

# Resilient Distributed Datasets: Spark

*CS 675: Distributed Systems (Spring 2020)*

Lecture 6

Yue Cheng

Some material taken/derived from:

- Matei Zaharia's NSDI'12 talk slides.
- Utah CS6450 by Ryan Stutsman.

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

# Announcement

- Deadline of the project report gets extended to 11:59pm next Friday, 04/03
- Doodle poll for Lab 2 demo and proposal discussion meetings (Thursday and Friday)
  - <https://doodle.com/poll/tbskp9hqaqsn7ysi>

# What's good with MapReduce

- Scaled analytics to thousands of machines
- Eliminated fault-tolerance as a concern

# Problems with MapReduce

- 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...

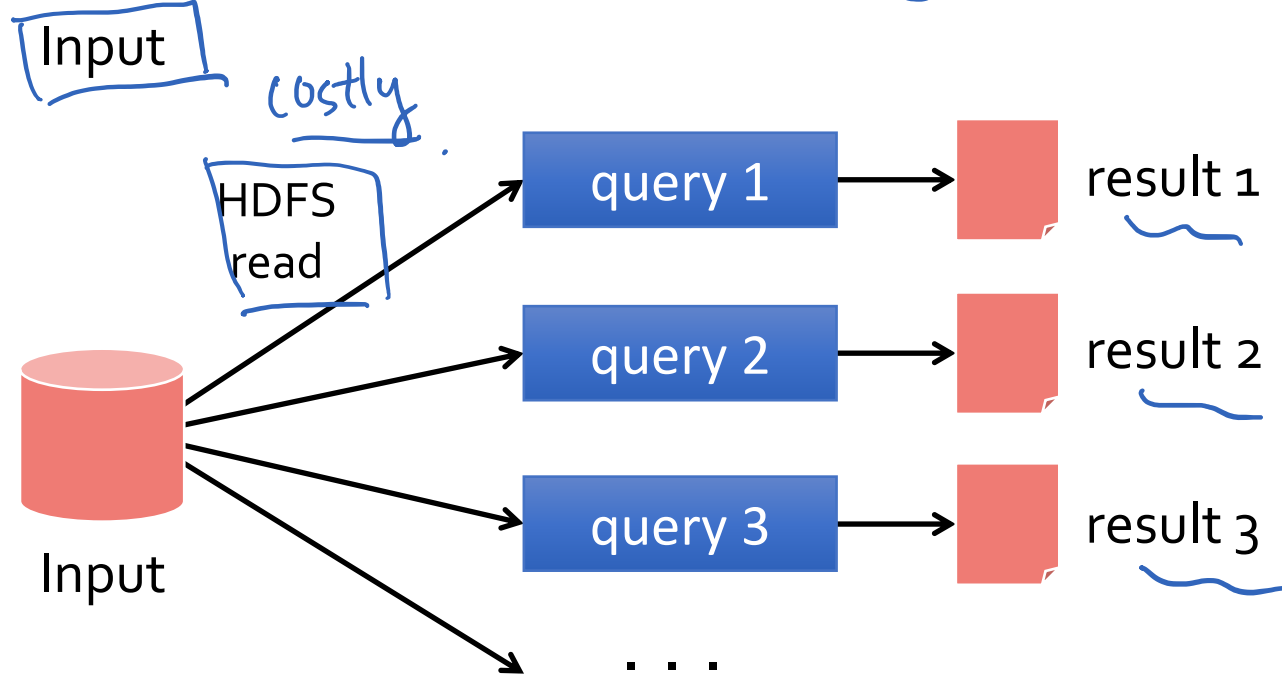
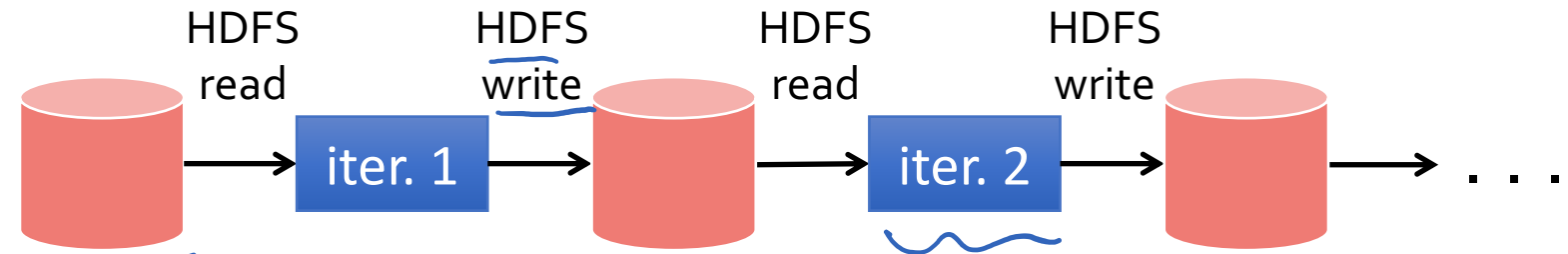


# Sharing data between iterations/ops

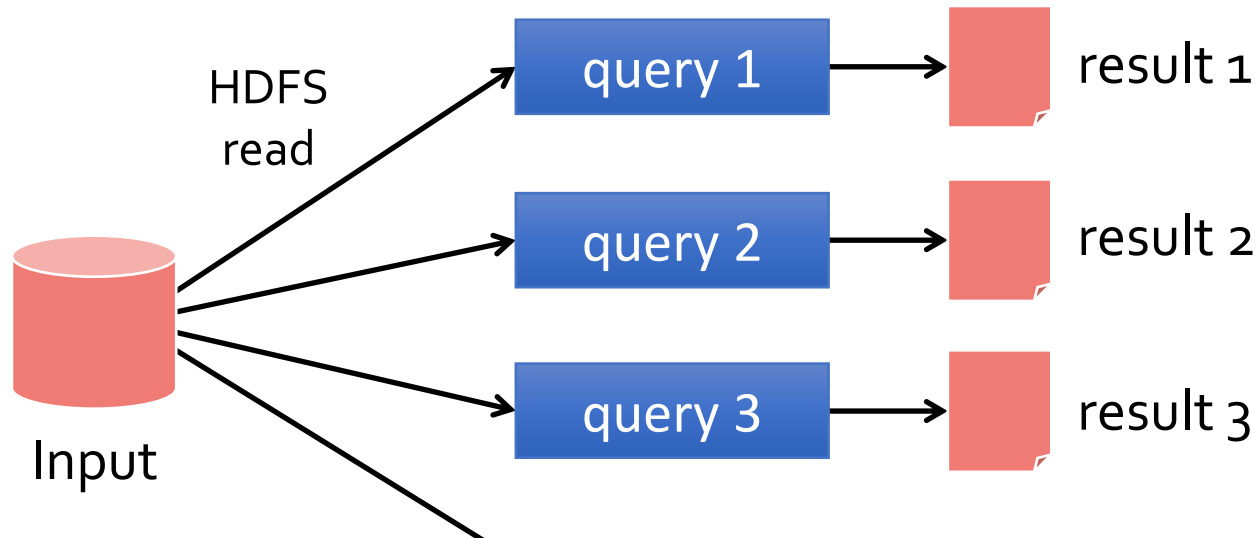
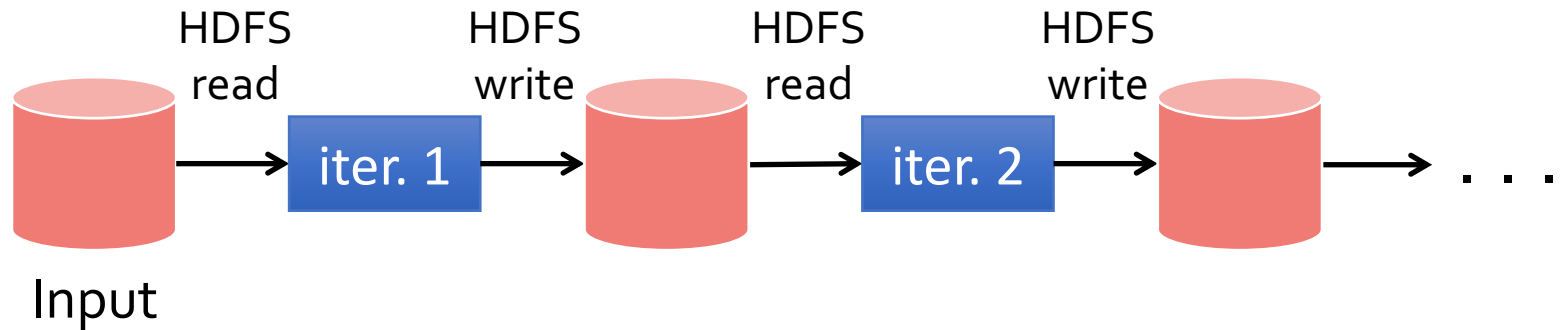
iter1 → output → iter 2.

