

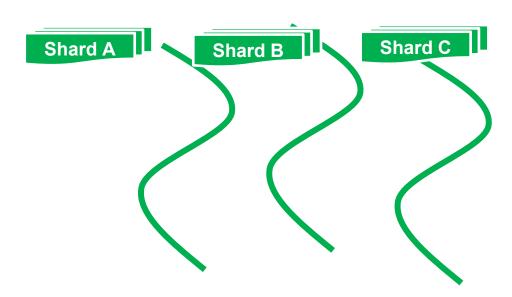
Apache Spark and RDDs

Ken Birman (with help) CS4414/5416 Lecture 23

Start with a close look at the MapReduce pattern: **Sharded data set**

Leader Worker threads





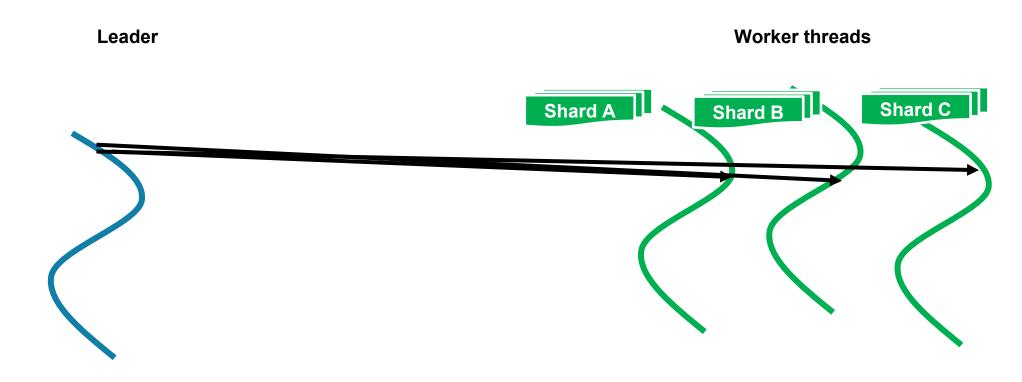
MapReduce: Map step

The leader **maps** some task over the n workers. This can be done in any way that makes sense for the application.

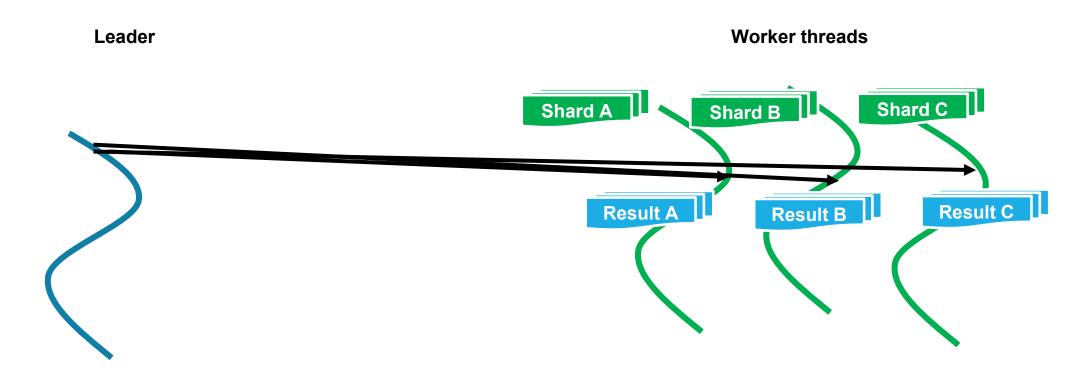
Each worker performs its share of the work by applying the requested function to the data in its shard.

When finished, each worker will have a list of new (key,value) pairs as its share of the result.

MapReduce pattern: Sharded data set



MapReduce pattern: Map (first step)



MapReduce: Shuffle exchange

Each worker breaks its key-value result set into n parts by applying the sharding rule to the keys.

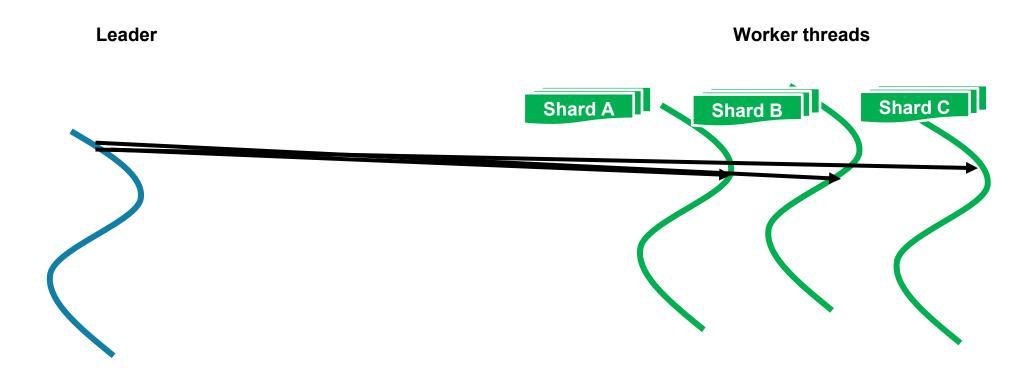
- Now it has one subset (perhaps empty) for each other worker.
- It hands that subset to corresponding worker.

Every worker waits until it has its one message from each worker.

