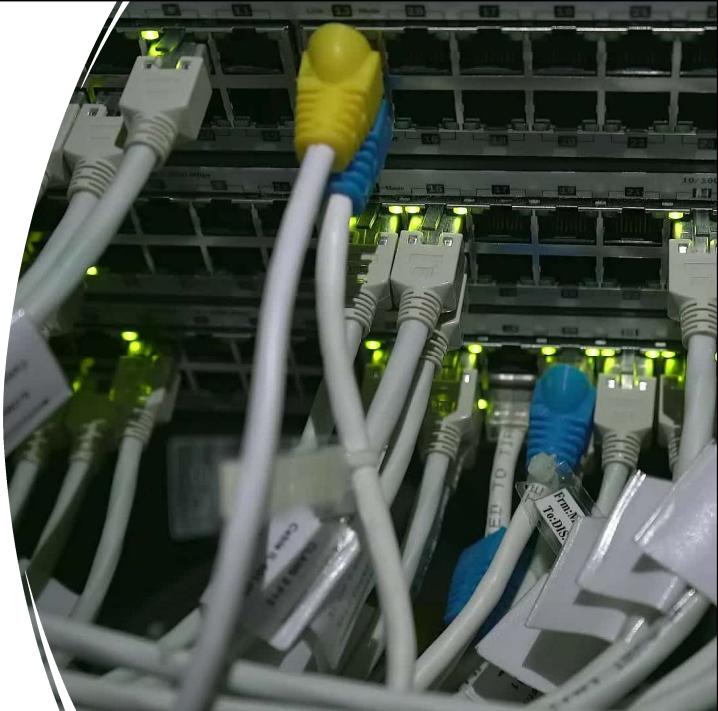


Data-Intensive Distributed Computing

CS431/451/631/651

Module 3 – From
MapReduce to Spark

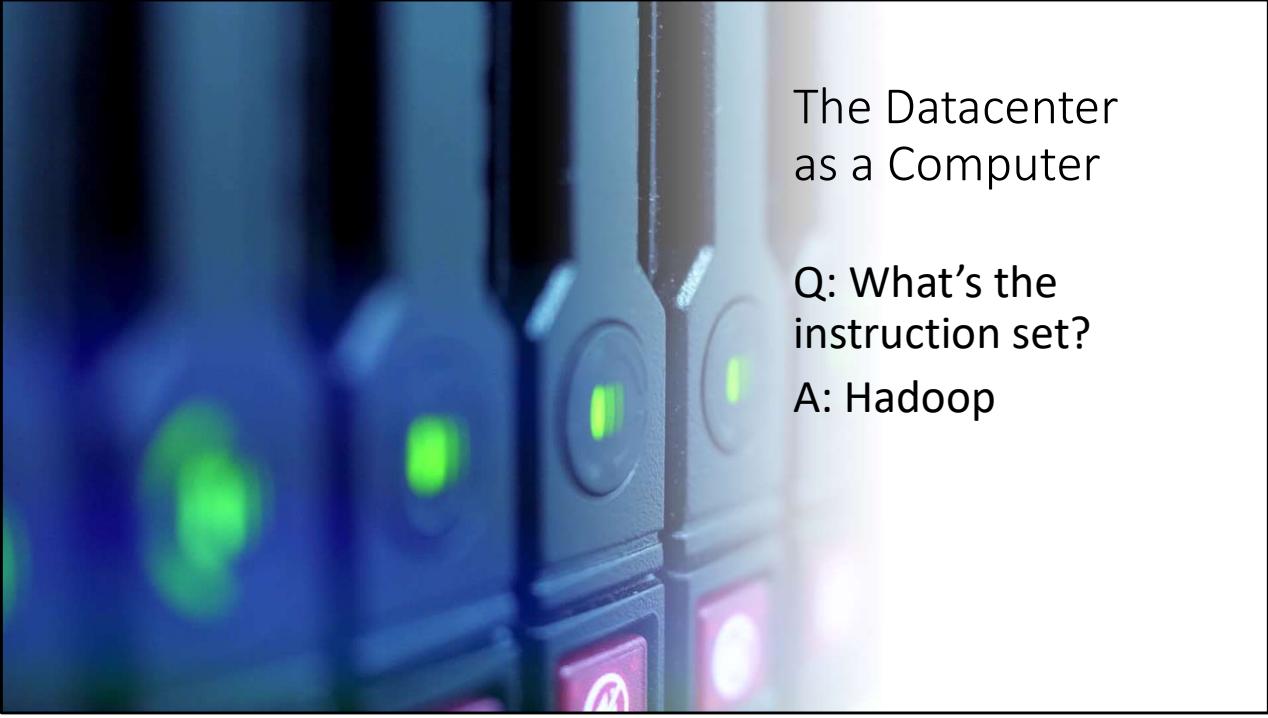


This Module's Agenda

Higher-Level Programming

Spark

Algorithm Design



The Datacenter as a Computer

Q: What's the
instruction set?

A: Hadoop

Layers of Abstraction

Higher Level Language (e.g. Python)

Lower Level Language (e.g. C)

Assembly

Machine Code

Instruction Set Architecture

Micro-Architecture

Gates, Adders, Registers, Etc.

Electronics (Transistors)

Physics

Data Center Abstraction

??? <TODAY'S TOPIC>

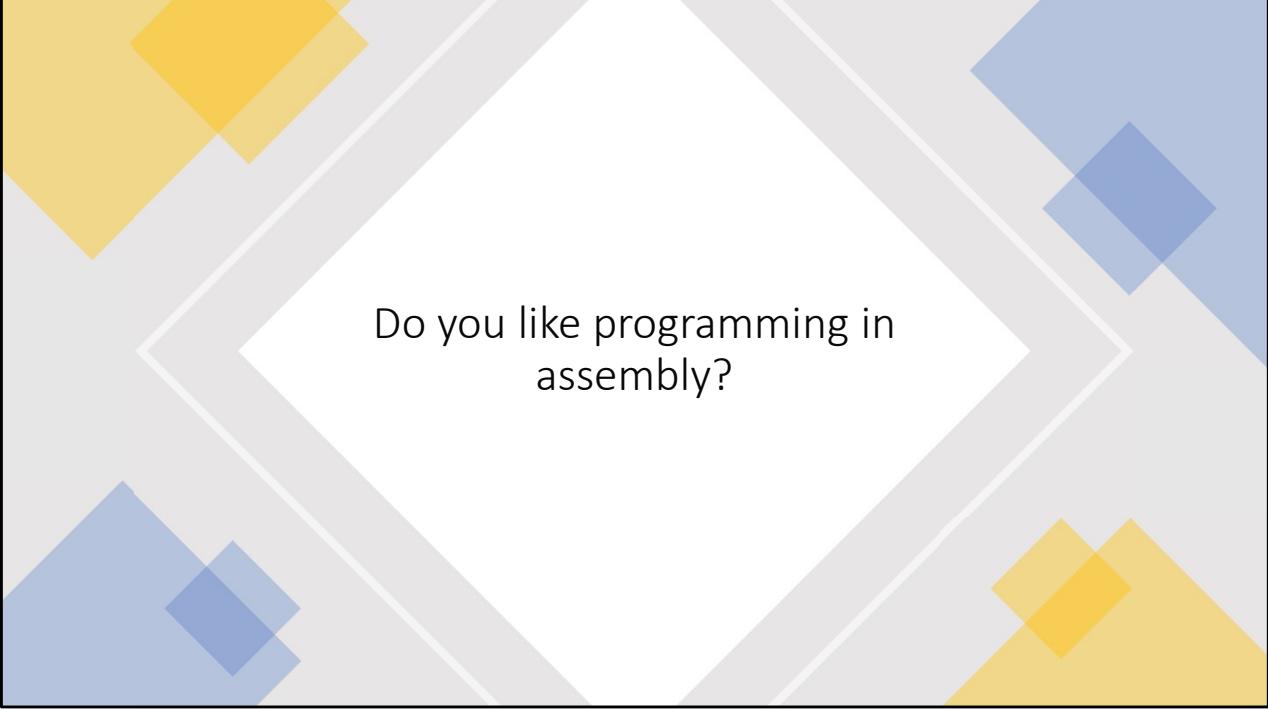
Hadoop Task

HDFS / Hadoop Framework

Cluster of Computers (Networking)

Individual Servers



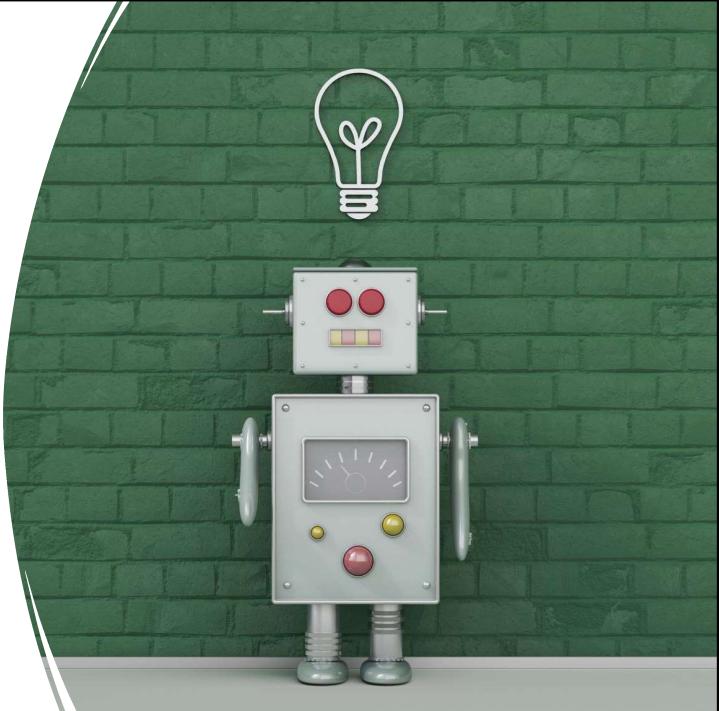


Do you like programming in
assembly?

It's OK if the answer is yes, there's no judgement here

What's the alternative?

- Hadoop is great, but has a lot of boilerplate and repetition
- It's also tedious to program
- Can we create a Distributed C (or Python) to Hadoop's Assembly?



Yes We* Can

facebook

What we really need
is SQL!