- Only way to share data between iterations / phases is through shared storage
  - **Slow!**
- 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 and Map or Reduce task since idempotent

# Examples

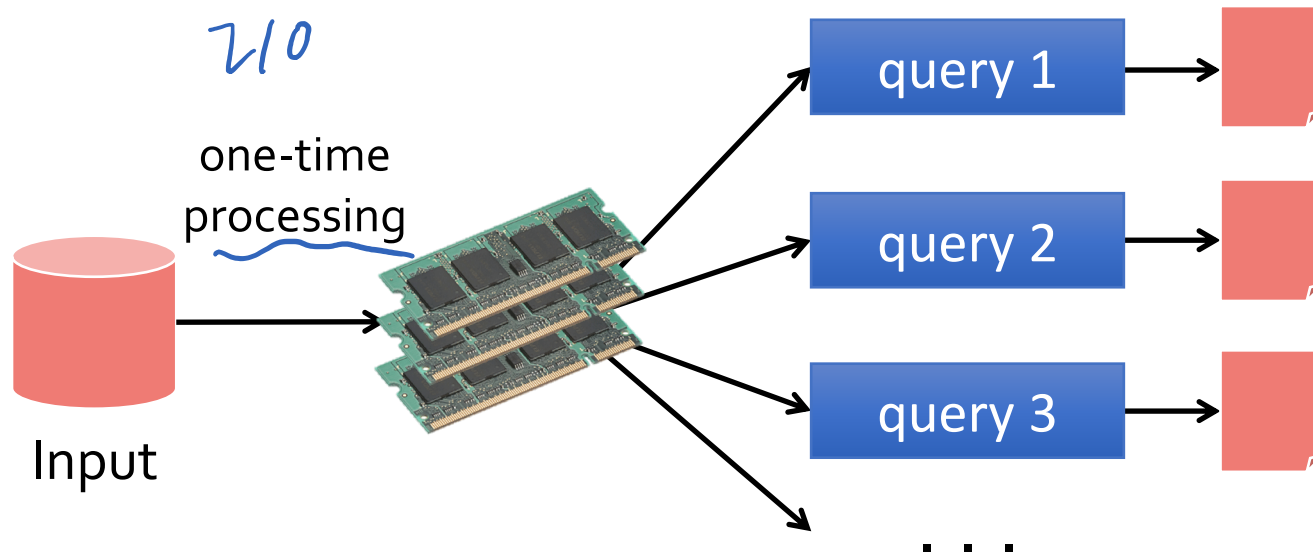
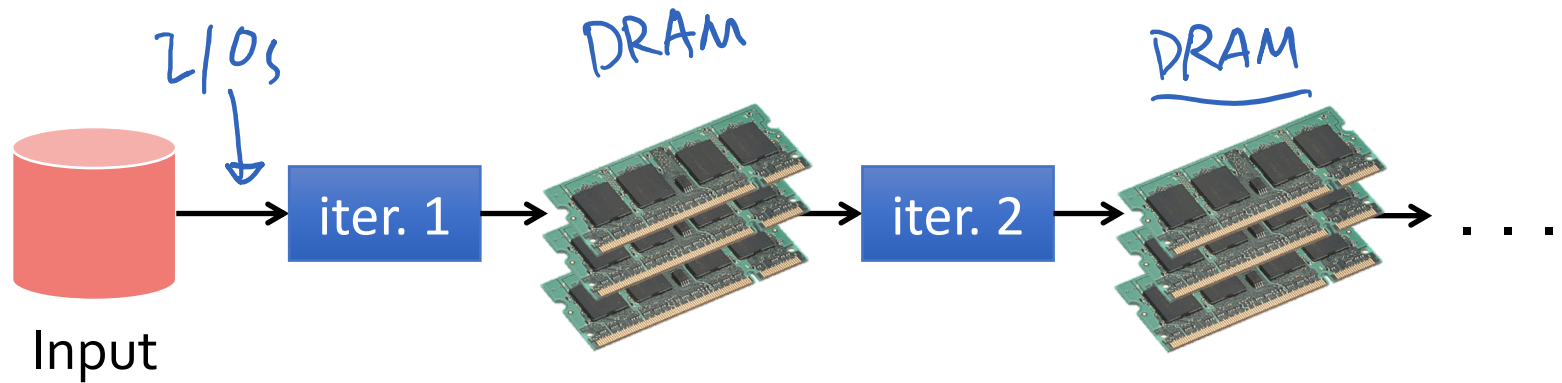


# Examples

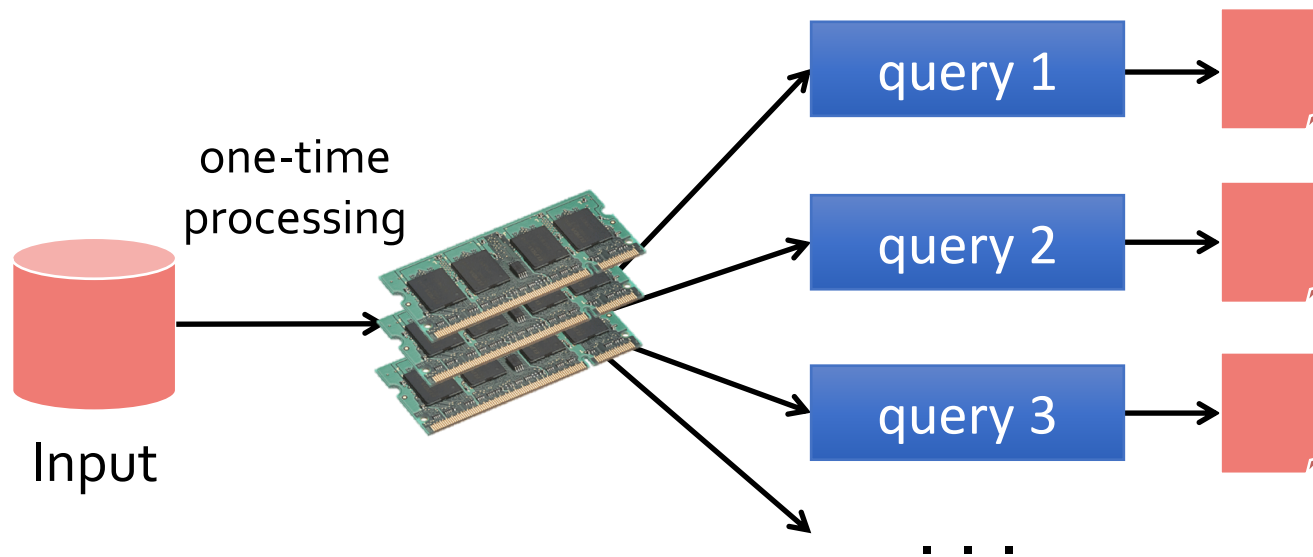
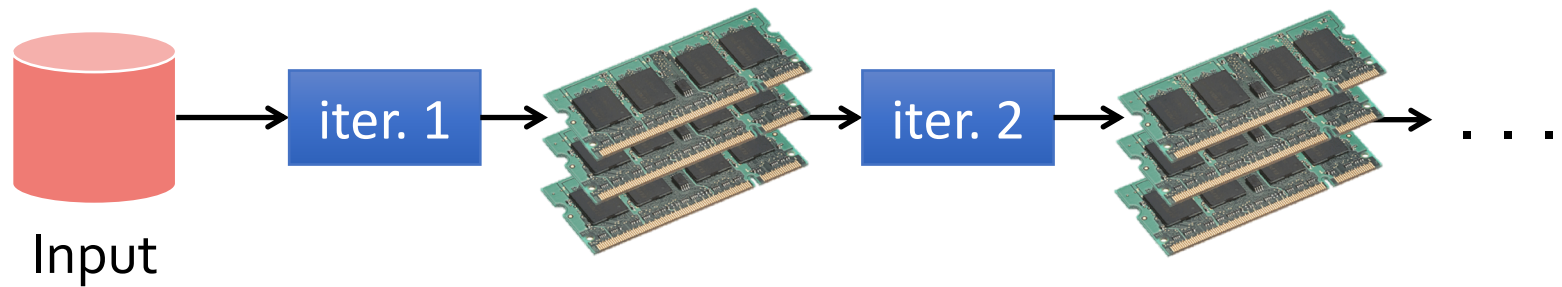


Slow due to replication and disk I/O,  
but necessary for fault tolerance

# Goal: In-memory data sharing



# Goal: In-memory data sharing



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

# Challenges

shared

- How to design a distributed <sup>shared</sup> memory abstraction that is both **fault-tolerant** and **efficient**?

# Challenges

- How to design a distributed memory abstraction that is both **fault-tolerant** and **efficient**?

- Existing <sup>in-memory</sup> storage systems allow **fine-grained** mutation to state
  - In-memory key-value stores → IMKV
  - Requires replicating data or logs across nodes for fault tolerance
    - Costly for data-intensive apps
    - 10-100x slower than memory write write-intensive
  - They also require costly on-the-fly replication for mutations

# Challenges

MR → batch processing  
IMKV → fine-grained

- How to design a distributed memory abstraction that is both **fault-tolerant** and **efficient**?

Spark: → batch-processing + Iterative + Interactive.



- Existing storage systems allow **fine-grained** mutation to state

**Insight:** leverage similar coarse-grained approach that transforms whole dataset per operation, like MapReduce (batch processing)

- 10-100x slower than memory write
- They also require costly on-the-fly replication for mutations



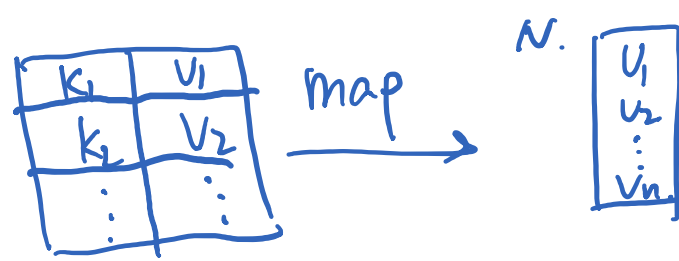
# Solution: Resilient Distributed Datasets (RDDs)

- 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<sub>1</sub> → RDD<sub>2</sub>)
- Actions on RDDs (RDD → output)
- Control over RDD partitioning (how items are split over nodes) *partitioner func.*
- Control over RDD persistence (in memory, on disk, or recompute on loss) *flag*

# Transformations



<p>Transformations (define a new RDD)</p>	<p><u>map</u> <u>filter</u> sample groupByKey reduceByKey sortByKey</p>	<p>flatMap union join cogroup cross mapValues</p>
---	---	---

RDDs in terms of Scala types → Scala semantics at workers

Transformations are lazy “thunks”; cause no cluster action

# Actions

Directed Acyclic Graph. - DAG.

