

SHIJIE ZHANG

SPARK CORE

OUTLINE

- ▶ Intro
- ▶ Standalone-view RDD concepts
 - ▶ Basic RDD
 - ▶ Key-Pair RDD
- ▶ Distributed-view RDD concepts
 - ▶ Lazy evaluation
 - ▶ Cache
 - ▶ Shuffle/Partition
 - ▶ Broadcast
 - ▶ Accumulator
- ▶ Conclusion

INTRO - LEGACY NUMBERS

Latency numbers “every programmer should know:”¹

L1 cache reference	0.5 ns	
Branch mispredict	5 ns	
L2 cache reference	7 ns	
Mutex lock/unlock	25 ns	
Main memory reference	100 ns	
Compress 1K bytes with Zippy	3,000 ns	= 3 µs
Send 2K bytes over 1 Gbps network	20,000 ns	= 20 µs
SSD random read	150,000 ns	= 150 µs
Read 1 MB sequentially from memory	250,000 ns	= 250 µs
Round trip within same datacenter	500,000 ns	= 0.5 ms
Read 1 MB sequentially from SSD*	1,000,000 ns	= 1 ms
Disk seek	10,000,000 ns	= 10 ms
Read 1 MB sequentially from disk	20,000,000 ns	= 20 ms
Send packet CA->Netherlands->CA	150,000,000 ns	= 150 ms

(Assuming ~1GB/sec SSD.)

INTRO

- ▶ To get a better intuition about differences of these numbers, humanize these durations

Method: multiply all these durations by a billion.

- ▶ Then, we can map each latency number to a human activity

Humanized durations grouped by magnitude:

INTRO

Minute:

L1 cache reference	0.5 s	One heart beat (0.5 s)
Branch mispredict	5 s	Yawn
L2 cache reference	7 s	Long yawn
Mutex lock/unlock	25 s	Making a coffee

Hour:

Main memory reference	100 s	Brushing your teeth
Compress 1K bytes with Zippy	50 min	One episode of a TV show

Day:

Send 2K bytes over 1 Gbps network	5.5 hr	From lunch to end of work day
-----------------------------------	--------	-------------------------------

Week:

SSD random read	1.7 days	A normal weekend
Read 1 MB sequentially from memory	2.9 days	A long weekend
Round trip within same datacenter	5.8 days	A medium vacation
Read 1 MB sequentially from SSD	11.6 days	Waiting for almost 2 weeks

Year:

Disk seek	16.5 weeks	A semester in university
Read 1 MB sequentially from disk	7.8 months	Almost producing a new human being
The above 2 together	1 year	

Decade:

Send packet CA->Netherlands->CA	4.8 years	Average time it takes to complete a bachelor's degree
---------------------------------	-----------	---

Humanized numbers

Shared Memory

Seconds

L1 cache reference.....0.5s
L2 cache reference.....7s
Mutex lock/unlock.....25s

Minutes

Main memory reference.....1m 40s

Distributed

Days

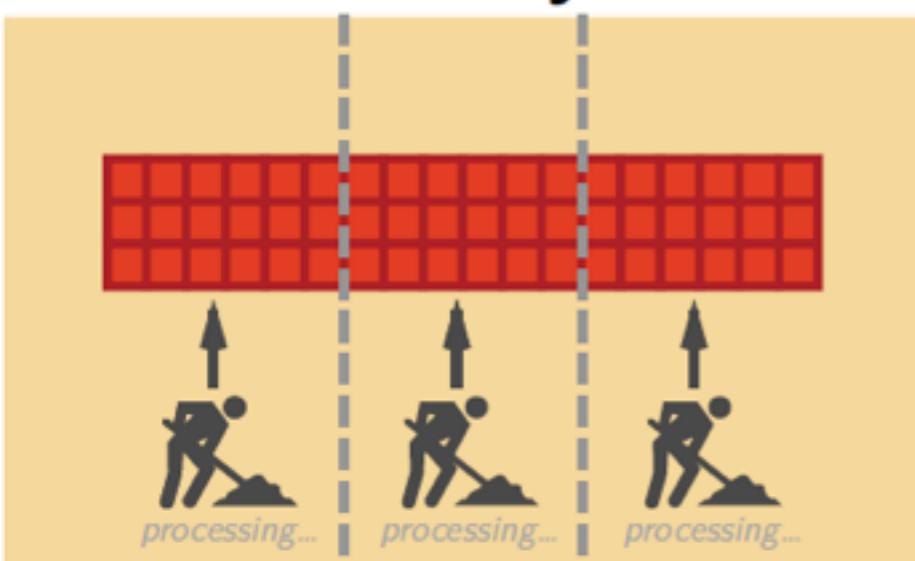
Roundtrip within
same datacenter.....5.8 days

Years

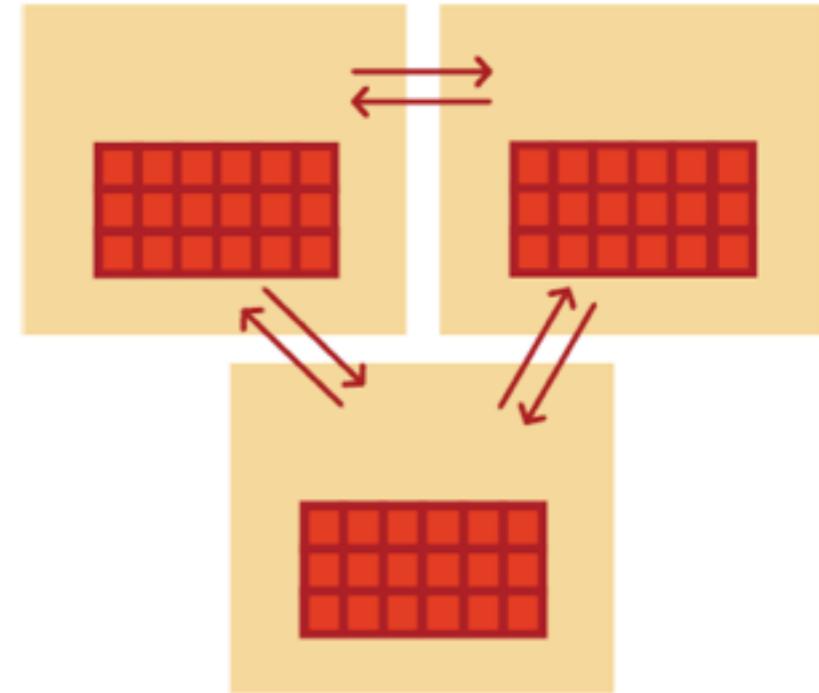
Send packet
CA->Netherlands->CA....4.8 years

What does **distributed** data-parallel look like?

Shared memory:



Distributed:



Shared memory case: Data-parallel programming model. Data partitioned in memory and operated upon in parallel.

Distributed case: Data-parallel programming model. Data partitioned between machines, network in between, operated upon in parallel.

Distributed data parallel programming, two more worries

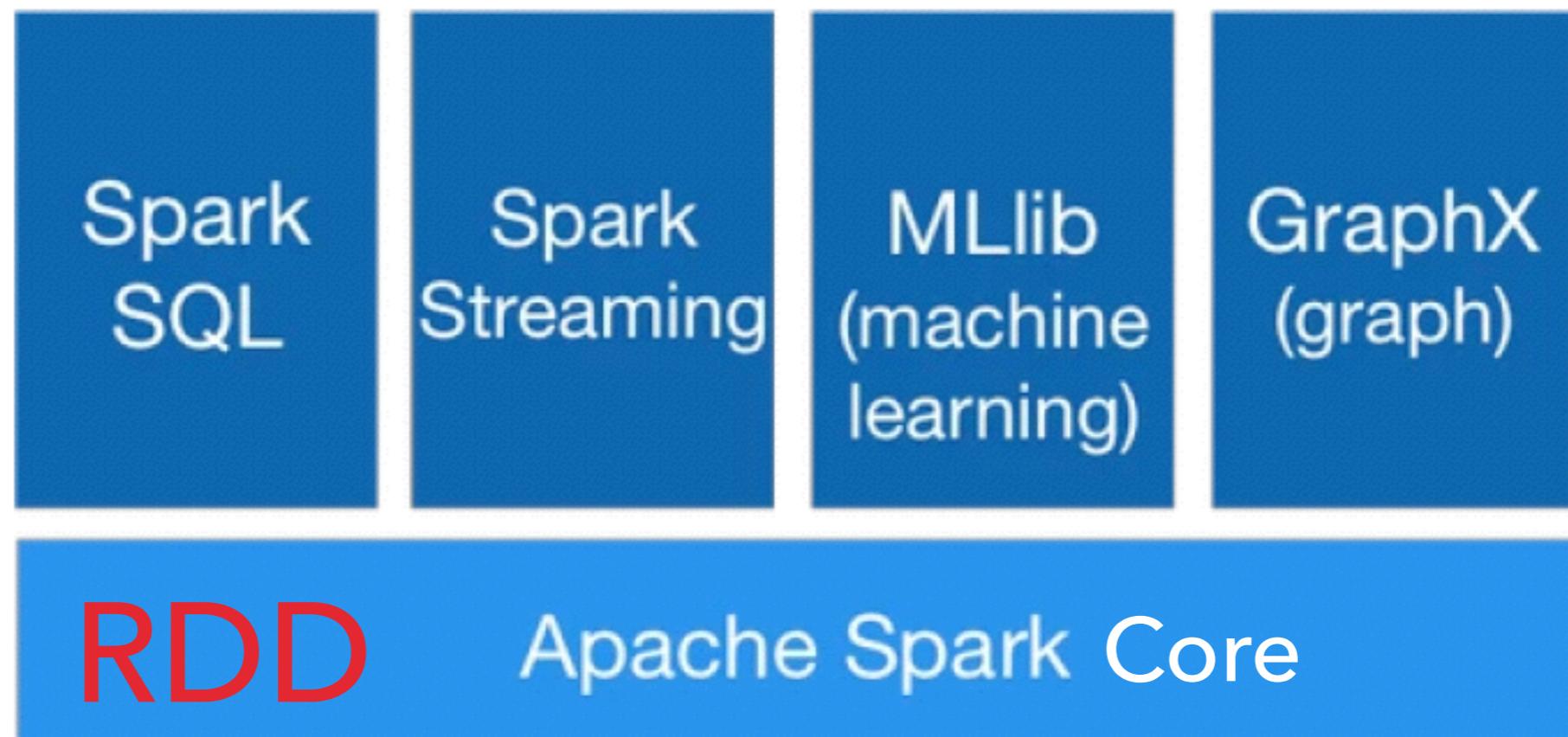
- ▶ Partial failure: a subset of machine crash
- ▶ Latency: certain operations have higher latency



Spark

- ▶ A distributed data parallel programming framework
- ▶ Implements a data parallel model called RDD

Spark components



WHAT'S RDD

- ▶ Just looks like Scala collections (but distributed across machines)

