

<https://www.youtube.com/watch?v=7ooZ4S7Ay6Y>



Algorithms
Machines
People

amp
lab

@

Berkeley

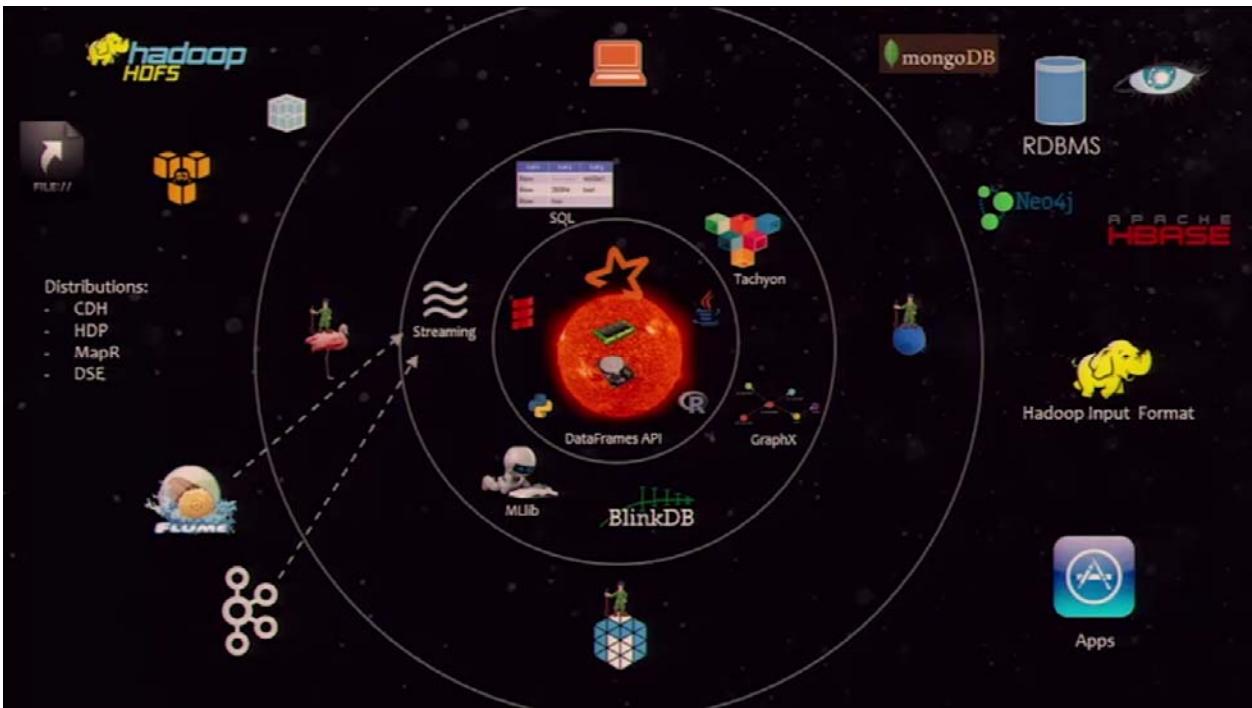
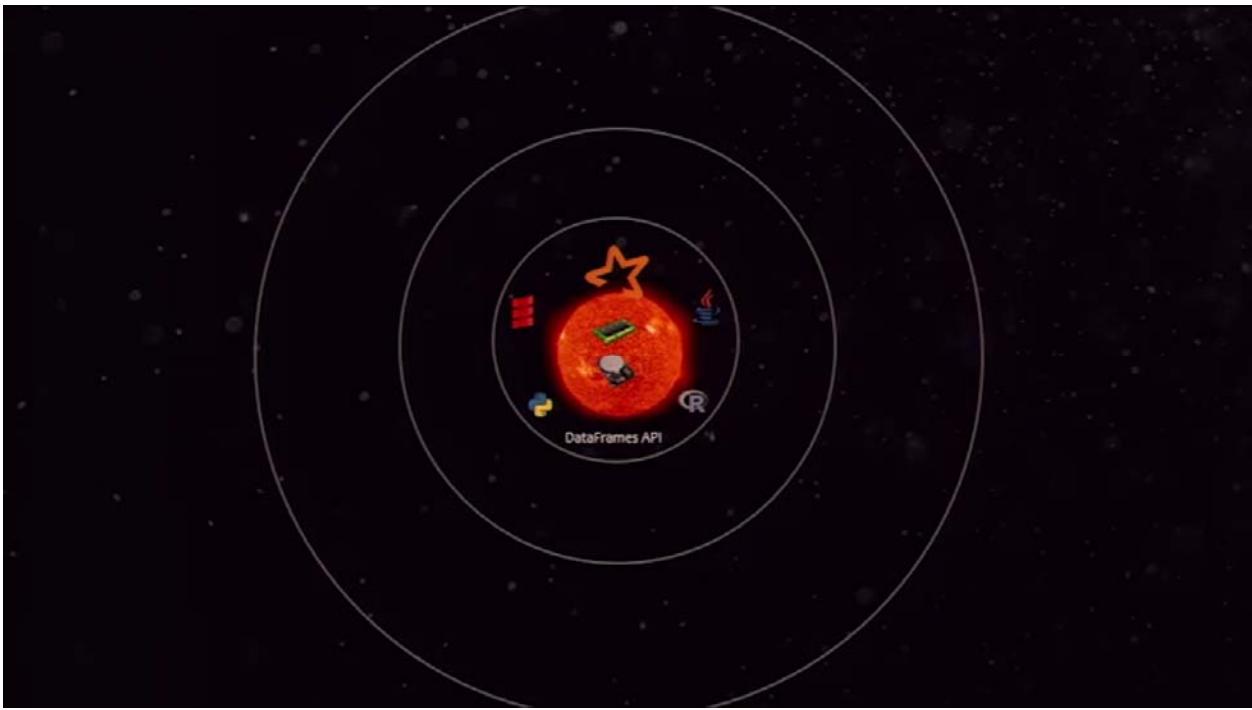
University of California

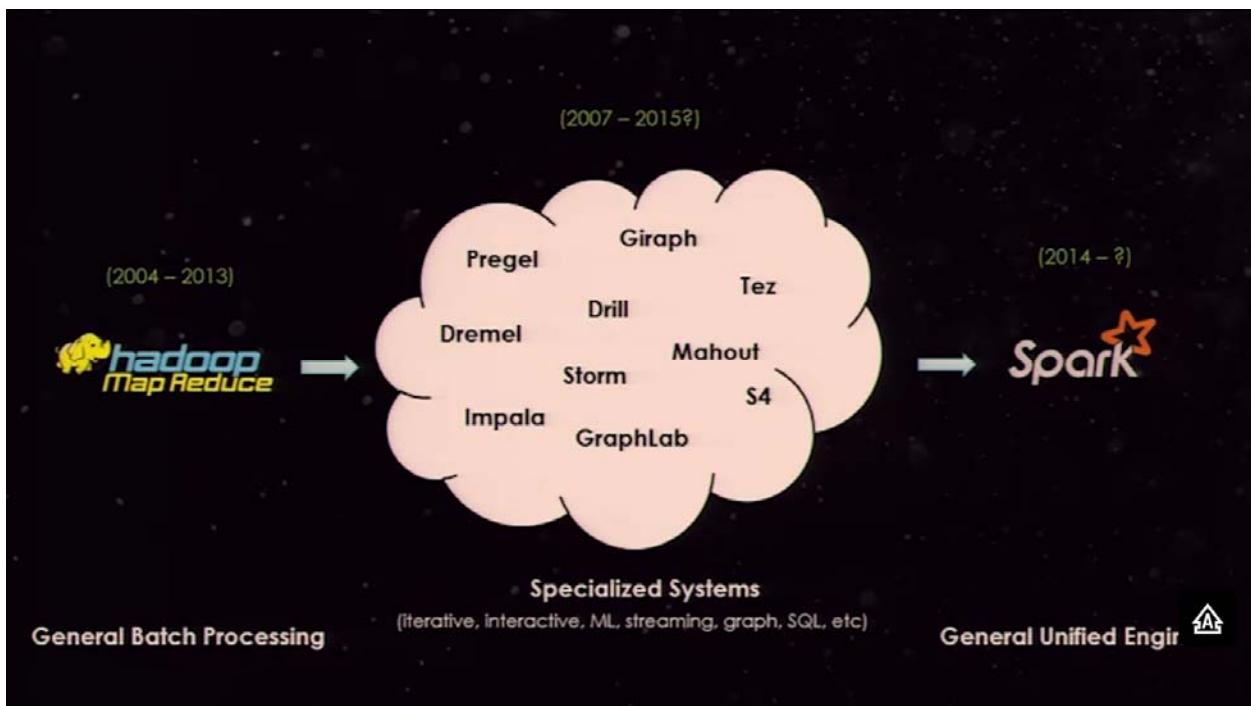
- AMPLab project was launched in Jan 2011, 6 year planned duration
- Personnel: ~65 students, postdocs, faculty & staff
- Funding from Government/Industry partnership, NSF Award, Darpa, DoE, 20+ companies
- Created BDAS, Mesos, SNAP. Upcoming projects: Succinct & Velox.

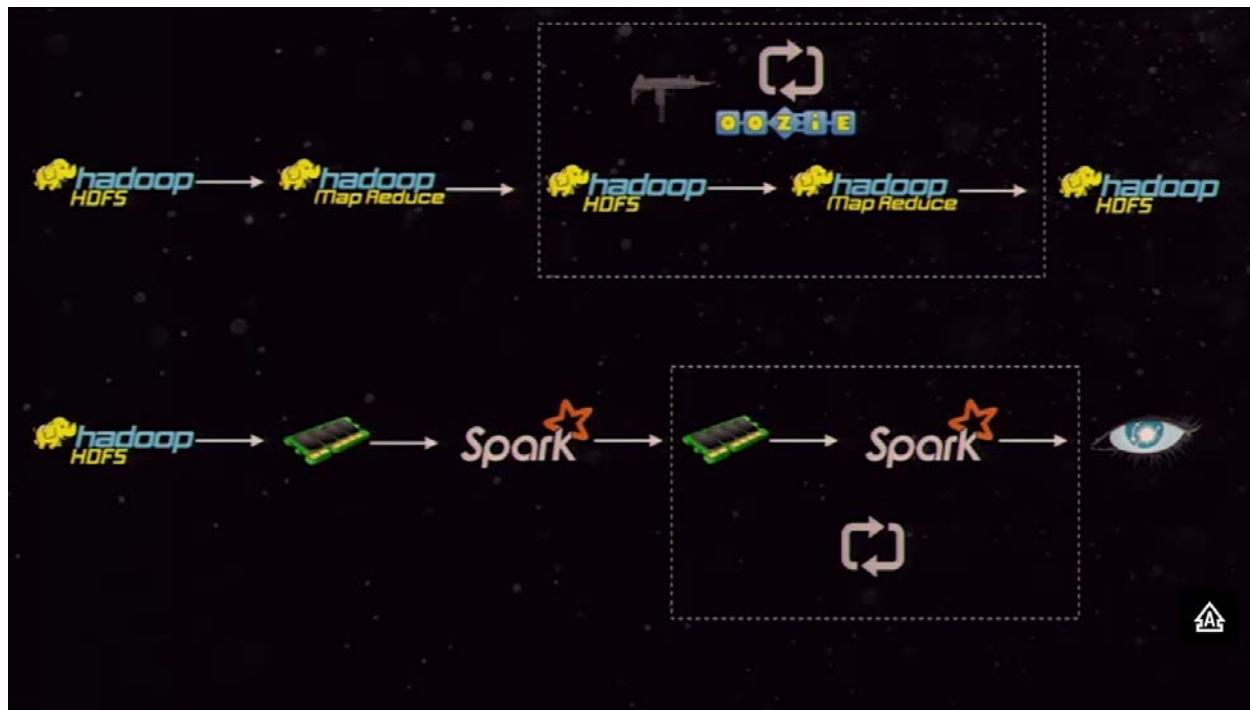
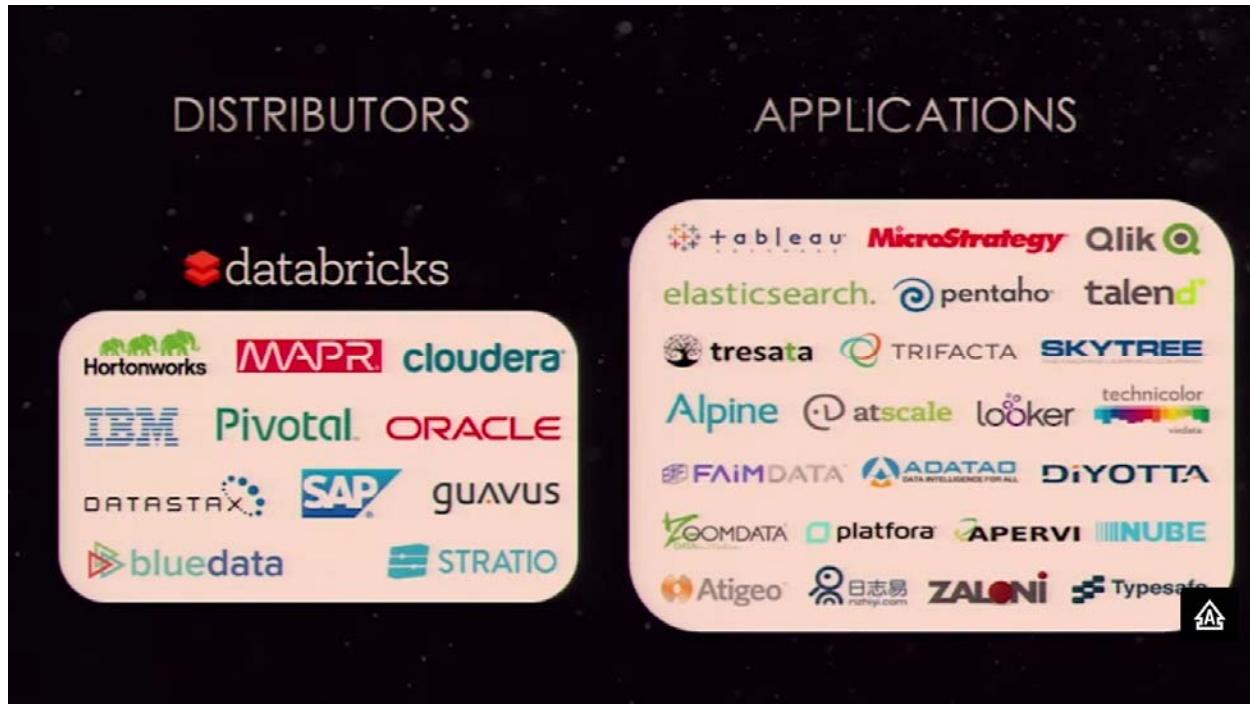
"Unknown to most of the world, the University of California, Berkeley's AMPLab has already left an indelible mark on the world of information technology, and even the web. But we haven't yet experienced the full impact of the group [...] Not even close"