<p><b>Actions</b> (return a result to driver program)</p>	<p>collect reduce <u>count</u> save <u>lookupKey</u></p>
---	--

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 pieces of the computation)

# Interactive debugging

RDD<sub>1</sub>

```
lines = textFile("hdfs://foo.log")
```

RDD<sub>2</sub> - 

```
errors = lines.filter(  
    _.startsWith("ERROR")
```

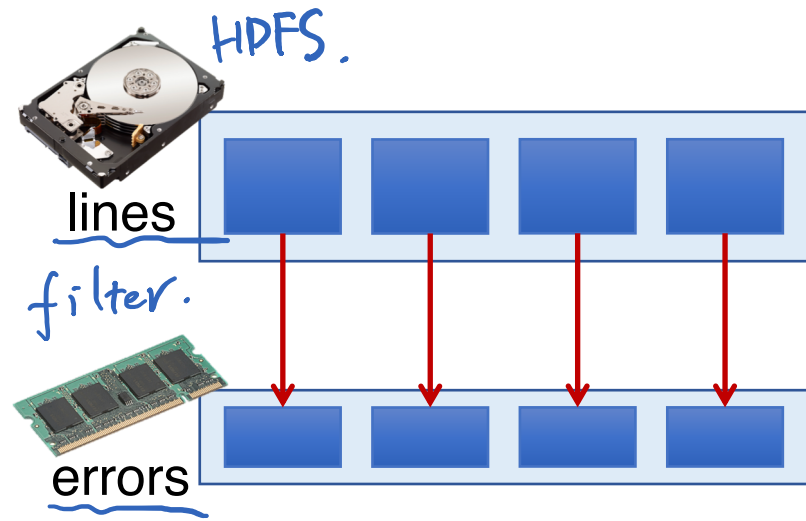
```
errors.persist()
```



# Interactive debugging

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

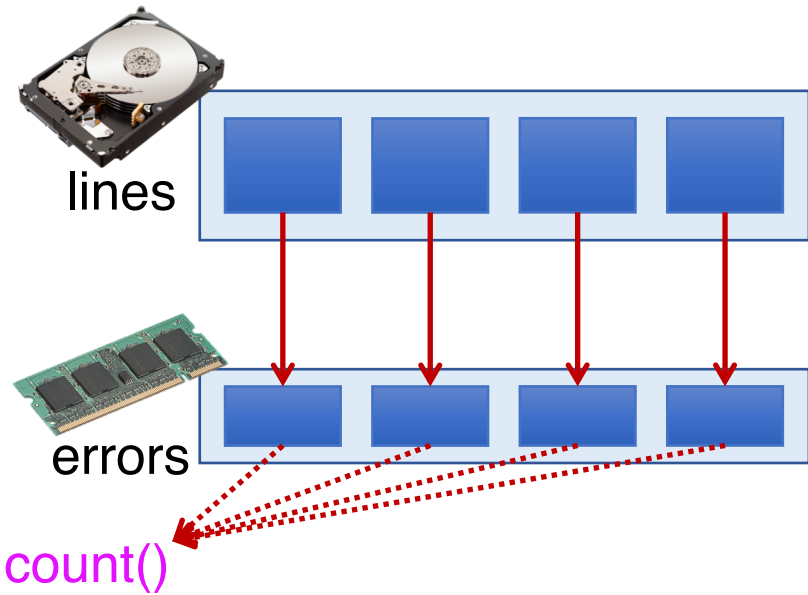
`errors.count()` ← Action.



# Interactive debugging

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

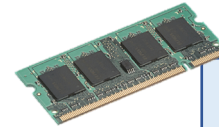
```
errors.count()
```



# Interactive debugging



lines



errors

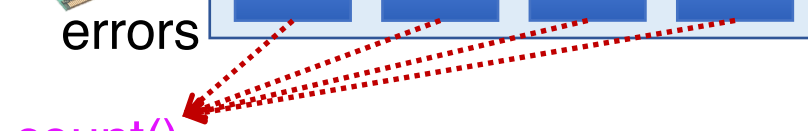


count()

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

```
errors.count()
```

```
errors.filter(  
    _.contains("MySQL"))
```



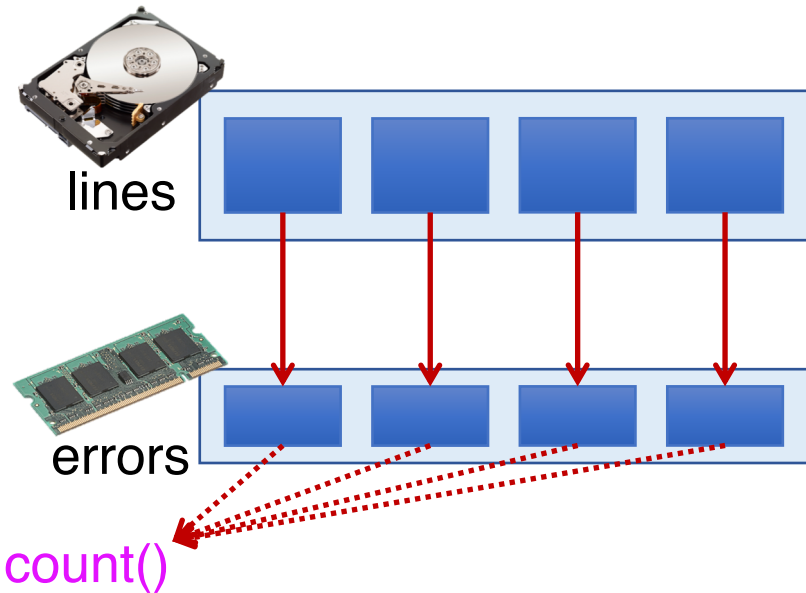


# Interactive debugging

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

```
errors.count()
```

```
errors.filter(
    _.contains("MySQL")).count()
↳
```



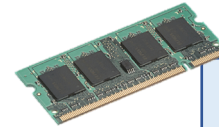
# Interactive debugging



lines



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



errors



count()

```
errors.count()
```

```
errors.filter(
    _.contains("MySQL")).count()
```



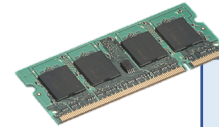
# Interactive debugging



lines



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



errors



```
errors.count()
```

count()

```
errors.filter(
    _.contains("MySQL")).count()
```

count()

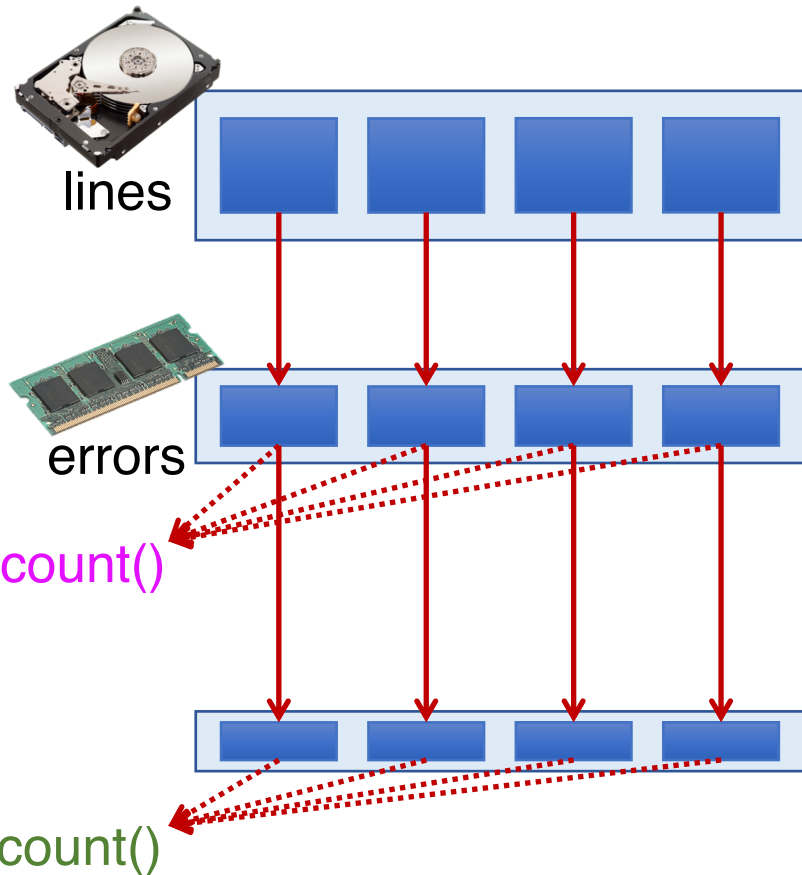


# Interactive debugging

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

```
errors.count()
```

```
errors.filter(
    _.contains("MySQL")).count()
errors.filter(
    _.contains("HDFS"))
```



# Interactive debugging

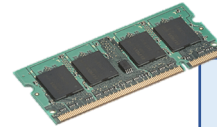
```
lines = textFile("hdfs://foo.log")
errors = lines.filter(
    _.startsWith("ERROR"))
errors.persist()
```

```
errors.count()
```

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



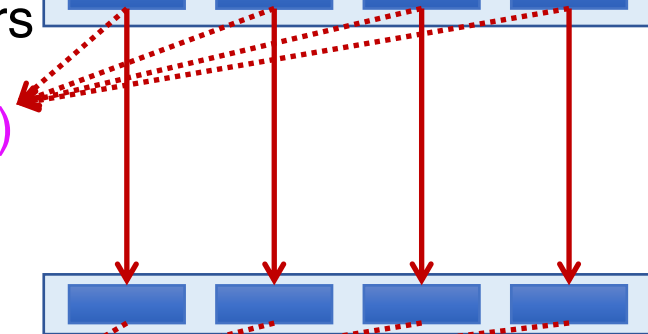
lines



errors



count()



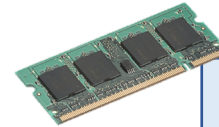
count()



# Interactive debugging



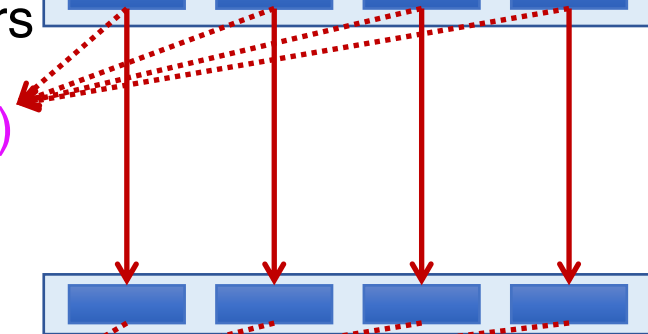
lines



errors



count()



count()

```
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()
```

Action

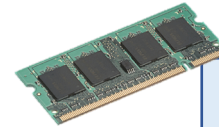
# Interactive debugging



lines



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



errors



```
errors.count()
```

count()

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

count()

filter



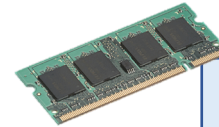
# Interactive debugging



lines



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



errors



```
errors.count()
```

count()