**Combinators on Scala
parallel/sequential collections:**

map
flatMap
filter
reduce
fold
aggregate

Combinators on RDDs:

map
flatMap
filter
reduce
fold
aggregate

WHAT'S RDD

- ▶ While signatures differ a bit, their semantics are the same

```
map[B](f: A => B): List[B] // Scala List
```

```
map[B](f: A => B): RDD[B] // Spark RDD
```

```
flatMap[B](f: A => TraversableOnce[B]): List[B] // Scala List
```

```
flatMap[B](f: A => TraversableOnce[B]): RDD[B] // Spark RDD
```

```
filter(pred: A => Boolean): List[A] // Scala List
```

```
filter(pred: A => Boolean): RDD[A] // Spark RDD
```

```
reduce(op: (A, A) => A): A // Scala List
```

```
reduce(op: (A, A) => A): A // Spark RDD
```

```
fold(z: A)(op: (A, A) => A): A // Scala List
```

```
fold(z: A)(op: (A, A) => A): A // Spark RDD
```

```
aggregate[B](z: => B)(seqop: (B, A) => B, combop: (B, B) => B): B // Scala
```

```
aggregate[B](z: B)(seqop: (B, A) => B, combop: (B, B) => B): B // Spark RDD
```

WHAT'S RDD

- ▶ Scientific Def: An interface containing five properties
 - ▶ Set of partitions
 - ▶ List of dependencies
 - ▶ Function to compute a partition given its parents
 - ▶ (Optional) Partition method
 - ▶ (Optional) Preferred locations

Spark Core components - code base (2012)

RDD Spark core: 16,000 LOC

Operators: 2000

Scheduler: 2500

Block manager: 2700

Networking: 1200

Accumulators: 200

Broadcast: 3500

Interpreter:
3300 LOC

Hadoop I/O:
400 LOC

Mesos backend:
700 LOC

Standalone backend:
1700 LOC

OUTLINE

- ▶ Intro
 - ▶ Standalone-view RDD concepts
 - ▶ Basic RDD
 - ▶ Key-Pair RDD
 - ▶ Distributed-view RDD concepts
 - ▶ Lazy evaluation
 - ▶ Cache
 - ▶ Partition/Shuffle
 - ▶ Broadcast
 - ▶ Accumulator
 - ▶ Conclusion
-
- ```
graph LR; Intro[] --> SV_RDD[]; SV_RDD --> BasicRDD[]; SV_RDD --> KeyPairRDD[]; DV_RDD[] --> LazyEvaluation[]; DV_RDD --> Cache[]; DV_RDD --> PartitionShuffle[]; DV_RDD --> Broadcast[]; DV_RDD --> Accumulator[]; Conclusion[] --> SV_RDD; Conclusion[] --> DV_RDD; Operators[Operators] --- BasicRDD; Operators --- KeyPairRDD; Scheduler[Scheduler] --- LazyEvaluation; Scheduler --- Cache; Scheduler --- PartitionShuffle; Scheduler --- Broadcast; Scheduler --- Accumulator; BlockManager[Block manager] --- LazyEvaluation; BlockManager --- Cache; BlockManager --- PartitionShuffle; BlockManager --- Broadcast; BlockManager --- Accumulator; Networking[Networking] --- LazyEvaluation; Networking --- Cache; Networking --- PartitionShuffle; Networking --- Broadcast; Networking --- Accumulator; Broadcast[Broadcast] --- LazyEvaluation; Broadcast --- Cache; Broadcast --- PartitionShuffle; Broadcast --- Broadcast; Broadcast --- Accumulator; Accumulator[Accumulator] --- LazyEvaluation; Accumulator --- Cache; Accumulator --- PartitionShuffle; Accumulator --- Broadcast; Accumulator --- Accumulator;
```

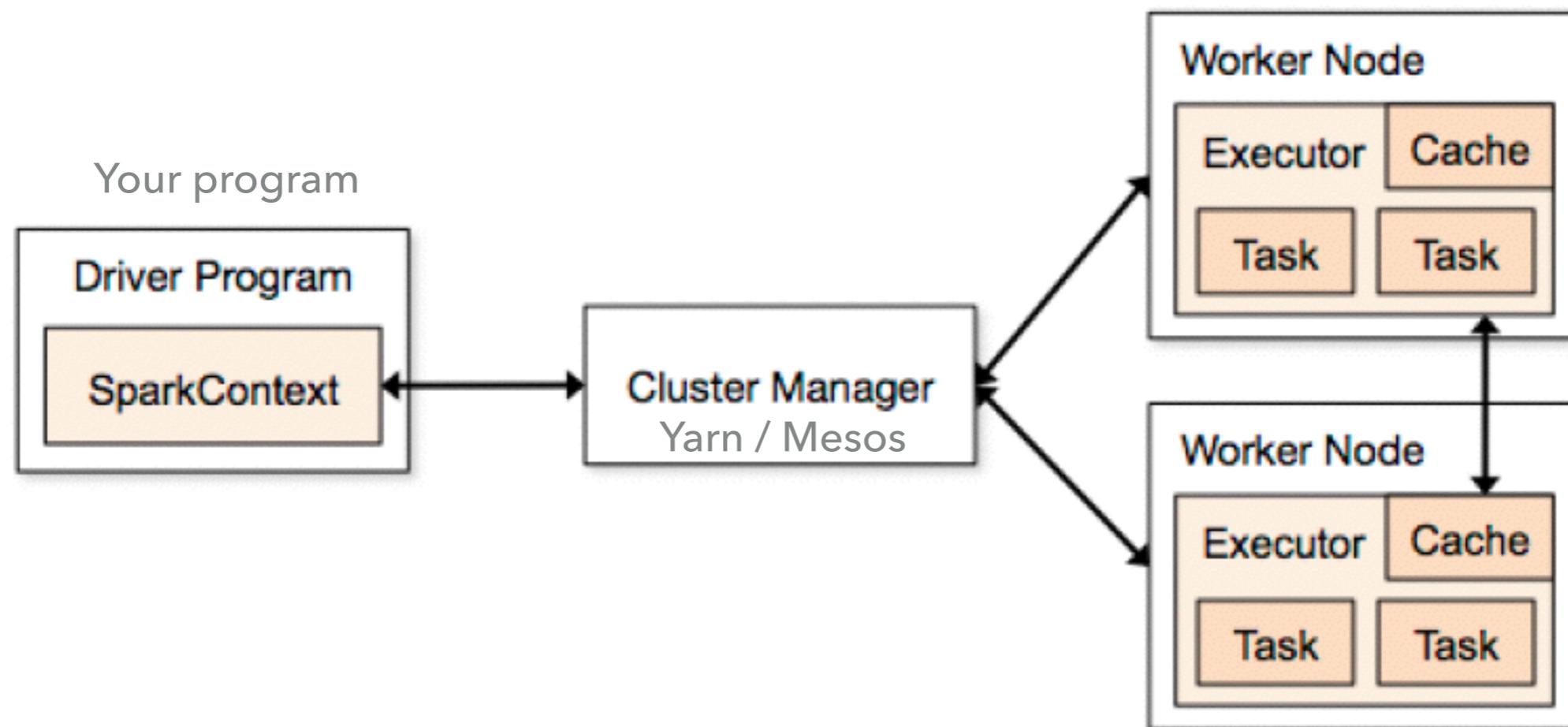
# OUTLINE

---

- ▶ Intro
  - ▶ **Standalone-view RDD concepts**
    - ▶ Basic RDD
    - ▶ Key-Pair RDD
  - ▶ **Distributed-view RDD concepts**
    - ▶ Lazy evaluation
    - ▶ Cache
    - ▶ Partition/Shuffle
    - ▶ Broadcast
    - ▶ Accumulator
  - ▶ Conclusion
- 
- ```
graph LR; Intro[ ] --> Standalone[Standalone-view RDD concepts]; Standalone --> BasicRDD[Basic RDD]; Standalone --> KeyPairRDD[Key-Pair RDD]; Distributed[ ] --> LazyEvaluation[Lazy evaluation]; Distributed --> Cache[Cache]; Distributed --> PartitionShuffle[Partition/Shuffle]; Distributed --> Broadcast[Broadcast]; Distributed --> Accumulator[Accumulator]; Operators[Operators] --- BasicRDD; Operators --- KeyPairRDD; Scheduler[Scheduler] --- LazyEvaluation; Scheduler --- Cache; Scheduler --- PartitionShuffle; Scheduler --- Broadcast; Scheduler --- Accumulator; BlockManager[Block manager] --- LazyEvaluation; BlockManager --- Cache; BlockManager --- PartitionShuffle; BlockManager --- Broadcast; BlockManager --- Accumulator; Networking[Networking] --- LazyEvaluation; Networking --- Cache; Networking --- PartitionShuffle; Networking --- Broadcast; Networking --- Accumulator; Broadcast[Broadcast] --- LazyEvaluation; Broadcast --- Cache; Broadcast --- PartitionShuffle; Broadcast --- Accumulator; Accumulator[Accumulator] --- LazyEvaluation; Accumulator --- Cache; Accumulator --- PartitionShuffle; Accumulator --- Broadcast;
```

Execution of Spark program

- Execution mode: Standalone / Cluster



OUTLINE

- ▶ Intro
 - ▶ **Standalone-view RDD concepts**
 - ▶ Basic RDD
 - ▶ Key-Pair RDD
 - ▶ Distributed-view RDD concepts
 - ▶ Lazy evaluation
 - ▶ Cache
 - ▶ Partition/Shuffle
 - ▶ Broadcast
 - ▶ Accumulator
 - ▶ Conclusion
-
- ```
graph TD; Intro[] --> Standalone[Standalone-view RDD concepts]; Standalone --> BasicRDD[Basic RDD]; Standalone --> KeyPairRDD[Key-Pair RDD]; Standalone --- Operators{Operators}; Operators --- BasicRDD; Operators --- KeyPairRDD; Standalone --> Distribute[]; Distribute --> LazyEvaluation[Lazy evaluation]; Distribute --> Cache[Cache]; Distribute --> PartitionShuffle[Partition/Shuffle]; Distribute --> Broadcast[Broadcast]; Distribute --> Accumulator[Accumulator]; Distribute --- Scheduler{Scheduler}; Scheduler --- LazyEvaluation; Distribute --- BlockManager{Block manager}; BlockManager --- Cache; Distribute --- Networking{Networking}; Networking --- PartitionShuffle; Distribute --- Broadcast{Broadcast}; Broadcast --- Broadcast; Distribute --- Accumulator{Accumulator}; Accumulator --- Accumulator;
```

## STANDALONE-VIEW

---



Note: actions/transformations will not modify input RDD, but will always create a new one

# RDD CREATION

- ▶ Loading an external dataset

*Example 3-5. parallelize() method in Python*

```
lines = sc.parallelize(["pandas", "i like pandas"])
```

- ▶ Parallelizing a collection

*Example 3-8. textFile() method in Python*

```
lines = sc.textFile("/path/to/README.md")
```

