
- But, fundamentally a shuffle on reduceByKey
- Optimization: custom partitioner on domain

# PageRank Performance

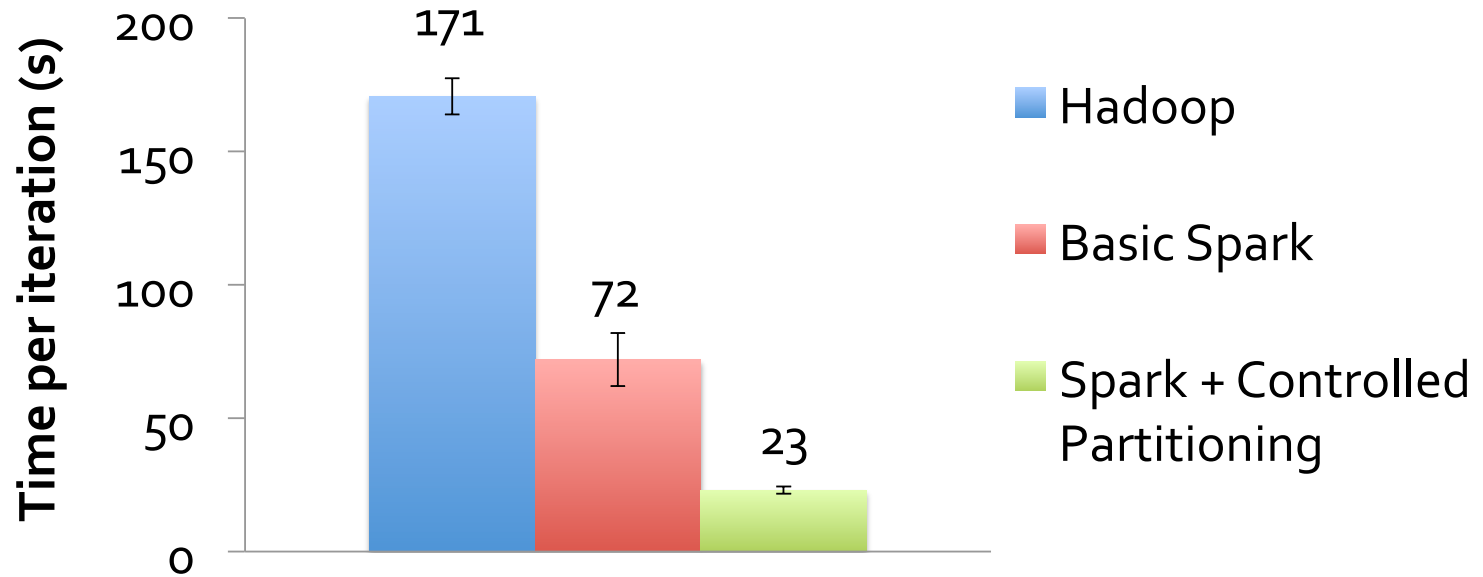
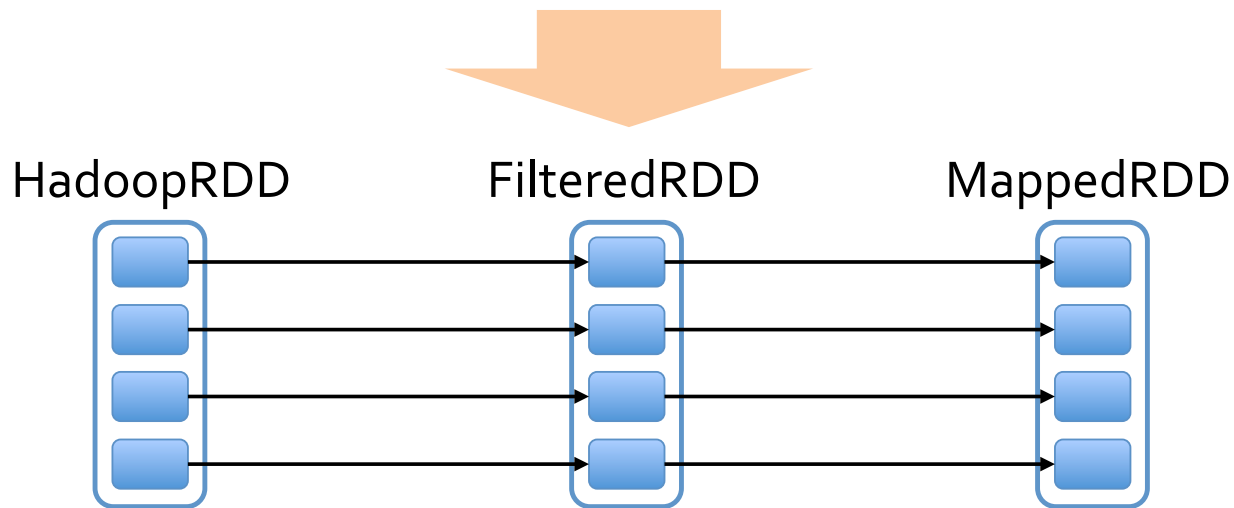