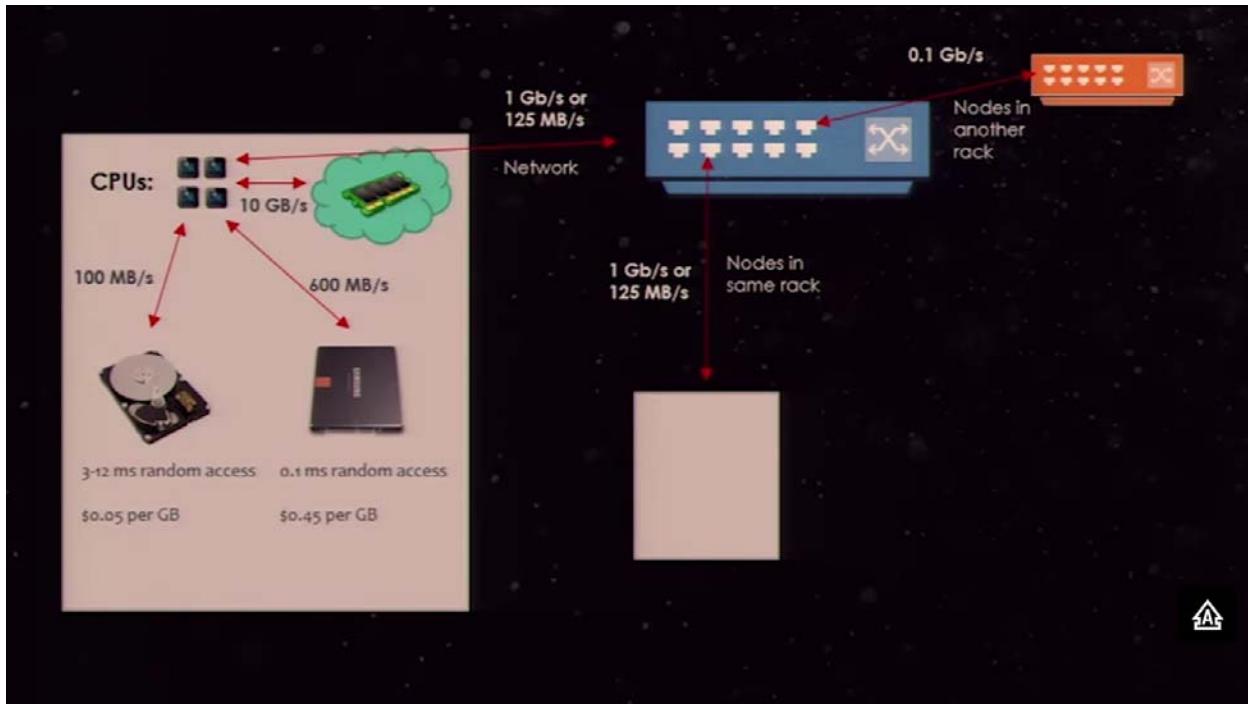
- Derrick Harris, [GigaOm, Aug 2014](#)











Resilient Distributed Datasets: A Fault-Tolerant Abstraction for In-Memory Cluster Computing

Mati Zuckerman, Mosharaf Chowdhury, Tatogata Das, Asikar Dasu, Justin Ma, Murphy McCauley, Michael J. Franklin, Scott Shenker, Ion Stoica
University of California, Berkeley

Abstract
We present Resilient Distributed Datasets (RDDs), a distributed memory abstraction that provides fault-tolerant in-memory computation for large clusters in a fault-tolerant manner. RDDs are motivated by two types of applications that current computing frameworks handle inefficiently: iterative algorithms and interactive data mining tools. In both cases, keeping data in memory can improve performance by an order of magnitude. To handle these applications, RDDs provide a fault-tolerant abstraction of shared memory based on constrained transformations rather than fine-grained updates and guarantees that RDDs are consistent and adaptive enough to capture a wide range of applications, including recent specialized programming models for iterative jobs, such as Pregel. We also show how RDDs can be used to support distributed machine learning. We have implemented RDDs as a system called Spark, which we evaluate through a variety of user applications and benchmarks.

1 Introduction
Current computing frameworks, like MapReduce [10] and Dryad [19], have been widely adopted for large-scale data analysis. These systems fit users' write parallel computations using a set of high-level operators, without forcing them to think about the underlying details.

While current frameworks provide useful abstractions for managing a cluster's computational resources, they lack abstractions for managing distributed memory. In particular, they do not support the class of computing applications that need intensive results across multiple computations. This occurs in many domains, including machine learning, streaming algorithms including Pregel [2], K-means clustering, and logistic regression. Another compelling use case is answering complex queries, such as ad-hoc queries or arbitrary queries on the same subset of the data. Until recently, in most current frameworks, the only way to reuse intermediate results was to write it to an external storage system, e.g., a distributed file system. This incurs substantial overhead due to data replication, disk I/O, and serialization.

tion, which can shorten application execution times.

Resilient distributed datasets have been developed as distributed frameworks for some applications that require iterative group computation that keeps intermediate data in memory. Examples include MapReduce [10] and Pregel [2]. However, these frameworks only support specific computation patterns (e.g., mapping a series of functions to a set of keys) and are not suitable for these patterns. They do not provide abstractions for more general cases, e.g., to fit a more broad set of distributed memory and update applications. RDDs are an attempt to address this gap. With RDDs, programmers can use a simple interface to define transformations that let them efficiently process intermediate results in memory, control their placement to optimize data placement, and manage their durability to ensure fault tolerance.

The main challenge in designing RDDs is defining a programming interface that can provide fault tolerance guarantees. This is particularly challenging for clusters, such as distributed shared memory [25], key-value stores [27], databases, and Pregel [2], offer an interface based on sequential updates to manage state consistency, i.e., to maintain consistency. While it is possible to build fault tolerance by replicating the data across machines or to log updates across machines, both approaches are not suitable for distributed memory applications that require copying large amounts of data over the cluster network, whose bandwidth is far lower than that of RAM memory.

In contrast to these systems, RDDs provide an interface based on coarse-grained transformations (e.g., map, flatMap, join) that apply the transformation to all data items. This abstraction is sufficiently powerful to provide fault tolerance by logging the transformations used to build a dataset, along with the actual data.¹ If a particular transformation fails, RDDs can recompute it using the information about how it was derived from other RDDs to reconstruct

¹Compressing the logs in some RDDs may be costly since it is a log-based approach, however, we defer this to [27].

"We present Resilient Distributed Datasets (RDDs), a distributed memory abstraction that lets programmers perform in-memory computations on large clusters in a fault-tolerant manner."

RDDs are motivated by two types of applications that current computing frameworks handle inefficiently: iterative algorithms and interactive data mining tools.

In both cases, keeping data in memory can improve performance by an order of magnitude."

"Best Paper Award and Honorable Mention for Community Award"
- NSDI 2012

- Cited 392 times!

April 2012

http://www.cs.berkeley.edu/~mattel/papers/2012/nsdi_spark.pdf

Spark BLINKDB

BlinkDB: Queries with Bounded Errors and Bounded Response Times on Very Large Data

Sameer Agarwal*, Barzan Morafizi*, Aurojit Panda*, Henry Milner*, Samuel Madden*, Jon Stoica*

*University of California, Berkeley *Massachusetts Institute of Technology *Conviva Inc.
 {sameerag, spanda, henrym, madden}@cs.berkeley.edu, {barzan, madden}@csail.mit.edu

Abstract

In this paper, we present BlinkDB, a massively parallel, approximate query engine for running interactive SQL queries on large volumes of data. BlinkDB allows users to trade off query accuracy for response time. It executes interactive queries on large amounts of data by running queries on data partitions and presenting results annotated with meaningful error bars. To achieve this, BlinkDB uses two key ideas: (1) an adaptive optimization framework that builds and maintains a set of multi-dimensional stratified samples from original data over time, and (2) a dynamic sample selection strategy that selects an appropriately sized sample based on a query's accuracy or response time requirements. We evaluate BlinkDB against the well-known TPC-H benchmarks and a real-world analytic workload derived from Conviva Inc., a company that manages video distribution over the Internet. Our experiments on a 100 node cluster show that BlinkDB can answer queries on up to 17 TiB of data in less than 2 seconds (over 20x faster than Hive), within an error of 2-10%.

1. Introduction

Modern data analytics applications involve computing aggregates over a large number of records to roll-up web clicks,

costing of large amounts of data by trading result accuracy for response time and space. These techniques include sampling [1, 14], sketches [13], and online aggregation [3]. To illustrate the utility of such techniques, consider the following simple query that computes the average SessionTime over all users originating in New York:

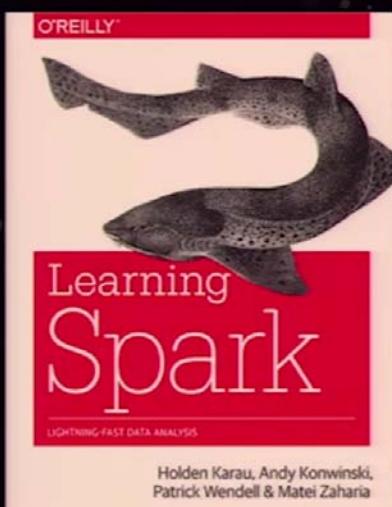
```
SELECT avg(sessionTime)
FROM Table
WHERE city='San Francisco'
WITHIN 2 SECONDS
```

Suppose the Sessions table contains 100 million tuples for New York, and cannot fit in memory. In that case, the above query may take a long time to execute, since disk reads are expensive, and such a query would need multiple disk accesses to stream through all the tuples. Suppose we instead executed the same query on a sample containing only 10,000 New York tuples, such that the entire sample fits in memory. This would be orders of magnitude faster, while still providing an approximation within a few percent of the actual value, assuming good enough for most analytical purposes. Using sampling theory we could even provide confidence bounds on the accuracy of the answer [16]. Previously described approximation techniques make different trade-offs between efficiency and the generality of the

SELECT avg(sessionTime)
FROM Table
WHERE city='San Francisco'
WITHIN 2 SECONDS

SELECT avg(sessionTime)
FROM Table
WHERE city='San Francisco'
ERROR 0.1 CONFIDENCE 95.0%

<https://www.cs.berkeley.edu/~sameerag/blinkdb-eurosys13.pdf> 

 O'REILLY

<http://shop.oreilly.com/product/0636920028512.do>

eBook: \$33.99 PDF, ePub, Mobi, DAISY
 Print: \$39.99 Shipping now!