# TRANSFORMATIONS

Table 3-2. Basic RDD transformations on an RDD containing {1, 2, 3, 3}

| Function name                                          | Purpose                                                                                                                               | Example                                   | Result             |
|--------------------------------------------------------|---------------------------------------------------------------------------------------------------------------------------------------|-------------------------------------------|--------------------|
| <code>map()</code>                                     | Apply a function to each element in the RDD and return an RDD of the result.                                                          | <code>rdd.map(x =&gt; x + 1)</code>       | {2, 3, 4, 4}       |
| <code>flatMap()</code>                                 | Apply a function to each element in the RDD and return an RDD of the contents of the iterators returned. Often used to extract words. | <code>rdd.flatMap(x =&gt; x.to(3))</code> | {1, 2, 3, 2, 3, 3} |
| <code>filter()</code>                                  | Return an RDD consisting of only elements that pass the condition passed to <code>filter()</code> .                                   | <code>rdd.filter(x =&gt; x != 1)</code>   | {2, 3, 3}          |
| <code>distinct()</code>                                | Remove duplicates.                                                                                                                    | <code>rdd.distinct()</code>               | {1, 2, 3}          |
| <code>sample(withReplacement, fraction, [seed])</code> | Sample an RDD, with or without replacement.                                                                                           | <code>rdd.sample(false, 0.5)</code>       | Nondeterministic   |

# TRANSFORMATIONS

Table 3-3. Two-RDD transformations on RDDs containing {1, 2, 3} and {3, 4, 5}

| Function name  | Purpose                                                      | Example                 | Result                      |
|----------------|--------------------------------------------------------------|-------------------------|-----------------------------|
| union()        | Produce an RDD containing elements from both RDDs.           | rdd.union(other)        | {1, 2, 3, 3, 4, 5}          |
| intersection() | RDD containing only elements found in both RDDs.             | rdd.intersection(other) | {3}                         |
| subtract()     | Remove the contents of one RDD (e.g., remove training data). | rdd.subtract(other)     | {1, 2}                      |
| cartesian()    | Cartesian product with the other RDD.                        | rdd.cartesian(other)    | {(1, 3), (1, 4), ... (3,5)} |

# ACTIONS

*Table 3-4. Basic actions on an RDD containing {1, 2, 3, 3}*

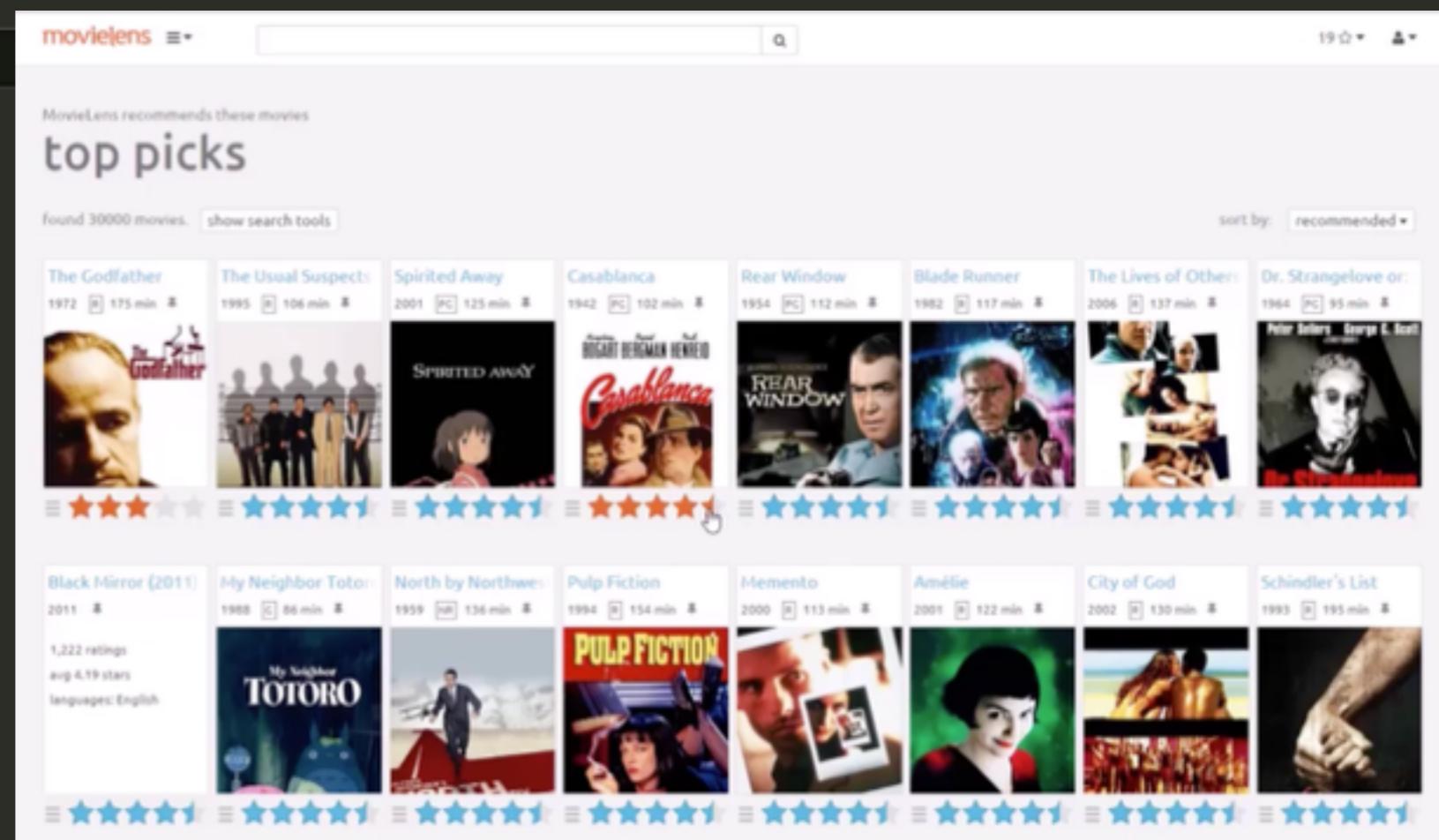
| Function name  | Purpose                                         | Example            | Result                         |
|----------------|-------------------------------------------------|--------------------|--------------------------------|
| collect()      | Return all elements from the RDD.               | rdd.collect()      | {1, 2, 3, 3}                   |
| count()        | Number of elements in the RDD.                  | rdd.count()        | 4                              |
| countByValue() | Number of times each element occurs in the RDD. | rdd.countByValue() | {(1, 1),<br>(2, 1),<br>(3, 2)} |

|                                                       |                                                                       |                                                                                                                |                  |
|-------------------------------------------------------|-----------------------------------------------------------------------|----------------------------------------------------------------------------------------------------------------|------------------|
| <code>take(num)</code>                                | Return num elements from the RDD.                                     | <code>rdd.take(2)</code>                                                                                       | {1, 2}           |
| <code>top(num)</code>                                 | Return the top num elements the RDD.                                  | <code>rdd.top(2)</code>                                                                                        | {3, 3}           |
| <code>takeOrdered(num)(ordering)</code>               | Return num elements based on provided ordering.                       | <code>rdd.takeOrdered(2)(myOrdering)</code>                                                                    | {3, 3}           |
| <code>takeSample(withReplacement, num, [seed])</code> | Return num elements at random.                                        | <code>rdd.takeSample(false, 1)</code>                                                                          | Nondeterministic |
| <code>reduce(func)</code>                             | Combine the elements of the RDD together in parallel (e.g., sum).     | <code>rdd.reduce((x, y) =&gt; x + y)</code>                                                                    | 9                |
| <code>fold(zero)(func)</code>                         | Same as <code>reduce()</code> but with the provided zero value.       | <code>rdd.fold(0)((x, y) =&gt; x + y)</code>                                                                   | 9                |
| <code>aggregate(zeroValue)(seqOp, combOp)</code>      | Similar to <code>reduce()</code> but used to return a different type. | <code>rdd.aggregate((0, 0))((x, y) =&gt; (x._1 + y, x._2 + 1), (x, y) =&gt; (x._1 + y._1, x._2 + y._2))</code> | (9, 4)           |
| <code>foreach(func)</code>                            | Apply the provided function to each element of the RDD.               | <code>rdd.foreach(func)</code>                                                                                 | Nothing          |

# PYTHON EXAMPLE - COUNT RATING NUM FOR EACH STAR

```
Ratings-Counter.py x
1 from pyspark import SparkConf, SparkContext
2 import collections
3
4 conf = SparkConf().setMaster("local").setAppName("RatingsHistogram")
5 sc = SparkContext(conf = conf)
6
7 lines = sc.textFile("file:///SparkCourse/ml-100k/u.data")
8 ratings = lines.map(lambda x: x.split()[2])
9 result = ratings.countByValue()
10
11 sortedResults = collections.OrderedDict(sorted(result.items()))
12 for key, value in sortedResults.iteritems():
13 print "%s %i" % (key, value)
14
15
```

```
u.data x README x
1 # movie rating data from grouplens.org
2 # user id | item id | rating | timestamp.
3
4 196 242 3 881250949
5 186 302 3 891717742
6 22 377 1 878887116
7 244 51 2 880606923
8 166 346 1 886397596
9 298 474 4 884182806
10 115 265 2 881171488
11 253 465 5 891628467
12 305 451 3 886324817
13 6 86 3 883603013
14 62 257 2 879372434
15 286 1014 5 879781125
16 200 222 5 876042340
17 210 40 3 891035994
18 224 29 3 888104457
19 303 785 3 879485318
20 122 387 5 879270459
21 194 274 2 879539794
22 291 1042 4 874834944
23 234 1184 2 892079237
24 119 392 4 886176814
25 167 486 4 892738452
26 299 144 4 877881320
```



For more detailed step, see:

<https://www.udemy.com/taming-big-data-with-apache-spark-hands-on/learn/v4/overview>

## PYTHON EXAMPLE - COUNT RATING NUM FOR EACH STAR

---

```
1 from pyspark import SparkConf, SparkContext → import statements
2 import collections
3
4 conf = SparkConf().setMaster("local").setAppName("RatingsHistogram")
5 sc = SparkContext(conf = conf)
6
7 lines = sc.textFile("file:///SparkCourse/ml-100k/u.data")
8 ratings = lines.map(lambda x: x.split()[2])
9 result = ratings.countByValue()
10
11 sortedResults = collections.OrderedDict(sorted(result.items()))
12 for key, value in sortedResults.iteritems():
13 print "%s %i" % (key, value)
14
15
```

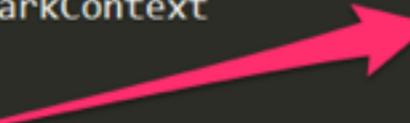
## PYTHON EXAMPLE - COUNT RATING NUM FOR EACH STAR