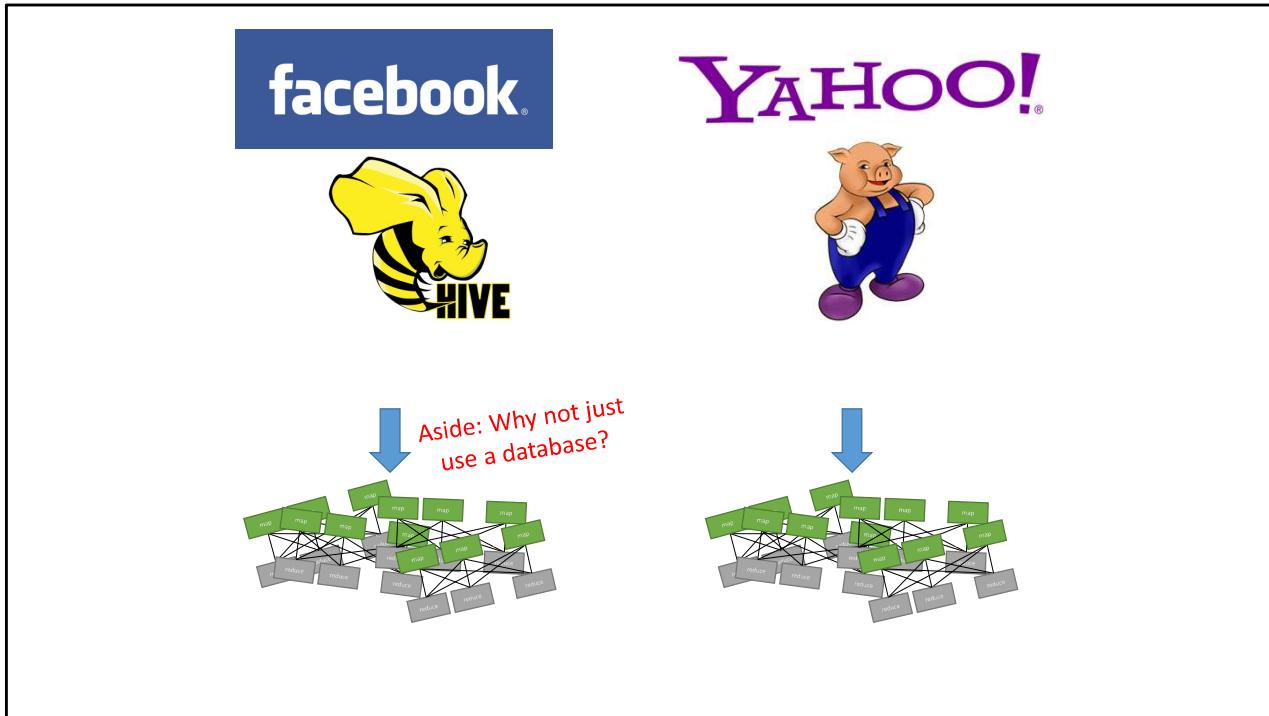


YAHOO!

What we really need
is a scripting
language!



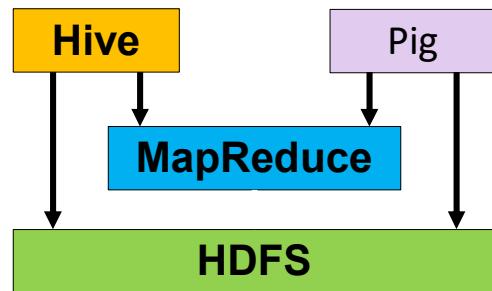
* - not me personally, but it has been done. Several times.



Both have their place. Hive is on top of MapReduce. It's good for huge datasets that are accessed in a linear fashion. One read, one write. SQL requires lots of read/write access to the data.

SQL – You need OLTP and/or low latency. Less-complicated data sets that need frequent updates

Hive – You don't care about latency, or have huge amounts of data (which means it doesn't matter whether or not you care, you're going to have latency). Batch processing of complicated data sets



10

Pig and Hive programs are converted to MapReduce jobs at the end of the day.

Pig Examples



Pig: Example

Task: Find the top 10 most visited pages in each category

Visits			URL Info		
User	Url	Time	Url	Category	PageRank
Amy	cnn.com	8:00	cnn.com	News	0.9
Amy	bbc.com	10:00	bbc.com	News	0.8
Amy	flickr.com	10:05	flickr.com	Photos	0.7
Fred	cnn.com	12:00	espn.com	Sports	0.9
	•		•		
	•		•		
	•		•		

12

Pig Slides adapted from Olston et al. (SIGMOD 2008)

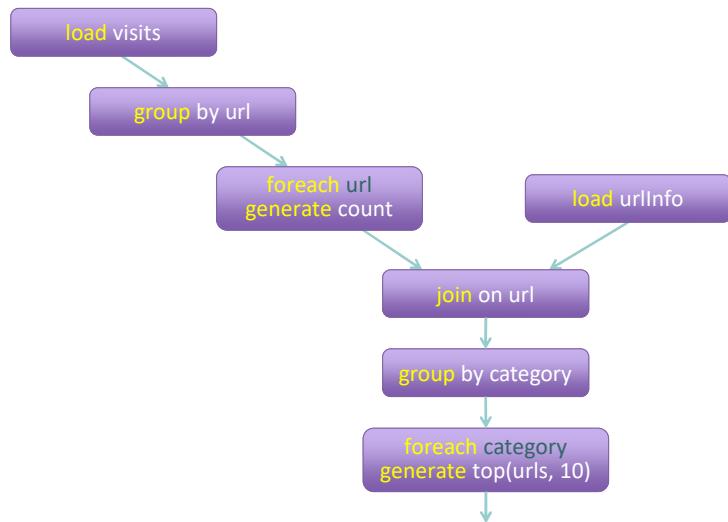
Pig: Example Script

```
visits = load '/data/visits' as (user, url, time);
gVisits = group visits by url;
visitCounts = foreach gVisits generate url, count(visits);
urlInfo = load '/data/urlInfo' as (url, category, pRank);
visitCounts = join visitCounts by url, urlInfo by url;
gCategories = group visitCounts by category;
topUrls = foreach gCategories generate
    top(visitCounts,10);

store topUrls into '/data/topUrls';
```

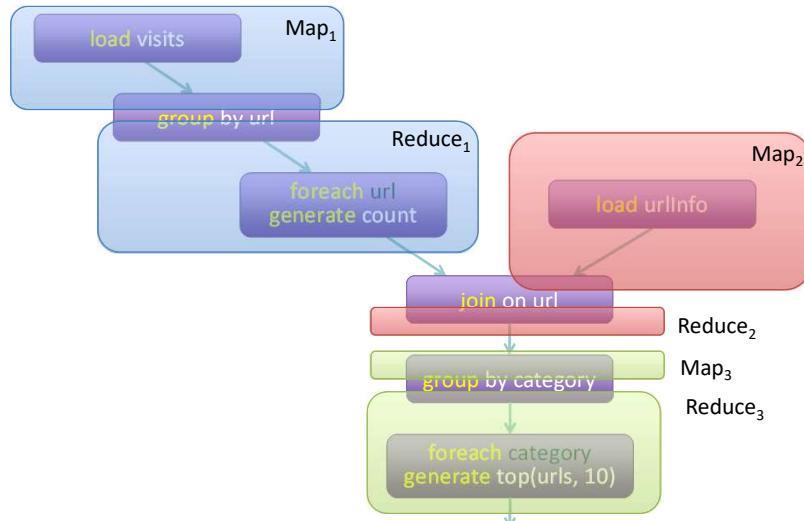
Pig Slides adapted from Olston et al. (SIGMOD 2008)

Pig Query Plan



Pig Slides adapted from Olston et al. (SIGMOD 2008)

Pig: MapReduce Execution



Pig Slides adapted from Olston et al. (SIGMOD 2008)

YUP, you can do a map that takes multiple inputs. Neato!

- Which would you rather:
 - Read
 - Write
 - Debug (!)

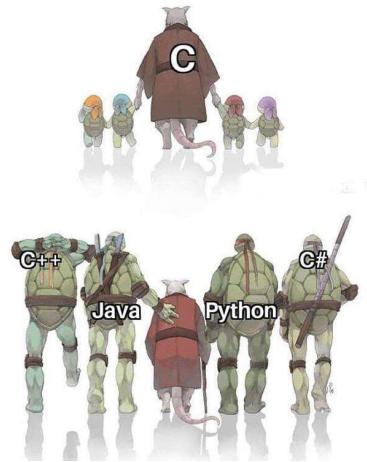
Pig Slides adapted from Olston et al. (SIGMOD 2008)

Isn't Pig Slower than Hadoop?

Potentially.

Isn't C slower than assembly?

Isn't Python slower than C?



The Data Center as a Computer

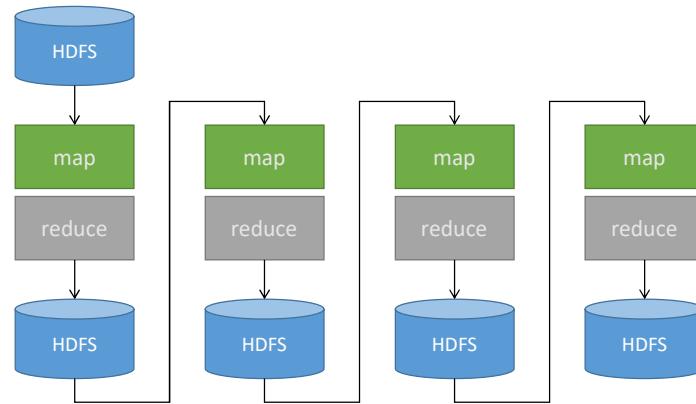
So Hadoop is the Instruction Set,
right?

What if I need two reduce passes.
Do I really need two jobs?

(On A1 yes, you do)



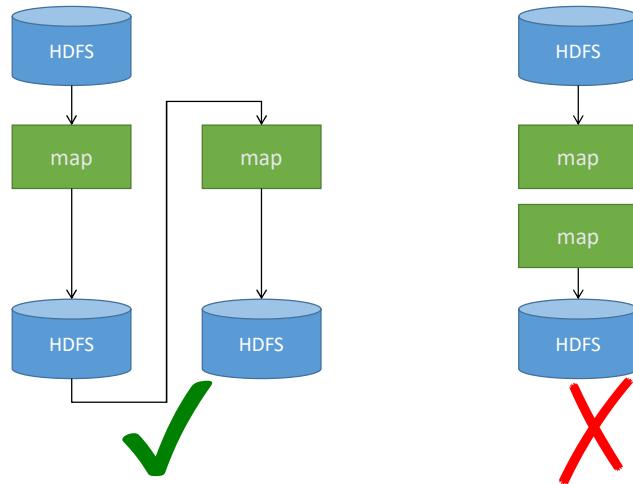
MapReduce Workflows



What's wrong?