\$30 @ Amazon:
<http://www.amazon.com/Learning-Spark-Lightning-Fast-Data-Analysis/dp/1449358624>

The screenshot shows the Apache Spark community page. At the top, there's a navigation bar with links for Download, Libraries, Documentation, Examples, Community, and FAQ. Below the navigation is a section titled "Latest News" which lists several recent releases and events. A red arrow points to a section titled "Mailing Lists" where it says: "Get help using Spark or contribute to the project on our mailing lists:" followed by two bullet points: "user@spark.apache.org is for usage questions, help, and announcements. ([subscribe](#)) ([unsubscribe](#)) ([archives](#))" and "dev@spark.apache.org is for people who want to contribute code to Spark. ([subscribe](#)) ([unsubscribe](#)) ([archives](#))".

The screenshot shows a Databricks lab document titled "Lab: Intro to Spark 1.1 on DSE 4.6". The document includes a logo for "Cassandra + Spark" and a link to "http://tinyurl.com/dsesparklab". It provides details about the lab creation date (Sept 2, 2014), last update (Dec 9, 2014), and contributors (Piotr Kolasowski, Holden Karau, Pat McDonough, Patrick Wendell, Matei Zaharia). It also specifies an estimated completion time of 2.5 hours and a Creative Commons license. The "Objective" section states that the lab will introduce users to using Apache Spark 1.1 on DataStax Enterprise Edition 4.6.0 in the Amazon cloud, assuming the audience is a beginner to both Cassandra and Spark. The "High Level Steps" section lists four steps: connecting via SSH, creating a new keyspace and table in C*, adding data, starting the Scala-based Spark shell, and importing fresh data into a Spark RDD.

102 pages

DevOps style

For complete beginners

Includes:

- Spark Streaming
- Dangers of GroupByKey vs. ReduceByKey

Labs: Intro to HDFS/YARN & Apache Spark on CDH 5.2



Lab created on: Dec, 2014
(please send edits and corrections to): sameer@databricks.com

Estimated lab completion time: 2 hours (spread throughout the day)

License: 

Objective:
This lab will introduce you to using 3 Hadoop ecosystem components in Cloudera's distribution: HDFS, Spark 1.1.0 and YARN. The lab will first walk you through the Cloudera Manager installation on a VM in AWS, followed by a CDH 5.2 binaries deployment on the same node. Then the lab will introduce students to Hadoop in a DevOps manner: experimenting with the distributed file system; looking at the XML config files; running a batch analytics workload with Spark from disk and from memory; writing some simple scala Spark code; running SQL commands with Spark SQL; breaking things and troubleshooting issues; etc.

The following high level steps are in the initial part of this lab:

- Connect via SSH to your Amazon instance
- Install Cloudera Manager and CDH 5.2
- Create a new folder in HDFS and add data files to it
- Start the scala based **Spark shell**
- Import the fresh data into Spark a RDD

<http://tinyurl.com/cdhsparklab>

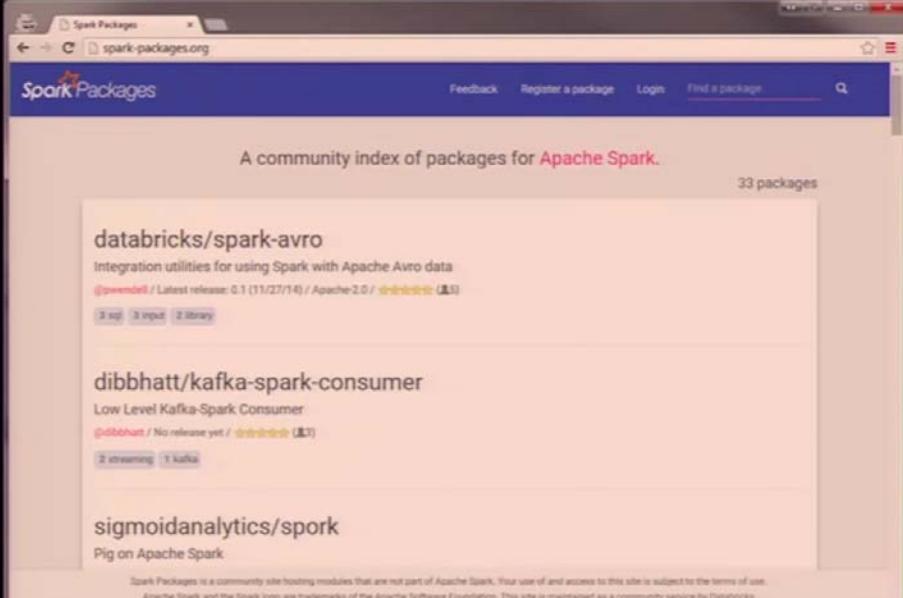
109 pages

DevOps style

For complete beginners

Includes:

- PySpark
- Spark SQL
- Spark-submit



A community index of packages for **Apache Spark**.

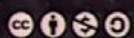
33 packages

databricks/spark-avro
Integration utilities for using Spark with Apache Avro data
[@pwendell](#) / Latest release: 0.1 (11/27/14) / Apache-2.0 / [View details](#) (1.5)

dibbhatt/kafka-spark-consumer
Low Level Kafka-Spark Consumer
[@dibbhatt](#) / No release yet / [View details](#) (1.0)

sigmoidanalytics/spork
Pig on Apache Spark

Spark Packages is a community site hosting modules that are not part of Apache Spark. Your use of and access to this site is subject to the terms of use. Apache Spork and the Spark logo are trademarks of the Apache Software Foundation. This site is maintained as a community service by Databricks.

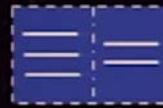


DEVOPS ADVANCED CLASS

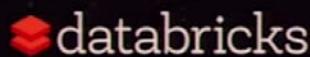
March 2015: Spark Summit East 2015

<http://spark-summit.org/east/training/devops>

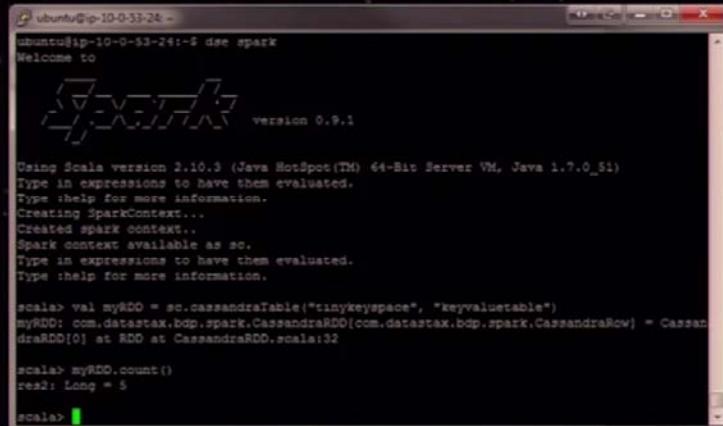
www.linkedin.com/in/blueplastic



RDD FUNDAMENTALS



INTERACTIVE SHELL



A screenshot of a terminal window showing a Scala interactive shell session. The session starts with the Spark logo and version information (version 0.9.1). It then creates a SparkContext and defines an RDD named myRDD from a CassandraTable. Finally, it prints the count of the RDD, which is 5.

```
ubuntu@ip-10-0-53-24:~$ spark
Welcome to
      ____          _ _ _   _ _ _ 
     / \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ 
    /   \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ 
   /     \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ 
  /       \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ 
 /         \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ 
version 0.9.1

Using Scala version 2.10.3 (Java HotSpot(TM) 64-Bit Server VM, Java 1.7.0_51)
Type in expressions to have them evaluated.
Type :help for more information.

Creating SparkContext...
Created spark context.
Spark context available as sc.
Type in expressions to have them evaluated.
Type :help for more information.

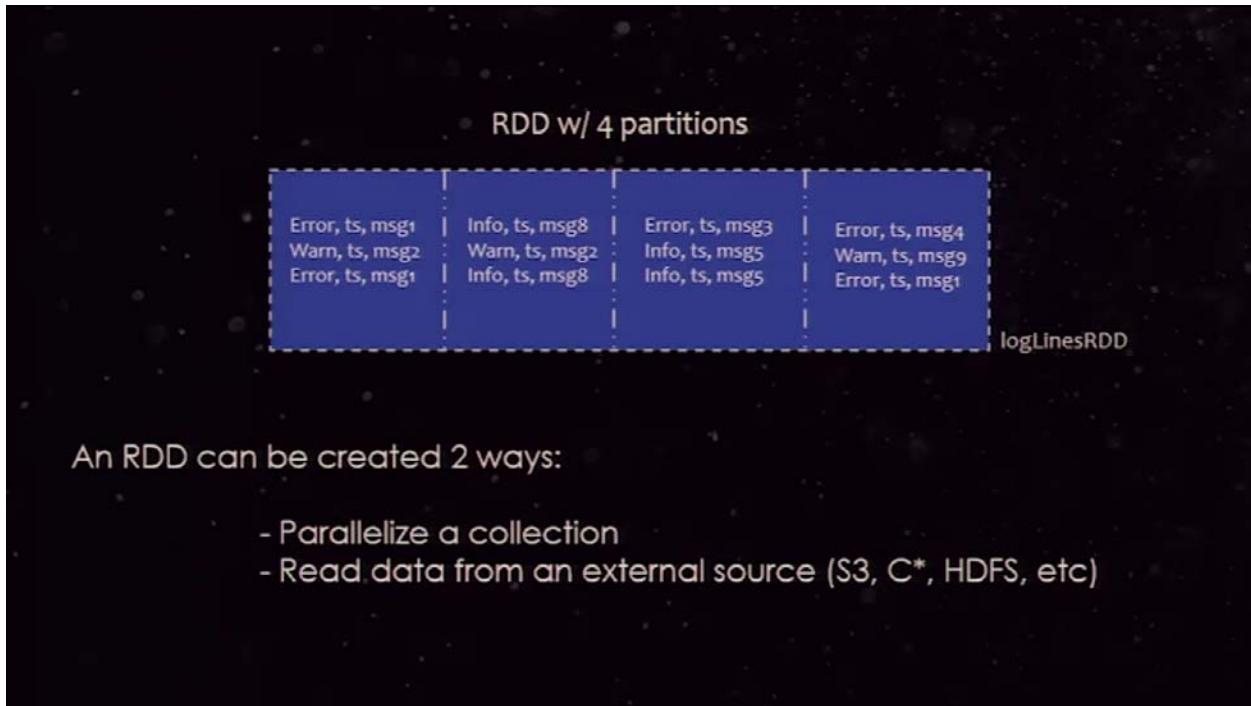
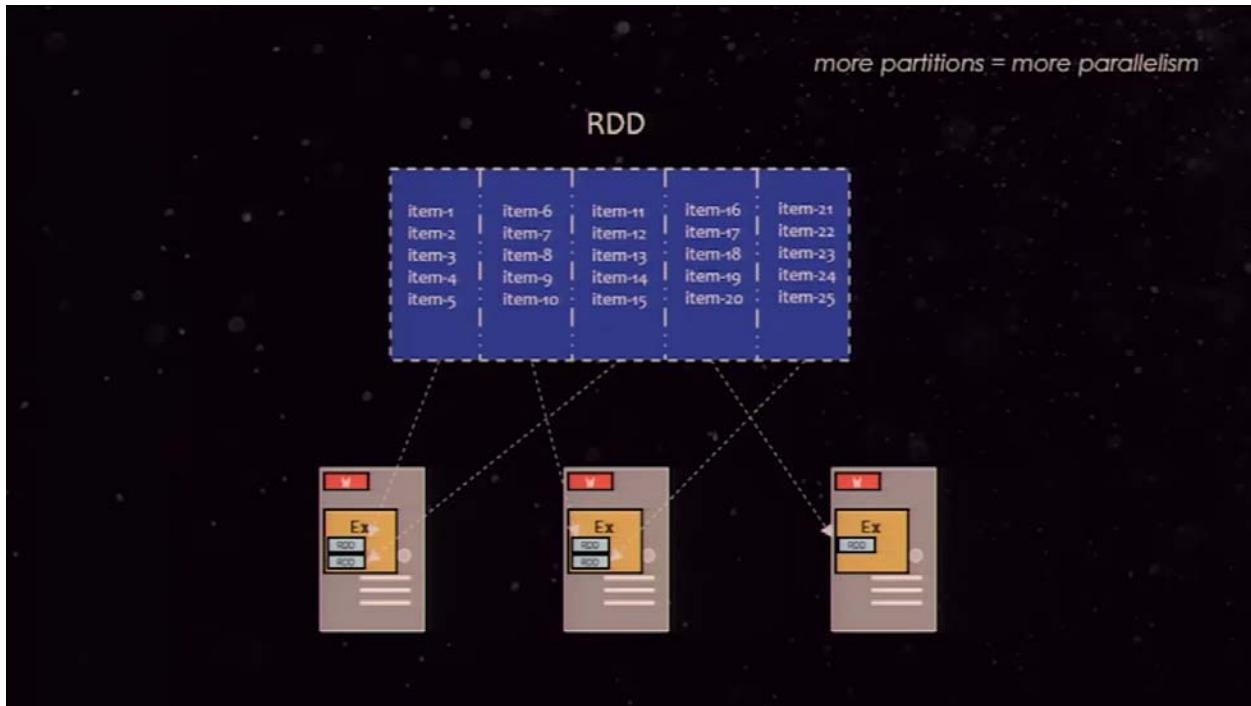
scala> val myRDD = sc.cassandraTable("tinykeyspace", "keyvaluetable")
myRDD: com.datastax.bdp.spark.CassandraRDD[com.datastax.bdp.spark.CassandraRow] = Cassan
draRDD(0) at RDD at CassandraRDD.scala:32

scala> myRDD.count()
res2: Long = 5

scala>
```