Figure 10a: 30 machines on 54 GB of Wikipedia data computing PageRank

# Fault Recovery

RDDs track the graph of transformations that built them (their *lineage*) to rebuild lost data

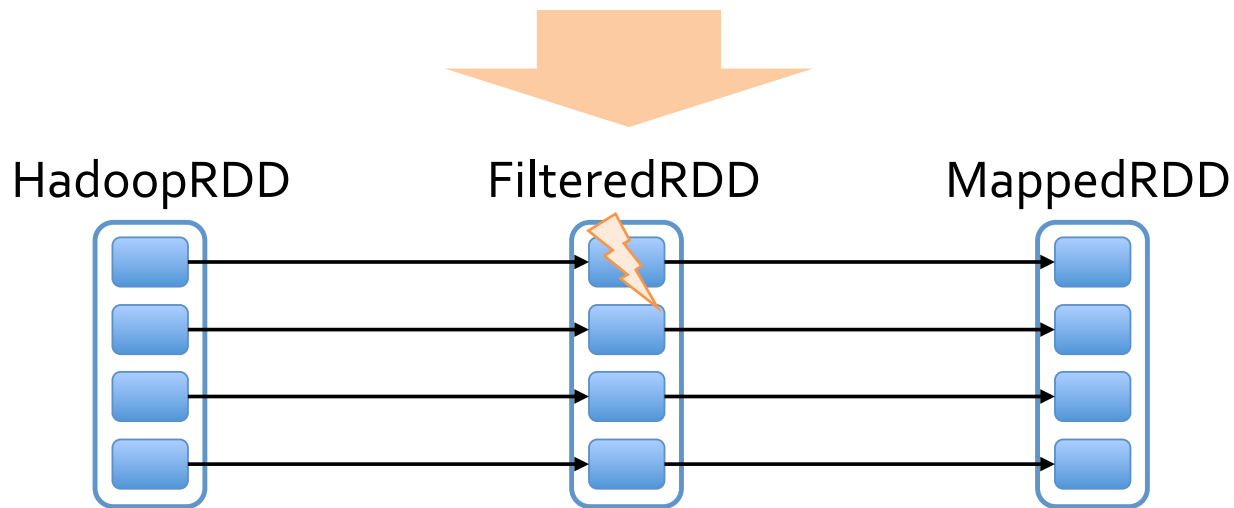
E.g.: `messages = textFile(...).filter(_.contains("error")).map(_.split('\t')(2))`



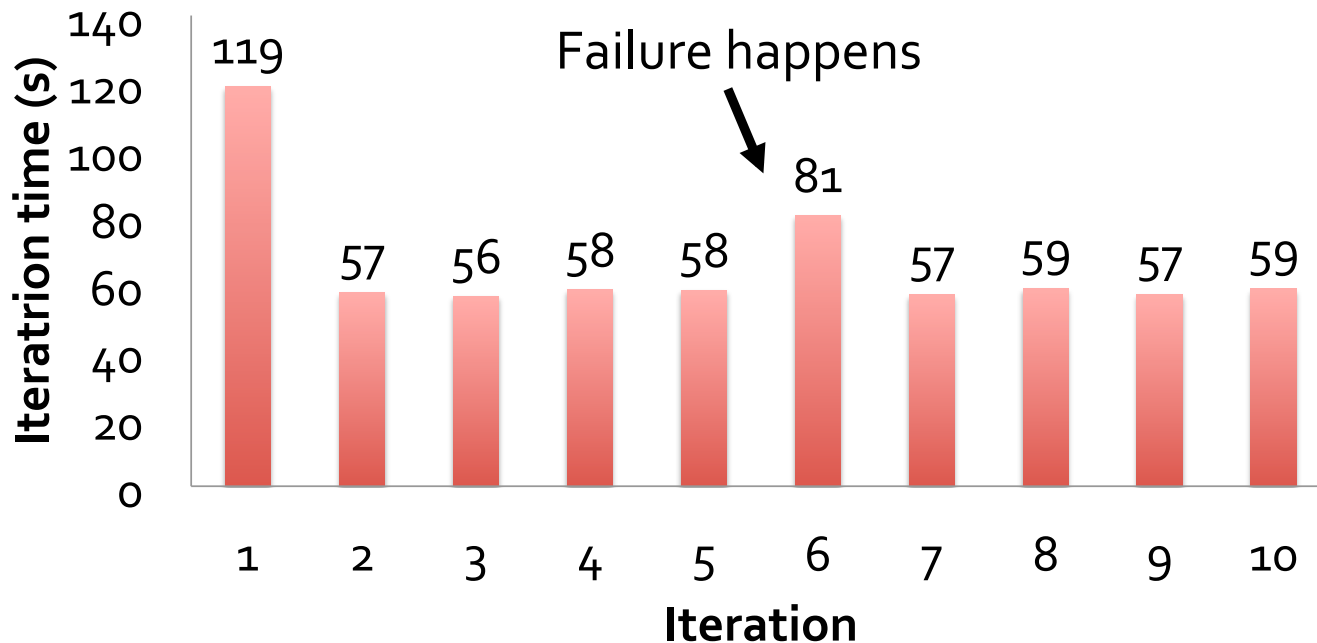
# Fault Recovery

RDDs track the graph of transformations that built them (their *lineage*) to rebuild lost data