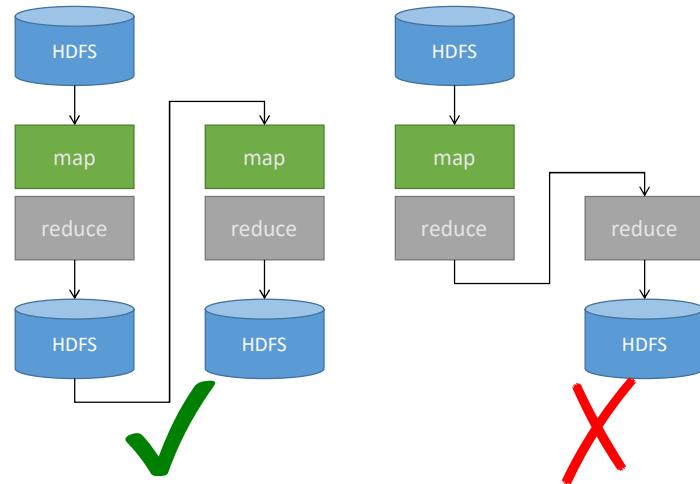
There is a lot of disk i/o involved which significantly reduces running MapReduce jobs like this.

Want Map-to-Map?



It's okay not to have reduce but the output of map cannot go to another map.

Want Map-to-Reduce-to-Reduce?



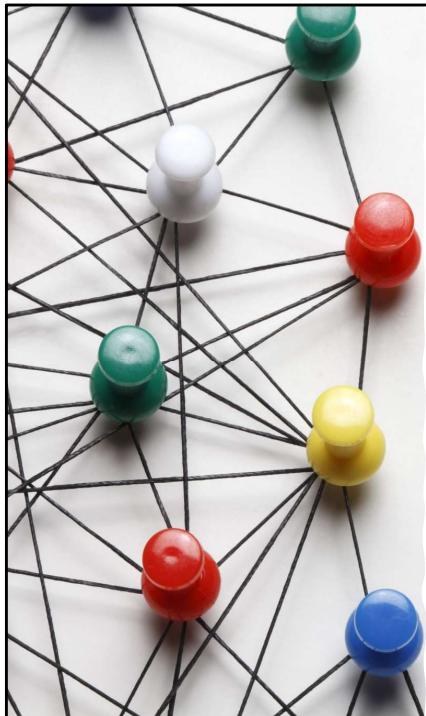
It's not OK to skip a map. (Strictly speaking, there's FILE access between map and reduce tasks, too, but then it doesn't fit as neatly onto the slides)

The Data Center as a Computer

Q: Is there a better
instruction set?

A: Hadoop 2





Hadoop 2.0

Nodes are now resource managers

Can do MapReduce the same as always

Can also do other things

Other Things? Like What



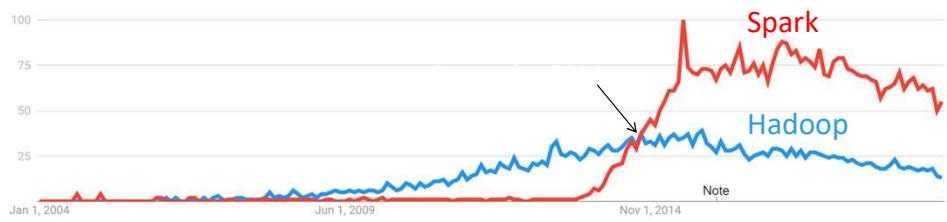
Brief history:

Developed at UC Berkeley AMPLab in 2009

Open-sourced in 2010

Became top-level Apache project in February 2014

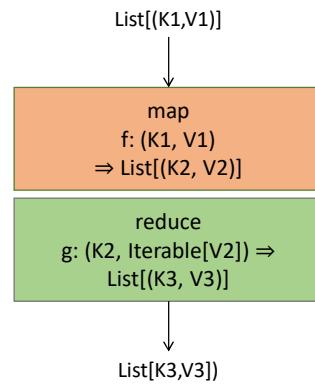
Spark vs. Hadoop



Spark is more popular than Hadoop today.

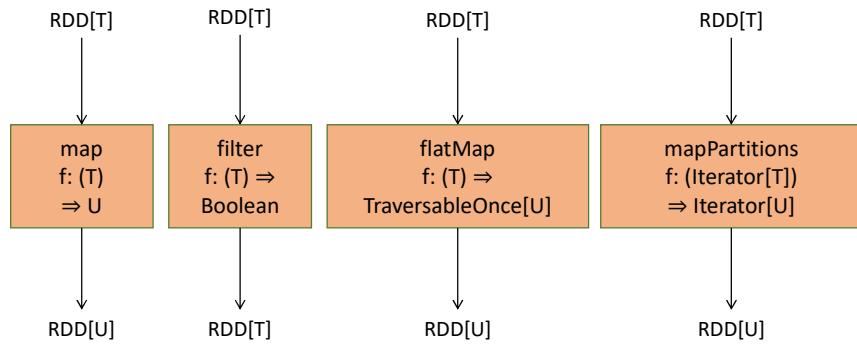


MapReduce



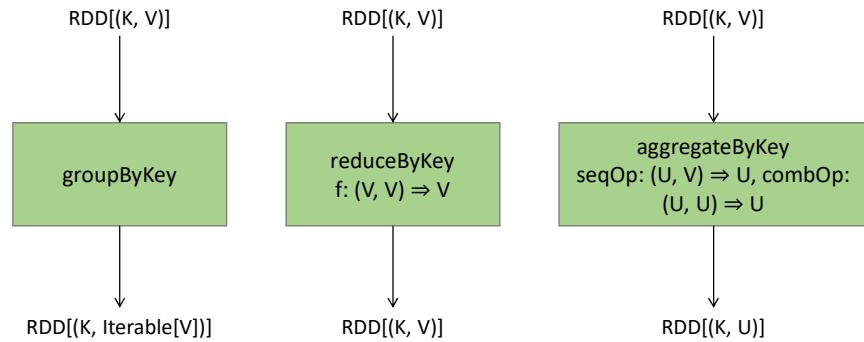
This is the only mechanism we had in MapReduce.

Map-like Operations



But Spark provides many more operations (enriched instruction set).

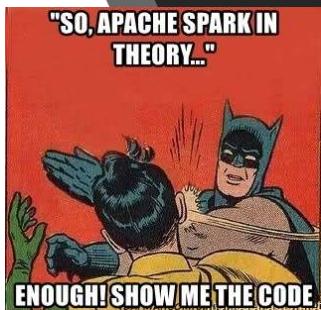
Reduce-like Operations





And many other operations!

Interactive Demo Time!



<Dan, showing off Spark Shell / PySpark>



Introduction to Apache Spark

Slides from: Patrick Wendell – Databricks

Memes from: Ali Abedi

What is Spark?

Fast and Expressive Cluster Computing
Engine Compatible with Apache Hadoop

Up to **10x** faster on disk,
100x in memory

Efficient

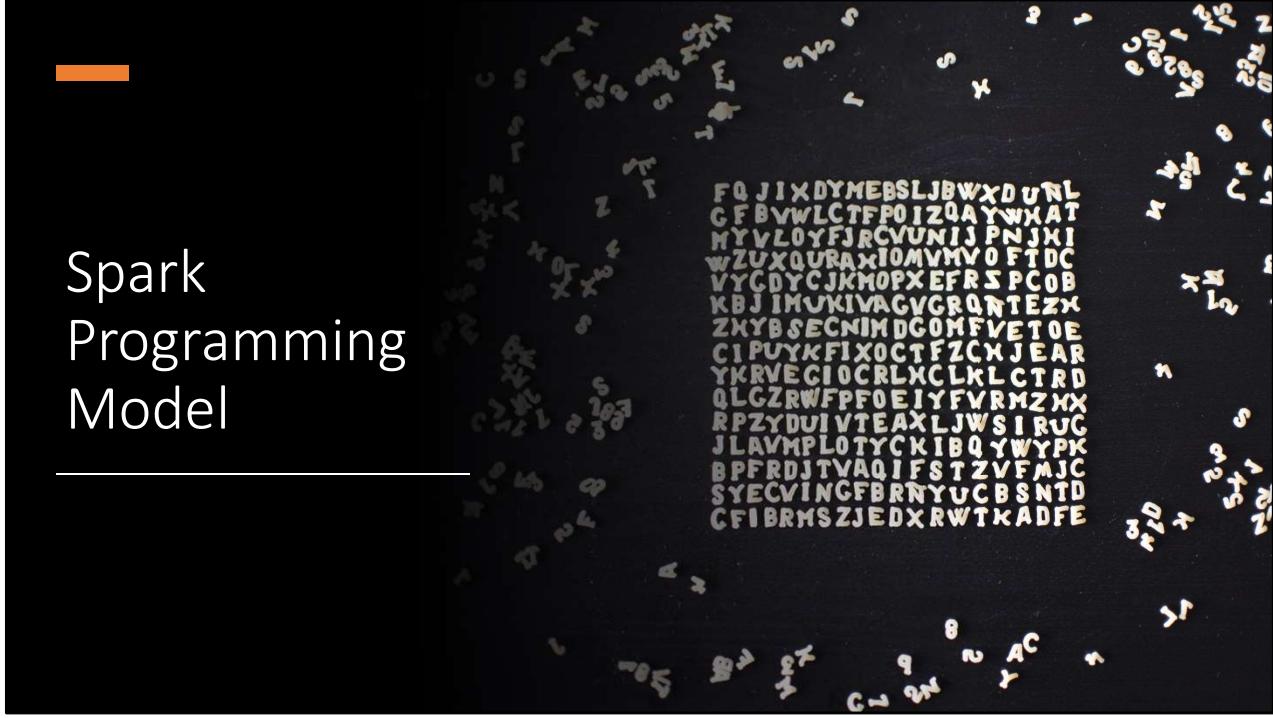
- General execution graphs
- In-memory storage