(Scala & Python only)





PARALLELIZE



```
# Parallelize in Python
wordsRDD = sc.parallelize(["fish", "cats", "dogs"])
```



```
// Parallelize in Scala
val wordsRDD = sc.parallelize(List("fish", "cats", "dogs"))
```



```
// Parallelize in Java
JavaRDD<String> wordsRDD = sc.parallelize(Arrays.asList("fish", "cats", "dogs"));
```

- Take an existing in-memory collection and pass it to SparkContext's parallelize method

- Not generally used outside of prototyping and testing since it requires entire dataset in memory on one machine

READ FROM TEXT FILE



```
# Read a local txt file in Python
linesRDD = sc.textFile("/path/to/README.md")
```

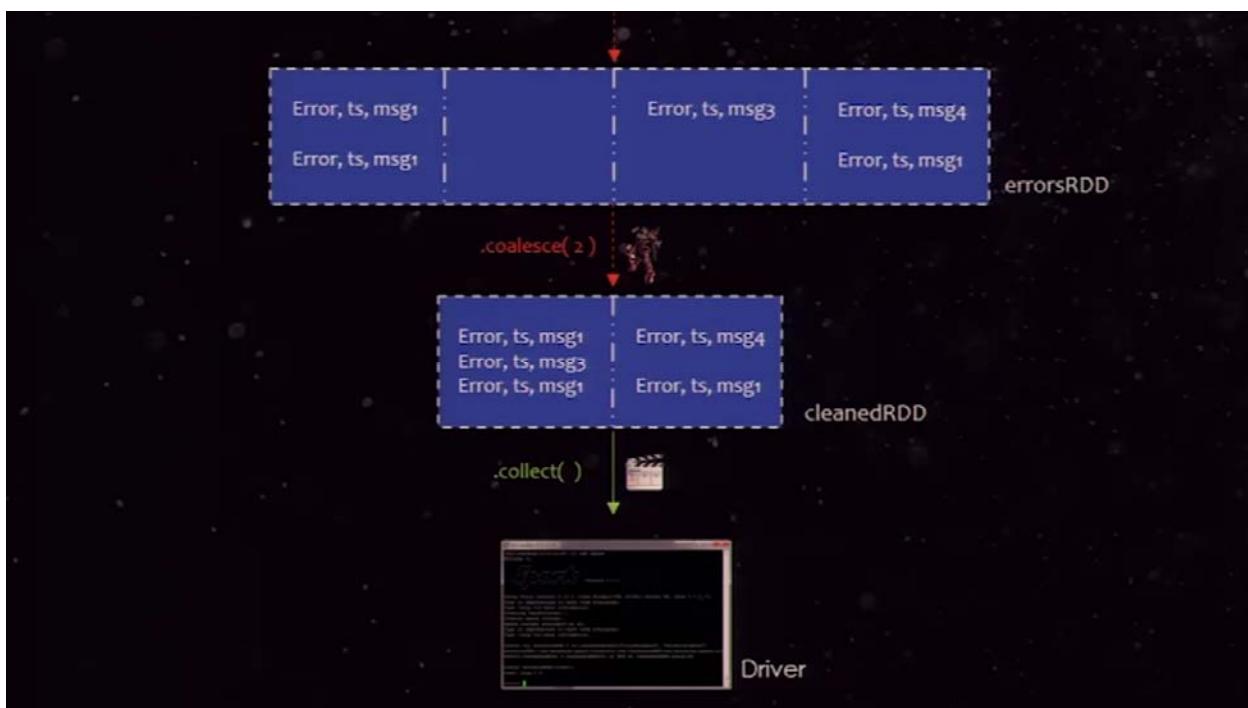
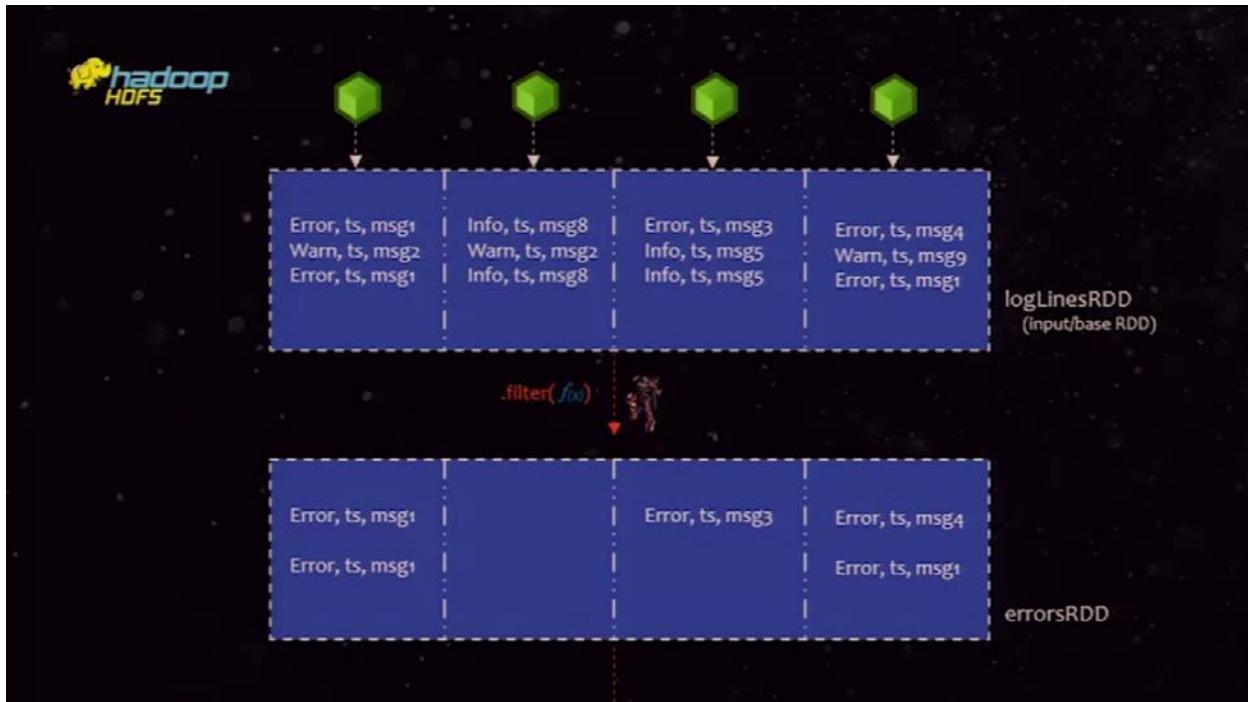
- There are other methods to read data from HDFS, C*, S3, HBase, etc.