E.g.: `messages = textFile(...).filter(_.contains("error")).map(_.split('\t')(2))`



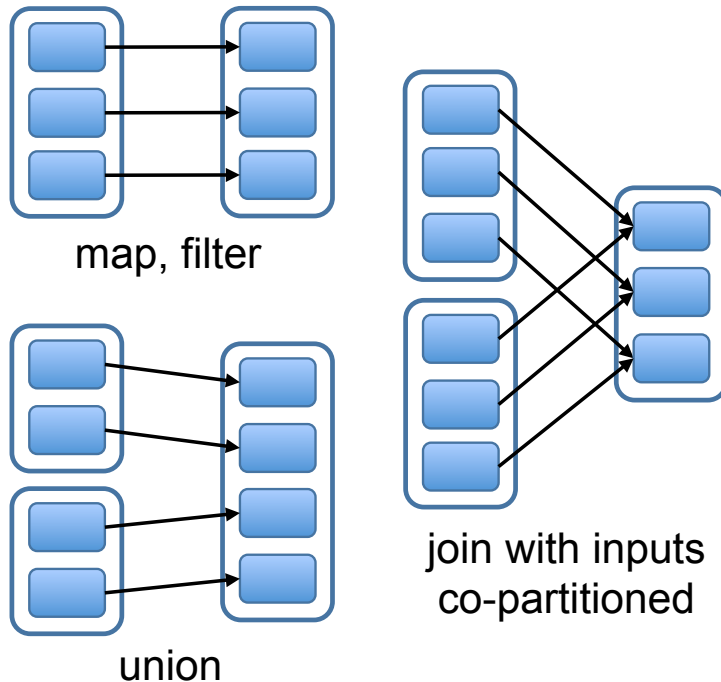
# Fault Recovery



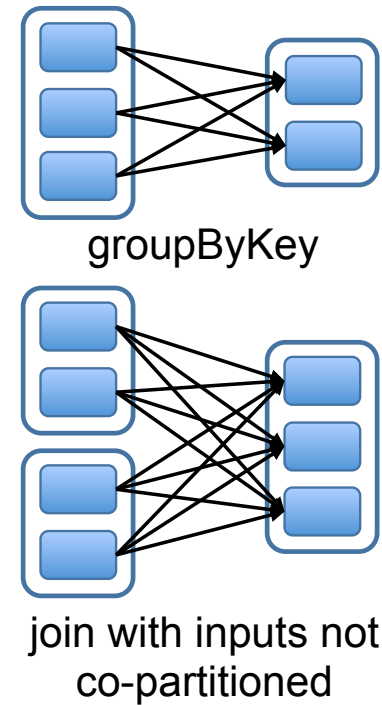
100 GB K-means Fig 11, partial reconstruction in step 6 much less than cost to write back results at each step with MR

# Narrow & Wide Dependencies

Narrow Dependencies:



Wide Dependencies:



**Wide:** multiple child partitions depend on partition

Must stall for all parent data, loss of child requires whole parent RDD (not just a small # of partitions)