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

count()

map





# Interactive debugging

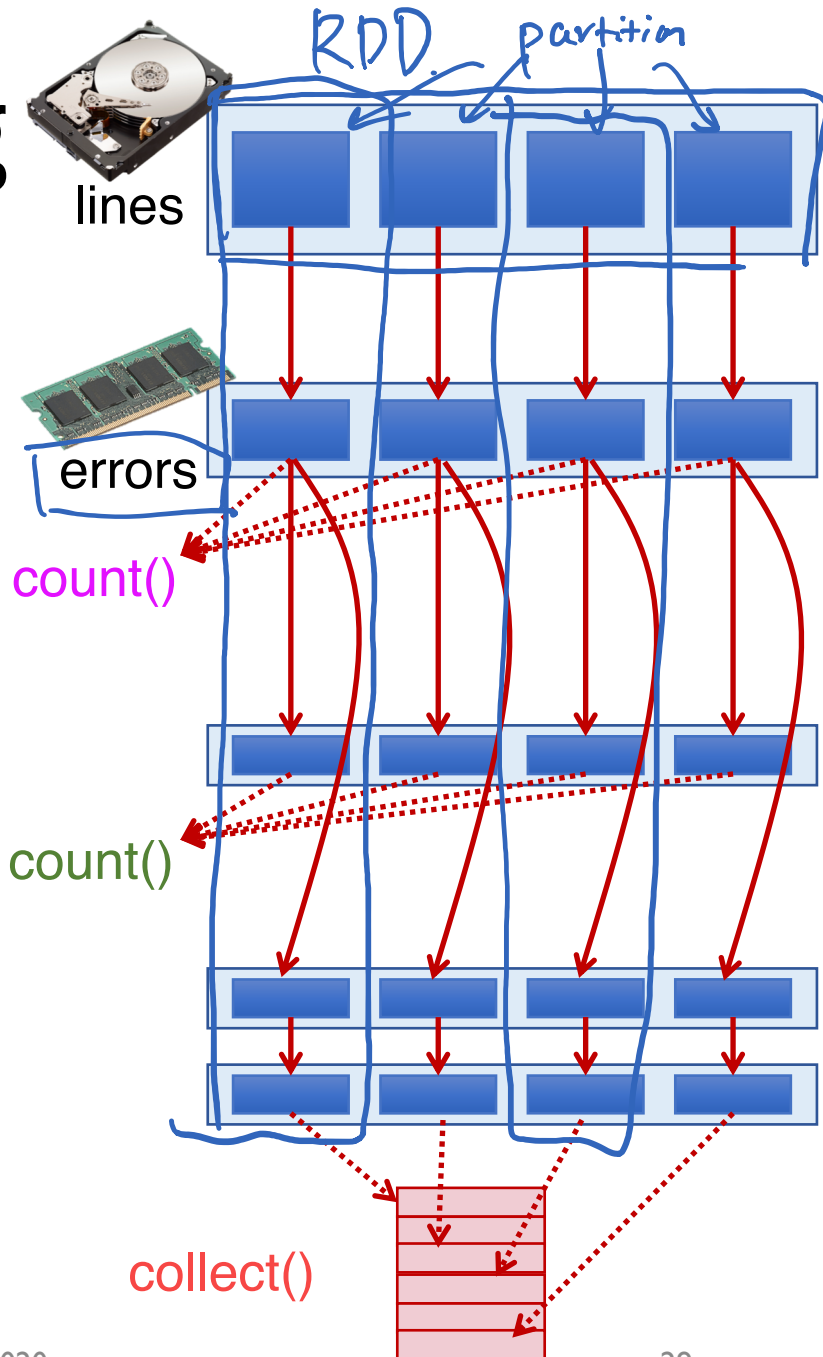
```
lines = textFile("hdfs://foo.log")
errors = lines.filter(
    _.startsWith("ERROR"))
errors.persist()
```

→ hint → sched.

```
errors.count()
```

```
errors.filter(
    _.contains("MySQL")).count()
```

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



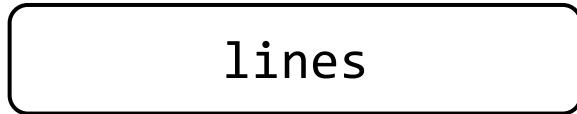
# persist()

- Not an action and not a transformation
- A scheduler hint
- Tells which RDDs the Spark scheduler should materialize and whether in memory or storage
- Gives the user control over reuse/recompute/recovery tradeoffs

# persist()

- Not an action and not a transformation
- A scheduler hint
- Tells which RDDs the Spark scheduler 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?

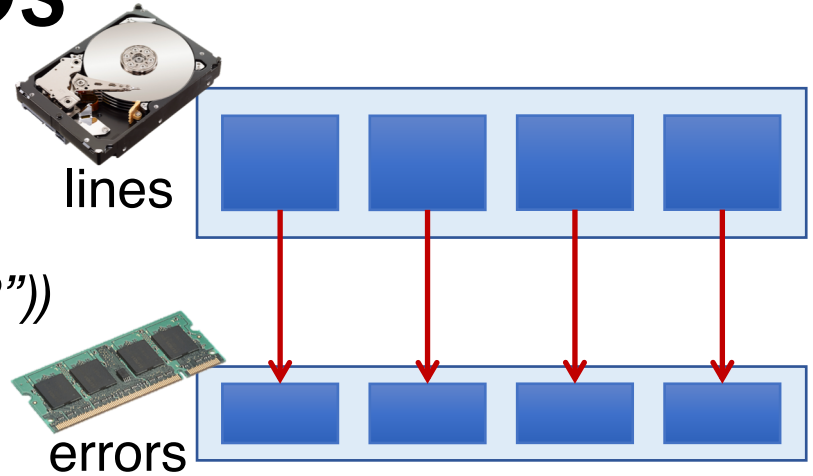
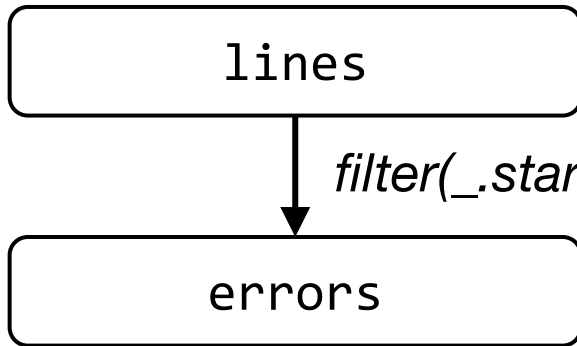
# Lineage graph of RDDs



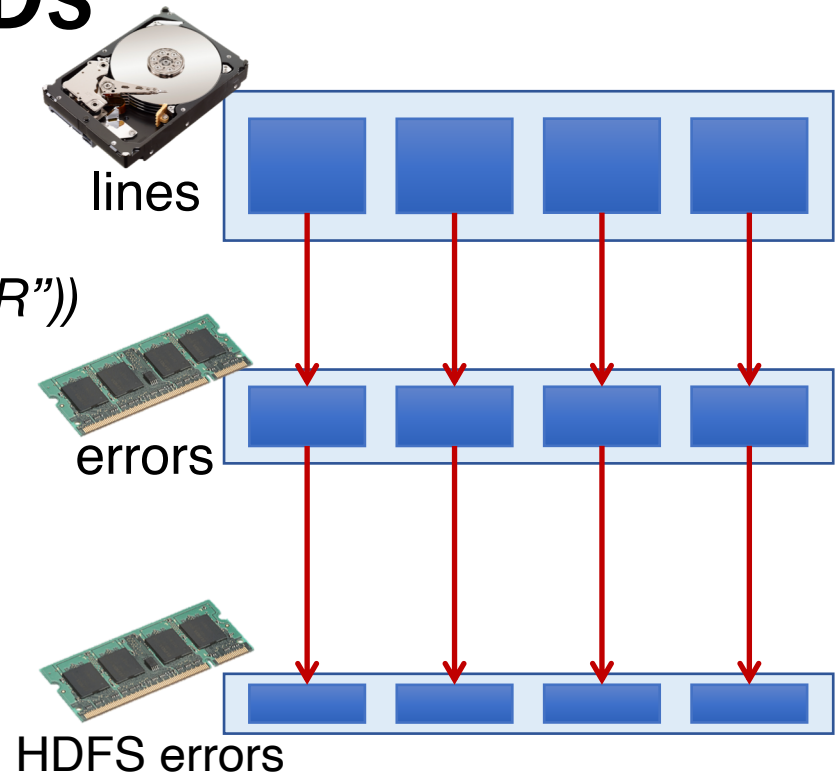
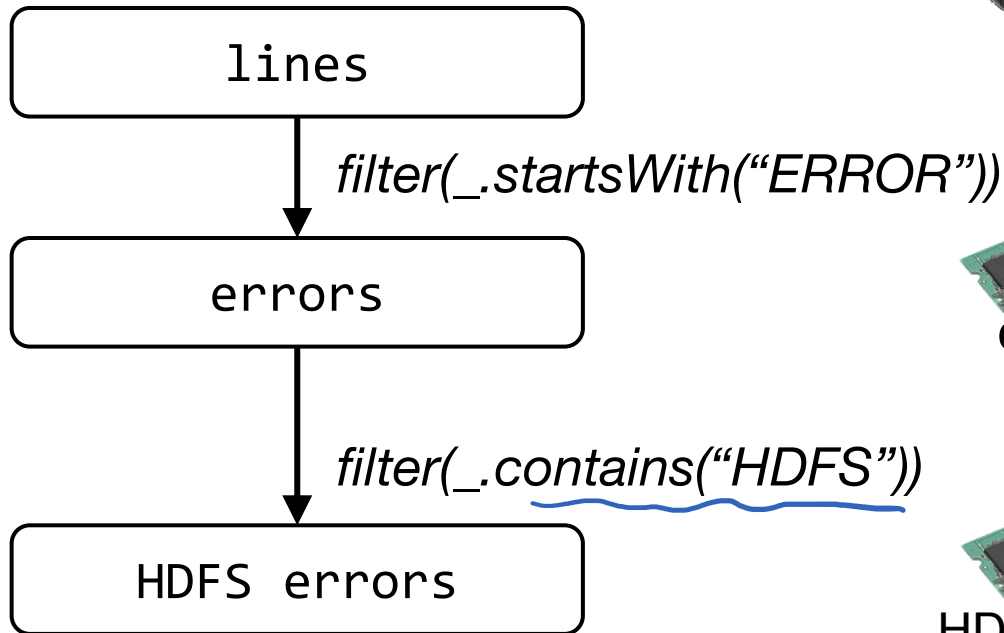
lines