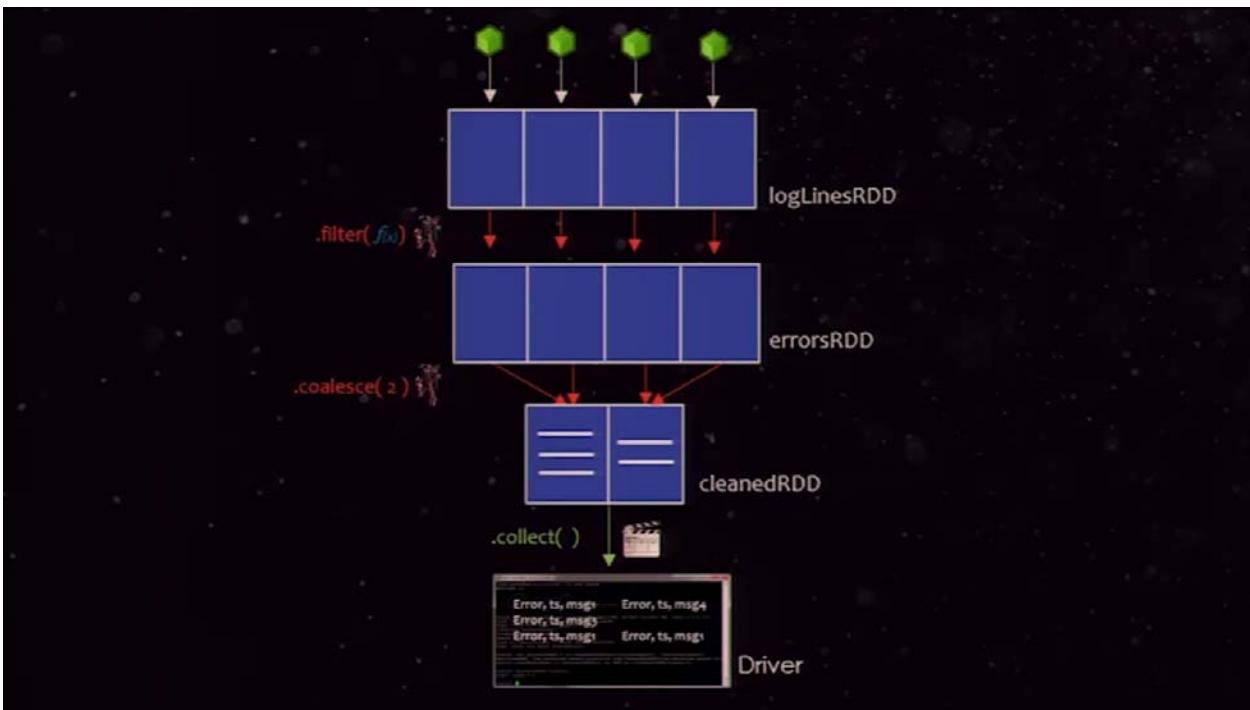


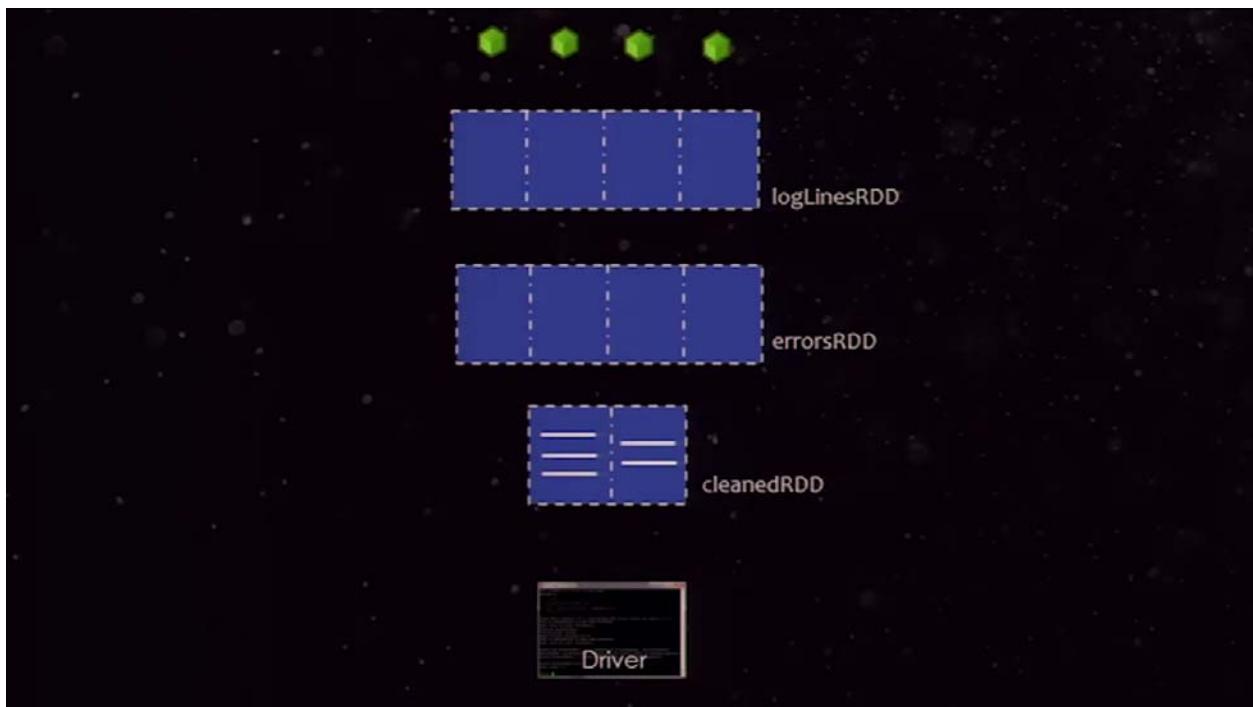
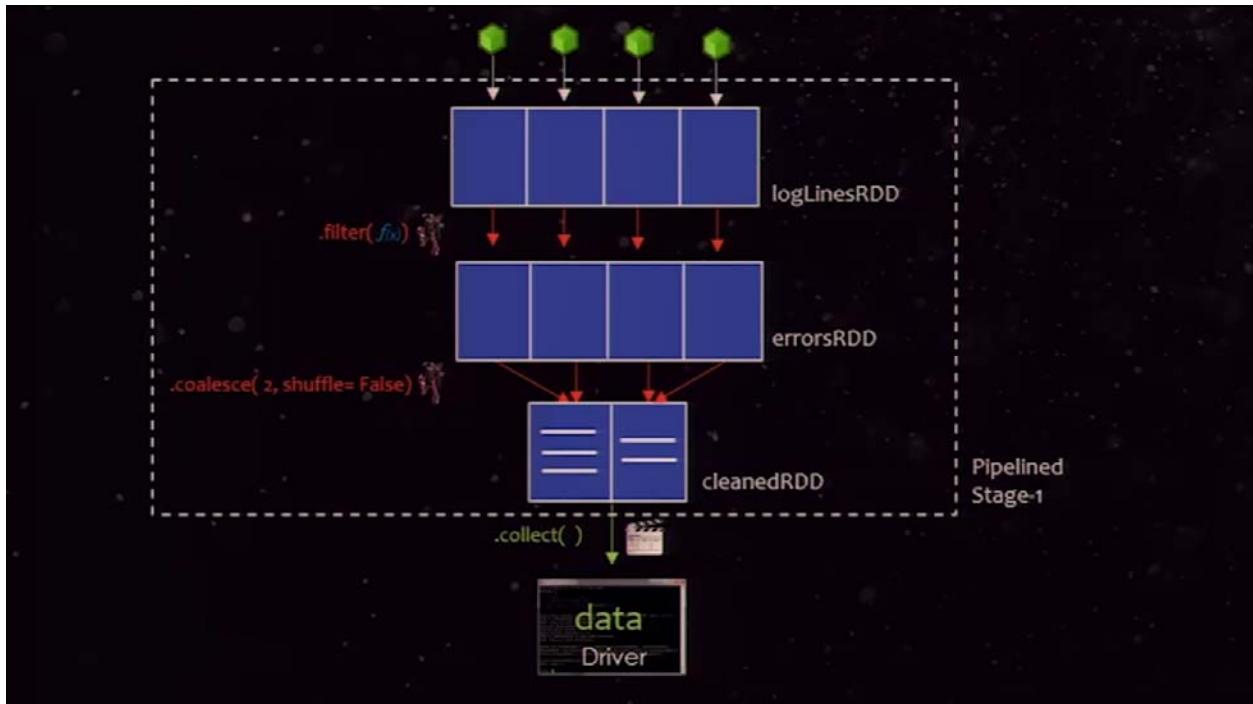
```
// Read a local txt file in Scala
val linesRDD = sc.textFile("/path/to/README.md")
```

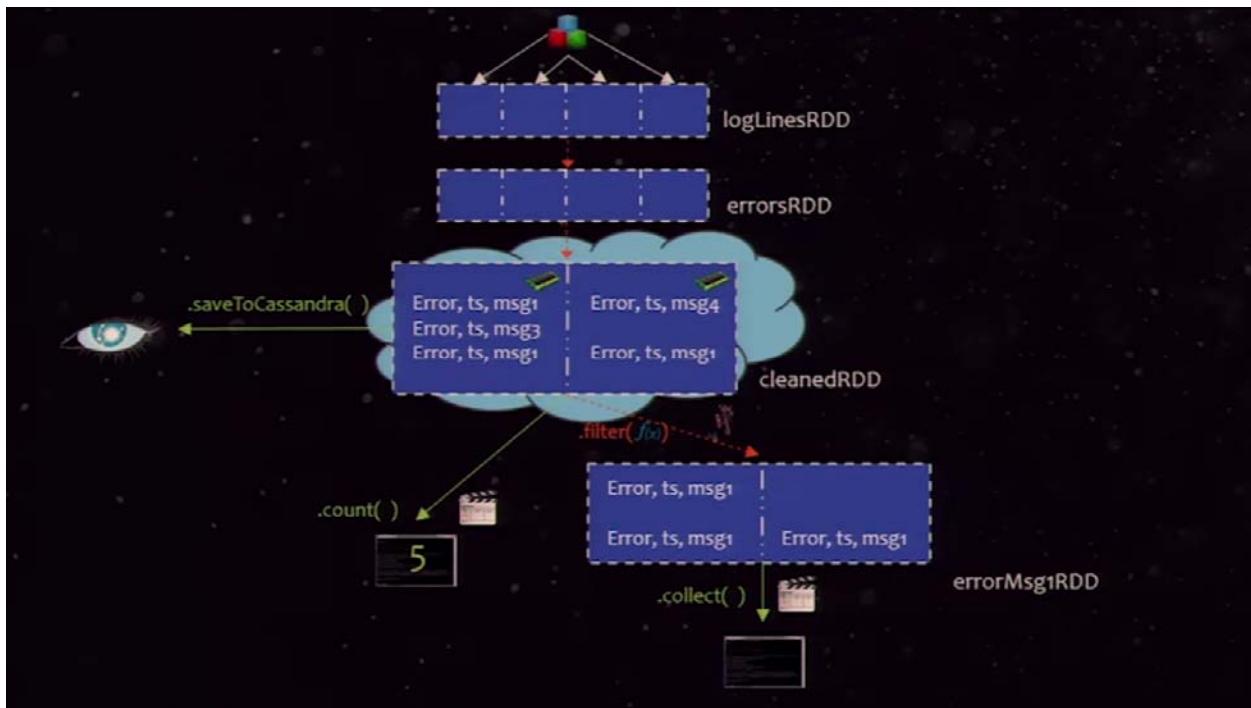
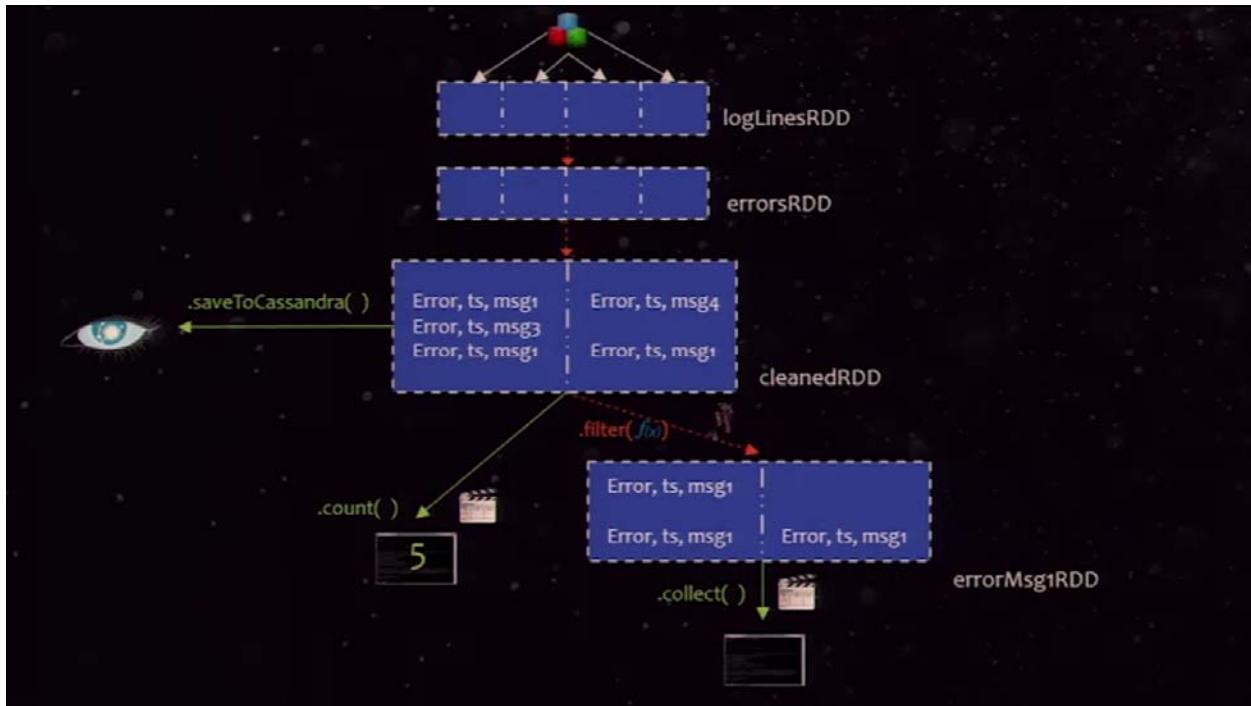


```
// Read a local txt file in Java
JavaRDD<String> lines = sc.textFile("/path/to/README.md");
```



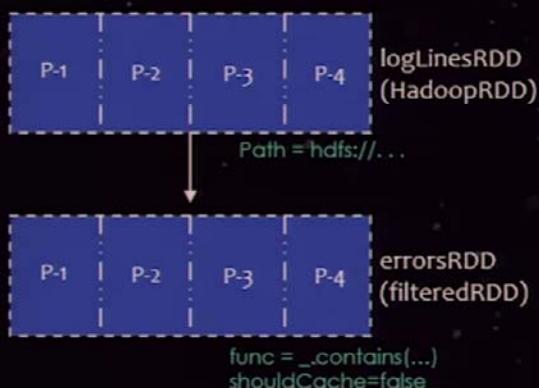




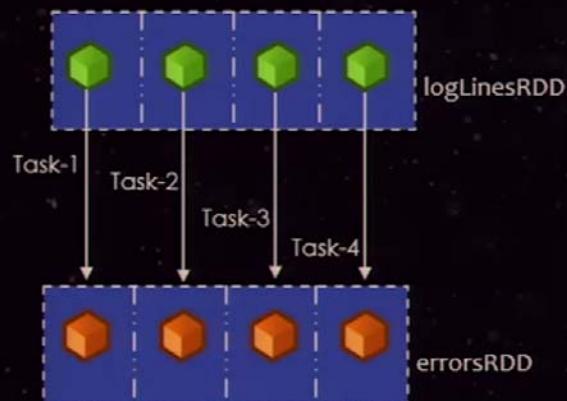


RDD GRAPH

Dataset-level view:



Partition-level view:



LIFECYCLE OF A SPARK PROGRAM

- 1) Create some input RDDs from external data or parallelize a collection in your driver program.
- 2) Lazily transform them to define new RDDs using transformations like `filter()` or `map()`.
- 3) Ask Spark to `cache()` any intermediate RDDs that will need to be reused.
- 4) Launch actions such as `count()` and `collect()` to kick off a parallel computation, which is then optimized and executed by Spark.

TRANSFORMATIONS (lazy)

map()	intersection()	cartesian()
flatMap()	distinct()	pipe()
filter()	groupByKey()	coalesce()
mapPartitions()	reduceByKey()	repartition()
mapPartitionsWithIndex()	sortByKey()	partitionBy()
sample()	join()	...
union()	cogroup()	...