---

- The Spark Framework can adopt several cluster managers

- ▶ Local Mode
- ▶ Standalone mode
- ▶ Apache Mesos
- ▶ Hadoop YARN

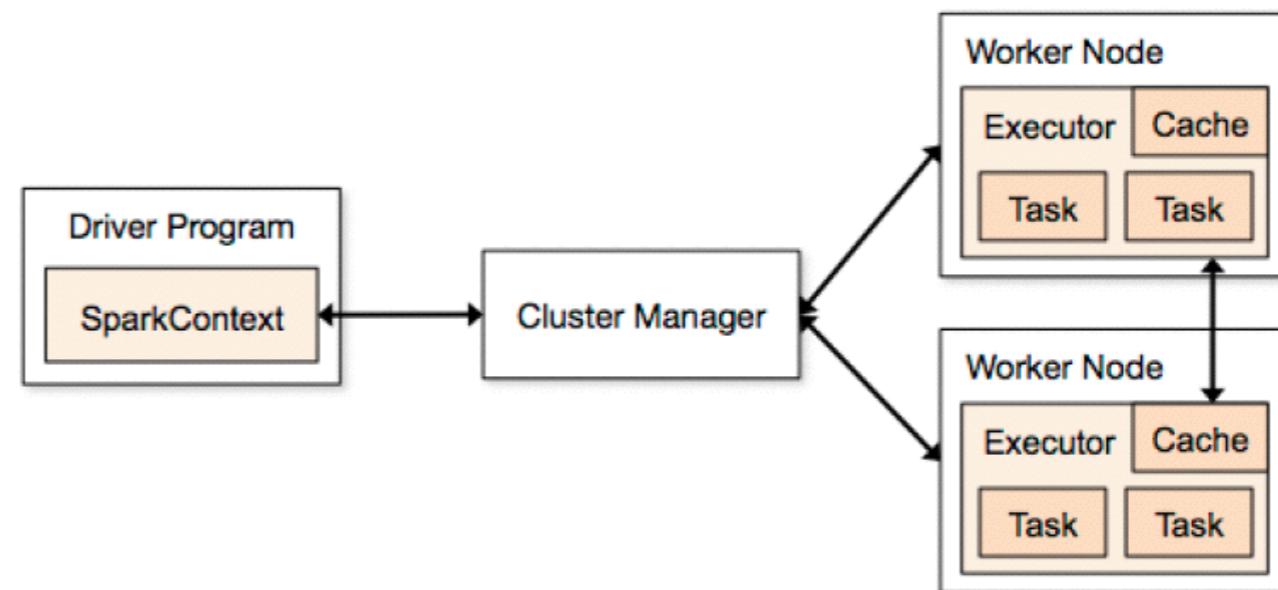
```
1 from pyspark import SparkConf, SparkContext
2 import collections
3
4 conf = SparkConf().setMaster("local").setAppName("RatingsHistogram")
5 sc = SparkContext(conf = conf)
6
7 lines = sc.textFile("file:///SparkCourse/ml-100k/u.data")
8 ratings = lines.map(lambda x: x.split()[2])
9 result = ratings.countByValue()
10
11 sortedResults = collections.OrderedDict(sorted(result.items()))
12 for key, value in sortedResults.iteritems():
13 print "%s %i" % (key, value)
14
15
```



## PYTHON EXAMPLE - COUNT RATING NUM FOR EACH STAR

---

```
1 from pyspark import SparkConf, SparkContext
2 import collections
3
4 conf = SparkConf().setMaster("local").setAppName("RatingsHistogram")
5 sc = SparkContext(conf = conf) → get a Spark Context, initialize the driver program
6
7 lines = sc.textFile("file:///SparkCourse/ml-100k/u.data")
8 ratings = lines.map(lambda x: x.split()[2])
9 result = ratings.countByValue()
10
11 sortedResults = collections.OrderedDict(sorted(result.items()))
12 for key, value in sortedResults.iteritems():
13 print "%s %i" % (key, value)
14
15
```



## PYTHON EXAMPLE - COUNT RATING NUM FOR EACH STAR

---

```
1 from pyspark import SparkConf, SparkContext
2 import collections
3
4 conf = SparkConf().setMaster("local").setAppName("RatingsHistogram")
5 sc = SparkContext(conf = conf)
6
7 lines = sc.textFile("file:///SparkCourse/ml-100k/u.data") → read a text file as a RDD
8 ratings = lines.map(lambda x: x.split()[2])
9 result = ratings.countByValue()
10
11 sortedResults = collections.OrderedDict(sorted(result.items()))
12 for key, value in sortedResults.iteritems():
13 print "%s %i" % (key, value)
14
15
```

## PYTHON EXAMPLE - COUNT RATING NUM FOR EACH STAR

---

```
1 from pyspark import SparkConf, SparkContext
2 import collections
3
4 conf = SparkConf().setMaster("local").setAppName("RatingsHistogram")
5 sc = SparkContext(conf = conf)
6
7 lines = sc.textFile("file:///SparkCourse/ml-100k/u.data")
8 ratings = lines.map(lambda x: x.split()[2])
9 result = ratings.countByValue()
10
11 sortedResults = collections.OrderedDict(sorted(result.items()))
12 for key, value in sortedResults.iteritems():
13 print "%s %i" % (key, value)
14
15
```



### transformation

| user id | item id | rating | timestamp |
|---------|---------|--------|-----------|
| 192     | 242     | 3      | 881250949 |
| 186     | 302     | 3      | 891717742 |

3  
3

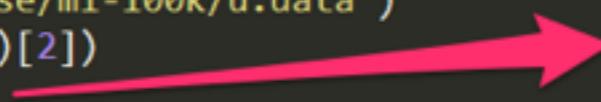
## PYTHON EXAMPLE - COUNT RATING NUM FOR EACH STAR

---

```
1 from pyspark import SparkConf, SparkContext
2 import collections
3
4 conf = SparkConf().setMaster("local").setAppName("RatingsHistogram")
5 sc = SparkContext(conf = conf)
6
7 lines = sc.textFile("file:///SparkCourse/ml-100k/u.data")
8 ratings = lines.map(lambda x: x.split()[2])
9 result = ratings.countByValue()
10
11 sortedResults = collections.OrderedDict(sorted(result.items()))
12 for key, value in sortedResults.iteritems():
13 print "%s %i" % (key, value)
14
15
```

Action which returns (value, count) pairs

3  
3  
3  
4  
4  
5  
**(3,3)**  
**(4,2)**  
**(5,1)**



## PYTHON EXAMPLE - COUNT RATING NUM FOR EACH STAR

---

```
1 from pyspark import SparkConf, SparkContext
2 import collections
3
4 conf = SparkConf().setMaster("local").setAppName("RatingsHistogram")
5 sc = SparkContext(conf = conf)
6
7 lines = sc.textFile("file:///SparkCourse/ml-100k/u.data")
8 ratings = lines.map(lambda x: x.split()[2])
9 result = ratings.countByValue()
10
11 sortedResults = collections.OrderedDict(sorted(result.items()))
12 for key, value in sortedResults.iteritems():
13 print "%s %i" % (key, value) → Loop through, printout python collection
14
15
```

```
(Canopy 64bit) C:\SparkCourse>spark-submit ratings-counter.py
15/09/29 11:39:35 WARN NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
15/09/29 11:39:37 WARN : Your hostname, HomeOfficePower resolves to a loopback/non-reachable address: fe80:0:0:0:5efe:c0a8:180%net2, but we couldn't find any external IP address!
15/09/29 11:39:38 WARN MetricsSystem: Using default name DAGScheduler for source because spark.app.id is not set.
1 6110
2 11370
3 27145
4 34174
5 21201
```

## KEY-PAIR RDD

---

- ▶ Def: RDDs containing key/value pairs
- ▶ Motivation: allows to act on each key

## KEY-PAIR RDD CREATION

---

- ▶ Create from RDDs by transforming to tuples

*Example 4-1. Creating a pair RDD using the first word as the key in Python*

```
pairs = lines.map(lambda x: (x.split(" ")[0], x))
```

- ▶ Load external data

# KEY-PAIR RDD TRANSFORMATIONS

---

## ► Key-Pair RDDs have all transformations available to standard RDDs

Table 4-1. Transformations on one pair RDD (example: `{(1, 2), (3, 4), (3, 6)}`)

| Function name                                                                                                                                          | Purpose                                                                                                                                                                        | Example                                          | Result                                                        |
|--------------------------------------------------------------------------------------------------------------------------------------------------------|--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|--------------------------------------------------|---------------------------------------------------------------|
| <code>reduceByKey(func)</code>                                                                                                                         | Combine values with the same key.                                                                                                                                              | <code>rdd.reduceByKey((x, y) =&gt; x + y)</code> | <code>{(1, 2), (3, 10)}</code>                                |
| <code>groupByKey()</code>                                                                                                                              | Group values with the same key.                                                                                                                                                | <code>rdd.groupByKey()</code>                    | <code>{(1, [2]), (3, [4, 6])}</code>                          |
| <code>combineByKey</code><br><code>Key(createCombiner,</code><br><code>mergeValue,</code><br><code>mergeCombiners,</code><br><code>partitioner)</code> | Combine values with the same key using a different result type.                                                                                                                | See Examples 4-12 through 4-14.                  |                                                               |
| <code>mapValues(func)</code>                                                                                                                           | Apply a function to each value of a pair RDD without changing the key.                                                                                                         | <code>rdd.mapValues(x =&gt; x+1)</code>          | <code>{(1, 3), (3, 5), (3, 7)}</code>                         |
| <code>flatMapValues(func)</code>                                                                                                                       | Apply a function that returns an iterator to each value of a pair RDD, and for each element returned, produce a key/value entry with the old key. Often used for tokenization. | <code>rdd.flatMapValues(x =&gt; (x to 5))</code> | <code>{(1, 2), (1, 3), (1, 4), (1, 5), (3, 4), (3, 5)}</code> |
| <code>keys()</code>                                                                                                                                    | Return an RDD of just the keys.                                                                                                                                                | <code>rdd.keys()</code>                          | <code>{1, 3, 3}</code>                                        |

## KEY-PAIR RDD TRANSFORMATIONS

---

- ▶ Key-Pair RDDs have all transformations available to standard RDDs

| Function name            | Purpose                           | Example                      | Result                   |
|--------------------------|-----------------------------------|------------------------------|--------------------------|
| <code>values()</code>    | Return an RDD of just the values. | <code>rdd.values()</code>    | {2, 4, 6}                |
| <code>sortByKey()</code> | Return an RDD sorted by the key.  | <code>rdd.sortByKey()</code> | {(1, 2), (3, 4), (3, 6)} |

*Table 4-2. Transformations on two pair RDDs (`rdd = {(1, 2), (3, 4), (3, 6)}` `other = {(3, 9)}`)*

| Function name               | Purpose                                                                         | Example                                | Result                                                  |
|-----------------------------|---------------------------------------------------------------------------------|----------------------------------------|---------------------------------------------------------|
| <code>subtractByKey</code>  | Remove elements with a key present in the other RDD.                            | <code>rdd.subtractByKey(other)</code>  | {(1, 2)}                                                |
| <code>join</code>           | Perform an inner join between two RDDs.                                         | <code>rdd.join(other)</code>           | {(3, (4, 9)), (3, (6, 9))}                              |
| <code>rightOuterJoin</code> | Perform a join between two RDDs where the key must be present in the first RDD. | <code>rdd.rightOuterJoin(other)</code> | {(3, (Some(4), 9)), (3, (Some(6), 9)))}                 |
| <code>leftOuterJoin</code>  | Perform a join between two RDDs where the key must be present in the other RDD. | <code>rdd.leftOuterJoin(other)</code>  | {(1, (2, None)), (3, (4, Some(9))), (3, (6, Some(9))))} |
| <code>cogroup</code>        | Group data from both RDDs sharing the same key.                                 | <code>rdd.cogroup(other)</code>        | {(1, ([2], [])), (3, ([4, 6], [9])))}                   |