**Narrow:** can pipeline on one machine

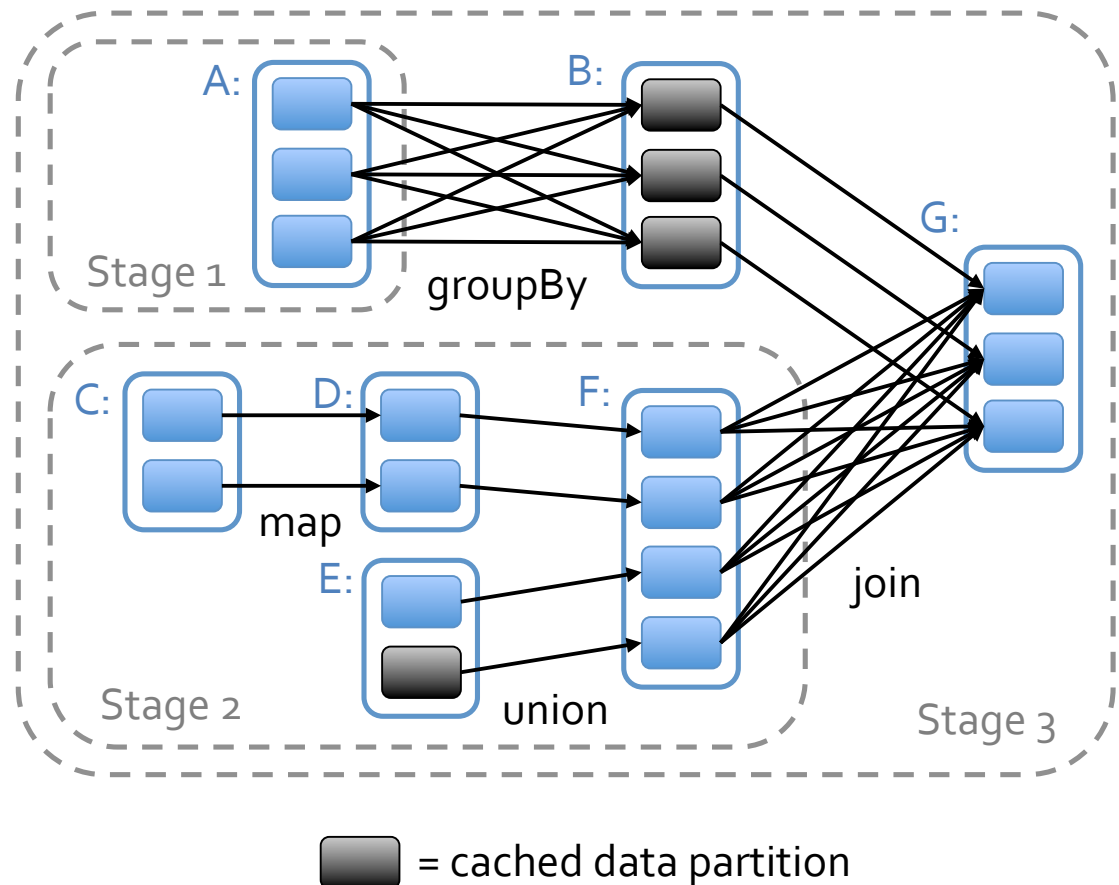
# Task Scheduler

Dryad-like DAGs

Pipelines functions within a stage

Locality & data reuse aware

Partitioning-aware to avoid shuffles

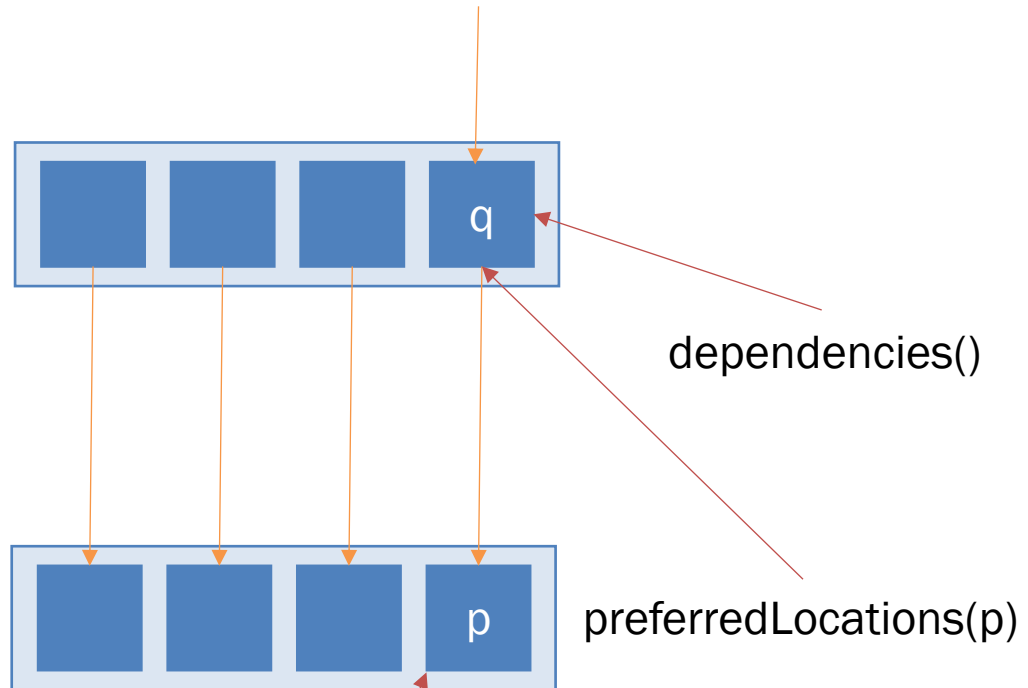




# RDD Implementation

- `partitions()`: set of partitions (ranges/hash range)
- `dependencies()`: set of parent RDDs
- closure for computing the transformation
- `preferredLocations(p)`: returns a set of locations where partition `p` can access parent data locally
- `partitioner()`: metadata about RDD partitioning scheme

# Volcano Model



parent.partitioner() ==  
child.partitioner()  
and no shuffle,  
so narrow

Iterator accesses are  
nested/recursive to pipeline  
work and avoid materialization  
while avoiding communication