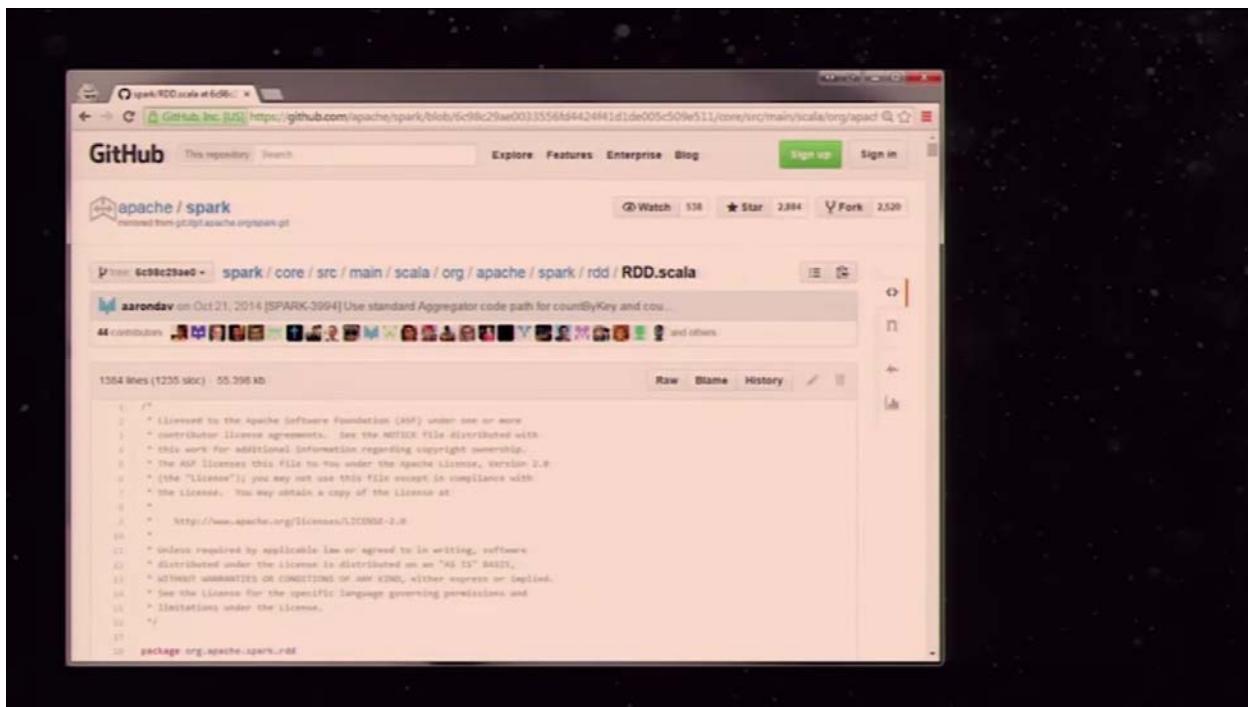
- Most transformations are element-wise (they work on one element at a time), but this is not true for all transformations

ACTIONS

reduce()	takeOrdered()
collect()	saveAsTextFile()
count()	saveAsSequenceFile()
first()	saveAsObjectFile()
take()	countByKey()
takeSample()	foreach()
saveToCassandra()	...

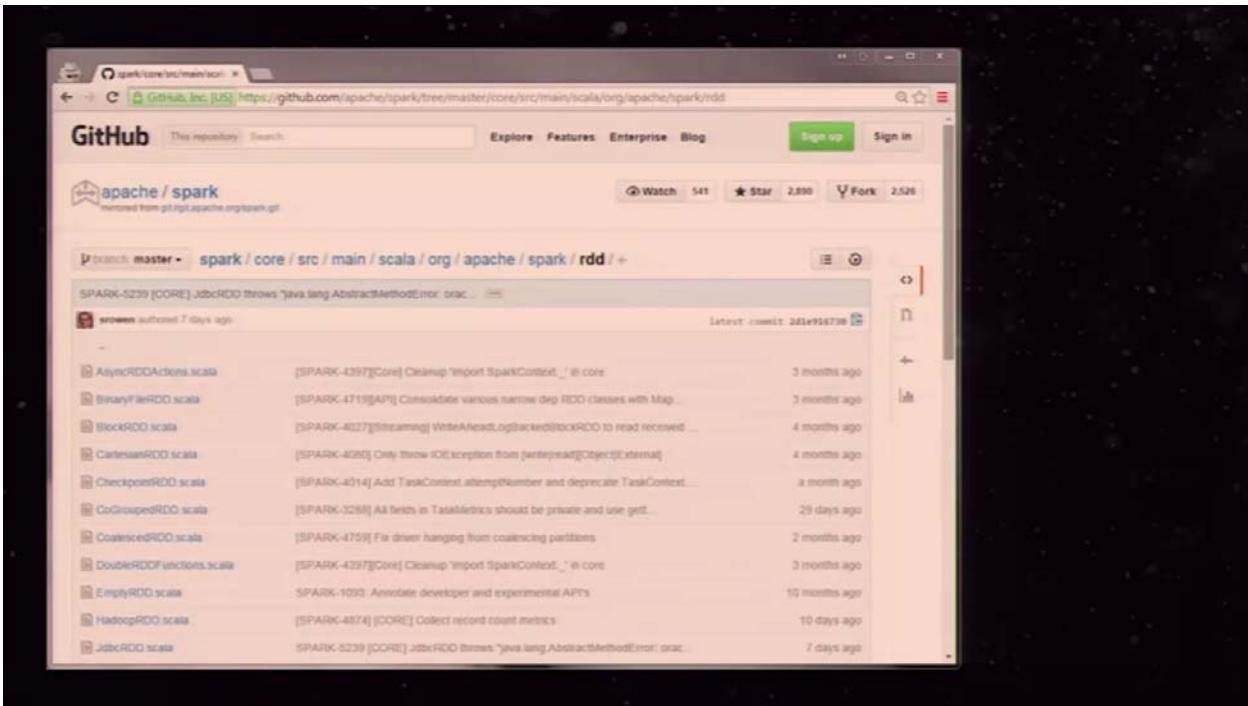
TYPES OF RDDS

- HadoopRDD
- FilteredRDD
- MappedRDD
- PairRDD
- ShuffledRDD
- UnionRDD
- PythonRDD
- DoubleRDD
- JdbcRDD
- JsonRDD
- SchemaRDD
- VertexRDD
- EdgeRDD
- CassandraRDD (*DataStax*)
- GeoRDD (*ESRI*)
- EsSpark (*ElasticSearch*)



The screenshot shows a GitHub repository page for the Apache Spark project. The URL is <https://github.com/apache/spark/blob/6c96c5c9ae0033556d4424ff1d1de005c509e511/core/src/main/scala/org/apache/rdd/RDD.scala>. The page displays the code for the RDD.scala file. The code is licensed under the Apache Software License (ASF) version 2.0. It includes a header with copyright information, a notice about the Apache Software Foundation, and a package declaration at the bottom.

```
/*
 * Licensed to the Apache Software Foundation (ASF) under one or more
 * contributor license agreements. See the NOTICE file distributed with
 * this work for additional information regarding copyright ownership.
 * The ASF licenses this file to you under the Apache License, Version 2.0
 * (the "License"); you may not use this file except in compliance with
 * the License. You may obtain a copy of the License at
 *
 *   http://www.apache.org/licenses/LICENSE-2.0
 *
 * Unless required by applicable law or agreed to in writing, software
 * distributed under the License is distributed on an "AS IS" BASIS,
 * WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied.
 * See the License for the specific language governing permissions and
 * limitations under the License.
 */
package org.apache.spark.rdd
```



RDD INTERFACE

- * 1) Set of **partitions** ("splits")
- * 2) List of **dependencies** on parent RDDs
- * 3) Function to **compute** a partition given parents
- * 4) Optional **preferred locations**
- * 5) Optional **partitioning info** for k/v RDDs (Partitioner)

This captures all current Spark operations!

EXAMPLE: HADOOPRDD

- * Partitions = one per HDFS block
- * Dependencies = none
- * Compute (partition) = read corresponding block

- * preferredLocations (part) = HDFS block location
- * Partitioner = none

EXAMPLE: FILTEREDRDD

- * Partitions = same as parent RDD
 - * Dependencies = "one-to-one" on parent
 - * Compute (partition) = compute parent and filter it
-
- * preferredLocations (part) = none (ask parent)
 - * Partitioner = none

EXAMPLE: JOINEDRDD

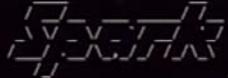
- * Partitions = One per reduce task
 - * Dependencies = "shuffle" on each parent
 - * Compute (partition) = read and join shuffled data
-
- * preferredLocations (part) = none
 - * Partitioner = HashPartitioner(numTasks)

READING DATA USING THE C* CONNECTOR

```
    Keyspace   Table  
↓           ↓  
val cassandraRDD = sc  
      .cassandraTable("ks", "mytable")  
  Server side column { .select("col-1", "col-3")  
& row selection { .where("col-5 = ?", "blue")  
}
```

INPUT SPLIT SIZE

(for dealing with wide rows)

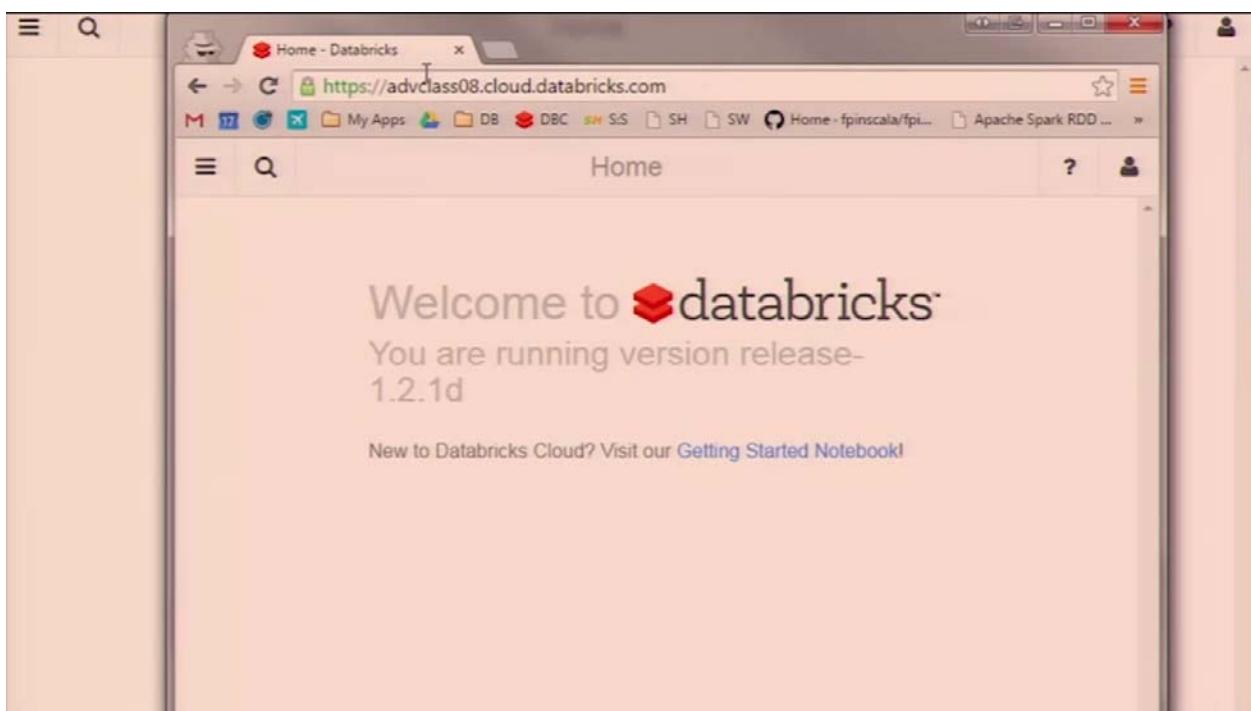
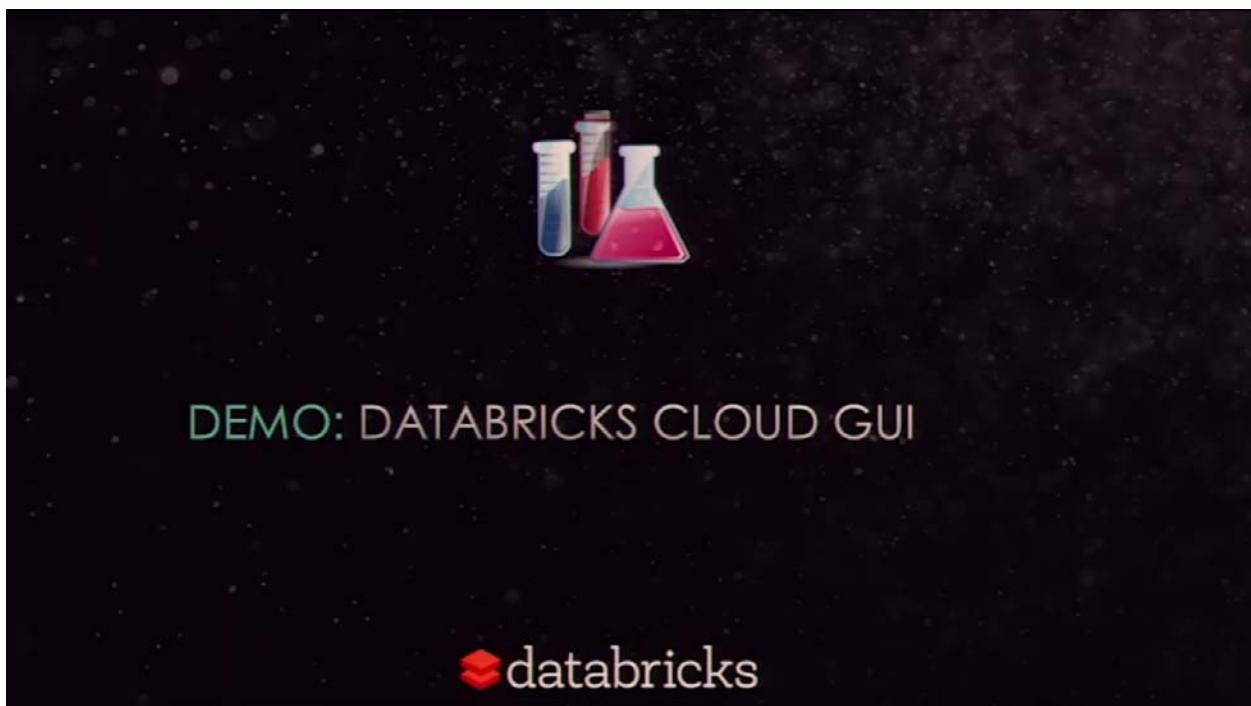
```
Start the Spark shell by passing in a custom cassandra.input.split.size:  
ubuntu@ip-10-0-53-24:~$ dse spark -Dspark.cassandra.input.split.size=2000  
Welcome to  
 version 0.9.1  
  
Using Scala version 2.10.3 (Java HotSpot(TM) 64-Bit Server VM, Java  
1.7.0_51)  
Type in expressions to have them evaluated.  
Type :help for more information.  
Creating SparkContext...  
Created spark context.  
Spark context available as sc.  
Type in expressions to have them evaluated.  
Type :help for more information.  
scala>
```