2-5x less code

Usable

- Rich APIs in Java, Scala, Python
- Interactive shell

Spark Programming Model



```
FQJIXDYM EBSLJBWXDU NL  
GFBVWLCTFP OIZQAYW XAT  
HYVLOYFJRCVUNI JPNJXI  
WZUXQURAXIOMVMV OF TDC  
VYCDYCJJKMOPXEFRS PCOB  
KBJIMU KIVAGVG RQNT EZX  
ZYBSEC NIMDGOMF VETOE  
CIPUYKFX OCTFZCX JEAR  
YKRVE CIOCRLXCLKLCTR D  
OLGZRWF PF OEIY FVRMZXX  
RPZYDUIVTEAXLJWSIRUC  
JLAVMPLOTYCKIBQYWYPK  
BPFRDJTVAQIFSTZVF MJC  
SYECVINGFBRNYUCBSNTD  
CFIBRMSZJEDXRWTKA DFE
```

Key Concept: RDD's

Write programs in terms of **operations** on
distributed datasets

Resilient Distributed Datasets

- Collections of objects spread across a cluster, stored in RAM or on Disk
- Built through parallel transformations
- Automatically rebuilt on failure

Operations

- Transformations
 - (e.g. map, filter, groupBy)
- Actions
 - (e.g. count, collect, save)

RDD structure



Example: Log Mining

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

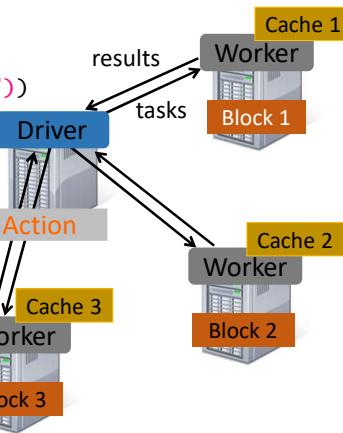
B Transformed RDD

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

```
messages.filter(lambda s: "mysql" in s).count()  
messages.filter(lambda s: "php" in s).count()  
....
```

Full-text search of Wikipedia

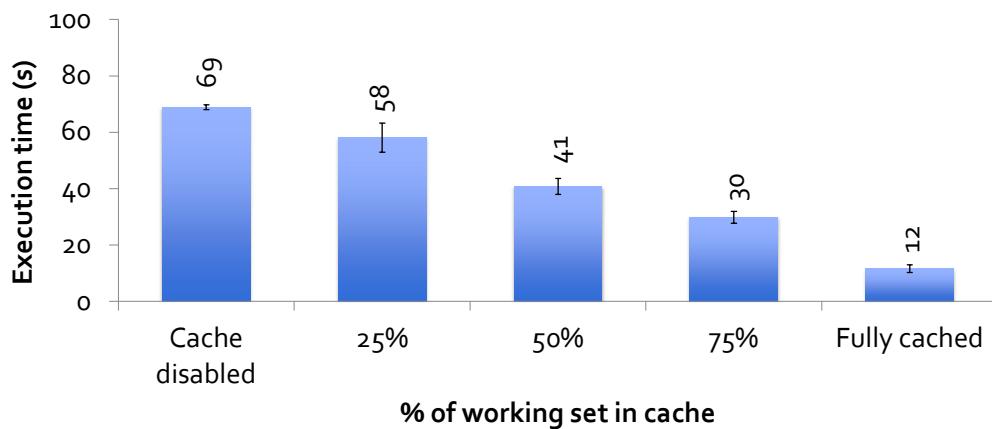
- 60GB on 20 EC2 machine
- 0.5 sec vs. 20s for on-disk



Lazy evaluation: Spark doesn't really do anything until it reaches an action! This helps Spark to optimize the execution and load only the data that is really needed for evaluation.

Dan adds: If you branch, then you cache!

Impact of Caching on Performance



Fault Recovery

RDDs track *lineage* information that can be used to efficiently recompute lost data

```
msgs = textFile.filter(lambda s: s.startswith("ERROR"))
    .map(lambda s: s.split("\t")[2])
```



Programming with RDD's

SparkContext

- Main entry point to Spark functionality
- Available in shell as variable `sc`
- In standalone programs, you'd make your own

Poor font choice I think? Lowercase “sc”

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")
```

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}  
  
# Map each element to zero or more others  
> nums.flatMap(lambda x: range(x))  
  > # => {0, 0, 1, 0, 1, 2}
```

Range object (sequence
of numbers 0, 1, ..., x-1)

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

# Write elements to a text file
> nums.saveAsTextFile("hdfs://file.txt")
```

Working with Key-Value Pairs

Spark's "distributed reduce" transformations operate on RDDs of key-value pairs

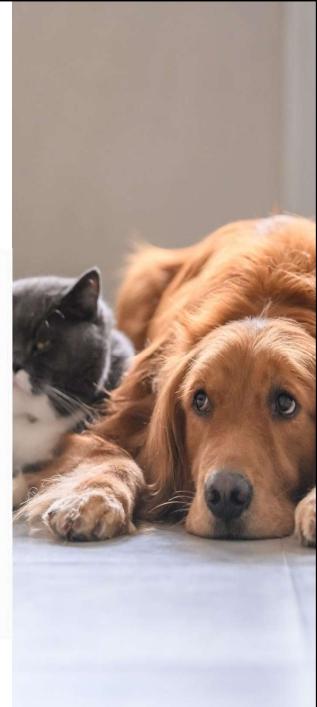
Python: pair = (a, b)
 pair[0] # => a
 pair[1] # => b

Scala: val pair = (a, b)
 pair._1 // => a
 pair._2 // => b

Java: Tuple2 pair = new Tuple2(a, b);
 pair._1 // => a
 pair._2 // => b

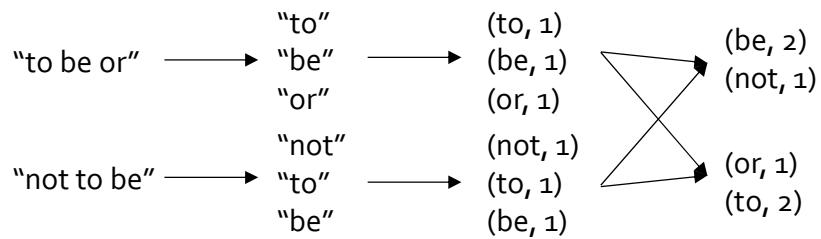
Some Key-Value Operations

```
> 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)}
```



Word Count (Python)

```
> lines = sc.textFile("hamlet.txt")
> counts = lines.flatMap(lambda line: line.split(" "))
    .map(lambda word: (word, 1))
    .reduceByKey(lambda x, y: x + y)
    .saveAsTextFile("results")
```



Word Count (Scala)

```
val textFile =  
  sc.textFile("hamlet.txt")  
  
textFile  
  .flatMap(line => line.split(" "))  
  .map(word => (word, 1))  
  .reduceByKey((x, y) => x + y)  
  .saveAsTextFile("results")
```

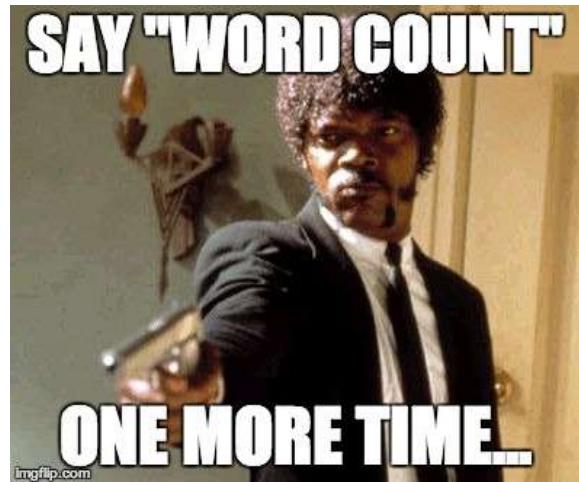
(Alternative Scala)

```
val textFile =  
  sc.textFile("hamlet.txt")  
  
textFile  
  .flatMap(_.split(" "))  
  .map((_, 1))  
  .reduceByKey(_ + _)  
  .saveAsTextFile("results")
```

In Scala, underscores mean “this expression is the body of an anonymous function”

“_ + _” means the same as “ $(x, y) \Rightarrow x + y$ ”

(_+_) looks like a butthole but we’re all going to just ignore that and be mature



I mean, word count, aka token frequency, is a building block for lots of text processing...just because it's easy doesn't mean it's not useful.

Other Key-Value Operations

```
> visits = sc.parallelize([ ("index.html", "1.2.3.4"),
   # ("about.html", "3.4.5.6"),
   # ("index.html", "1.3.3.1") ])
```



```
> pageNames = sc.parallelize([ ("index.html", "Home"),
   ("about.html", "About") ])
```



```
> visits.join(pageNames)
# ("index.html", ("1.2.3.4", "Home"))
# ("index.html", ("1.3.3.1", "Home"))
# ("about.html", ("3.4.5.6", "About"))
```



```
> visits.cogroup(pageNames)
# ("index.html", ([["1.2.3.4"], ["1.3.3.1"]], [[["Home"]]]))
# ("about.html", ([["3.4.5.6"]], [[["About"]]]))
```

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)
```

Dan adds: So does scc.textFile (and other base RDDs). However, for these this is “minimum number of tasks”. E.g. if a file is split into 16 blocks on HDFS, and you open it with textFile(PathString, 10), you’ll still get 16 partitions, not 10.

If you’re submitting a job with a total of 8 vCores, you should always have 8 partitions if you can manage it. Otherwise a core will be idle. (In fact, it’s usually better to have more tasks than cores, so that tasks bottle necked on reading will be able to share a single core).

What’s the default?

It uses the same number of partitions for destination as the source has. Eg a reduceByKey on an RDD with 8 partitions will result in another RDD with 8 partitions.

For joins, it’s the minimum of the LHS and RHS RDDs. Eg join an RD with 3 parts to one with 8, you will get 3.