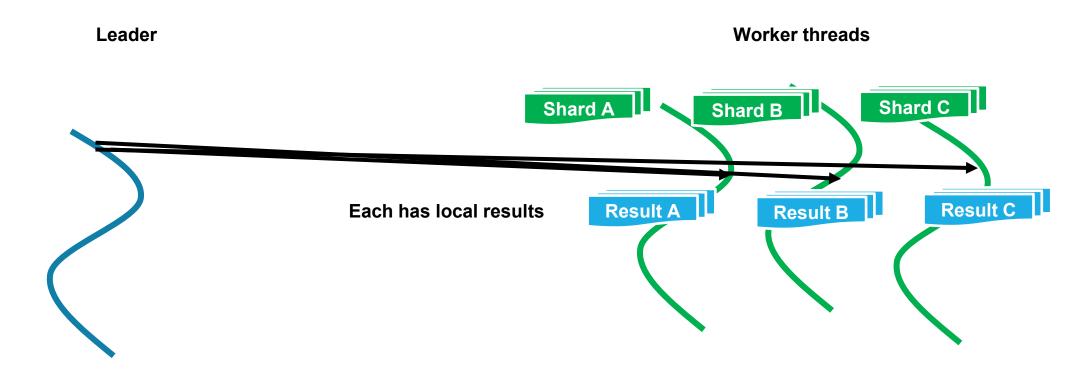
Now it can merge the n "pieces", sort them, group by key

It now has a list of (key, {set-of-values}) tuples. It calls reduce one by one on these.

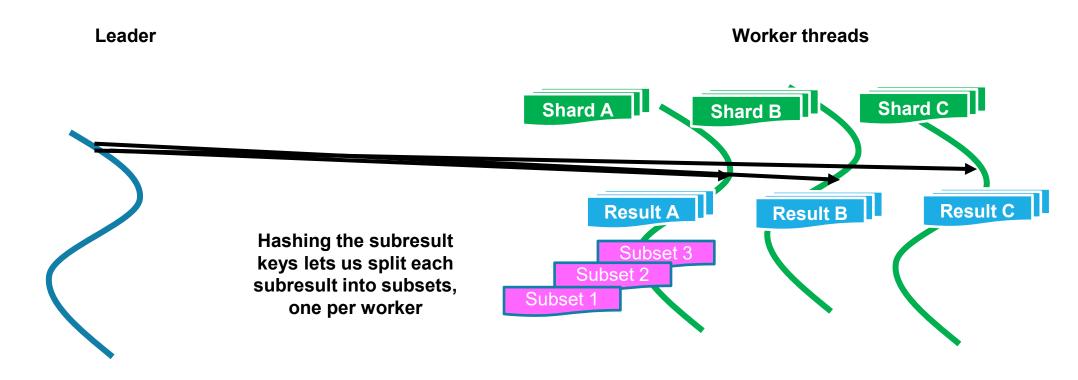
MapReduce pattern: Map (first step)



MapReduce pattern: Map (first step)

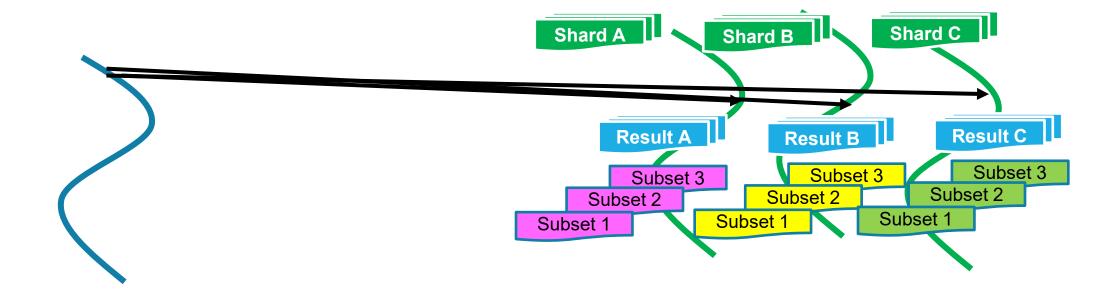


MapReduce pattern: Map (hash to split data)



MapReduce pattern: Parallel behavior, including the shuffle exchange

Leader Worker threads



MapReduce pattern: Each worker sorts its distinct share of the subresults.

to the worker that owns it

Leader

Shard C Shard A Shard B Result C Result A Result B After forwarding each subset to the worker that will own it, Subset 2 Subset 3 Subset 1 everyone sorts. Now they each Subset 2 Supset 3 Supset 1 have DIFFERENT data: each has Supset 3 Supset 2 OUDSEL a contribution from every other worker, but each slice is unique

Worker threads

MapReduce pattern: Each worker applies reduce on its sorted portion of data

of the AllReduce results

Leader

Shard C Shard A Shard B Result C Result A Result B The reducing function runs on the sorted data sets and each worker ends up with its own private share of the MapReduce results Reduced results B Reduced results C Reduced results A This is in contrast with AllReduce where everyone shares **EVERYTHING**, so they end up with identical "complete" copies

Worker threads

Example: Word Count

The use case scenario: Start with standard WC for one file.

We have a large file of documents (the input elements)

Documents are words separated by whitespace.

Count the number of times each distinct word appears in the file.

... with MapReduce we can extend this concept to huge numbers of files.

Example: Word Count

Why Do We Care About Counting Words?

- > NLP systems train on n-grams: counts of n-word sequences.
- Word or n-gram count is challenging over massive amounts of data
 - Using a single compute node would be too time-consuming
 - Using distributed nodes requires moving data
 - Number of unique words can easily exceed available memory -- would need to store to disk
- Many common tasks are very similar to word count, e.g., log file analysis where we might look for the storage devices with the highest error rates