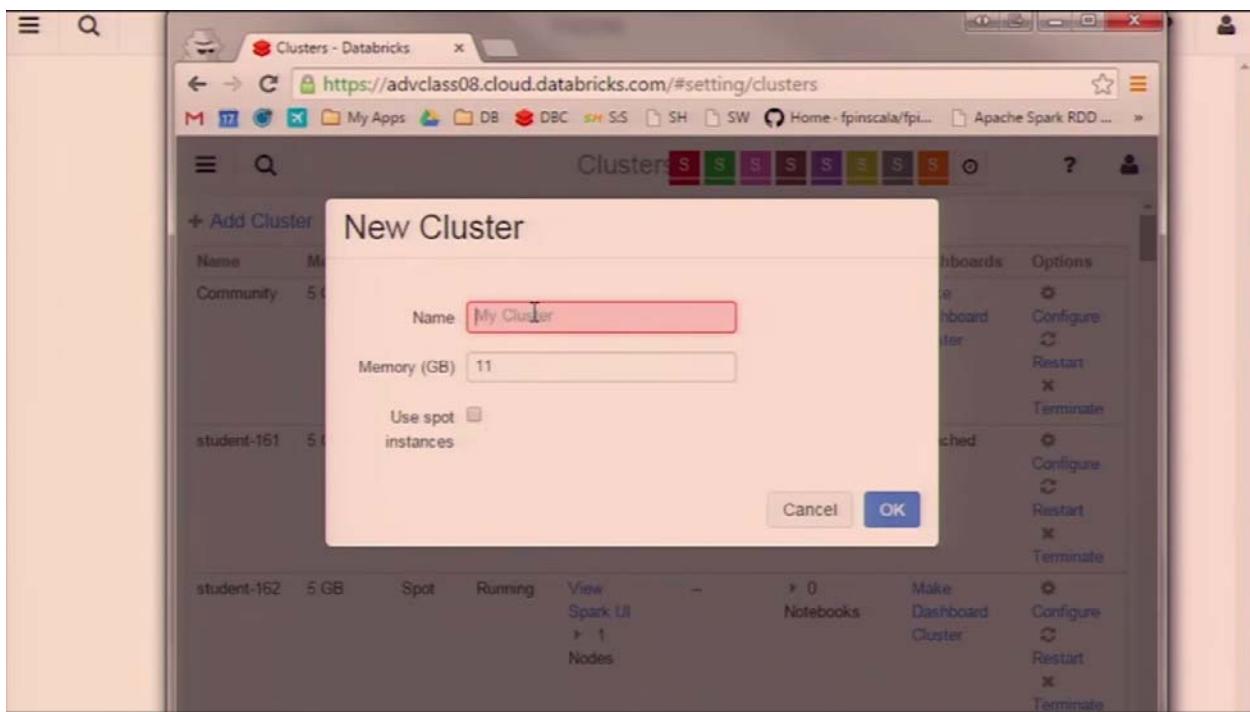
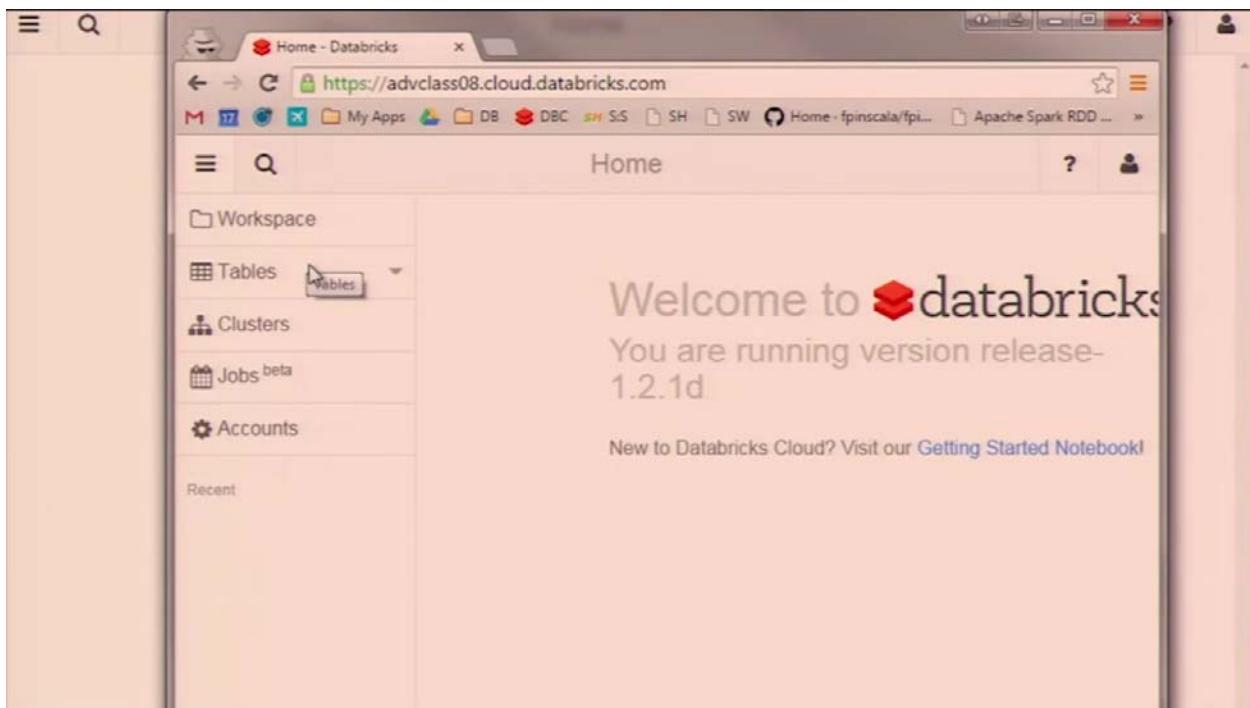
The `cassandra.input.split.size` parameter defaults to 100,000. This is the approximate number of physical rows in a single Spark partition. If you have really wide rows (thousands of columns), you may need to lower this value. The higher the value, the fewer Spark tasks are created. Increasing the value too much may limit the parallelism level."

<https://github.com/datastax/spark-cassandra-connector>

- Open Source
- Implemented mostly in Scala
- Scala + Java APIs
- Does automatic type conversions







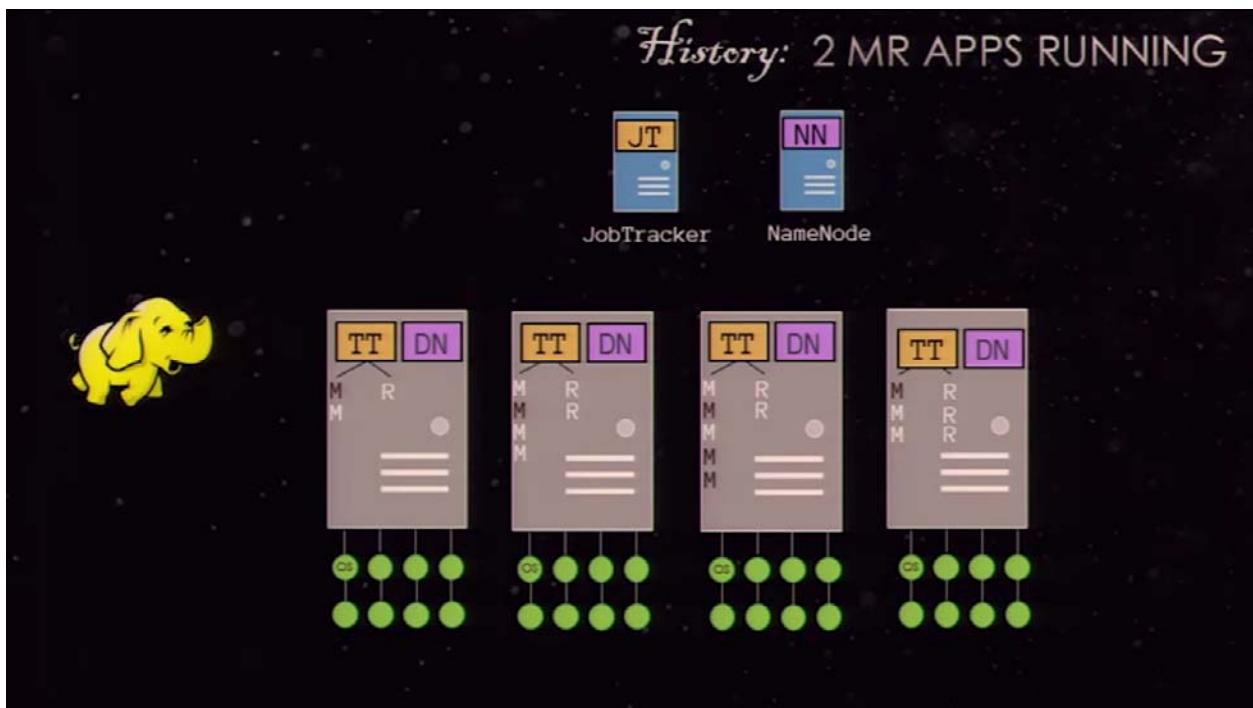
WAYS TO RUN SPARK

-  - Local 
-  - Standalone Scheduler 
-  - YARN 
-  - Mesos 

WAYS TO RUN SPARK

-  - Local 
 -  - Standalone Scheduler 
 -  - YARN 
 -  - Mesos 
- Static Partitioning
- Dynamic Partitioning

History: 2 MR APPS RUNNING





LOCAL MODE



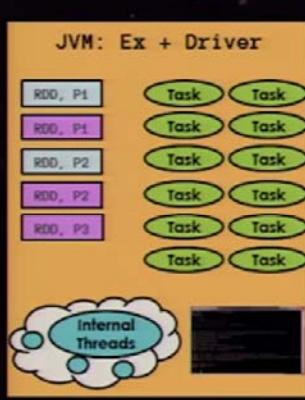
LOCAL MODE

3 options:
- local
- local[N]
- local[*]

> ./bin/spark-shell --master local[12]

> ./bin/spark-submit --name "MyFirstApp"
--master local[12] myApp.jar

```
val conf = new SparkConf()  
    .setMaster("local[12]")  
    .setAppName("MyFirstApp")  
    .set("spark.executor.memory", "3g")  
val sc = new SparkContext(conf)
```