If you specify spark.default.parallelism it will use this as the default instead! (For shuffles only, not parallelize textFile or other base RDDs)

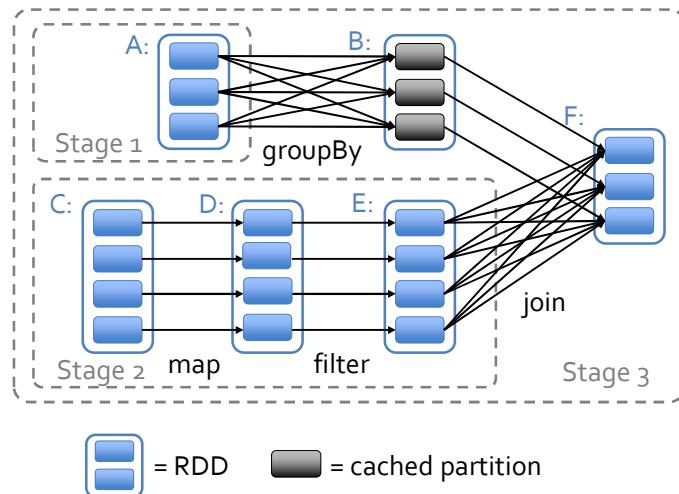


Please watch your jobs on datasci and kill things that seem to be stuck. Try on student.cs FIRST, only run on datasci when you're confident. If you need to make changes, rerun on student.cs first!!!

I've added a bonus slide at the end with some tips about viewing Spark jobs on the cluster. (A big reason for "runs forever" on datasci is a reducer that's $O(n)$ – usually caused by stripes being merged inefficiently.)

Under The Hood: DAG Scheduler

- General task graphs
- Automatically pipelines functions
- Data locality aware
- Partitioning aware to avoid shuffles



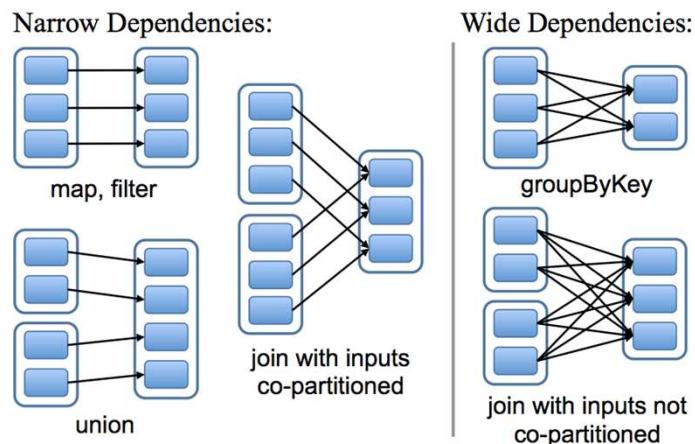
Directed Acyclic Graph (DAG)

A job is broken down to multiple stages that form a DAG.

You can get the DAG from an RDD using the `toDebugString` method. (print it, since it contains newlines and will be illegible as a string value)

It's also viewable through the Hadoop monitoring page.

Physical Operators



Narrow dependency is much faster than wide dependency because it does not require shuffling data between working nodes.

Also: `reduceByKey`, `groupByKey`, etc will also have narrow dependencies if the upstream RDD is already partitioned by key. Its less common but not unheard of.

More RDD Operators

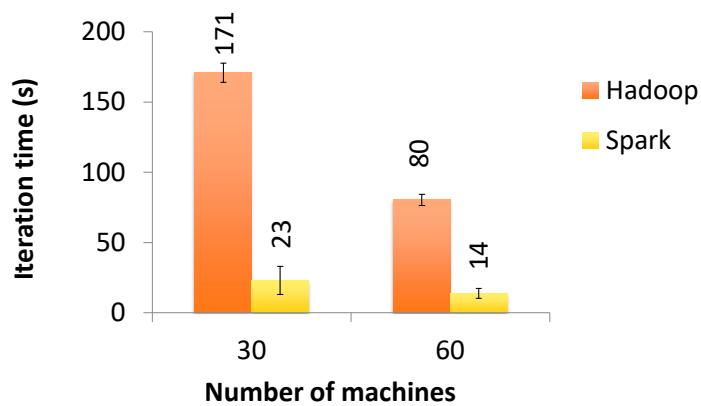
- map
- filter
- groupBy
- sort
- union
- join
- leftOuterJoin
- rightOuterJoin
- reduce
- count
- fold
- reduceByKey
- groupByKey
- cogroup
- cross
- zip
- sample
- take
- first
- partitionBy
- mapwith
- pipe
- save
- ...





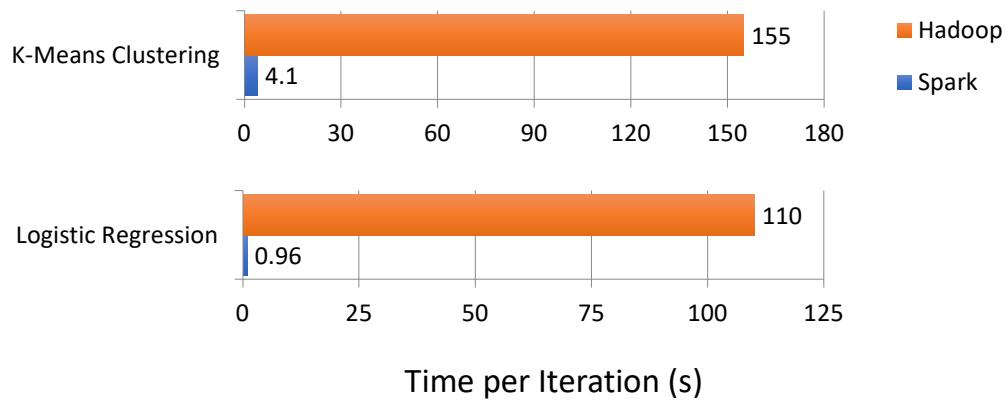
Performance

PageRank Performance



Since spark avoids heavy disk i/o, it significantly improves the performance.

Other Iterative Algorithms



Spark outperforms Hadoop in iterative programs because it tries to keep the data that will be used again in the next iteration in memory. In contrast with Hadoop which always read and write from/to disk.



Hadoop Ecosystem and Spark

YARN

Hadoop's (original) limitations:

Can only run MapReduce

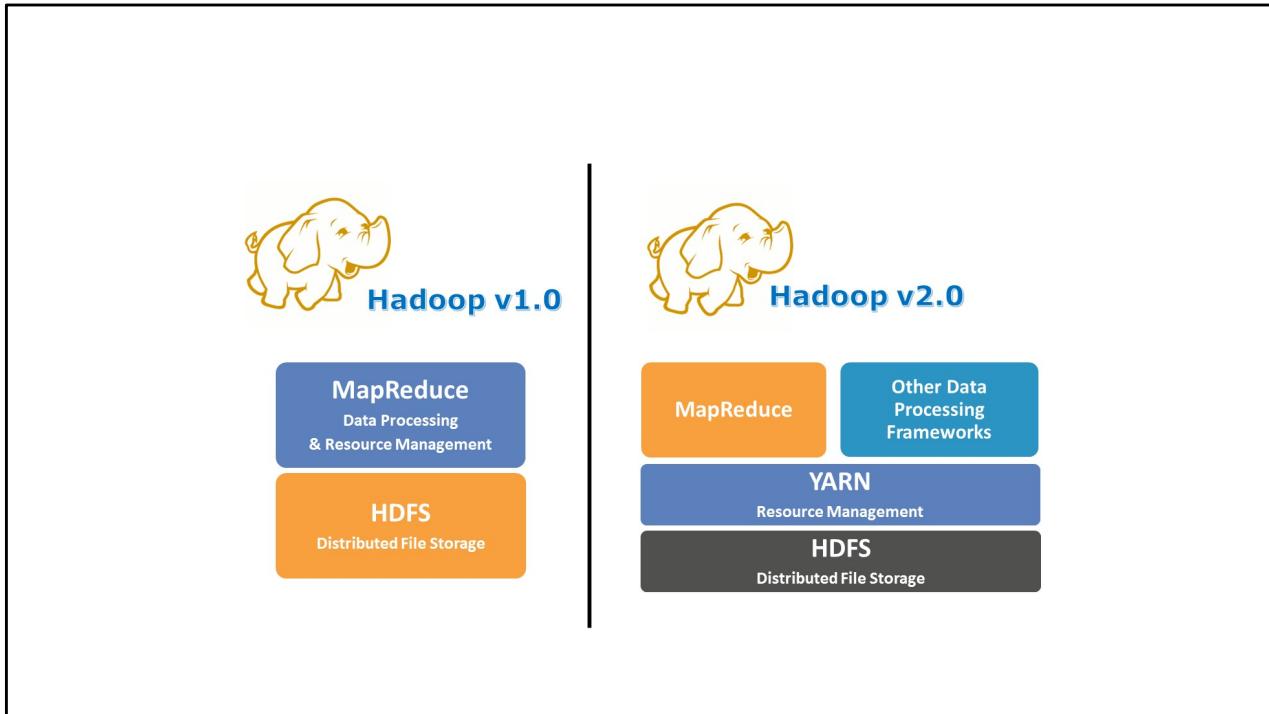
What if we want to run other distributed frameworks?

YARN = Yet-Another-Resource-Negotiator

Provides API to develop any generic distributed application

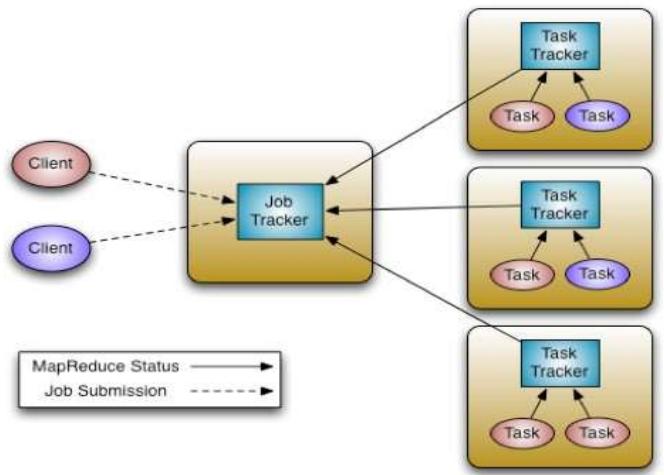
Handles scheduling and resource request

MapReduce (MR2) is one such application in YARN



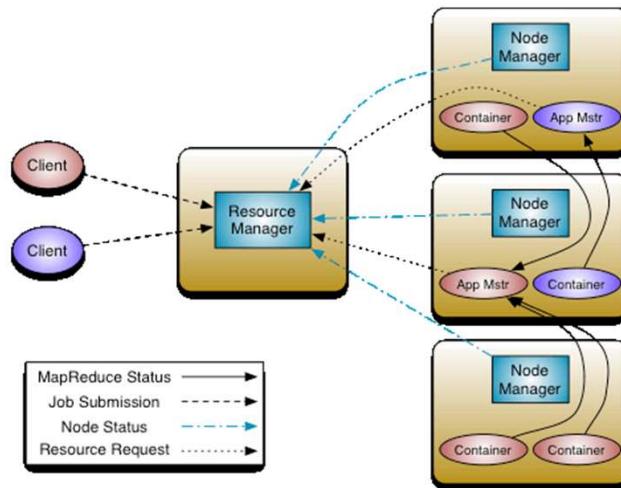
In Hadoop v1.0, the architecture was designed to support Hadoop MapReduce only. But later we realised that it is a good idea if other frameworks can also run on Hadoop cluster (rather than building a separate cluster for each framework). So in v2.0, YARN provides a general resource management system that can support different platforms on the same physical cluster.