Word Count Using MapReduce

Word Count Using MapReduce

Input

the cat sat on the mat the aardvark sat on the sofa Map & Reduce

Result

aardvark 1

cat 1

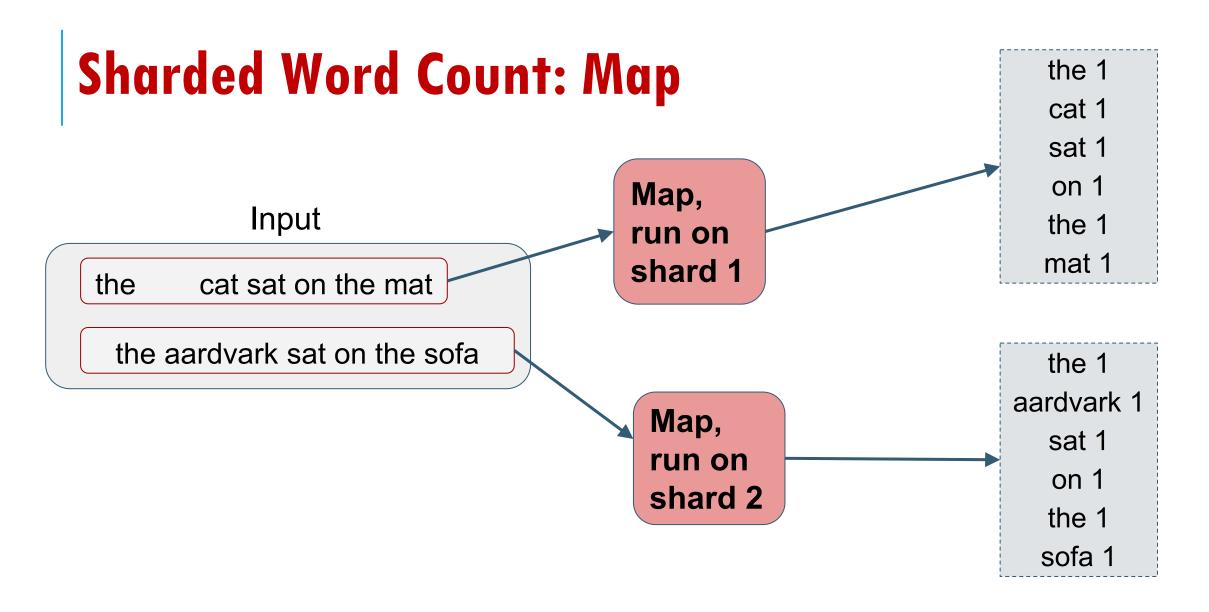
mat 1

on 2

sat 2

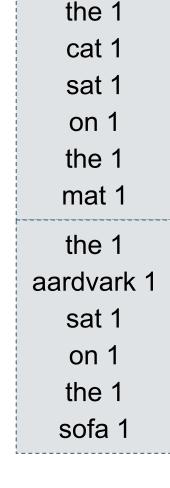
sofa 1

the 4



Shuffle & Sort

Mapper Output



Intermediate Data

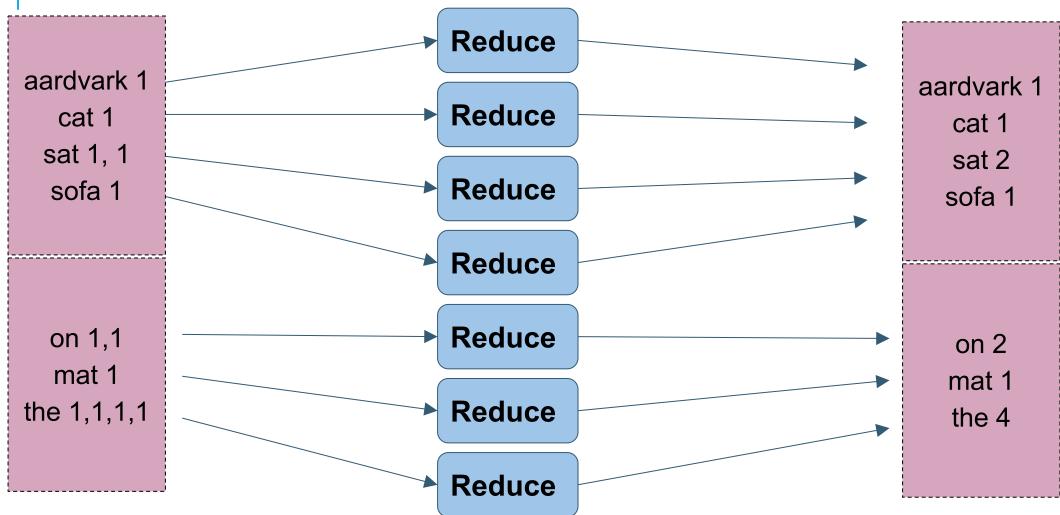


Keys that mapped to shard 1 are still on shard 1. The sort was internal to shard 1

on 1,1 mat 1 the 1,1,1,1

Keys that mapped to shard 2

Word Count: Reducer



Notice that...

Data stays sharded at all times. At the start document names determined which document was on which shard. Now, after the shuffle exchange, the words themselves are the keys, and determine which shard owns that word and count

Keys are sorted and grouped shard-by-shard.

Reduce runs on (key, $\{v_1, \dots, v_k\}$) and outputs (key, reduced-value), once per key

Output is never collected to one place: **We never merge and sort the full data set.** AllReduce does that and is easier to understand, but the results can be too large to hold on a single machine, which forces use of MapReduce for big data

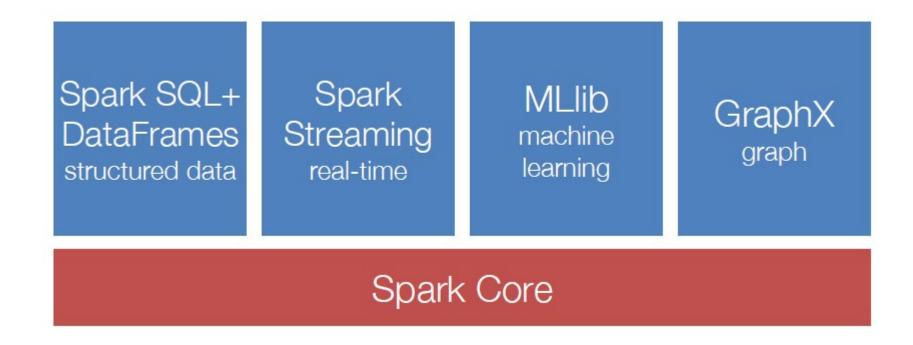
Spark Project (became launch point for the DataBricks company in San Francisco)

Undertaken at UC Berkeley

Goal was to standardize use of MapReduce and to speed up, focusing on efficient caching of reusable results

Part of the Berkeley "View from the clouds" vision for cloud computing research, authored by Ion Stoica

Spark Ecosystem: A Unified Pipeline



Note: Spark is <u>not</u> designed for IoT real-time. The streaming layer is used for continuous input streams like financial data from stock markets, where events occur steadily and must be processed as they occur. But there is no sense of direct I/O from sensors/actuators. For IoT use cases, Spark would not be suitable.