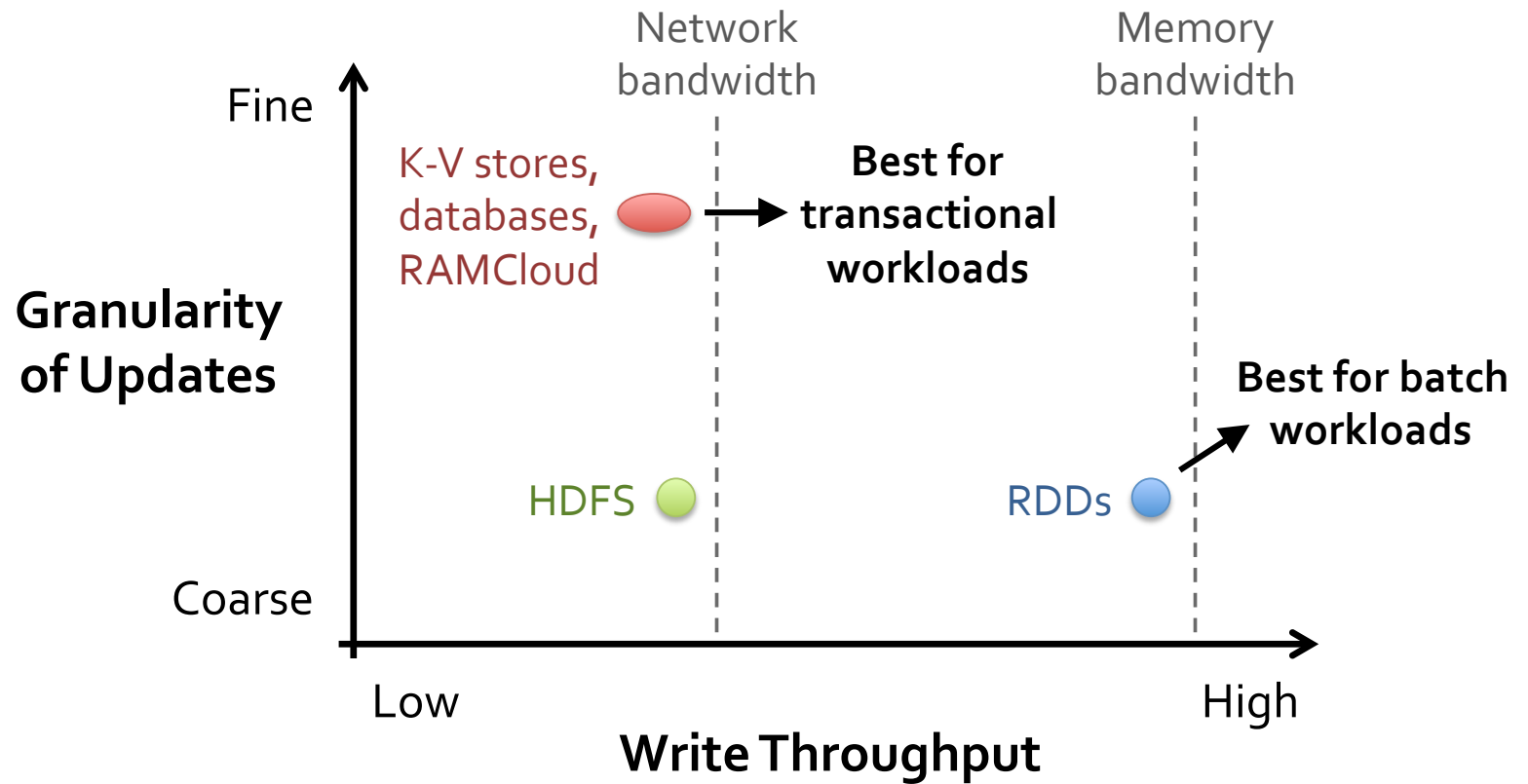
`i = iterator(p, qlter)`

`child.count() { return sum([row for row in i]) }`

# Generality of RDDs

- RDDs can express many parallel algorithms
  - They already apply the same operation to many items
- Unifies many programming models
  - Data flow models: MapReduce, Dryad, SQL, ...
  - Specialized models for iterative apps: BSP/Pregel, iterative MapReduce (Haloop), bulk incremental, ..
- Supports new applications that these models don't

# Tradeoff Space



# Programming Models as Libraries

RDDs can express many existing parallel models

- » **MapReduce, DryadLINQ**
- » **Pregel** graph processing [200 LOC]
- » **Iterative MapReduce** [200 LOC]
- » **SQL**: Hive on Spark (Shark) [in progress]