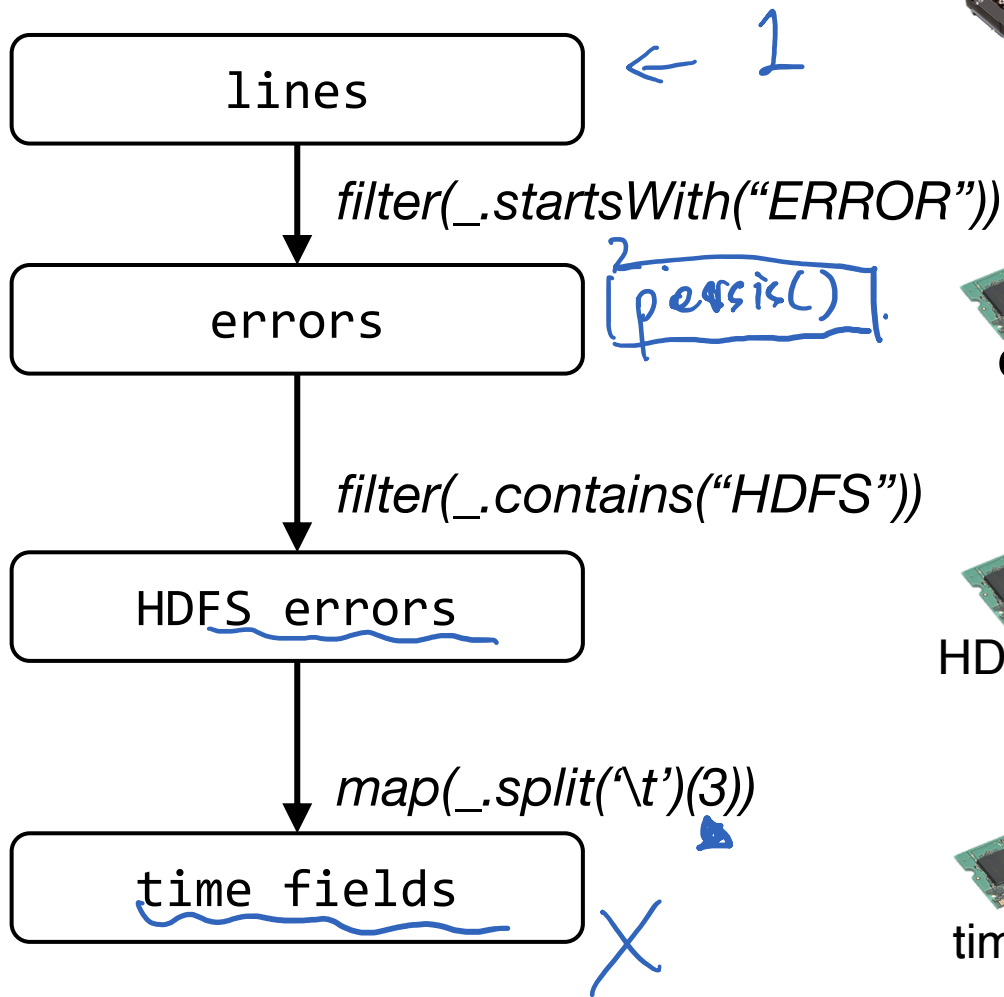
# Lineage graph of RDDs



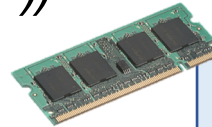
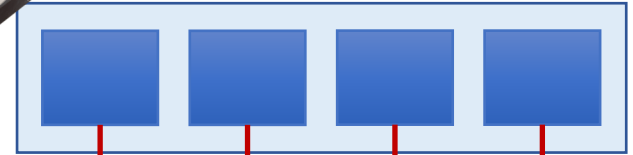
# Lineage graph of RDDs