Key ideas

In Hadoop, prior to Spark, each developer tended to invent their own mapping from data mining goal to MapReduce or AllReduce

With Spark, serious effort was invested to **standardize** around the idea that in data mining, parallel code often runs for many "cycles" or "iterations" in which a lot of reuse of information occurs. So caching can be a big win.

Spark centers on Resilient Distributed Dataset, RDDs, that capture the information being reused and compile to MapReduce

How this works

You express your application as a data flow graph of RDDs.

The graph is only evaluated as needed, and they only compute the RDDs actually needed for the output you have requested.

Then Spark can be told to cache the reuseable information either in memory, in SSD storage or even on disk, based on when it will be needed again, how big it is, and how costly it would be to recreate.

You write the RDD logic and control all of this via hints

Spark Basics

There are two ways to manipulate data in Spark

- Spark Shell:
 - Interactive for learning or data exploration
 - Python or Scala
- Spark Applications
 - For large scale data processing
 - Python, Scala, or Java

Spark Shell

The Spark Shell provides interactive data exploration (REPL)

Python Shell: pyspark

Scala Shell: spark-shell

REPL: Repeat/Evaluate/Print Loop

Spark Fundamentals

Example of an application:

```
val sc = new SparkContext("spark://...", "MyJob", home,
    jars)

val file = sc.textFile("hdfs://...") // This is an RDD

val errors = file.filter(_.contains("ERROR")) // This is an RDD

errors.cache()

errors.count() // This is an action
```

- Spark Context
- Resilient DistributedData
- Transformations
- Actions

Spark Context (1)

- Every Spark application requires a spark context: the main entry point to the Spark API
- Spark Shell provides a preconfigured Spark Context called "sc"

```
Using Python version 2.7.8 (default, Aug 27 2015 05:23:36)
SparkContext available as sc, HiveContext available as sqlCtx.

>>> sc.appName
u'PySparkShell'

Spark context available as sc.

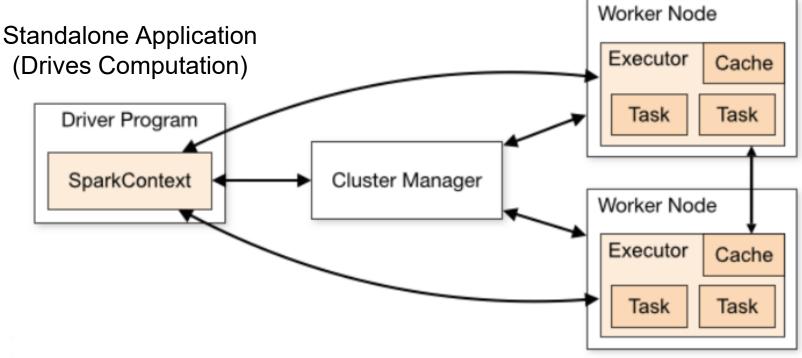
SQL context available as sqlContext.

Scala

scala> sc.appName
res0: String = Spark shell
```

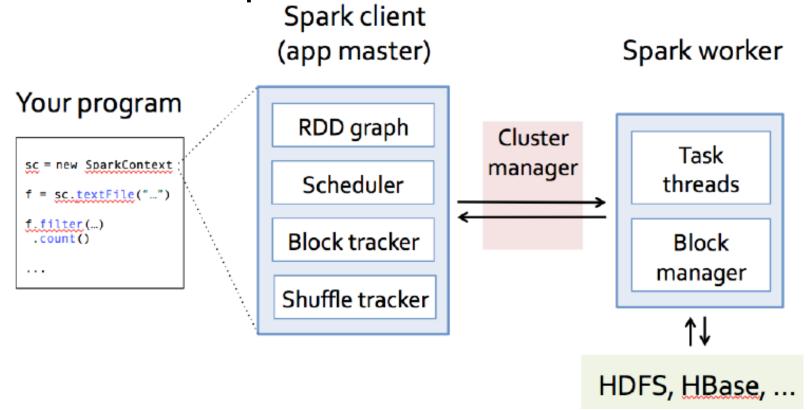
Spark Context (2)

- Standalone applications → Driver code → Spark Context
- Spark Context holds configuration information and represents connection to a Spark cluster



Spark Context (3)

Spark context works as a client and represents connection to a Spark cluster



Spark Fundamentals

Example of an application:

```
val sc = new SparkContext("spark://...", "MyJob", home,
    jars)

val file = sc.textFile("hdfs://...") // This is an RDD

val errors = file.filter(_.contains("ERROR")) // This is
    an RDD

errors.cache()

errors.count() // This is an action
```

- Spark Context
- Resilient DistributedData
- Transformations
- Actions

Resilient Distributed Dataset (RDD)

The RDD (Resilient Distributed Dataset) is the fundamental unit of data in Spark: An *Immutable* collection of objects (or records, or elements) that can be operated on "in parallel" (spread across a cluster)

Resilient -- if data in memory is lost, it can be recreated

- Recover from node failures
- An RDD keeps its lineage information → it can be recreated from parent RDDs

Distributed -- processed across the cluster

Each RDD is composed of one or more partitions → (more partitions – more parallelism)

Dataset -- initial data can come from a file or be created

RDDs

Key Idea: Write applications in terms of transformations on distributed datasets. One RDD per transformation.