## KEY-PAIR RDD ACTIONS

---

*Table 4-3. Actions on pair RDDs (example  $\{(1, 2), (3, 4), (3, 6)\}$ )*

| Function       | Description                                         | Example            | Result                      |
|----------------|-----------------------------------------------------|--------------------|-----------------------------|
| countByKey()   | Count the number of elements for each key.          | rdd.countByKey()   | $\{(1, 1), (3, 2)\}$        |
| collectAsMap() | Collect the result as a map to provide easy lookup. | rdd.collectAsMap() | Map{(1, 2), (3, 4), (3, 6)} |
| lookup(key)    | Return all values associated with the provided key. | rdd.lookup(3)      | [4, 6]                      |

# OUTLINE

---

- ▶ Intro
  - ▶ Standalone-view RDD concepts
    - ▶ Basic RDD
    - ▶ Key-Pair RDD
  - ▶ Distributed-view RDD concepts
    - ▶ Lazy evaluation
    - ▶ Cache
    - ▶ Partition/Shuffle
    - ▶ Broadcast
    - ▶ Accumulator
  - ▶ Conclusion
- 
- ```
graph LR; A[Operators] --- B[Basic RDD]; A --- C[Key-Pair RDD]; D[Scheduler] --- E[Lazy evaluation]; F[Block manager] --- G[Cache]; H[Networking] --- I[Partition/Shuffle]; J[Broadcast] --- K[Broadcast]; L[Accumulator] --- M[Accumulator]
```

OUTLINE

- ▶ Intro
 - ▶ Standalone-view RDD concepts
 - ▶ Basic RDD
 - ▶ Key-Pair RDD
 - ▶ Distributed-view RDD concepts
 - ▶ **Lazy evaluation**
 - ▶ Cache
 - ▶ Partition/Shuffle
 - ▶ Broadcast
 - ▶ Accumulator
 - ▶ Conclusion
-
- ```
graph TD; A[Operators] --- B[Basic RDD]; A --- C[Key-Pair RDD]; D[Scheduler] --- E[Lazy evaluation]; F[Block manager] --- G[Cache]; F --- H[Partition/Shuffle]; F --- I[Broadcast]; J[Networking] --- K[Accumulator]; L[Broadcast] --- M[Accumulator];
```

## LAZY-EVALUATION

---

- ▶ Def:
  - ▶ Transformations are **lazy**. Their result RDD is not immediately computed.
  - ▶ Actions are **eager**. Their result is immediately computed.
- ▶ Example:

```
lines = sc.textFile("file:///SparkCourse/ml-100k/u.data") → loading data operation, not executed at this point
ratings = lines.map(lambda x: x.split()[2])
result = ratings.countByValue()
```

## LAZY-EVALUATION

---

- ▶ Def:
  - ▶ Transformations are **lazy**. Their result RDD is not immediately computed.
  - ▶ Actions are **eager**. Their result is immediately computed.
- ▶ Example:

```
lines = sc.textFile("file:///SparkCourse/ml-100k/u.data")
ratings = lines.map(lambda x: x.split()[2])
result = ratings.countByValue()
```



**map transformation operation, not executed at this point**

## LAZY-EVALUATION

---

- ▶ Def:
  - ▶ Transformations are **lazy**. Their result RDD is not immediately computed.
  - ▶ Actions are **eager**. Their result is immediately computed.
- ▶ Example:

```
lines = sc.textFile("file:///SparkCourse/ml-100k/u.data")
ratings = lines.map(lambda x: x.split()[2])
result = ratings.countByValue()
```



**countByValue action operation, executed immediately**

## LAZY-EVALUATION

---

- ▶ Why:
  - ▶ Enable computation optimization

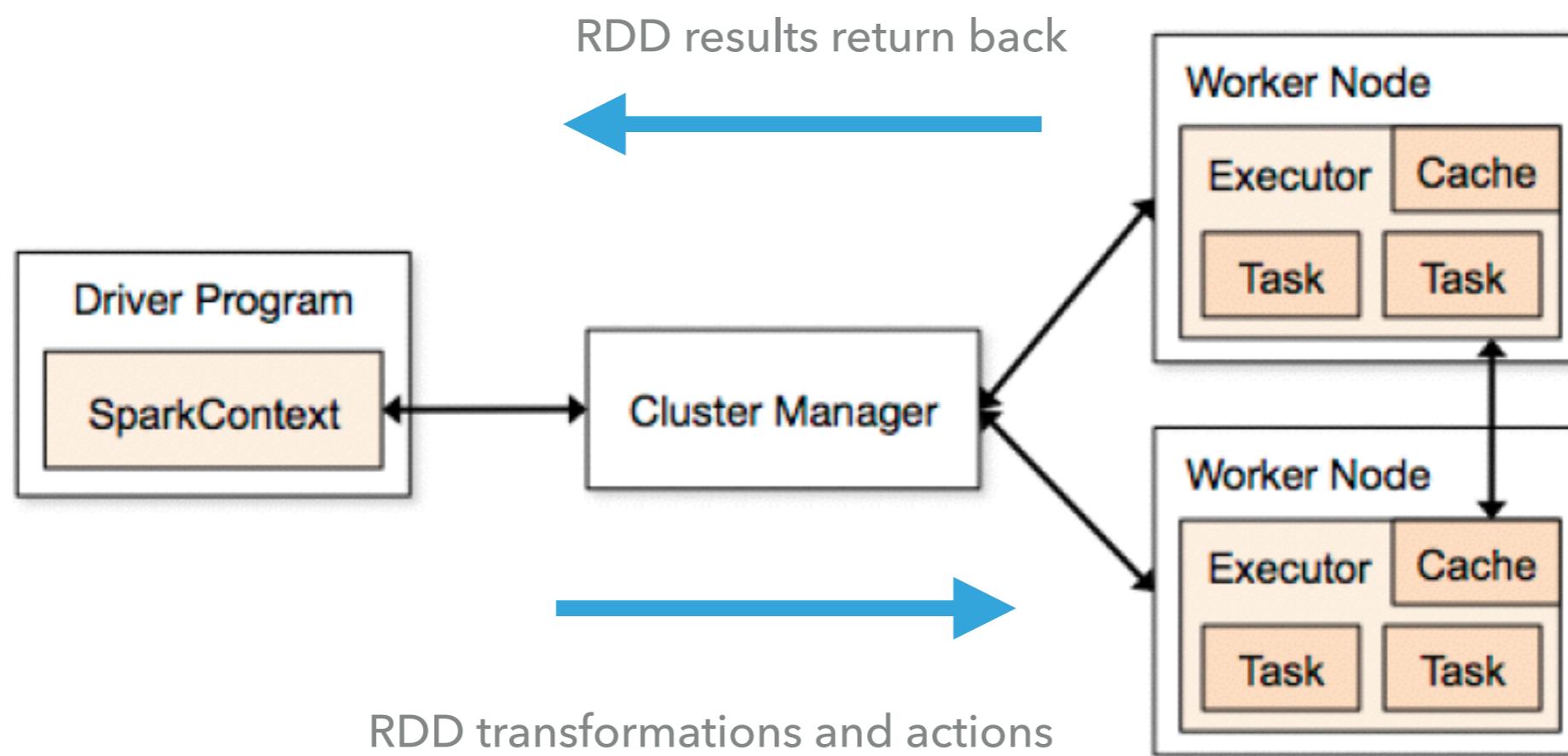
```
val lastYearsLogs: RDD[String] = ...
```

```
val firstLogsWithErrors = lastYearsLogs.filter(_.contains("ERROR")).take(10)
```

## LAZY-EVALUATION

---

- ▶ Why:
  - ▶ Limit network communication
    - ▶ If you perform an action, its result returns back to the driver program



## LAZY-EVALUATION

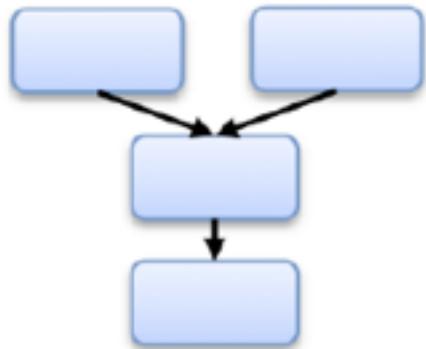
---

- ▶ What does it do when loading/transformation are performed:
  - ▶ Record RDD dependencies with a directed acyclic graph (DAG)
    - ▶ For failure recovery
    - ▶ Pipeline operations execution order
  - ▶ Optimize by utilizing network traffic based on partition and shuffling

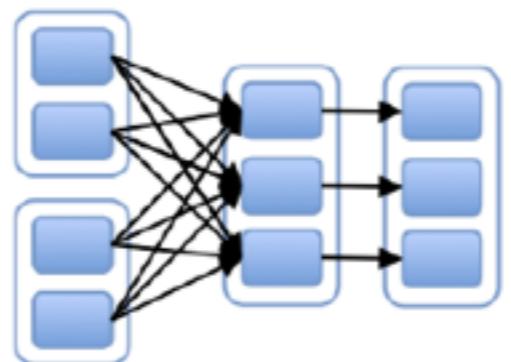
# LAZY EVALUATION

## Life time of a job in Spark

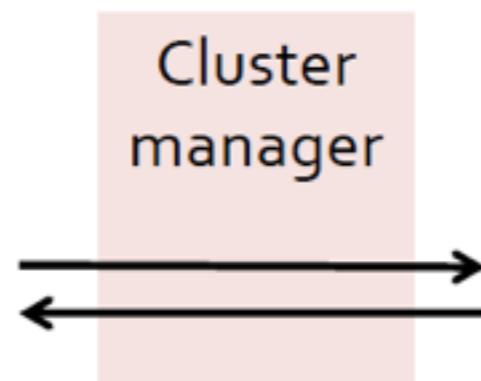
RDD Objects



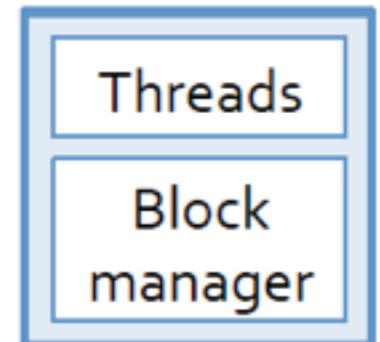
DAG Scheduler



Task Scheduler



Worker



`rdd1.join(rdd2)  
.groupBy(...)  
.filter(...)`

Build the operator DAG

Split the DAG into  
*stages of tasks*

Submit each stage and  
its tasks as ready

Launch tasks via Master

Retry failed and strag-  
gler tasks

Execute tasks

Store and serve blocks

# OUTLINE

---

- ▶ Intro
  - ▶ Standalone-view RDD concepts
    - ▶ Basic RDD
    - ▶ Key-Pair RDD
  - ▶ Distributed-view RDD concepts
    - ▶ Lazy evaluation
    - ▶ Cache
    - ▶ Partition/Shuffle
    - ▶ Broadcast
    - ▶ Accumulator
  - ▶ Conclusion
- 
- ```
graph TD; A[Operators] --- B[Basic RDD]; A --- C[Key-Pair RDD]; D[Scheduler] --- E[Lazy evaluation]; F[Block manager] --- G[Cache]; F --- H[Partition/Shuffle]; F --- I[Broadcast]; J[Networking] --- K[Accumulator]; L[Broadcast] --- M[Accumulator];
```