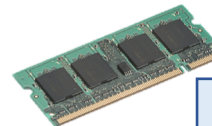
# Lineage graph of RDDs



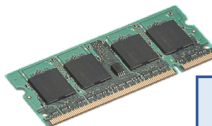
lines



errors



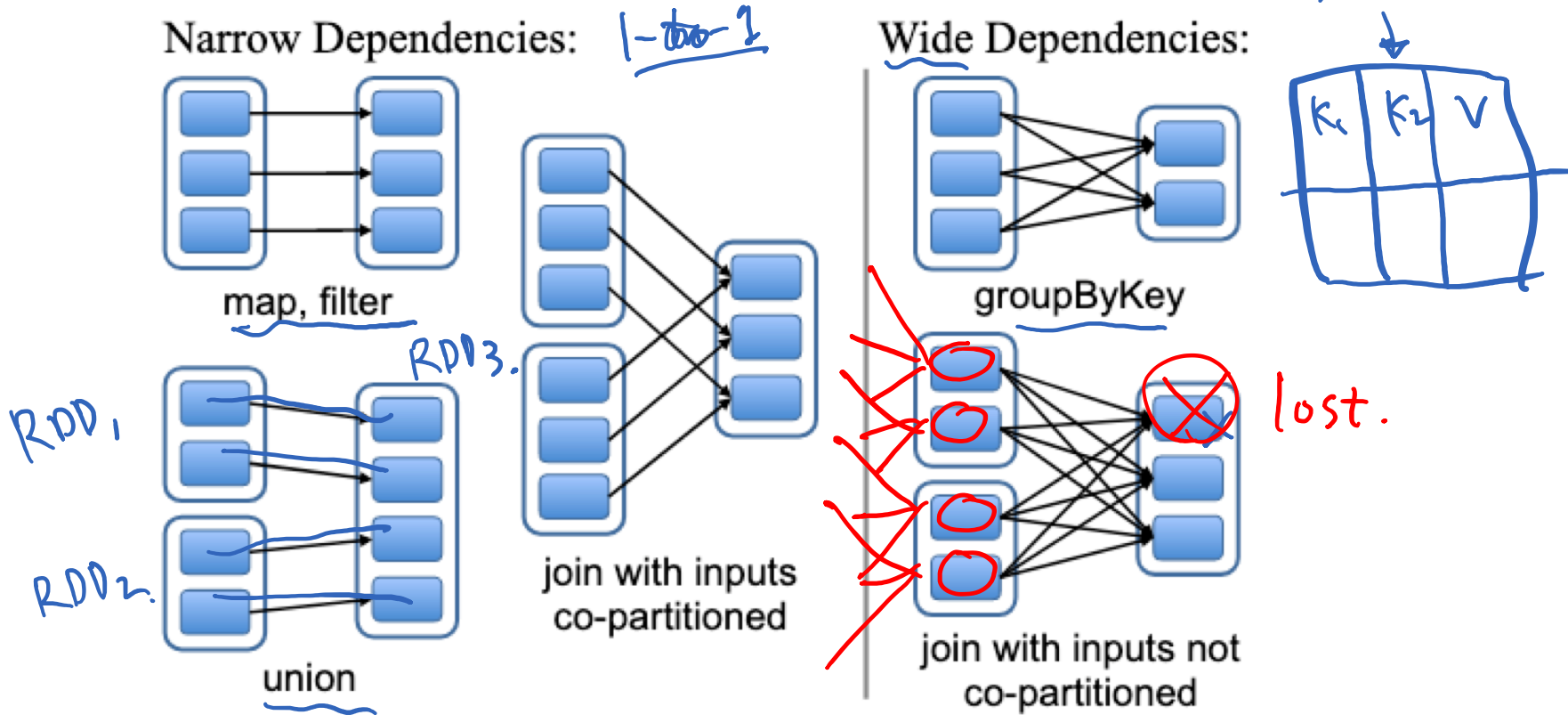
HDFS errors



time fields



# Narrow & wide dependencies



**Narrow:** each parent partition used by at most one child partition (can partition on one machine)

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

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



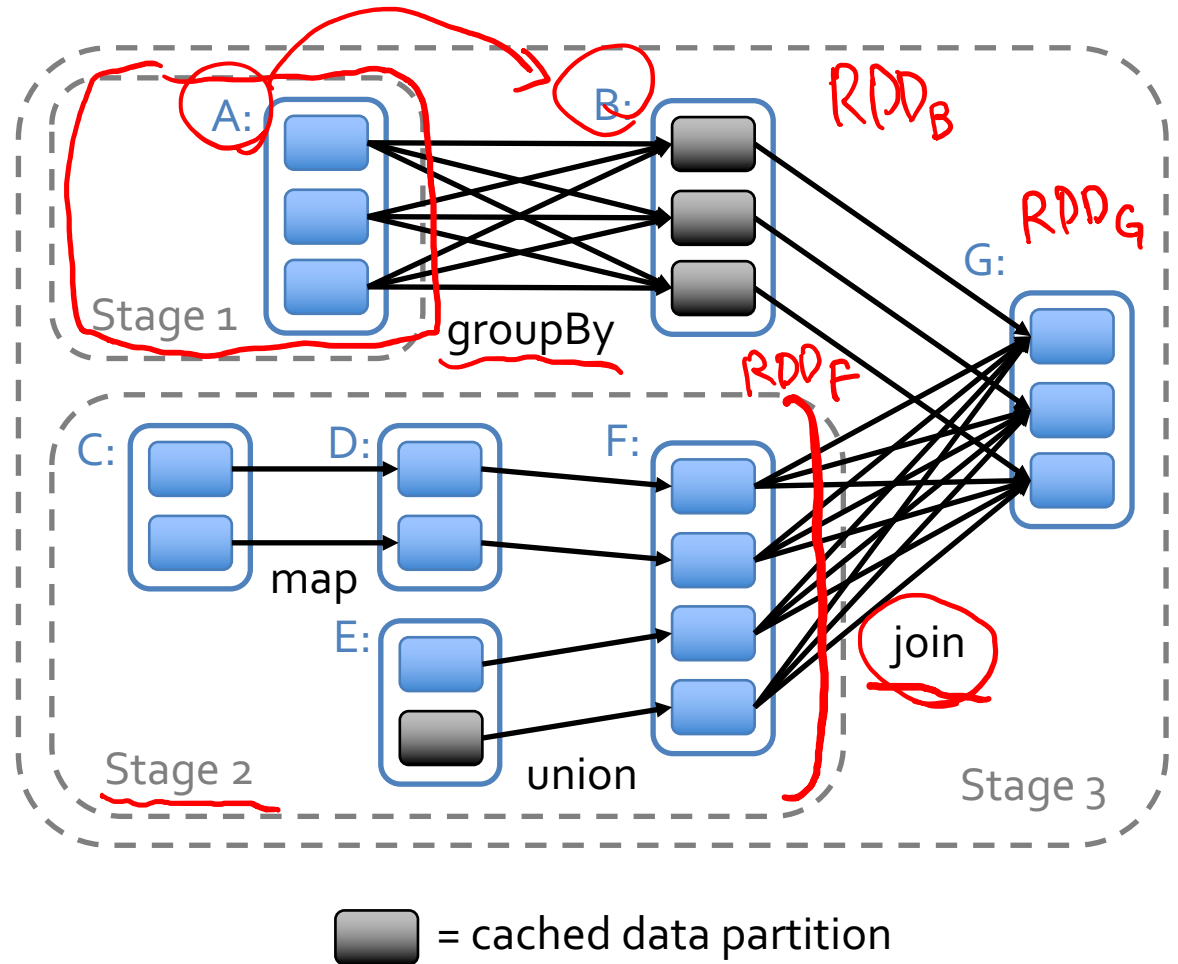
# Task scheduler

Dryad-like DAGs

Pipelines functions within a stage

Locality & data reuse aware

Partitioning-aware to avoid shuffles



# Interactive debugging (control and data flow)

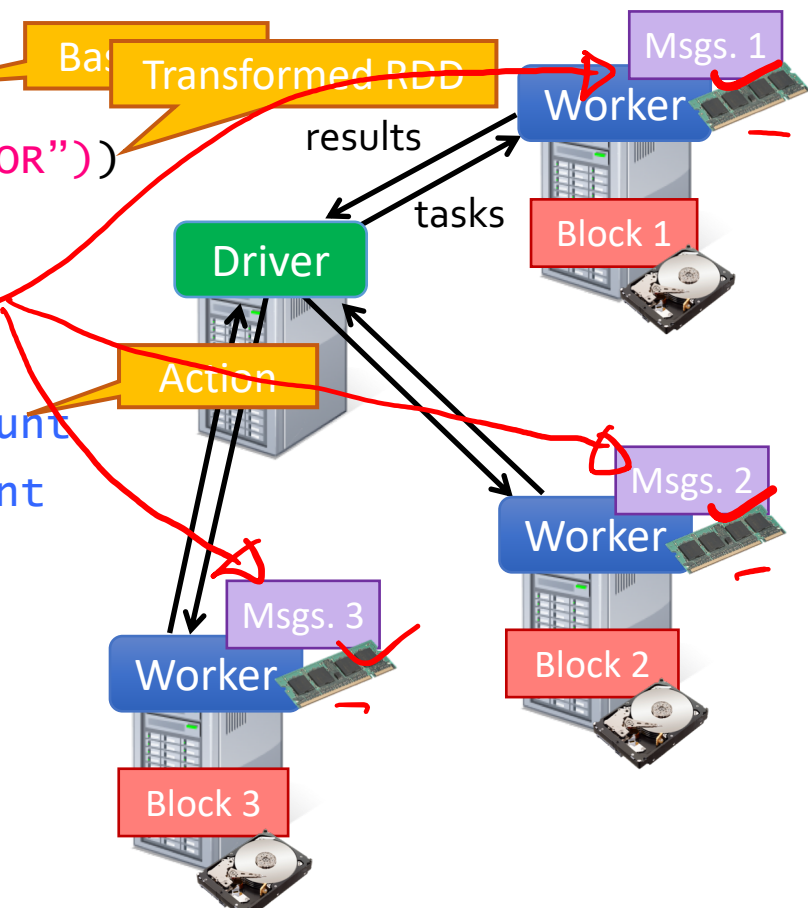
Load error messages from a log into memory, then interactively search for various patterns

```
lines = spark.textFile("hdfs://...")
errors = lines.filter(_.startsWith("ERROR"))
messages = errors.map(_.split('\t')(2))
messages.persist()
messages.filter(_.contains("MySQL")).count
messages.filter(_.contains("HDFS")).count
```

*transf*

*Action*

**Result:** scaled to 1 TB data in 5-7 sec  
(vs 170 sec for on-disk data)



# Fault recovery

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

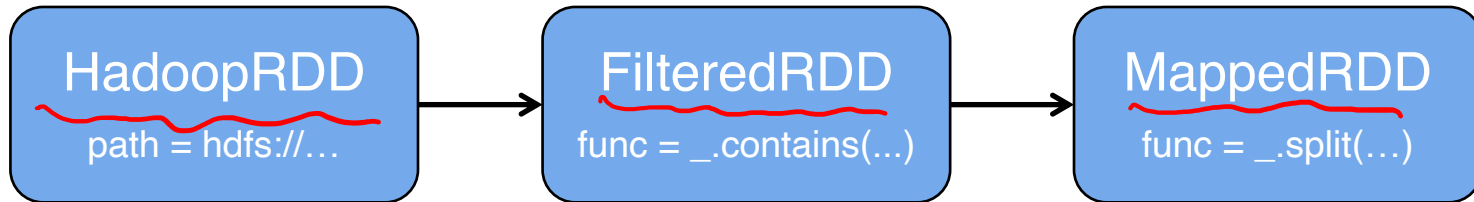
E.g.:

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

pattern



2nd col

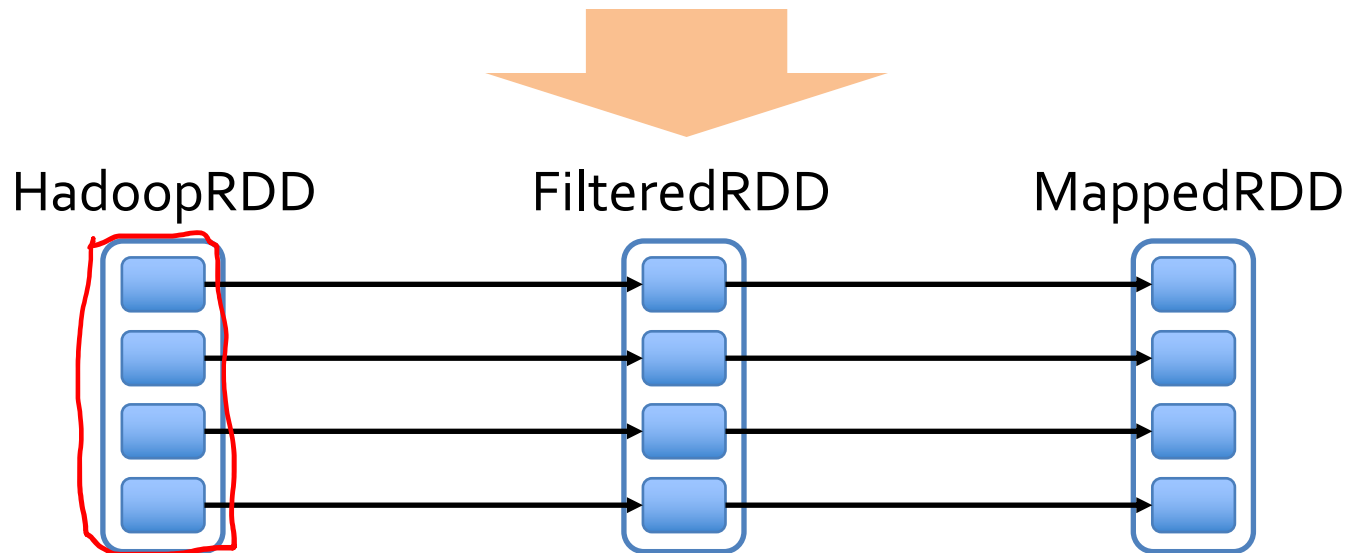


# Fault recovery

- 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))
```

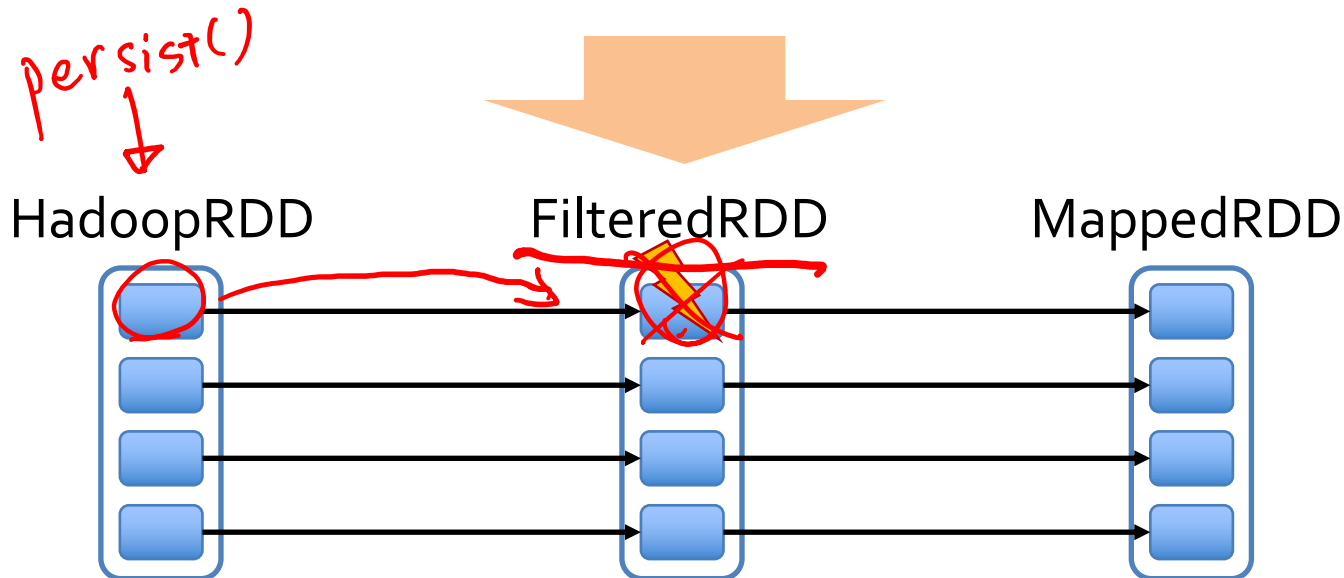


# Fault recovery

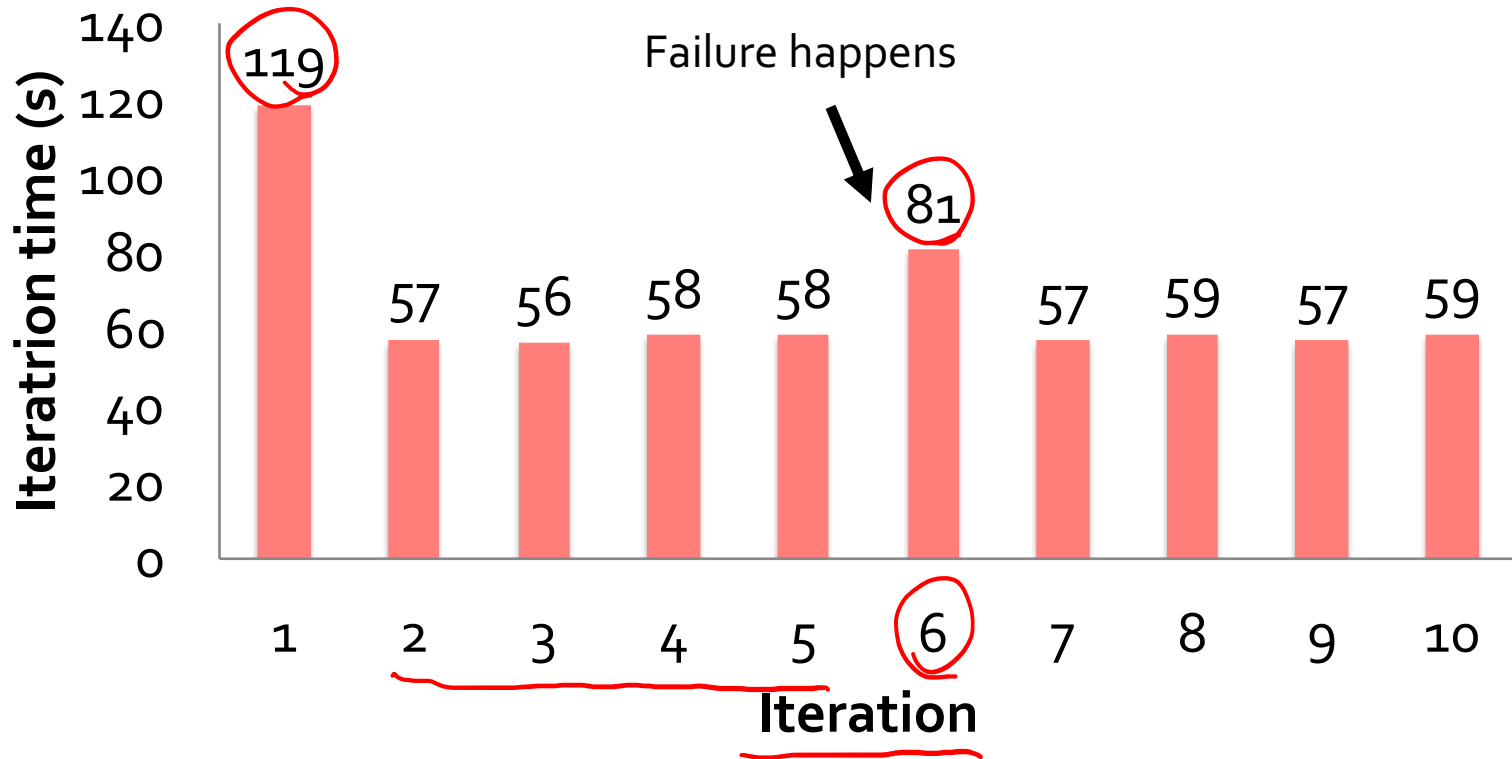
- RDDs track the graph of transformations that built them (their *lineage*) to rebuild lost data

E.g.:

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



# Fault recovery results



# Example: PageRank

foo.com



1. Start each page with a rank of 1

2. On each iteration, update each page's rank to

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

$$\frac{4}{|\text{neighbors}_{\text{foo}}|} = \frac{4}{2} = 2$$

links = // RDD of (url, neighbors) pairs

ranks = // RDD of (url, rank) pairs

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

Contrib factor

# 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(_ + _)
}
```

For each neighbor in links emits (URL, RankContrib)

Reduce to RDD[(URL, Rank)]



# Join ( $\bowtie$ )

Age

Alice	5
Bob	6
Claire	4



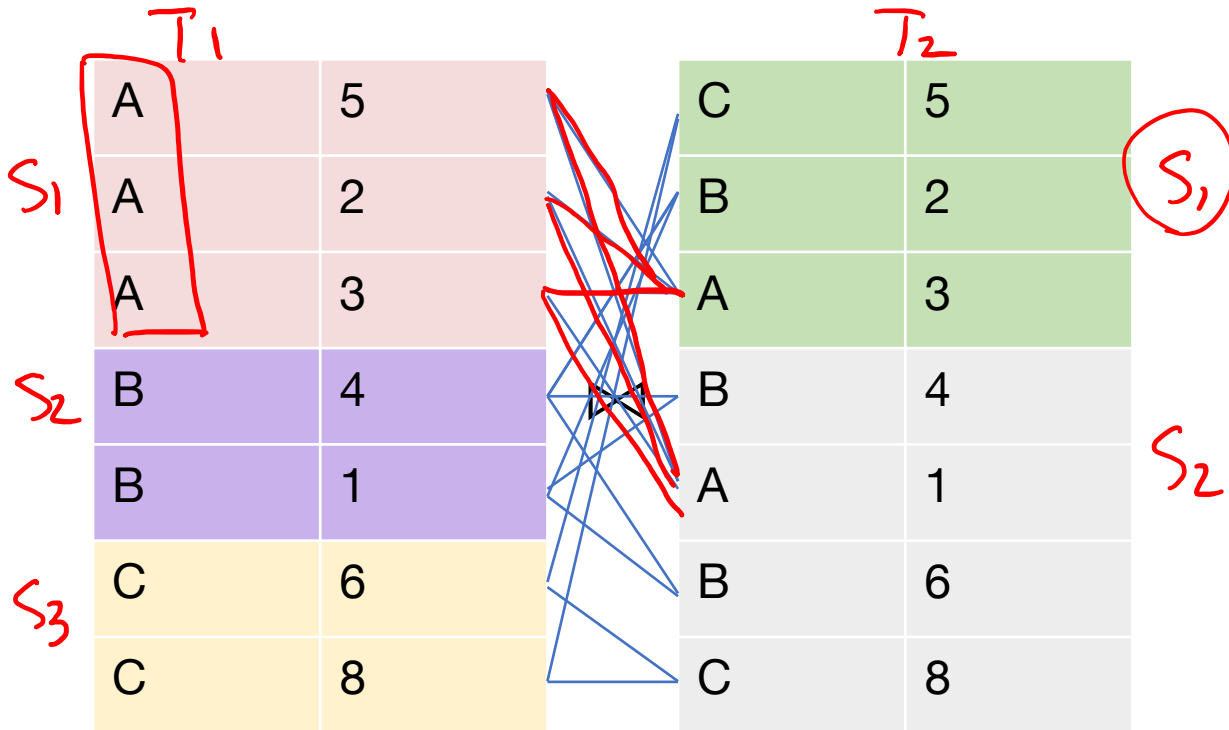
Gender.

Alice	F
Bob	M
Claire	F

=

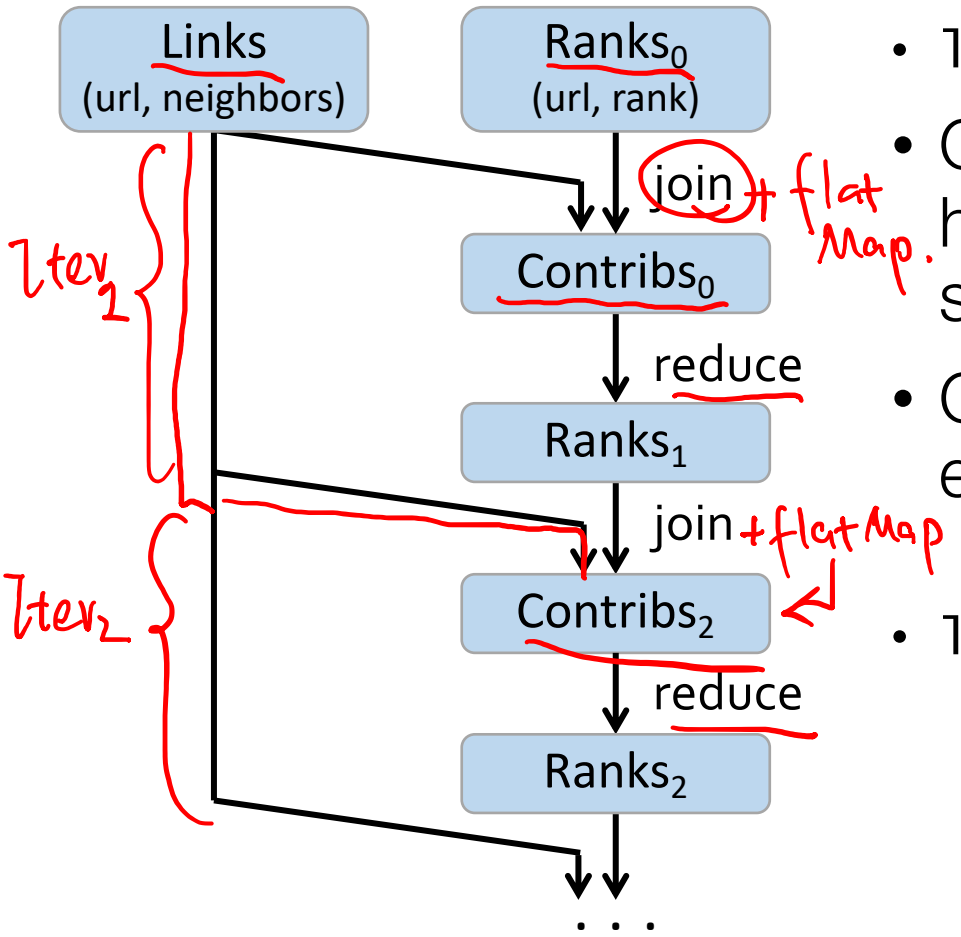
Age. Gender.