All are based on  
coarse-grained  
operations

Enables apps to efficiently *intermix* these models

# Memory Reuse Impact

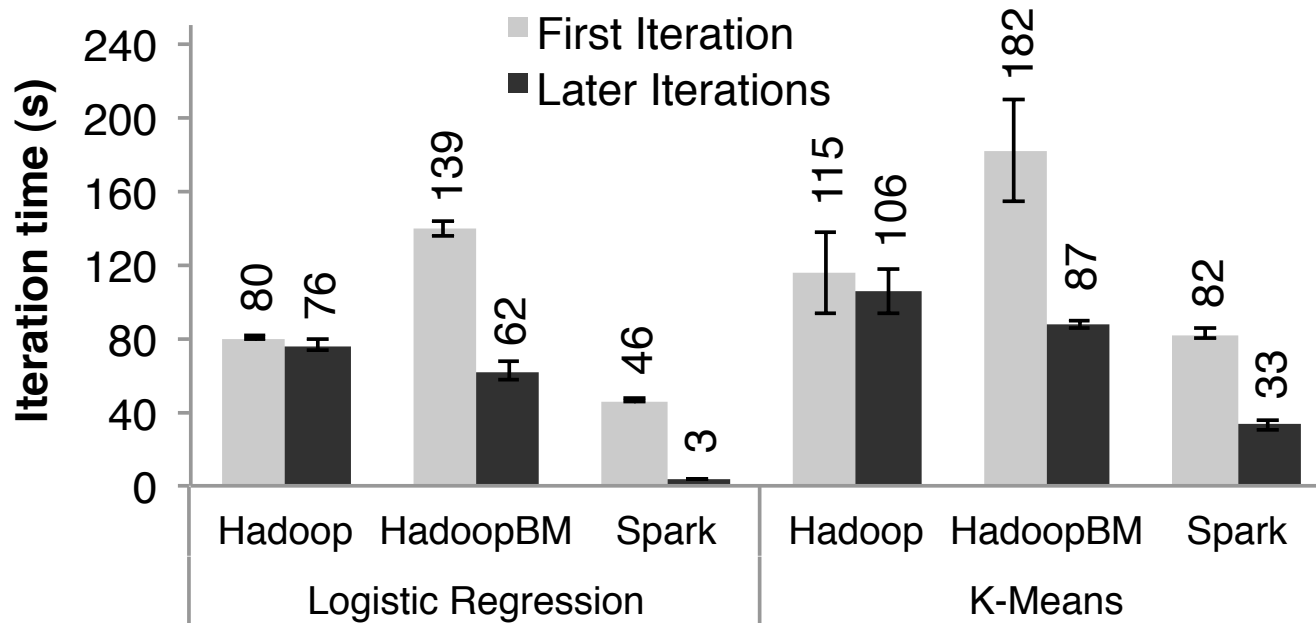


Figure 7: Duration of the first and later iterations in Hadoop, HadoopBinMem and Spark for logistic regression and k-means using 100 GB of data on a 100-node cluster.