- ▶ Def
 - ▶ By default, RDDs are recomputed each time you run an action on them

Example 3-39. Double execution in Scala

```
val result = input.map(x => x*x)
println(result.count())
println(result.collect().mkString(","))
```

Example 3-40. persist() in Scala

```
val result = input.map(x => x * x)
result.persist(StorageLevel.DISK_ONLY)
println(result.count())
println(result.collect().mkString(","))
```

CACHE

▶ Cache options

Table 3-6. Persistence levels from org.apache.spark.storage.StorageLevel and pyspark.StorageLevel; if desired we can replicate the data on two machines by adding _2 to the end of the storage level

Level	Space used	CPU time	In memory	On disk	Comments
MEMORY_ONLY	High	Low	Y	N	
MEMORY_ONLY_SER	Low	High	Y	N	
MEMORY_AND_DISK	High	Medium	Some	Some	Spills to disk if there is too much data to fit in memory.
MEMORY_AND_DISK_SER	Low	High	Some	Some	Spills to disk if there is too much data to fit in memory. Stores serialized representation in memory.
DISK_ONLY	Low	High	N	Y	

CACHE

- ▶ Why
 - ▶ Avoid network IO and recomputation
- ▶ How
 - ▶ Implemented with Least Recently Used (LRU) cache policy
 - ▶ Users do not have to worry about caching too much data

OUTLINE

- ▶ Intro
 - ▶ Standalone-view RDD concepts
 - ▶ Basic RDD
 - ▶ Key-Pair RDD
 - ▶ Distributed-view RDD concepts
 - ▶ Lazy evaluation
 - ▶ Cache
 - ▶ **Shuffle/Partition**
 - ▶ Broadcast
 - ▶ Accumulator
 - ▶ Conclusion
-
- ```
graph TD; Operators[Operators] --- BasicRDD[Basic RDD]; Operators --- KeyPairRDD[Key-Pair RDD]; Scheduler[Scheduler] --- LazyEvaluation[Lazy evaluation]; Scheduler --- Cache[Cache]; BlockManager[Block manager] --- ShufflePartition[Shuffle/Partition]; BlockManager --- Broadcast[Broadcast]; BlockManager --- Accumulator[Accumulator]; Networking[Networking];
```

# OUTLINE

---

- ▶ Intro
  - ▶ Standalone-view RDD concepts
    - ▶ Basic RDD
    - ▶ Key-Pair RDD
  - ▶ Distributed-view RDD concepts
    - ▶ Lazy evaluation
    - ▶ Cache
    - ▶ **Shuffle/Partition**
    - ▶ Broadcast
    - ▶ Accumulator
  - ▶ Conclusion
- 
- ```
graph TD; Operators[Operators] --- BasicRDD[Basic RDD]; Operators --- KeyPairRDD[Key-Pair RDD]; Scheduler[Scheduler] --- LazyEvaluation[Lazy evaluation]; Scheduler --- Cache[Cache]; BlockManager[Block manager] --- ShufflePartition[Shuffle/Partition]; BlockManager --- Broadcast[Broadcast]; BlockManager --- Accumulator[Accumulator]; Networking[Networking];
```

SHUFFLE

Let's start with an example. Given:

```
case class CFFPurchase(customerId: Int, destination: String, price: Double)
```

Assume we have an RDD of the purchases that users of the CFF mobile app have made in the past month.

```
val purchasesRdd: RDD[CFFPurchase] = sc.textFile(...)
```

Goal: calculate how many trips, and how much money was spent by each individual customer over the course of the month.

SHUFFLE

Goal: calculate how many trips, and how much money was spent by each individual customer over the course of the month.

```
val purchasesRdd: RDD[CFFPurchase] = sc.textFile(...)

// Returns: Array[(Int, (Int, Double))]

val purchasesPerMonth =
  purchasesRdd.map(p => (p.customerId, p.price)) // Pair RDD
    .groupByKey() // groupByKey returns RDD[K, Iterable[V]]
    .map(p => (p._1, (p._2.size, p._2.sum)))
    .collect()
```

SHUFFLE

Let's start with an example dataset:

```
val purchases = List(CFFPurchase(100, "Geneva", 22.25),  
                     CFFPurchase(300, "Zurich", 42.10),  
                     CFFPurchase(100, "Fribourg", 12.40),  
                     CFFPurchase(200, "St. Gallen", 8.20),  
                     CFFPurchase(100, "Lucerne", 31.60),  
                     CFFPurchase(300, "Basel", 16.20))
```

What might the cluster look like with this data distributed over it?

SHUFFLE

What might the cluster look like with this data distributed over it?

Starting with purchasesRdd:

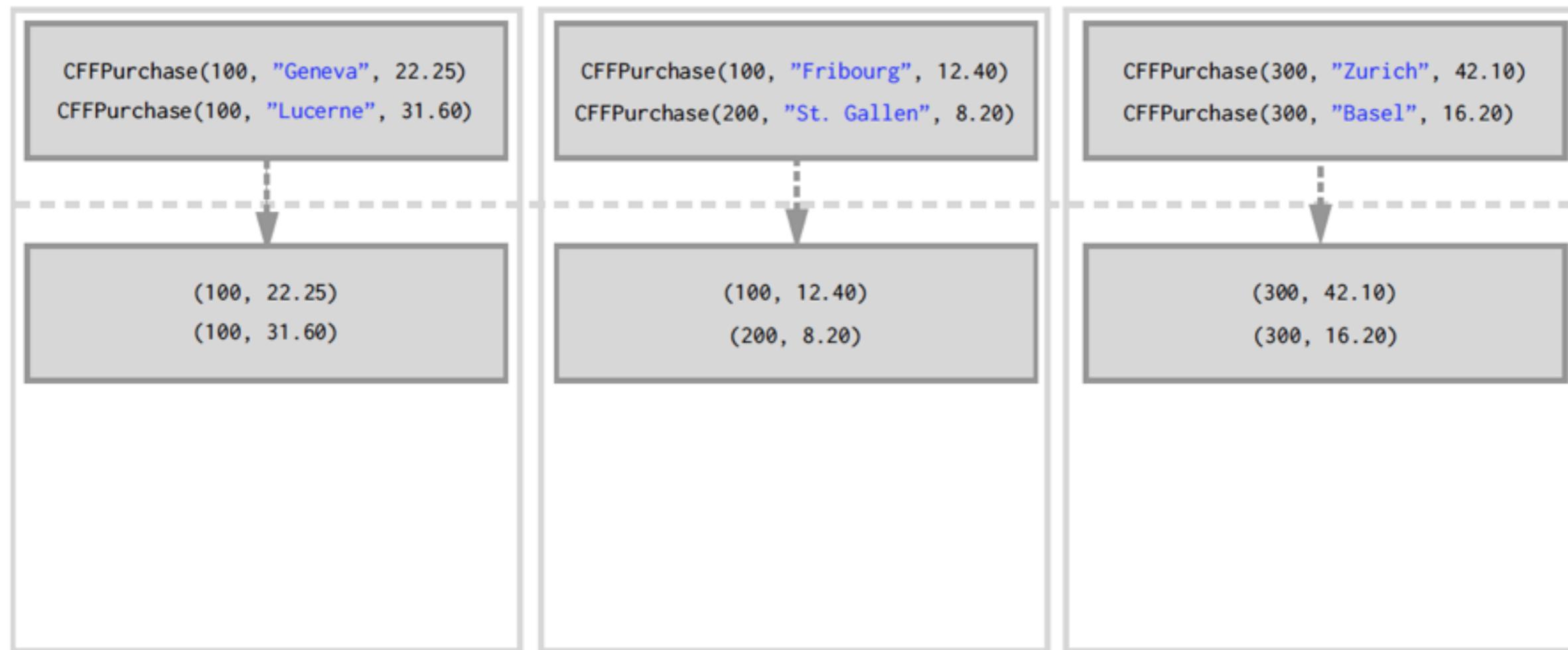
```
CFFPurchase(100, "Geneva", 22.25)  
CFFPurchase(100, "Lucerne", 31.60)
```

```
CFFPurchase(100, "Fribourg", 12.40)  
CFFPurchase(200, "St. Gallen", 8.20)
```

```
CFFPurchase(300, "Zurich", 42.10)  
CFFPurchase(300, "Basel", 16.20)
```

SHUFFLE

What might this look like on the cluster?



SHUFFLE

Goal: calculate how many trips, and how much money was spent by each individual customer over the course of the month.

```
val purchasesRdd: RDD[CFFPurchase] = sc.textFile(...)
```

```
val purchasesPerMonth =  
    purchasesRdd.map(p => (p.customerId, p.price)) // Pair RDD  
        .groupByKey() // groupByKey returns RDD[K, Iterable[V]]
```

Note: groupByKey results in one key-value pair per key. And this single key-value pair cannot span across multiple worker nodes.

SHUFFLE

```
CFFPurchase(100, "Geneva", 22.25)  
CFFPurchase(100, "Lucerne", 31.60)
```

```
CFFPurchase(100, "Fribourg", 12.40)  
CFFPurchase(200, "St. Gallen", 8.20)
```

```
CFFPurchase(300, "Zurich", 42.10)  
CFFPurchase(300, "Basel", 16.20)
```

```
(100, 22.25)  
(100, 31.60)
```

```
(100, 12.40)  
(200, 8.20)
```

```
(300, 42.10)  
(300, 16.20)
```

SHUFFLE
"Shuffles" data
across network

```
(100, [22.25, 12.40, 31.60])
```

```
(200, [8.20])
```

```
(300, [42.10, 16.20])
```

map

groupByKey

REMINDER: LATENCY MATTERS (HUMANIZED)