Alice	5	F
Bob	6	M
Claire	4	F



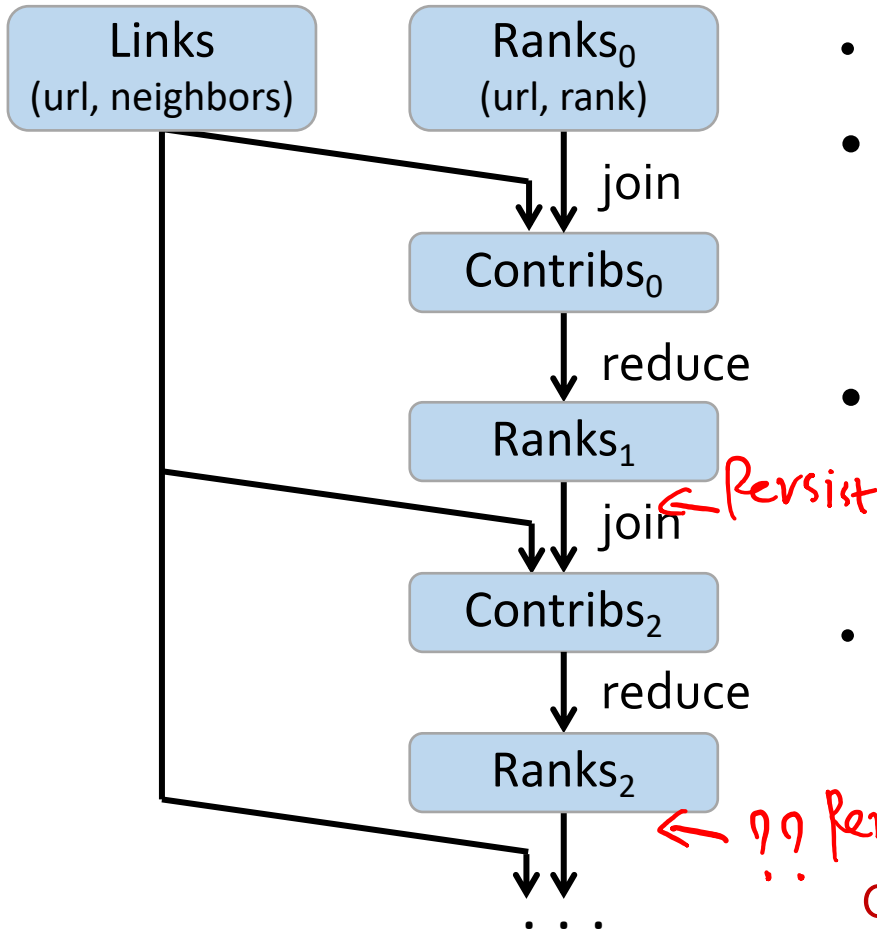
If partitioning doesn't match, then need to reshuffle to match pairs. Same problem in `reduce()` for MapReduce.

# 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())`

# 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()`?

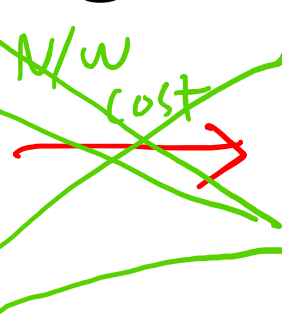
flat Map .

reduceByKey

# Co-partitioning example

ranks  
↓

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



bar.	2
bad.	0
foo.	5+4+2=11 → 1/
widget	3

5  
1

4  
1

link:

ranks,

flatten →

rank  
neighbor  
key

bar	foo
bad	foo
foo	widget
widget	bar, foo



bar	5
bad	4
foo	3
widget	4

bar	foo ✓	5
bad	foo ✓	4
foo	widget ✓	3
widget	bar, foo	4

3  
1

4  
2

Co-partitioning can avoid shuffle on join

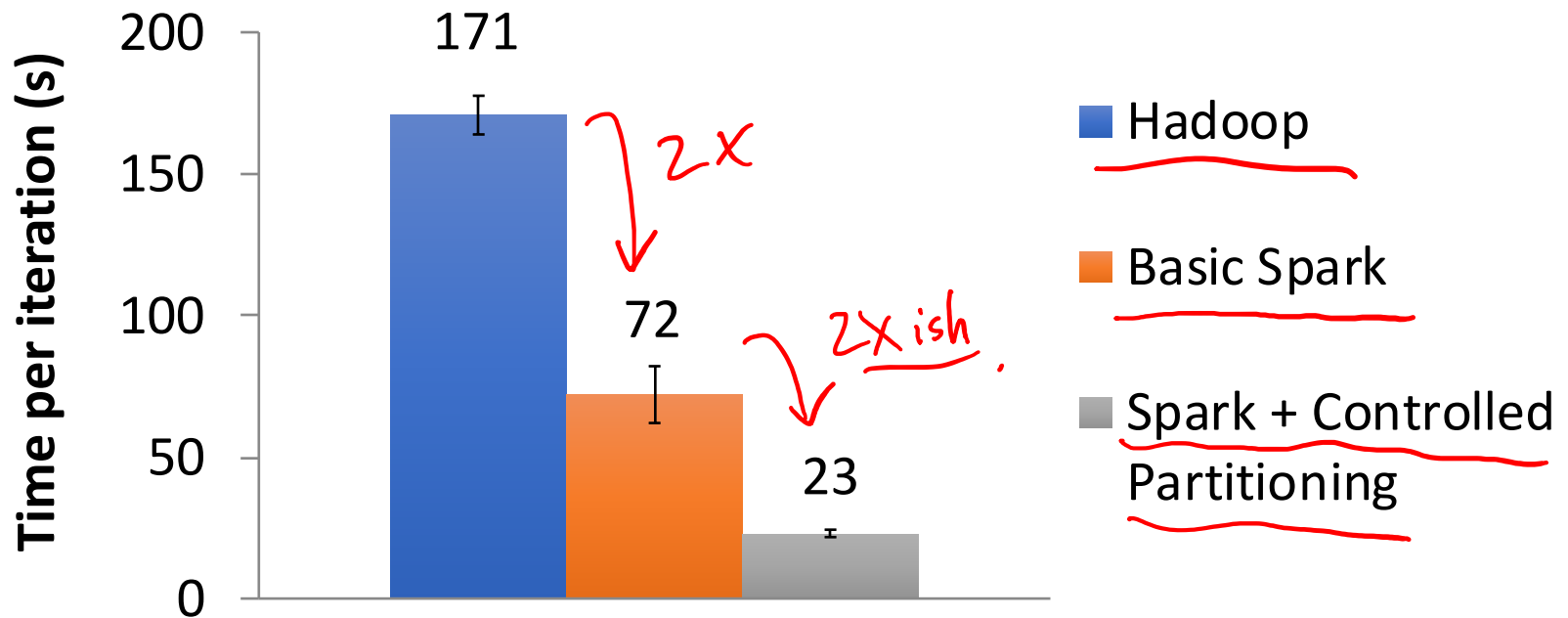
But, fundamentally a shuffle on reduceByKey

Optimization: custom partitioner on domain

4  
2

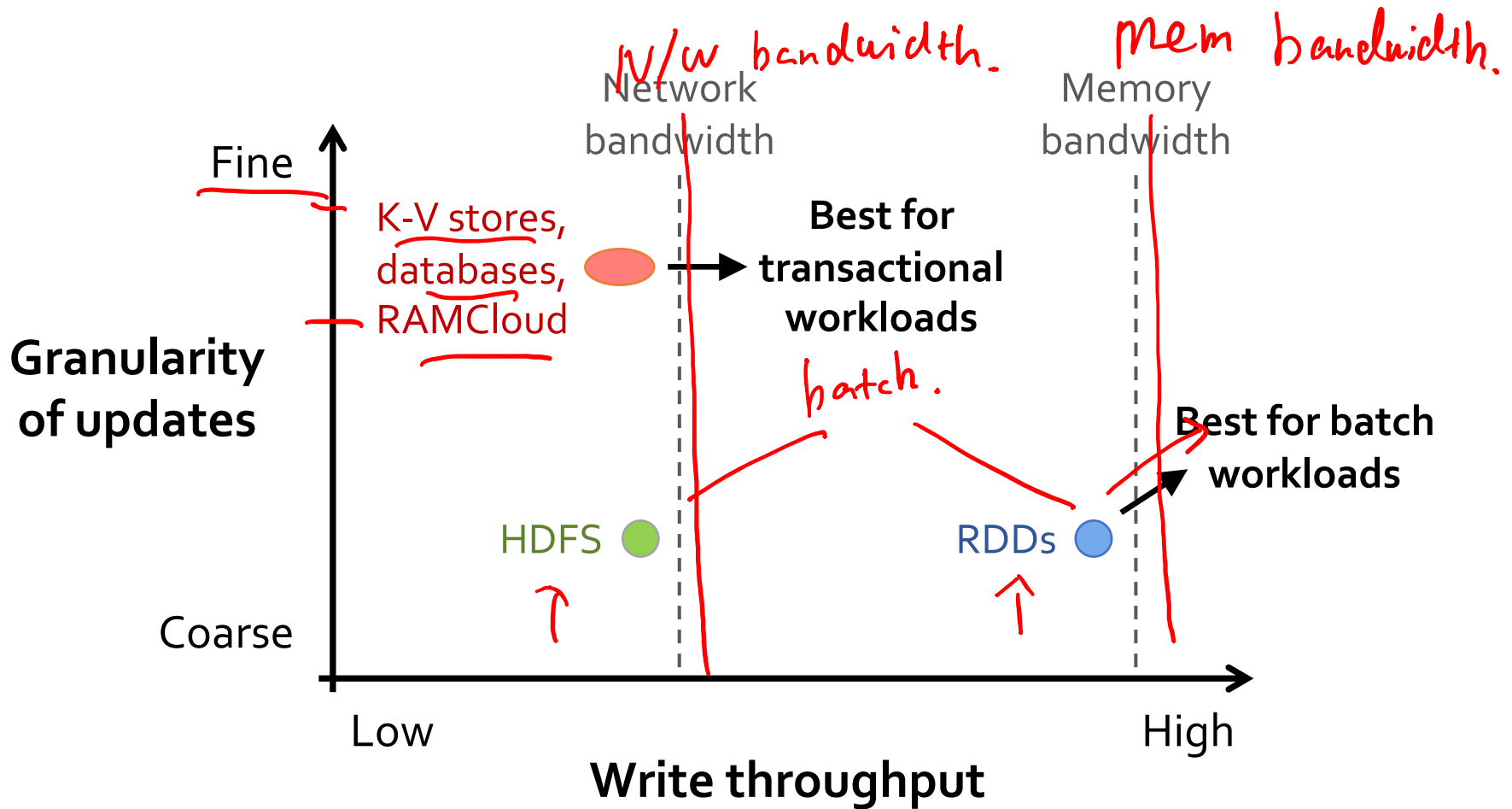
.com .edu

# PageRank performance



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

# Tradeoff space



# Discussion & wrap-up