# Analysis of Speedup

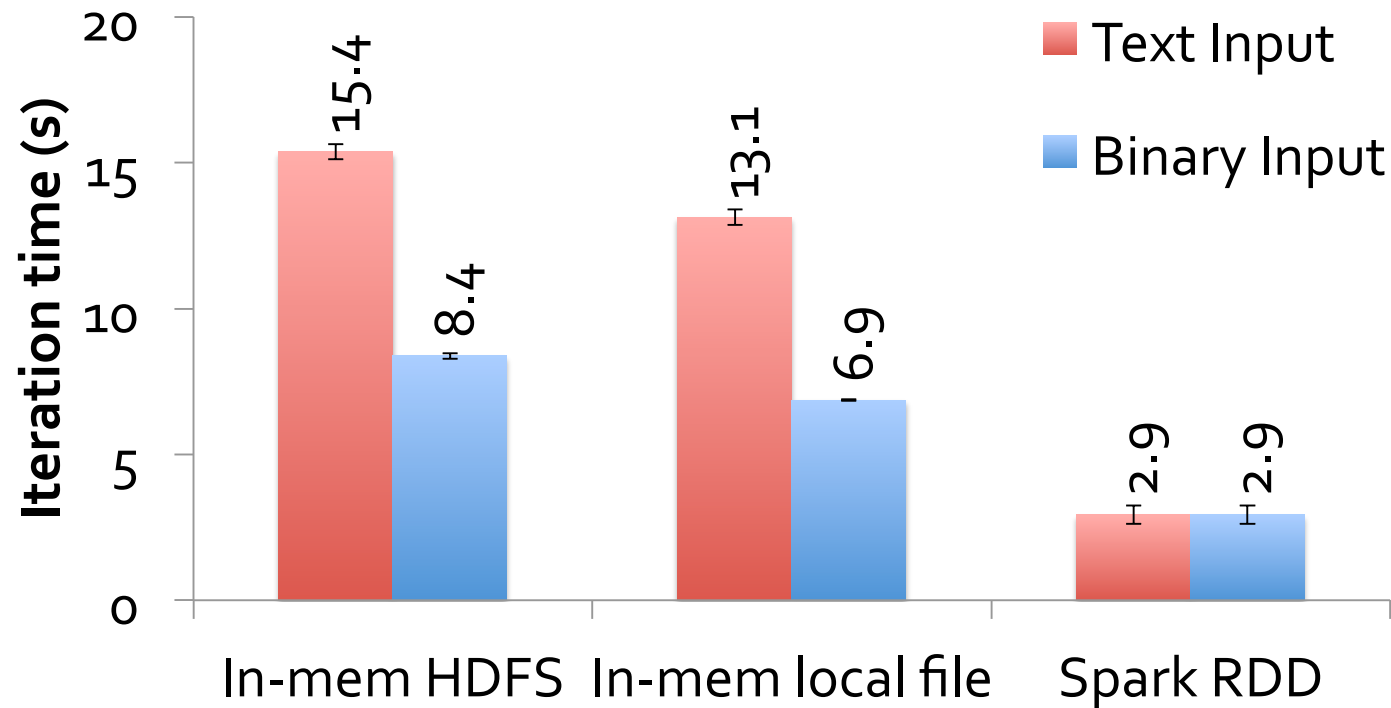


Fig 9: 256 GB on a single machine

# Eviction and Working Sets

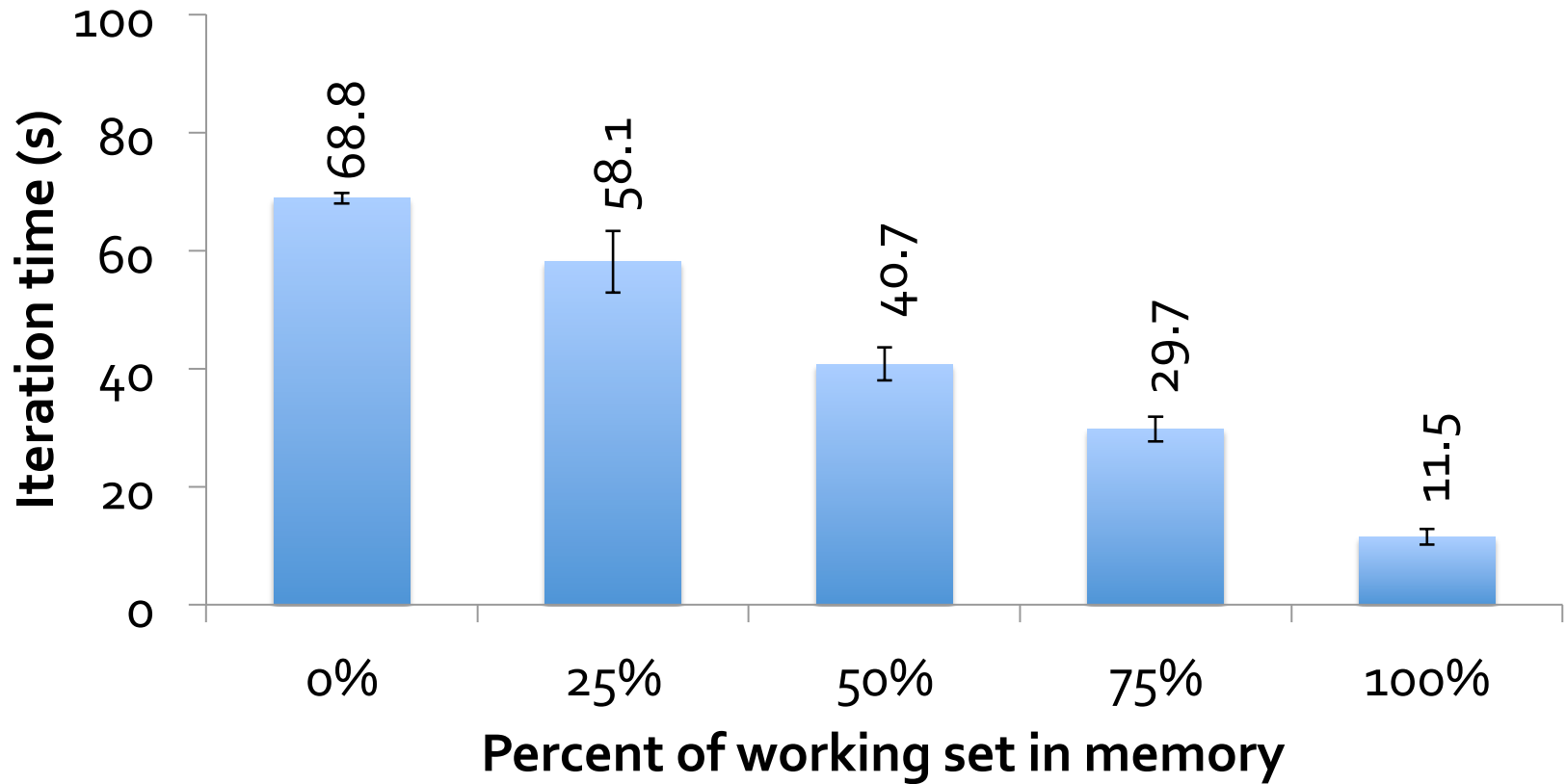
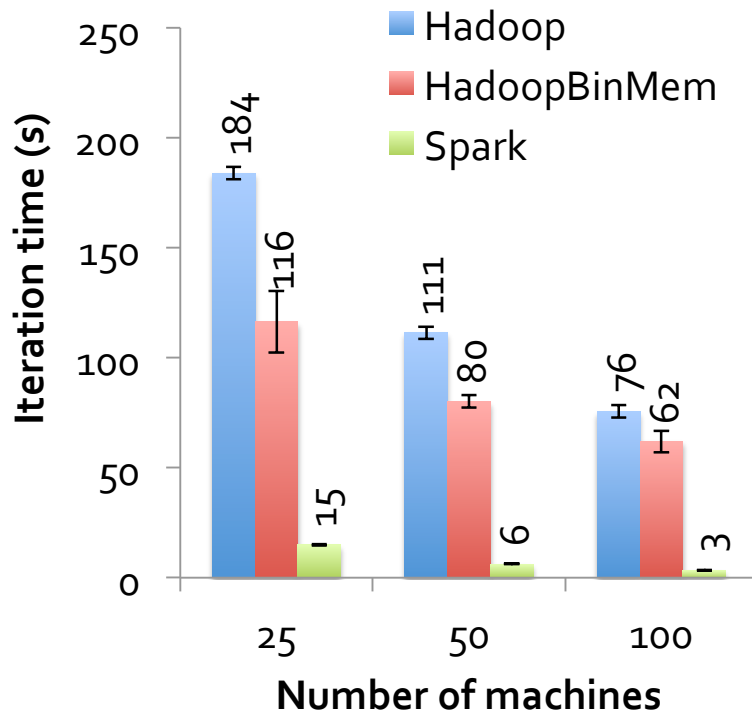