Hadoop v1.0



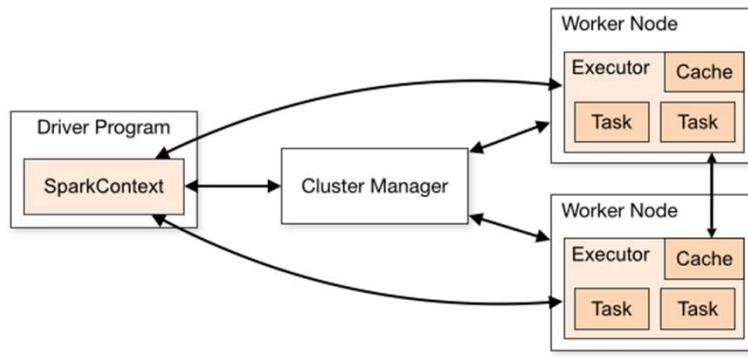
The Job tracker in v1.0 was specific to Hadoop jobs.

Hadoop v2.0



But the resource manager in v2.0 can support different types of jobs (e.g., Hadoop, Spark,...).

Spark Architecture



Important –
There are
multiple
tasks per
executor

Why is this important?

To work, the Spark driver must send relevant code (Scala or Python) to run each **task**.

```
thresh = 5
```

```
myRdd.map(lambda x: x >= thresh)
```

The lambda “captures” thresh, so it gets packaged up too. (That’s bad if it’s large)



Broadcast

If you Broadcast a value, then Spark only sends one copy per Executor (worker machine) not per Task

```
thresh = sc.broadcast(5)  
myRdd.filter(lambda x: x > thresh.value)
```

(It makes no difference here, but would if broadcasting a lookup table)



Constant means Constant

Broadcast variables are read-only

```
thresh = sc.broadcast(5)
thresh.value = 6
Error: value is not a member of ...Broadcast[int]
Error: value is not assignable
```

(Global variables are too, but will silently fail)

The errors are what you'd see in Scala or Python

Accumulators

A Broadcast variable carries information from Driver to Executor

What if we want communication from Executor back to Driver?

A: Accumulator



Counter Accumulators (Python)

```
lineCounter = sc.accumulator(0)

def split_and_count(line):
    lineCounter.add(1)
    return line.split()

myRdd.map(split_and_count). ...
lineCounter.value()
```

Counter Accumulators (Scala)

```
val lineCounter = sc.longAccumulator

def split_and_count(line : String) = {
    lineCounter.add(1)
    line.split()
}

myRdd.map(split_and_count). ...
lineCounter.value
```

Types of Accumulator

longAccumulator, doubleAccumulator

(In Python, they're just called accumulator)

Used for accumulating numerical values

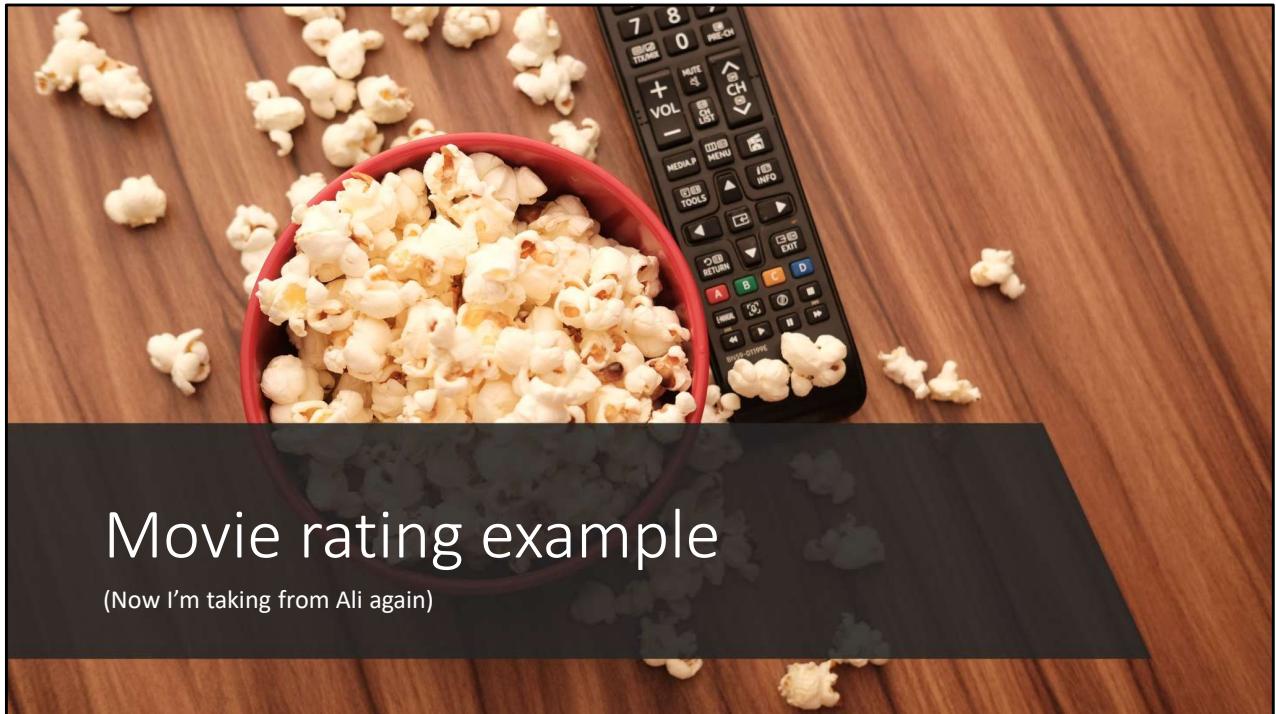
Driver can inspect the value (and take average of values accumulated)

Workers can only write

Partitioners

By default Spark shuffles use a hash partitioner (just like MapReduce)

Also like MapReduce, can override.



This example in particular is not very helpful in slide-only form. I alt-tab and do some goofing around in spark-shell or pyspark

Input Format

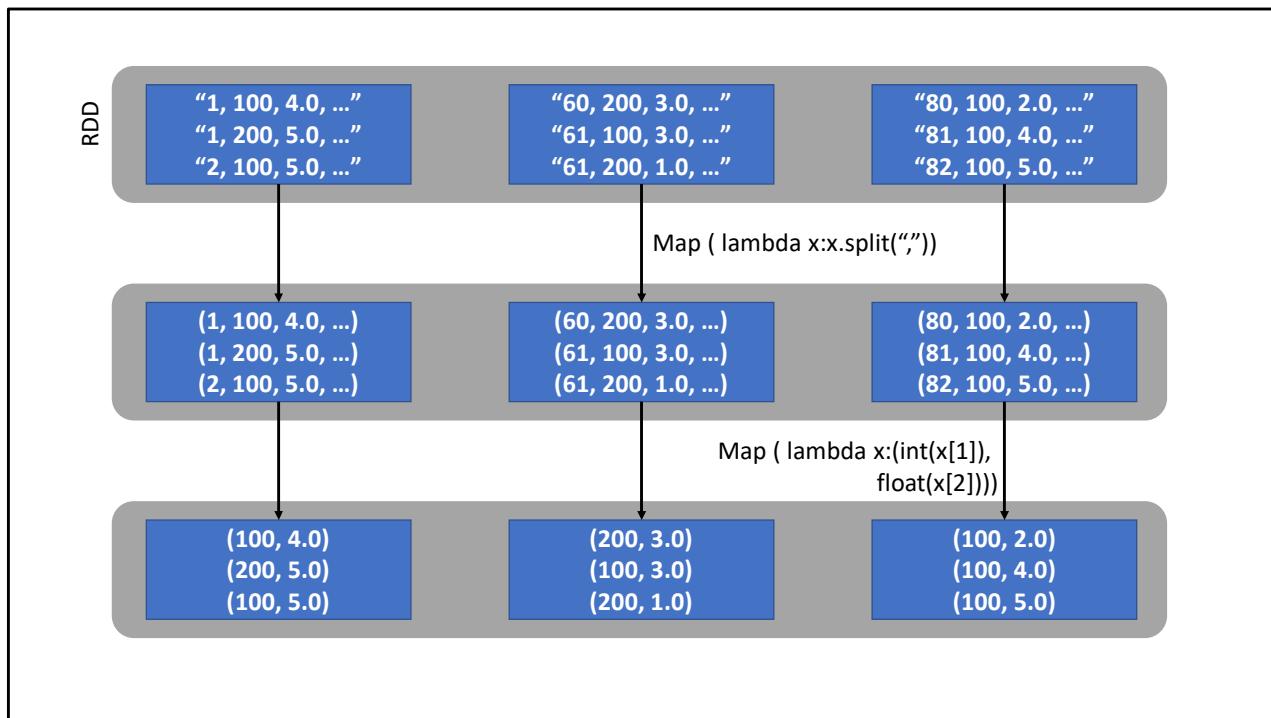
CSV file

Fields:

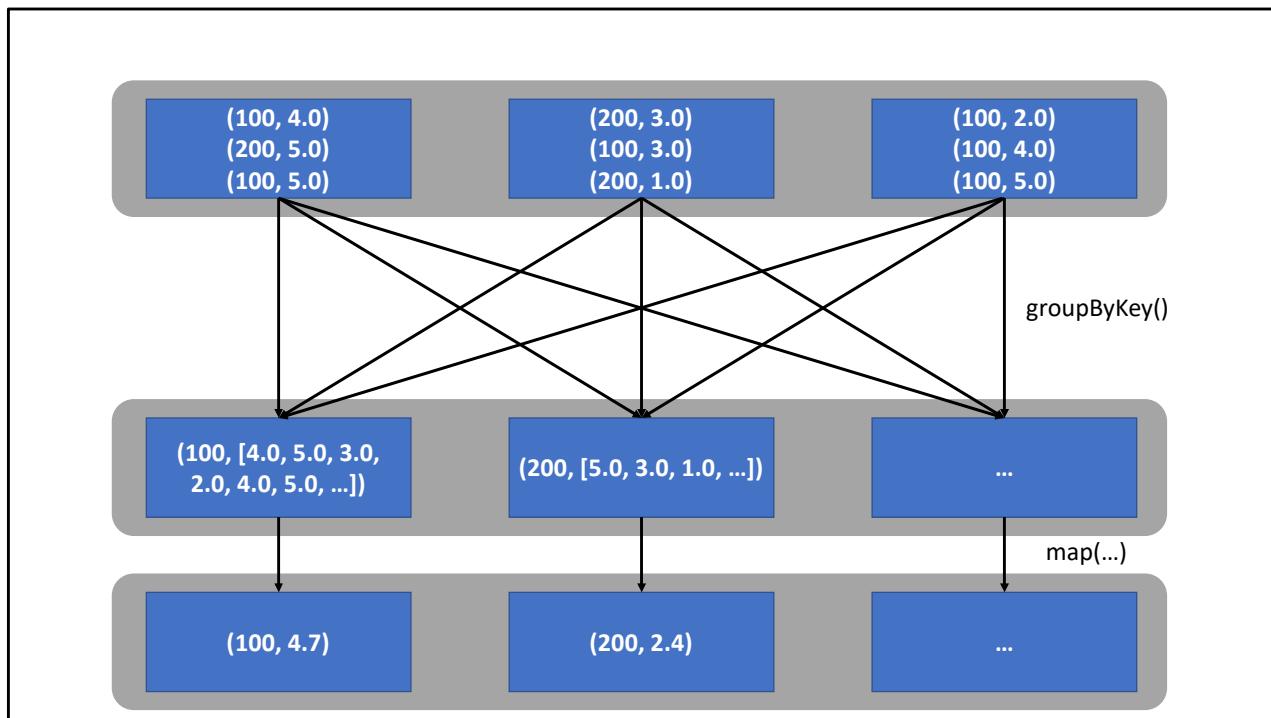
- User ID (unique key per user)
- Movie ID (unique key per movie)
- Rating (1-5 stars)
- Text of review (optional)

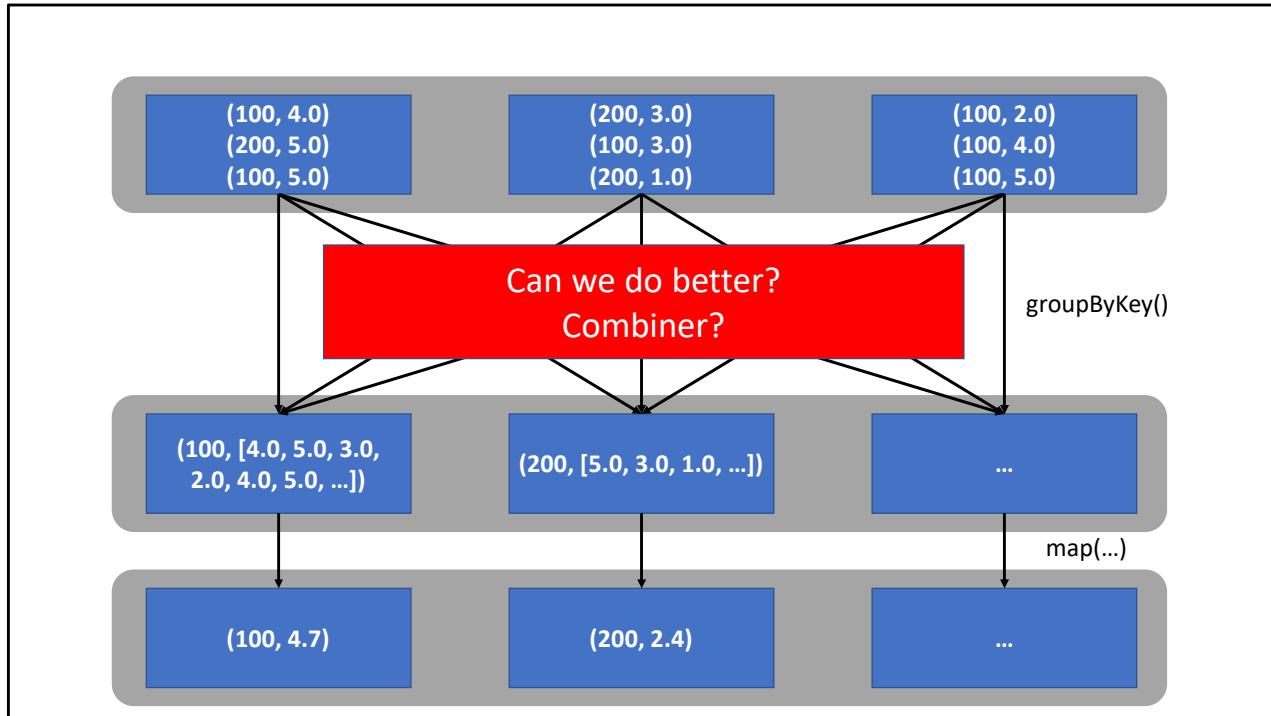
e.g.

“1, 100, 3.5, ‘s aight”



The “” are missing from the tuples in the middle because I’m not going to type that many “”! It’d be really hard to read I think





Avoid `groupByKey` if you can – MapReduce (without combiner) in Spark is essentially – flatmap -> `groupByKey` -> flatmap. We're trying to do better than MapReduce though.

Wait...converting MapReduce to Spark doesn't use the `reduceByKey` function??? That's right. MapReduce's reduce is more flexible.

Reduce vs reduceByKey

Reduce (MapReduce)

- $(K_2, V_2) \Rightarrow \text{List}[K_3, V_3]$
- KVP are partitioned and shuffled by Partitioner
- Reduce job calls reduce on keys in sorted order

reduceByKey (Spark)

- $V \Rightarrow V$
- $\text{RDD}[(K, V)] \Rightarrow \text{RDD}[(K, V)]$
- Less flexible
 - But does what reduce should normally be used for
- Reduces before shuffle (combiner)
- Reduces after shuffle (reducer)

reduceByKey vs combineByKey

combineByKey gives more fine-grained control (if needed)

$\text{RDD}[(\text{K}, \text{V})].\text{combineByKey}(\text{create}, \text{append}, \text{merge}) \Rightarrow \text{RDD}[(\text{K}, \text{C})]$

create – make a C from a V

append – take a C and add a V to it

merge – combine two C

reduceByKey(reduce) calls combineByKey(identity, reduce, reduce)

reduceByKey vs aggregateByKey

aggregateByKey is between reduceByKey and combineByKey

`RDD[(K, V)].aggregateByKey(init, append, merge) => RDD[(K, C)]`

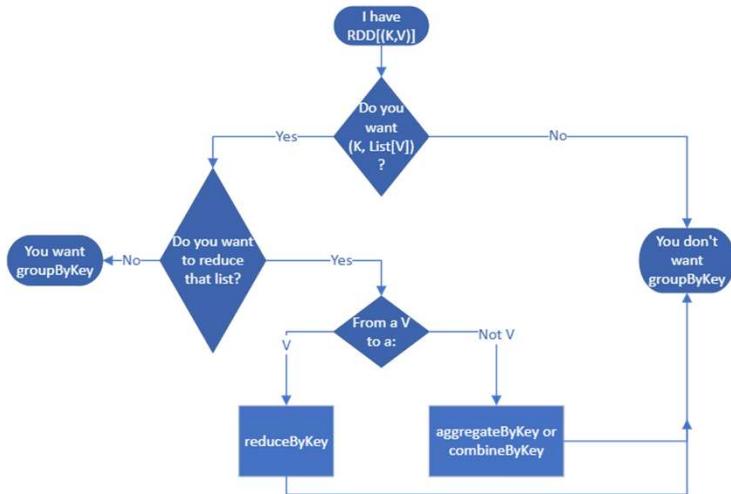
init – initial (or zero) value [type C]