Shared Memory

Seconds

L1 cache reference.....0.5s
L2 cache reference.....7s
Mutex lock/unlock.....25s

Minutes

Main memory reference.....1m 40s

Distributed

Days

Roundtrip within
same datacenter.....5.8 days

Years

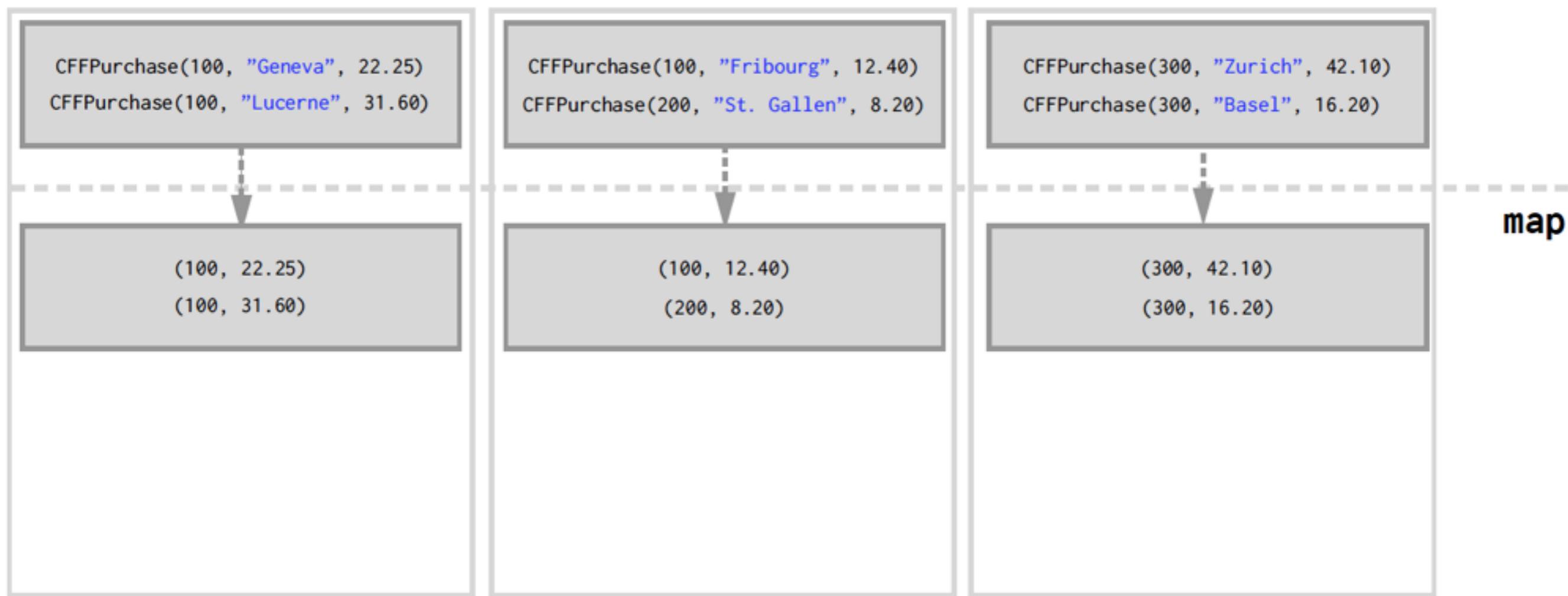
Send packet
CA->Netherlands->CA....4.8 years

We don't want to be sending all of our data over the network if it's not absolutely required. Too much network communication kills performance.

SHUFFLE

- ▶ Can we do a better job

Perhaps we don't need to send all pairs over the network.



Perhaps we can reduce before we shuffle. This could greatly reduce the amount of data we have to send over the network.

SHUFFLE

We can use `reduceByKey`.

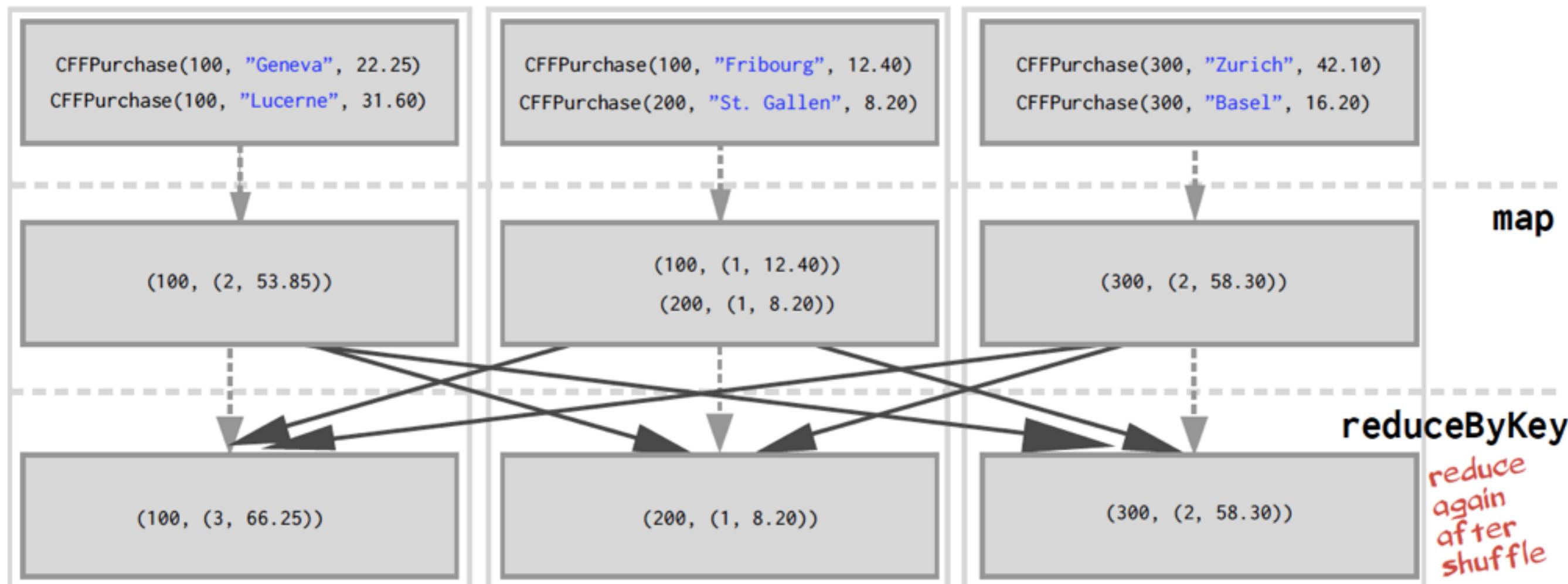
Conceptually, `reduceByKey` can be thought of as a combination of first doing `groupByKey` and then `reduce-ing` on all the values grouped per key. It's more efficient though, than using each separately. We'll see how in the following example.

Signature:

```
def reduceByKey(func: (V, V) => V): RDD[(K, V)]
```

SHUFFLE

What might this look like on the cluster?



What are the benefits of this approach?

By reducing the dataset first, the amount of data sent over the network during the shuffle is greatly reduced.

This can result in non-trivial gains in performance!

Let's benchmark on a real cluster.

SHUFFLE

```
> val purchasesPerMonthSlowLarge = purchasesRddLarge.map(p => (p.customerId, p.price))
    .groupByKey()
    .map(p => (p._1, (p._2.size, p._2.sum)))
    .count()

purchasesPerMonthSlowLarge: Long = 100000
Command took 15.48s
```



```
> val purchasesPerMonthFastLarge = purchasesRddLarge.map(p => (p.customerId, (1, p.price)))
    .reduceByKey((v1, v2) => (v1._1 + v2._1, v1._2 + v2._2))
    .count()

purchasesPerMonthFastLarge: Long = 100000
Command took 4.65s
```

Full example with 20 million element RDD can be found in the lecture2-apr2 notebook on our Databricks Cloud installation.

SHUFFLE

Recall our example using `groupByKey`:

```
val purchasesPerCust =  
    purchasesRdd.map(p => (p.customerId, p.price)) // Pair RDD  
        .groupByKey()
```

Grouping all values of key-value pairs with the same key requires collecting all key-value pairs with the same key on the same machine.

But how does Spark know which key to put on which machine?

- ▶ By default, Spark uses *hash partitioning* to determine which key-value pair should be sent to which machine.

OUTLINE

- ▶ Intro
 - ▶ Standalone-view RDD concepts
 - ▶ Basic RDD
 - ▶ Key-Pair RDD
 - ▶ Distributed-view RDD concepts
 - ▶ Lazy evaluation
 - ▶ Cache
 - ▶ Shuffle/**Partition**
 - ▶ Broadcast
 - ▶ Accumulator
 - ▶ Conclusion
-
- ```
graph TD; A[Operators] --- B[Basic RDD]; A --- C[Key-Pair RDD]; D[Scheduler] --- E[Lazy evaluation]; F[Block manager] --- G[Cache]; H[Networking] --- I[Shuffle/Partition]; J[Broadcast] --- K[Broadcast]; L[Accumulator] --- M[Accumulator]
```

## PARTITION

---

The data within an RDD is split into several *partitions*.

### Properties of partitions:

- ▶ Partitions never span multiple machines, i.e., tuples in the same partition are guaranteed to be on the same machine.
- ▶ Each machine in the cluster contains one or more partitions.
- ▶ The number of partitions to use is configurable. By default, it equals the *total number of cores on all executor nodes*.

### Two kinds of partitioning available in Spark:

- ▶ Hash partitioning
- ▶ Range partitioning

**Customizing a partitioning is only possible on Pair RDDs.**

## HASH PARTITION

---

Back to our example. Given a Pair RDD that should be grouped:

```
val purchasesPerCust =
 purchasesRdd.map(p => (p.customerId, p.price)) // Pair RDD
 .groupByKey()
```

## HASH PARTITION

---

Back to our example. Given a Pair RDD that should be grouped:

```
val purchasesPerCust =
 purchasesRdd.map(p => (p.customerId, p.price)) // Pair RDD
 .groupByKey()
```

groupByKey first computes per tuple  $(k, v)$  its partition  $p$ :

```
p = k.hashCode() % numPartitions
```

Then, all tuples in the same partition  $p$  are sent to the machine hosting  $p$ .

**Intuition: hash partitioning attempts to spread data evenly across partitions *based on the key*.**

## RANGE PARTITION

---

Consider a Pair RDD, with keys [8, 96, 240, 400, 401, 800], and a desired number of partitions of 4.

## RANGE PARTITION

---

Consider a Pair RDD, with keys [8, 96, 240, 400, 401, 800], and a desired number of partitions of 4.

Furthermore, suppose that hashCode() is the identity ( $n.hashCode() == n$ ).

## RANGE PARTITION

---

Consider a Pair RDD, with keys [8, 96, 240, 400, 401, 800], and a desired number of partitions of 4.

Furthermore, suppose that hashCode() is the identity (`n.hashCode() == n`).

In this case, hash partitioning distributes the keys as follows among the partitions:

- ▶ partition 0: [8, 96, 240, 400, 800]
- ▶ partition 1: [401]
- ▶ partition 2: []
- ▶ partition 3: []

The result is a very unbalanced distribution which hurts performance.

## RANGE PARTITION

---

Pair RDDs may contain keys that have an *ordering* defined.

- ▶ Examples: Int, Char, String, ...

For such RDDs, *range partitioning* may be more efficient.

Using a range partitioner, keys are partitioned according to:

1. an *ordering* for keys
2. a set of *sorted ranges* of keys

*Property:* tuples with keys in the same range appear on the same machine.

## RANGE PARTITION

---

Using *range partitioning* the distribution can be improved significantly:

- ▶ Assumptions: (a) keys non-negative, (b) 800 is biggest key in the RDD.
- ▶ Set of ranges: [1, 200], [201, 400], [401, 600], [601, 800]

## RANGE PARTITION

---

Using *range partitioning* the distribution can be improved significantly:

- ▶ Assumptions: (a) keys non-negative, (b) 800 is biggest key in the RDD.
- ▶ Set of ranges: [1, 200], [201, 400], [401, 600], [601, 800]

In this case, range partitioning distributes the keys as follows among the partitions:

- ▶ partition 0: [8, 96]
- ▶ partition 1: [240, 400]
- ▶ partition 2: [401]
- ▶ partition 3: [800]

The resulting partitioning is much more balanced.