- Organize the RDDs into a DAG showing how data flows.
- RDD can be saved and reused or recomputed. Spark can save it to disk if the dataset does not fit in memory
- Built through parallel transformations (map, filter, group-by, join, etc). Automatically rebuilt on failure
- Controllable persistence (e.g. caching in RAM)

RDDs are designed to be "immutable"

- Create once, then reuse without changes. Spark knows lineage → can be recreated at any time → Fault-tolerance
- Avoids data inconsistency problems (no simultaneous updates) → Correctness
- Easily live in memory as on disk → Caching → Safe to share across processes/tasks → Improves performance
- Tradeoff: (Fault-tolerance & Correctness) vs (Disk Memory & CPU)

Creating a RDD

Three ways to create a RDD

- From a file or set of files
- From data in memory
- From another RDD

RDD COMPILES TO MAPREDUCE!

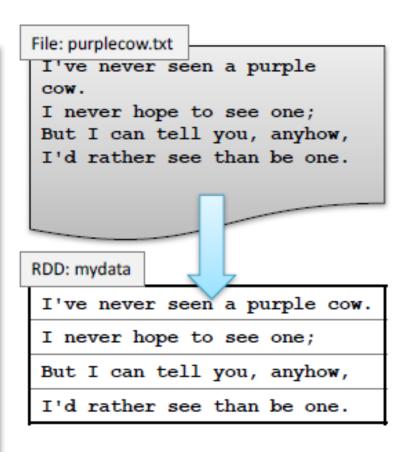
The RDD language is designed around primitives that all can compile to the MapReduce pattern

In effect, you are doing high-level MapReduce programming, and your code will automatically be parallelized.

By adding hints about caching, performance can be very high!

Example: A File-based RDD

```
> val mydata = sc.textFile("purplecow.txt")
15/01/29 06:20:37 INFO storage.MemoryStore:
  Block broadcast 0 stored as values to
  memory (estimated size 151.4 KB, free 296.8
> mydata.count()
15/01/29 06:27:37 INFO spark.SparkContext: Job
  finished: take at <stdin>:1, took
  0.160482078 s
4
```



Spark Fundamentals

Example of an application:

```
val sc = new SparkContext("spark://...", "MyJob", home,
    jars)

val file = sc.textFile("hdfs://...") // This is an RDD

val errors = file.filter(_.contains("ERROR")) // This is an RDD

errors.cache()

errors.count() // This is an action
```

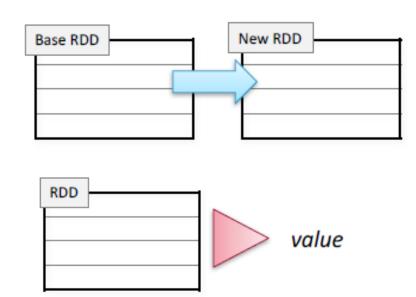
- Spark Context
- Resilient Distributed
 Data
- Transformations
- Actions

RDD Operations

Two types of operations

Transformations: Define a new RDD based on current RDD(s)

Actions: return values



```
val sc = new SparkContext("spark://...", "MyJob", home,
    jars)

val file = sc.textFile("hdfs://...") // This is an RDD

val errors = file.filter(_.contains("ERROR")) // This is an RDD

errors.cache()
errors.count() // This is an action
```

RDD Transformations

- Set of operations on a RDD that define how they should be transformed
- As in relational algebra, the application of a transformation to an RDD yields a new RDD (because RDD are immutable)
- Transformations are lazily evaluated, which allow for optimizations to take place before execution
- Examples: map(), filter(), groupByKey(), sortByKey(), etc.

Example: map and filter Transformations

```
I've never seen a purple cow.
                               I never hope to see one;
                               But I can tell you, anyhow,
                               I'd rather see than be one.
         map(lambda line: line.upper())
                                                    map(line => line.toUpperCase)
                               I'VE NEVER SEEN A PURPLE COW.
                               I NEVER HOPE TO SEE ONE;
                               BUT I CAN TELL YOU, ANYHOW,
                               I'D RATHER SEE THAN BE ONE.
filter(lambda line: line.startswith('I'))
                                                  filter(line => line.startsWith('I'))
                               I'VE NEVER SEEN A PURPLE COW.
                               I NEVER HOPE TO SEE ONE;
                               I'D RATHER SEE THAN BE ONE.
```

RDD Actions

- Apply transformation chains on RDDs, eventually performing some additional operations (e.g., counting)
- Some actions only store data to an external data source (e.g. HDFS), others fetch data from the RDD (and its transformation chain) upon which the action is applied, and convey it to the driver
- Some common actions
 - >count() return the number of elements
 - >take(n) return an array of the first n elements
 - >collect()- return an array of all elements
 - >saveAsTextFile(file) save to text file(s)

Graph of RDDs

- A collection of RDDs can be understood as a graph
- Nodes in the graph are the RDDs, which means the code but also the actual data object that could would create at runtime when executed on specific parameters + data. Reminder: Hadoop is a "read only" model, so we can "materialize" an RDD any time we like.
- Edges represent how data objects are accessed: RDD B might consume the object created by RDD A. This gives us a directed edge $A \rightarrow B$

Lazy Execution of RDDs (1)

Data in RDDs is not processed until an action is performed



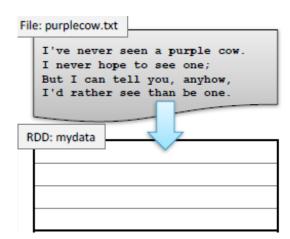
File: purplecow.txt

I've never seen a purple cow. I never hope to see one; But I can tell you, anyhow, I'd rather see than be one.