append- take a C and add a V to it

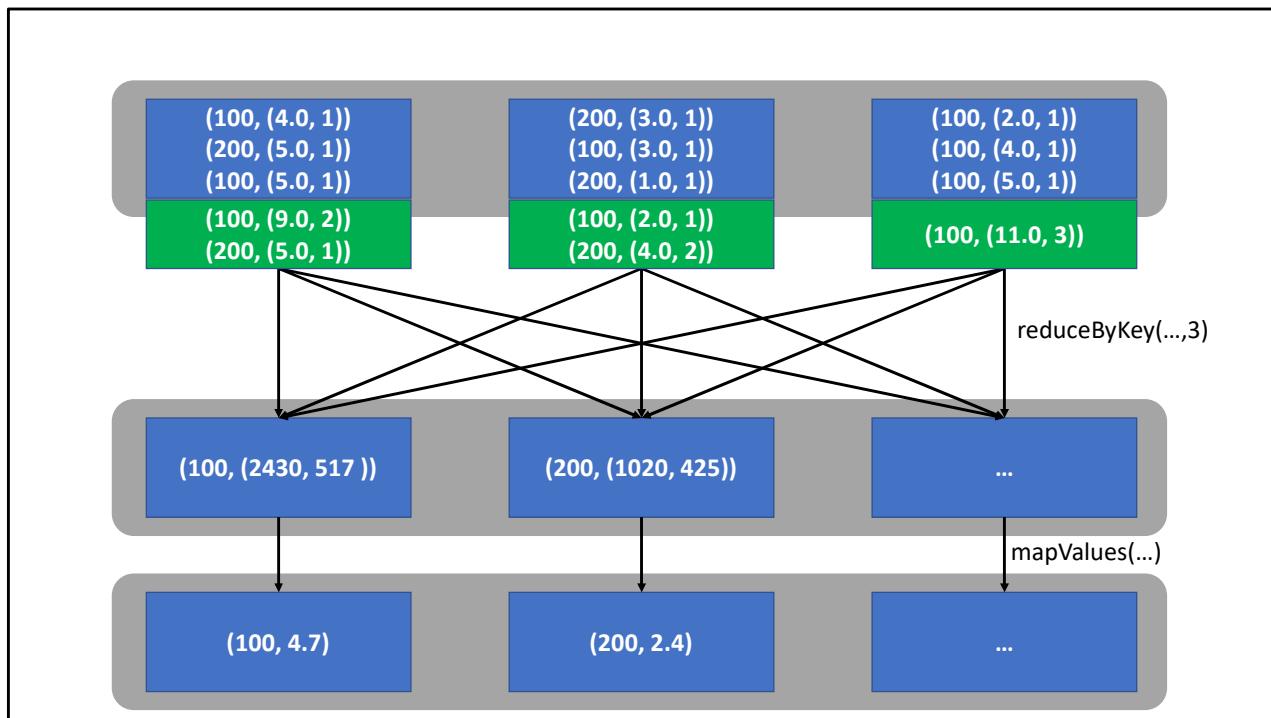
merge- combine two C

groupByKey

You (probably) don't want to use it



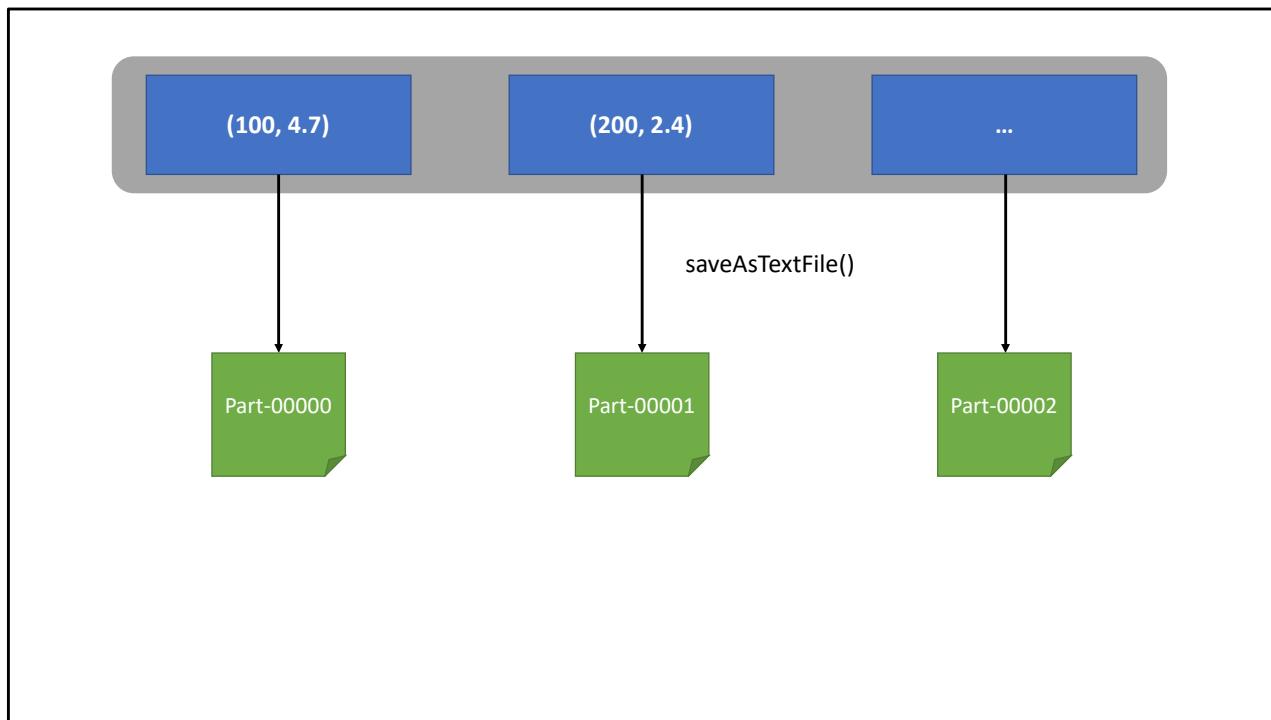
This is incomplete! If your reduce action needs to know what the key is (meaning, if some keys need to be treated differently) then groupByKey -> map or mapPartitions might be what you want.

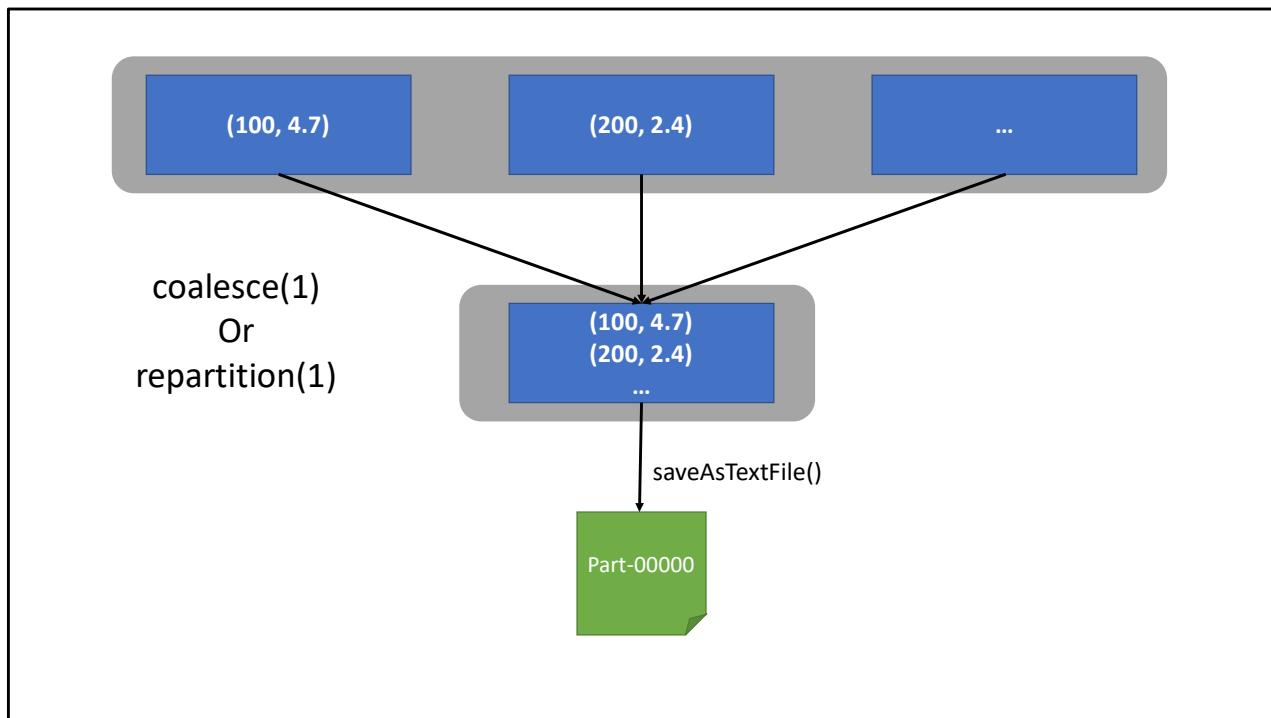


Spark's `reduceByKey` is NOT like the Reduce phase of MapReduce!

`reduceByKey` – partitions the RDD, then reduces each partition, THEN shuffles for a final reduce.

The second parameter here is optional (the default number of partitions is a Spark configuration option)





Repartition triggers shuffling but it gives more balanced partitions. It can be used to increase or decrease the number of partitions.

Coalesce can be used to only reduce the number of partitions. It avoids full shuffling so it is faster than repartition but it may give unbalanced partitions.

Just the Code (Scala)

```
sc.textFile("movies.csv").  
  map(_.split(",")).  
  map(lst => (lst(1).toInt, (lst(2).toDouble,1))).  
  reduceByKey({case ((s1,c1), (s2,c2)) =>  
    (s1 + s2, c1 + c2)}).  
  mapValues({case (sum, cnt) =>  
    sum / cnt }).  
  coalesce(1).  
  saveAsTextFile("averages")
```

Behold the power of pattern matching anonymous functions! Pattern matching is one of several reasons to love Scala

I could have written `(p1, p2) => (p1._1 + p2._1, p1._2 + p2._2)` but that's ugly!

Also...pro tip for live coding in front of an audience. Names like "count" and "cnt" are easy to typo. There was some scandalized gasps one lecture, let me tell you...

Just the Code (Python)

```
sc.textFile("movies.csv").\
    map(lambda line: line.split(",")).\
    map(lambda lst:
        (int(lst[1]), (float(lst[2]),1))).\
    reduceByKey(lambda p1, p2:
        (p1[0] + p2[0], p1[1] + p2[1])).\
    mapValues(lambda pair: pair[0] / pair[1]).\
    coalesce(1).\
    saveAsTextFile("averages")
```

D'ya like DAGs?

```
print(rdd.toDebugString())
(1) CoalescedRDD[13] at coalesce at ...
|  MapPartitionsRDD[12] at map at ...
|  ShuffledRDD[11] at reduceByKey at ...
+- (2) MapPartitionsRDD[10] at map at ...
   |  MapPartitionsRDD[9] at map at ...
   |  movies.csv MapPartitionsRDD[8] at textFile ...
   |  movies.csv HadoopRDD[7] at textFile ...
```

Read bottom-to-top

1. Text file loaded and partitioned (like MapReduce in Hadoop, this will try to allocate the jobs to workers that already have that chunk of HDFS data)
2. Map is applied to existing partitions (split the lines)
3. Map is applied to existing partitions (extract useful fields, convert to appropriate types, convert rating to (rating, 1) for averages)
4. reduceByKey triggers a repartition based on the keys (movie IDs)
5. Map is applied to the new partitions (convert (sum , count) to sum / count)
6. Coalesce merges data into 1 partition

BONUS SLIDE, NEVER SEEN IN CLASS

- For CS451 students – the Hadoop cluster page you viewed on A0 is useful for figuring out what's going on with your Spark jobs!
- If you click “ApplicationManager” you can explore the DAG graphically, including seeing all of the individual tasks created.
- Caution – if your map / flatMap is slow...it might actually be the next stage that's inefficient:
 - RDD.flatMap(...).reduceByKey(...) – as the flatMap emits pairs, they'll be combined by reduceByKey's lambda (like a MapReduce combiner).
 - If this combiner is expensive, it'll look like flatMap is slow