### **How do we set a partitioning for our data?**

There are two ways to create RDDs with specific partitionings:

1. Call `partitionBy` on an RDD, providing an explicit Partitioner.
2. Using transformations that return RDDs with specific partitioners.

## PARTITIONING DATA:

---

Invoking `partitionBy` creates an RDD with a specified partitioner.

Example:

```
val pairs = purchasesRdd.map(p => (p.customerId, p.price))
val tunedPartitioner = new RangePartitioner(8, pairs)
val partitioned = pairs.partitionBy(tunedPartitioner).persist()
```

Creating a RangePartitioner requires:

1. Specifying the desired number of partitions.
2. Providing a Pair RDD with *ordered keys*. This RDD is *sampled* to create a suitable set of *sorted ranges*.

**Important: the result of `partitionBy` should be persisted. Otherwise, the partitioning is repeatedly applied (involving shuffling!) each time the partitioned RDD is used.**

## PARTITIONING DATA:

---

### **Partitioner from parent RDD:**

Pair RDDs that are the result of a transformation on a *partitioned* Pair RDD typically is configured to use the hash partitioner that was used to construct it.

### **Automatically-set partitioners:**

Some operations on RDDs automatically result in an RDD with a known partitioner – for when it makes sense.

For example, by default, when using `sortByKey`, a `RangePartitioner` is used. Further, the default partitioner when using `groupByKey`, is a `HashPartitioner`, as we saw earlier.

## PARTITIONING DATA:

---

Operations on Pair RDDs that hold to (and propagate) a partitioner:

- ▶ cogroup
- ▶ groupWith
- ▶ join
- ▶ leftOuterJoin
- ▶ rightOuterJoin
- ▶ groupByKey
- ▶ reduceByKey
- ▶ foldByKey
- ▶ combineByKey
- ▶ partitionBy
- ▶ sort
- ▶ mapValues (if parent has a partitioner)
- ▶ flatMapValues (if parent has a partitioner)
- ▶ filter (if parent has a partitioner)

**All other operations will produce a result without a partitioner.**

## PARTITIONING DATA:

---

...All other operations will produce a result without a partitioner.

### Why?

Consider the `map` transformation. Given have a hash partitioned Pair RDD, why would it make sense for `map` to lose the partitioner in its result RDD?

Because it's possible for `map` to change the key . *E.g.,:*

```
rdd.map((k: String, v: Int) => ("doh!", v))
```

In this case, if the `map` transformation preserved the partitioner in the result RDD, it no longer make sense, as now the keys are all different.

**Hence `mapValues`. It enables us to still do `map` transformations without changing the keys, thereby preserving the partitioner.**

## OPTIMIZING PREVIOUS EXAMPLE USING RANGE PARTITIONING

---

```
> val purchasesPerMonthSlowLarge = purchasesRddLarge.map(p => (p.customerId, p.price))
 .groupByKey()
 .map(p => (p._1, (p._2.size, p._2.sum)))
 .count()
```

purchasesPerMonthSlowLarge: Long = 100000

Command took 15.48s

```
> val purchasesPerMonthFastLarge = purchasesRddLarge.map(p => (p.customerId, (1, p.price)))
 .reduceByKey((v1, v2) => (v1._1 + v2._1, v1._2 + v2._2))
 .count()
```

purchasesPerMonthFastLarge: Long = 100000

Command took 4.65s

## On the range partitioned data:

```
> val purchasesPerMonthFasterLarge = partitioned.map(x => x)
 .reduceByKey((v1, v2) => (v1._1 + v2._1, v1._2 + v2._2))
 .count()
```

purchasesPerMonthFasterLarge: Long = 100000

Command took 1.79s

almost a 9x speedup over  
purchasePerMonthSlowLarge!

## HOW DO I KNOW A SHUFFLE WILL OCCUR?

---

**Rule of thumb:** a shuffle *can* occur when the resulting RDD depends on other elements from the same RDD or another RDD.

*Note: sometimes one can be clever and avoid much or all network communication while still using an operation like join via smart partitioning*

## HOW DO I KNOW A SHUFFLE WILL OCCUR?

---

You can also figure out whether a shuffle has been planned/executed via:

1. The return type of certain transformations, e.g.,

```
org.apache.spark.rdd.RDD[(String, Int)] = ShuffledRDD[366]
```

2. Using function `toDebugString` to see its execution plan:

```
partitioned.reduceByKey((v1, v2) => (v1._1 + v2._1, v1._2 + v2._2))
 .toDebugString
res9: String =
(8) MapPartitionsRDD[622] at reduceByKey at <console>:49 []
 | ShuffledRDD[615] at partitionBy at <console>:48 []
 | CachedPartitions: 8; MemorySize: 1754.8 MB; DiskSize: 0.0 B
```

## HOW DO I KNOW A SHUFFLE WILL OCCUR?

---

- ▶ Operations that might cause a shuffle

- ▶ cogroup
- ▶ groupWith
- ▶ join
- ▶ leftOuterJoin
- ▶ rightOuterJoin
- ▶ groupByKey
- ▶ reduceByKey
- ▶ combineByKey
- ▶ distinct
- ▶ intersection
- ▶ repartition
- ▶ coalesce

# OUTLINE

---

- ▶ Intro
  - ▶ Standalone-view RDD concepts
    - ▶ Basic RDD
    - ▶ Key-Pair RDD
  - ▶ Distributed-view RDD concepts
    - ▶ Lazy evaluation
    - ▶ Cache
    - ▶ Shuffle/Partition
    - ▶ Broadcast
    - ▶ Accumulator
  - ▶ Conclusion
- 
- ```
graph TD; Operators[Operators] --- BasicRDD[Basic RDD]; Operators --- KeyPairRDD[Key-Pair RDD]; Scheduler[Scheduler] --- LazyEvaluation[Lazy evaluation]; Scheduler --- Cache[Cache]; BlockManager[Block manager] --- ShufflePartition[Shuffle/Partition]; Networking[Networking] --- Broadcast[Broadcast]; Broadcast --- Accumulator[Accumulator];
```

SHARED VARIABLES

- ▶ Start with an example

```
val points: RDD[Point] = // ...
var clusterCenters = new Array[Point](k)
```

```
val closestCenter = points.map {
  p => findClosest(clusterCenters, p)
}
...
```

How should
this variable
be shipped?

BROADCAST VARIABLES

- ▶ Motivation:
 - ▶ Normally, Spark closures, including variables they use, are sent separately within each task.
 - ▶ But in some cases, the same variable will be used in multiple parallel operations
 - ▶ By default, Spark task launching mechanism is optimized for small task sizes.
 - ▶ But in some cases, a large read-only variable needs to be shared across tasks, or across operations
- ▶ Examples: large lookup tables

Example 6-6. Country lookup in Python

```
# Look up the locations of the call signs on the
# RDD contactCounts. We load a list of call sign
# prefixes to country code to support this lookup.
signPrefixes = loadCallSignTable()

def processSignCount(sign_count, signPrefixes):
    country = lookupCountry(sign_count[0], signPrefixes)
    count = sign_count[1]
    return (country, count)

countryContactCounts = (contactCounts
    .map(processSignCount)
    .reduceByKey((lambda x, y: x + y)))
```



Example 6-7. Country lookup with Broadcast values in Python

```
# Look up the locations of the call signs on the
# RDD contactCounts. We load a list of call sign
# prefixes to country code to support this lookup.
signPrefixes = sc.broadcast(loadCallSignTable())

def processSignCount(sign_count, signPrefixes):
    country = lookupCountry(sign_count[0], signPrefixes.value)
    count = sign_count[1]
    return (country, count)

countryContactCounts = (contactCounts
    .map(processSignCount)
    .reduceByKey((lambda x, y: x + y)))
```

OUTLINE

- ▶ Intro
 - ▶ Standalone-view RDD concepts
 - ▶ Basic RDD
 - ▶ Key-Pair RDD
 - ▶ Distributed-view RDD concepts
 - ▶ Lazy evaluation
 - ▶ Cache
 - ▶ Shuffle/Partition
 - ▶ Broadcast
 - ▶ **Accumulator**
 - ▶ Conclusion
-
- The diagram illustrates the relationship between RDD concepts and system components. Brackets group concepts into categories:
- Operators**: Basic RDD, Key-Pair RDD
 - Scheduler**: Lazy evaluation
 - Block manager**: Cache
 - Networking**: Shuffle/Partition
 - Broadcast**: Broadcast
 - Accumulator**: Accumulator

BROADCAST VARIABLES

```
# Look up the locations of the call signs on the  
# RDD contactCounts. We load a list of call sign  
# prefixes to country code to support this lookup.  
signPrefixes = sc.broadcast(loadCallSignTable())
```

Create with
`sparkContext.broadcast(initVal)`

```
def processSignCount(sign_count, signPrefixes):  
    country = lookupCountry(sign_count[0], signPrefixes.value)  
    count = sign_count[1]  
    return (country, count)
```

Access with `.value` inside tasks
First task on each node to fetch val

```
countryContactCounts = (contactCounts  
    .map(processSignCount)  
    .reduceByKey((lambda x, y: x+ y)))
```

!!! Cannot modify after creation
if you try, change only on one node

```
countryContactCounts.saveAsTextFile(outputDir + "/countries.txt")
```

ACCUMULATORS

- ▶ Motivation:
 - ▶ Normally, Spark closures, including variables they use, are sent separately within each task.
 - ▶ But values inside worker nodes are not propagated back to the driver
 - ▶ Normally, we could get the entire RDD back with reduce(),
 - ▶ But sometimes, we just need a simpler way
 - ▶ Examples: counters

Example 6-2. Accumulator empty line count in Python

OUTLINE

- ▶ Intro
 - ▶ Standalone-view RDD concepts
 - ▶ Basic RDD
 - ▶ Key-Pair RDD
 - ▶ Distributed-view RDD concepts
 - ▶ Lazy evaluation
 - ▶ Cache
 - ▶ Shuffle/Partition
 - ▶ Broadcast
 - ▶ Accumulator
 - ▶ Conclusion
-
- ```
graph TD; A[Operators] --- B[Basic RDD]; A --- C[Key-Pair RDD]; D[Scheduler] --- E[Lazy evaluation]; F[Block manager] --- G[Cache]; F --- H[Shuffle/Partition]; F --- I[Broadcast]; J[Networking] --- K[Accumulator]; L[Broadcast] --- M[Accumulator];
```

## CONCLUSION

---

- ▶ Spark core
  - ▶ basic Spark RDDs
  - ▶ Things to take care when programming in clusters
- ▶ What's next?
  - ▶ Spark SQL/Hive

## REFERENCES

---

- ▶ Udemy hands-on course
  - ▶ <https://www.udemy.com/taming-big-data-with-apache-spark-hands-on/learn/v4/overview>
- ▶ Several slides:
  - ▶ <http://heather.miller.am/teaching/cs212/slides>
  - ▶ <http://ampcamp.berkeley.edu/wp-content/uploads/2012/06/matei-zaharia-amp-camp-2012-advanced-spark.pdf>
  - ▶ <https://databricks-training.s3.amazonaws.com/slides/advanced-spark-training.pdf>
- ▶ <<Learning Spark>>