Lazy Execution of RDDs (2)

Data in RDDs is not processed until an action is performed

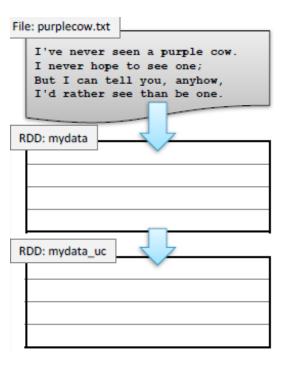
> val mydata = sc.textFile("purplecow.txt")



Lazy Execution of RDDs (3)

Data in RDDs is not processed until an action is performed

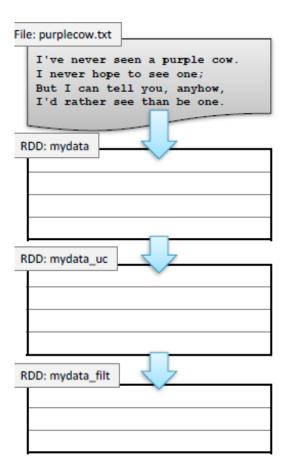
```
> val mydata = sc.textFile("purplecow.txt")
> val mydata_uc = mydata.map(line =>
    line.toUpperCase())
```



Lazy Execution of RDDs (4)

Data in RDDs is not processed until an action is performed

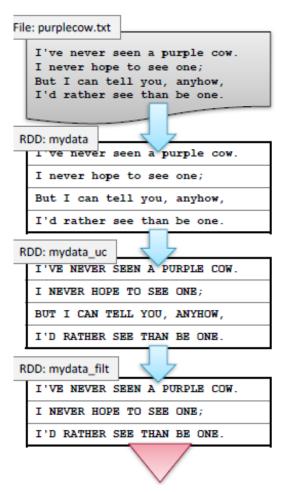
```
> val mydata = sc.textFile("purplecow.txt")
> val mydata_uc = mydata.map(line =>
    line.toUpperCase())
> val mydata_filt = mydata_uc.filter(line
    => line.startsWith("I"))
```



Lazy Execution of RDDs (5)

Data in RDDs is not processed until an action is performed

```
> val mydata = sc.textFile("purplecow.txt")
> val mydata_uc = mydata.map(line =>
    line.toUpperCase())
> val mydata_filt = mydata_uc.filter(line
    => line.startsWith("I"))
> mydata_filt.count()
3
```



Output Action "triggers" computation, pull model

Opportunities This Enables

- Automated compilation to MapReduce: This is the fundamental reason we use Spark. Instead of coding by hand we "script" parallel compute as a graph of RDDs
- On-demand optimization: Spark can behave like a compiler by first building a
 potentially complex RDD graph, but then trimming away unneeded computations that
 for today's purpose, won't be used.
- Caching for later reuse.
- Graph transformations: A significant amount of effort is going into this area. It is a lot like compiler-managed program transformation and aims at simplifying and speeding up the computation that will occur.
- Dynamic decisions about what to schedule and when. Concept: minimum adequate set of input objects: RDD can run if all its inputs are ready

Example: Mine error logs

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

```
lines = spark.textFile("hdfs://...")    HadoopRDD
errors = lines.filter(lambda s: s.startswith("ERROR"))    FilteredRDD
messages = errors.map(lambda s: s.split("\t")[2])
messages.cache()
messages.filter(lambda s: "foo" in s).count()
```

Result: full-text search of Wikipedia in 0.5 sec (vs 20 sec for on-disk data)

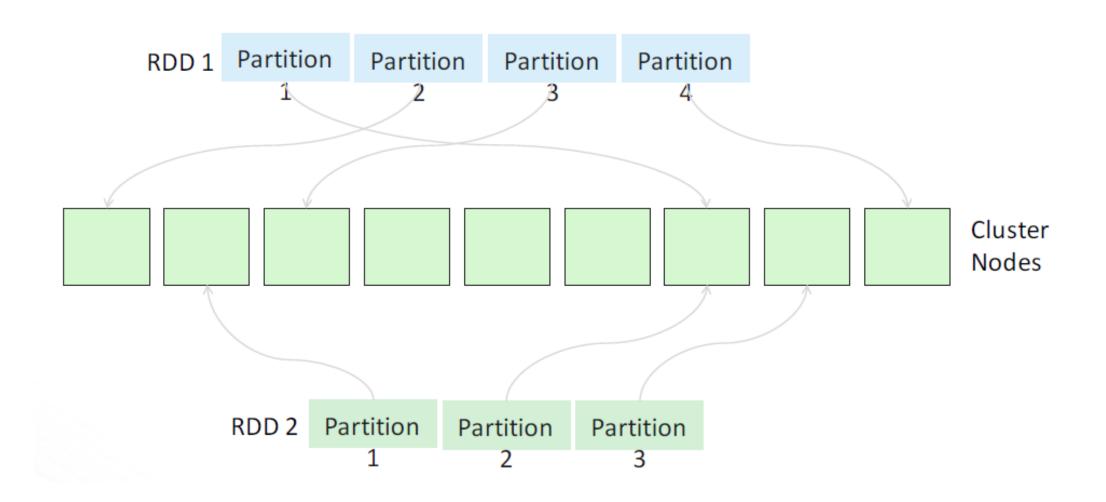
Key feature: Elastic parallelism

RDDs operations are designed to leverage embarrassing parallelism and are coded to view "how many workers?" as a runtime decision.

Spark will always spread the task over the full set of nodes it is allocated. Normally this is one worker per K shards, so that the full set of shards is mapped to the full set of workers. The resulting pattern is a highly concurrent execution that minimizes delays: a "partitioned computation".