Fig 12: 100 GB LR on 25 machines

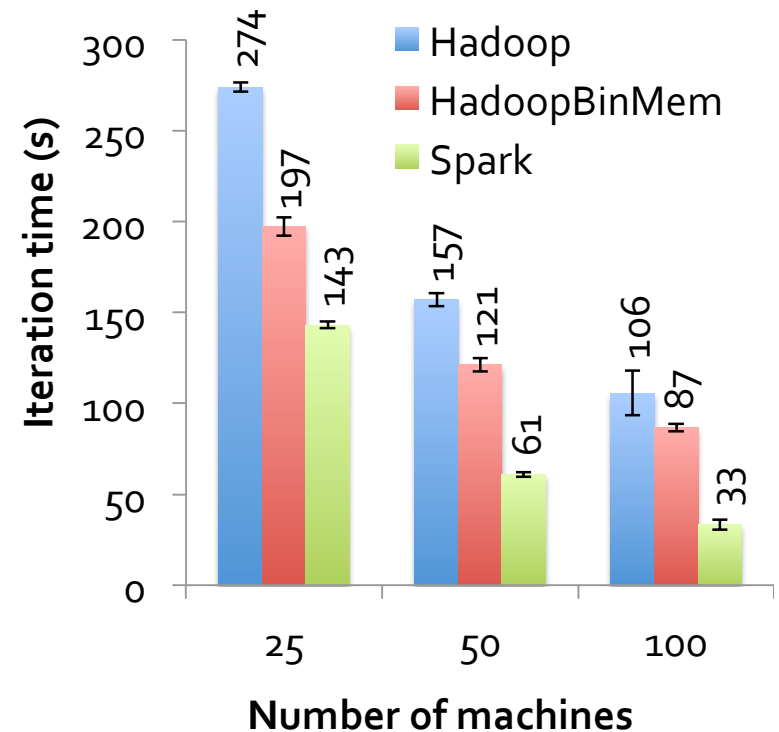


# Scalability

## Logistic Regression



## K-Means



100 GB datasets each

# Evolution of Spark

- Want to write Spark programs in different languages
  - Not everyone loves Scala, or JVM-based languages
- Problem: RDD semantics are bound to JVM
- Move toward Dataframes
  - Effectively, DB relations
  - Can manipulate representation w/o ser/des
  - Can use bindings in any language
- Capturing closures in many languages?  
Mismatching language/dataset semantics...

# Conclusions

- M-R expressivity and performance have been a central point of sadness
- Several attempts to make improvements
- Spark improves expressivity, which also improves performance since scheduler can “think” across the whole pipeline
- Still preserves a lot of M-R fault-tolerance
- Does force users to reason a bit about fault-tolerance, though through careful `persist()` calls
- Not dead yet: TensorFlow and many more ...