If some component crashes or even is just slow, Spark simply kills that task and launches a substitute.

RDD and Partitions (Parallelism example)



RDD Graph: Data Set vs Partition Views

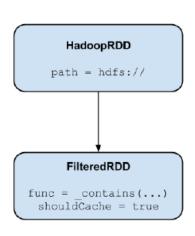
Much like in Hadoop MapReduce, each RDD is associated to (input) partitions

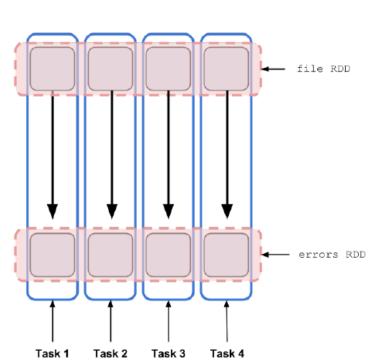
```
val sc = new SparkContext("spark://...", "MyJob", home,
    jars)

val file = sc.textFile("hdfs://...") // This is an RDD

val errors = file.filter(_.contains("ERROR")) // This is
    an RDD

errors.cache()
errors.count() // This is an action
```





Worker 1 Worker 2 Worker 3 Worker 4

RDDs: Data Locality

- Data Locality Principle
- Keep high-value RDDs precomputed, in cache or SDD
- Run tasks that need the specific RDD with those same inputs on the node where the cached copy resides.
- This can maximize in-memory computational performance.

Requires cooperation between your hints to Spark when you build the RDD, Spark runtime and optimization planner, and the underlying YARN resource manager.

RDDs -- Summary

RDD are partitioned, locality aware, distributed collections

RDD are immutable

RDD are data structures that:

- Either point to a direct data source (e.g. HDFS)
- Apply some transformations to its parent RDD(s) to generate new data elements

Computations on RDDs

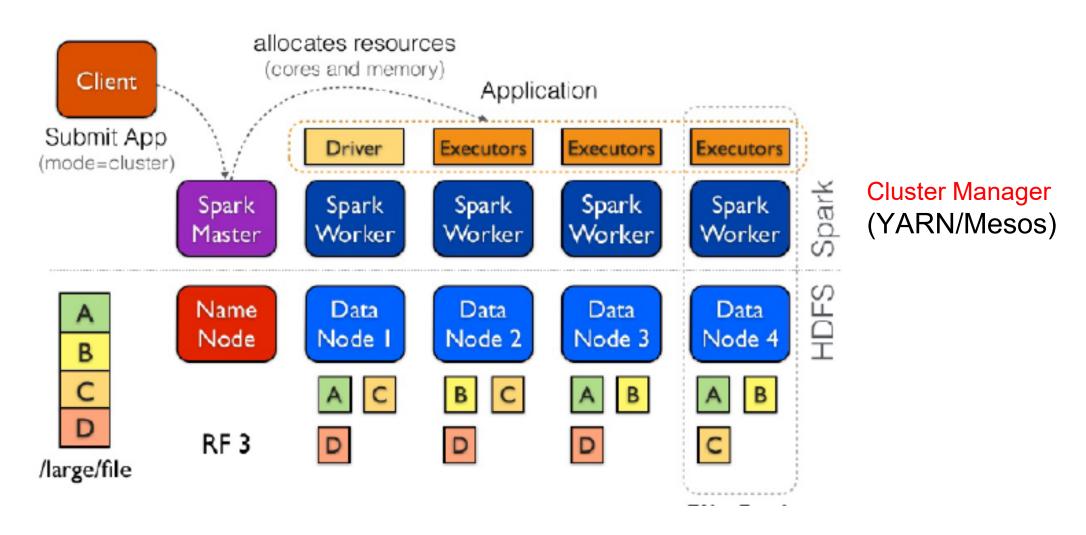
Represented by lazily evaluated lineage DAGs composed by chained RDDs

Lifetime of a Job in Spark

Build the operator DAG

RDD Objects DAG Scheduler Task Scheduler Worker Cluster Threads manager Block manager rdd1.join(rdd2) Launch tasks via Master Execute tasks Split the DAG into .groupBy(...) stages of tasks .filter(...) Retry failed and strag-Store and serve blocks gler tasks Submit each stage and its tasks as ready

Anatomy of a Spark Application



Typical RDD pattern of use

Instead of doing a lot of work in each RDD, developers split tasks into lots of small RDDs

These are then organized into a DAG.

Developer anticipates which will be costly to recompute and hints to Spark that it should cache those.

Why is this a good strategy?

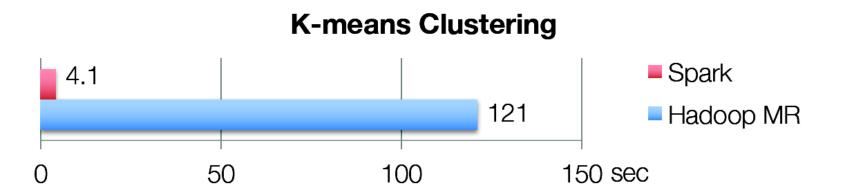
Spark tries to run tasks that will need the same intermediary data on the same nodes.

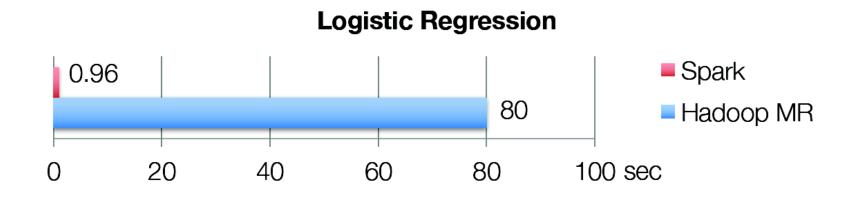
If MapReduce jobs were arbitrary programs, this wouldn't help because reuse would be very rare.

But in fact the MapReduce model is very repetitious and iterative, and often applies the same transformations again and again to the same input files.

- Those particular RDDs become great candidates for caching.
- MapReduce programmer may not know how many iterations will occur, but Spark itself is smart enough to evict RDDs if they don't actually get reused.

Iterative Algorithms: Spark vs MapReduce





Today's Topics

- Motivation
- Spark Basics
- Spark Programming

Spark Programming (1)

Creating RDDs

```
# Turn a Python collection into an RDD
sc.parallelize([1, 2, 3])
# Load text file from local FS, HDFS, or S3
sc.textFile("file.txt")
sc.textFile("directory/*.txt")
sc.textFile("hdfs://namenode:9000/path/file")
# Use existing Hadoop InputFormat (Java/Scala only)
sc.hadoopFile(keyClass, valClass, inputFmt, conf)
```

Spark Programming (2)

Basic Transformations

```
nums = sc.parallelize([1, 2, 3])
# Pass each element through a function
squares = nums.map(lambda x: x*x) // {1, 4, 9}
# Keep elements passing a predicate
even = squares.filter(lambda x: x % 2 == 0) // {4}
```

Spark Programming (3)

Basic Actions

```
nums = sc.parallelize([1, 2, 3])
# Retrieve RDD contents as a local collection
nums.collect() \# = [1, 2, 3]
# Return first K elements
nums.take(2) \# = [1, 2]
# Count number of elements
nums.count() \# => 3
# Merge elements with an associative function
nums.reduce(lambda x, y: x + y) # => 6
```

Spark Programming (4)

Working with Key-Value Pairs

pair. 2 // => b

```
Spark's "distributed reduce" transformations operate on RDDs of key-value pairs
Python: pair = (a, b)
           pair[0] # => a
           pair[1] # => b
Scala: val pair = (a, b)
           pair. 1 // => a
           pair. 2 // => b
Java: Tuple2 pair = new Tuple2(a, b);
           pair. 1 // => a
```

Spark Programming (5)

Some Key-Value Operations

```
pets = sc.parallelize([("cat", 1), ("dog", 1), ("cat", 2)])

pets.reduceByKey(lambda x, y: x + y)  # => {(cat, 3), (dog, 1)}

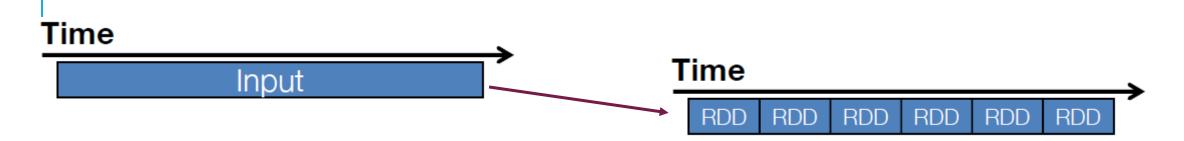
pets.groupByKey()  # => {(cat, [1, 2]), (dog, [1])}

pets.sortByKey()  # => {(cat, 1), (cat, 2), (dog, 1)}
```

Example: Word Count

```
lines = sc.textFile("hamlet.txt")
counts = lines.flatMap(lambda line: line.split(" "))
                  .map(lambda word: (word, 1))
                  .reduceByKey(lambda x, y: x + y)
                          "to"
                                        (to, 1)
                                                        (be, 2)
                          "be"
                                        (be, 1)
         "to be or"
                                                        (not, 1)
                          "or"
                                        (or, 1)
                          "not"
                                        (not, 1)
                                                        (or, 1)
                          "to"
                                        (to, 1)
         "not to be".
                                                        (to, 2)
                          "be"
                                        (be, 1)
```

Example: Spark Streaming



Represents streams as a series of RDDs over time (typically sub second intervals, but it is configurable)

```
val spammers = sc.sequenceFile("hdfs://spammers.seq")
sc.twitterStream(...)
    .filter(t => t.text.contains("Santa Clara University"))
    .transform(tweets => tweets.map(t => (t.user, t)).join(spammers))
    .print()
```

Spark: Combining Libraries (Unified Pipeline)

```
Load data using Spark SQL
points = spark.sql("select latitude, longitude from tweets")
# Train a machine learning model
model = KMeans.train(points, 10)
# Apply it to a stream
sc.twitterStream(...)
   .map(lambda t: (model.predict(t.location), 1))
   .reduceByWindow("5s", lambda a, b: a + b)
```

Spark: Setting the Level of Parallelism

All the pair RDD operations take an optional second parameter for number of tasks

```
words.reduceByKey(lambda x, y: x + y, 5)
words.groupByKey(5)
visits.join(pageViews, 5)
```

Summary

Spark is a powerful "manager" for big data computing.

It centers on a job scheduler for Hadoop (MapReduce) that is smart about where to run each task: co-locate task with data.

The data objects are "RDDs": a kind of recipe for generating a file from an underlying data collection. RDD caching allows Spark to run mostly from memory-mapped data, for speed.

Online tutorials: spark.apache.org/docs/latest

SELF-STUDY QUESTIONS

Spark RDDs look very similar to Python SQL or LINQ

What are some things Spark can do that would not automatically occur if you coded the same RDD query as a Python SQL query and just treated the entire data set as a single giant collection of objects?

SELF-STUDY QUESTIONS

Matei Zaharia thinks of Spark RDDs as a high level programming language that compiles to MapReduce.

Research his publications on Spark RDD "transformations," such as this paper. What are some PL optimization concepts that Spark is adopting when it maps RDDs to MapReduce on sharded data?

Self study questions

Is Spark a universal programming language for embarrassingly parallel computing?

Try to construct an example of a task that Spark cannot be used to solve, or at least not in a natural way.

Hint: If you know SQL, think about joins.