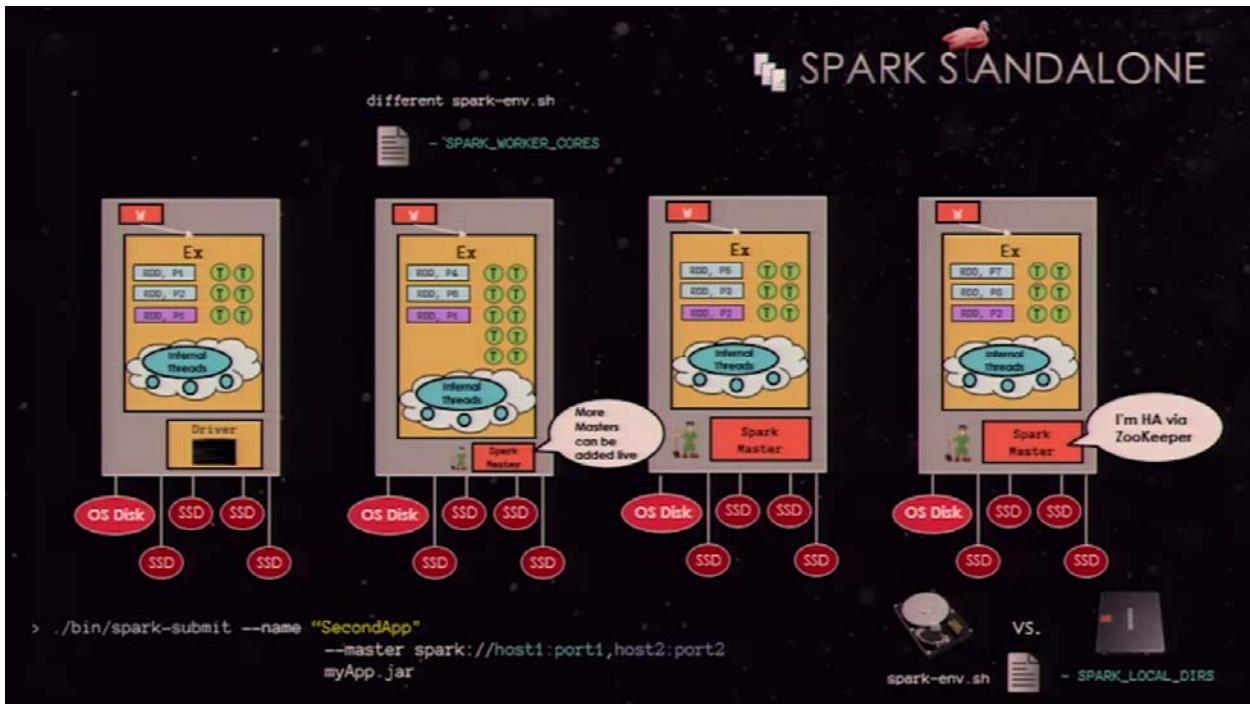
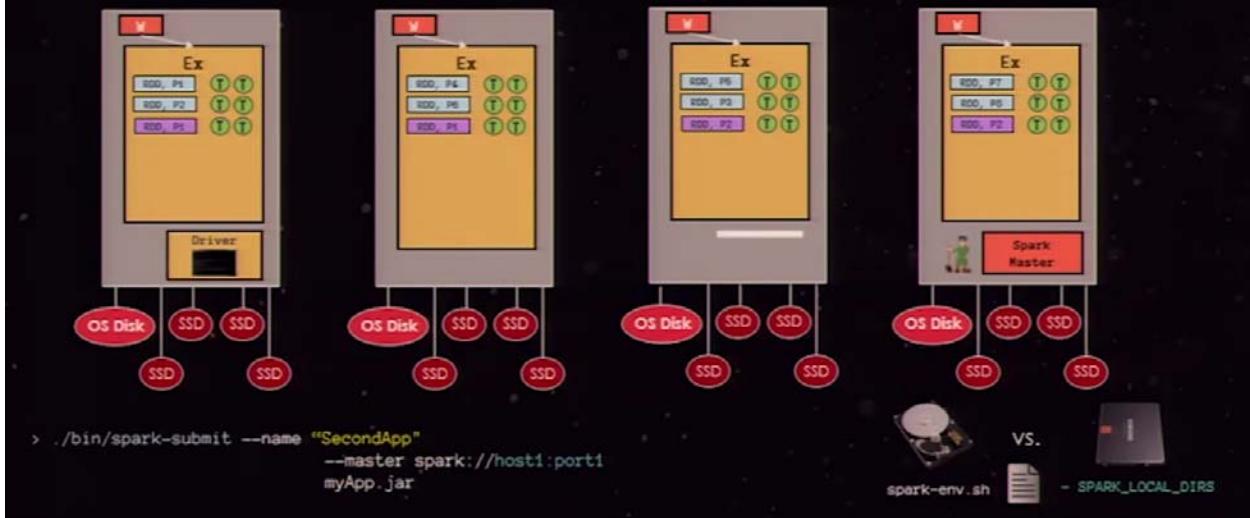
The diagram shows four identical host machines represented by grey rectangles. Each host has three horizontal white bars inside and a small grey circle at the top right. Below each host is a red oval containing the text "OS Disk".

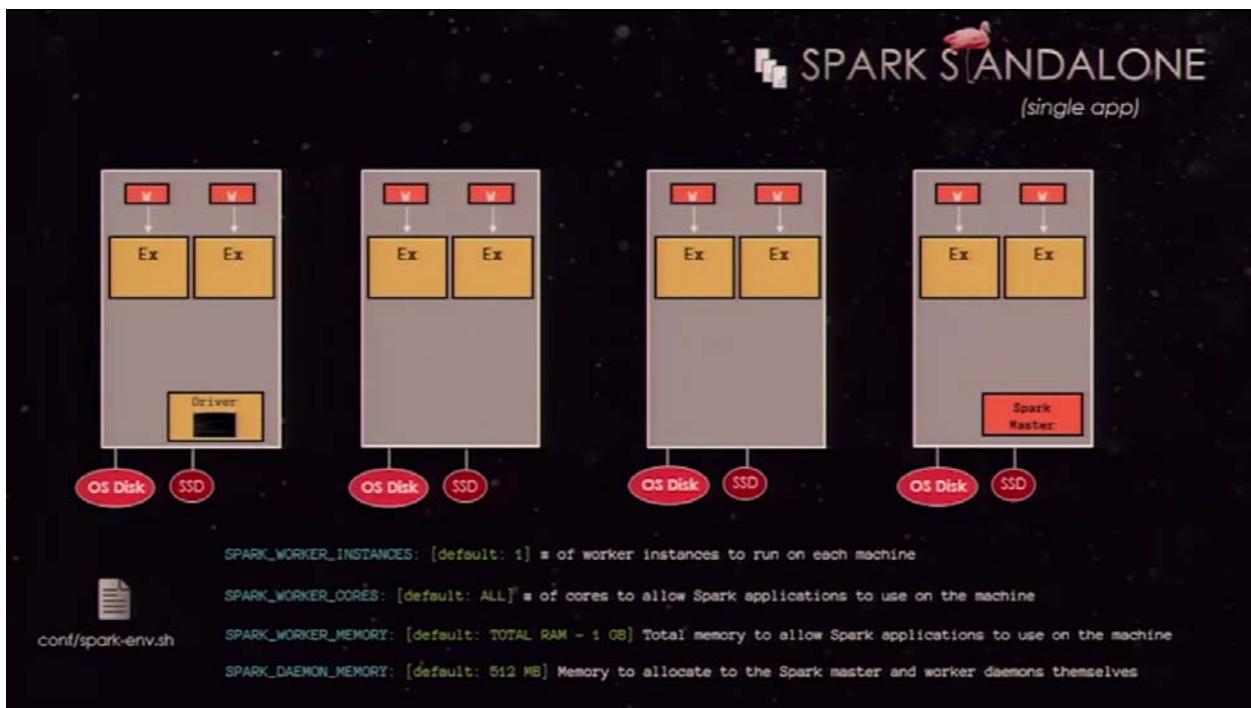
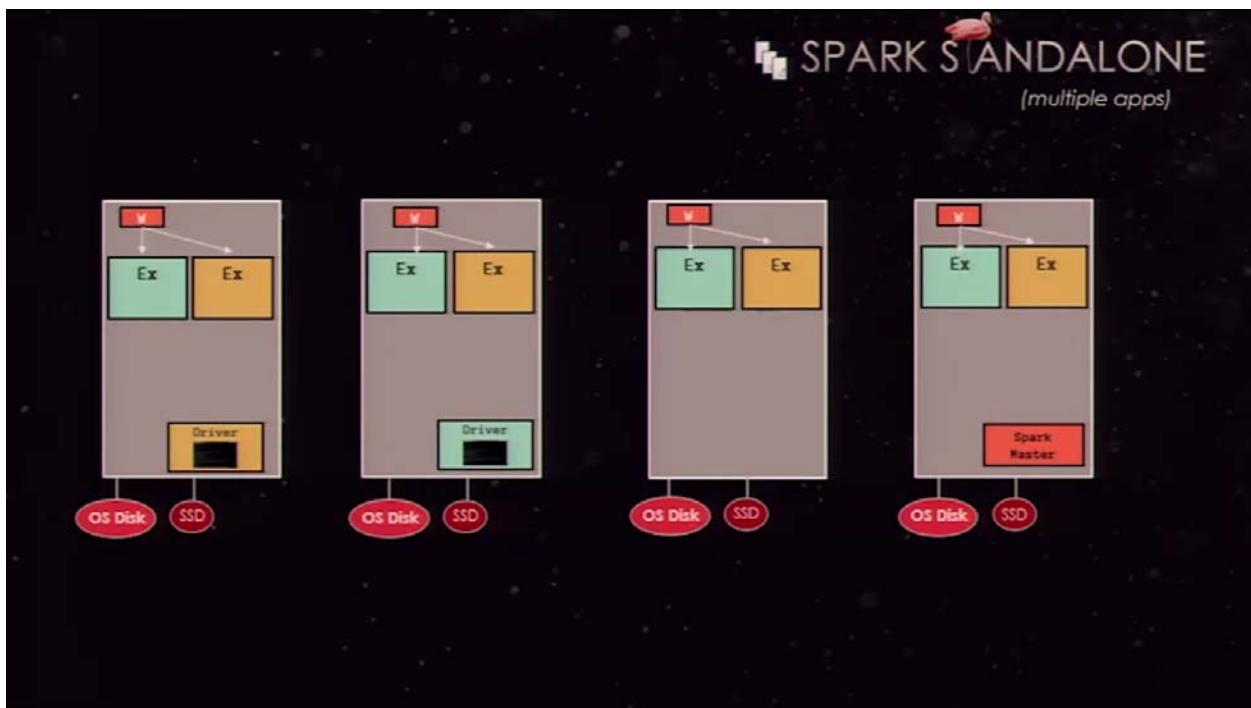
SPARK STANDALONE

```
> ./bin/spark-submit --name "SecondApp"
--master spark://host1:port1
myApp.jar
```

spark-env.sh - SPARK_LOCAL_DIRS

SPARK STANDALONE







Standalone settings

- Apps submitted will run in FIFO mode by default

`spark.cores.max`: maximum amount of CPU cores to request for the application from across the cluster

`spark.executor.memory`: Memory for each executor

The screenshot shows the Spark Master UI at `spark://10.0.64.177:7077`. The top section displays cluster statistics:

- URL: `spark://10.0.64.177:7077`
- Workers: 1
- Cores: 2 Total, 0 Used
- Memory: 4.0 GB Total, 0.0 B Used
- Applications: 0 Running, 0 Completed
- Drivers: 0 Running, 0 Completed

A callout bubble points to the total memory: "Total potential memory this Spark cluster has access to is 4 GB (aka sum of how much memory each Worker, below, has access to)".

Below the statistics is a table titled "Workers" showing one worker entry:

ID	Address	State	Cores	Memory
worker-20140905191420-10.0.64.177-33571	10.0.64.177:33571	ALIVE	2 (0 Used)	4.0 GB (0.0 B Used)

A callout bubble points to the worker's memory: "Amount of potential memory this particular Spark worker has access to".

Below the workers table are sections for "Running Applications" and "Completed Applications", each with a table header:

ID	Name	Cores	Memory per Node	Submitted Time	User	State	Duration
----	------	-------	-----------------	----------------	------	-------	----------

Spark Master at spark://10.0.12.60:7077

URL: spark://10.0.12.60:7077
Workers: 1
Cores: 3 Total, 3 Used
Memory: 7.7 GB Total, 512.0 MB Used
Applications: 1 Running, 0 Completed
Drivers: 0 Running, 0 Completed
Status: ALIVE

Workers

ID	Address	State	Cores	Memory
worker-20141110195851-10.0.12.60-35935	10.0.12.60:35935	ALIVE	3 (3 Used)	7.7 GB (512.0 MB Used)

Running Applications

ID	Name	Cores	Memory per Node	Submitted Time	User	State	Duration
app-20141110204831-0000	Spark shell	3	512.0 MB	2014/11/10 20:48:31	ec2-user	RUNNING	23 min

Completed Applications

ID	Name	Cores	Memory per Node	Submitted Time	User	State	Duration
----	------	-------	-----------------	----------------	------	-------	----------

Spark shell - Spark Jobs

Jobs Stages Storage Environment Executors Spark shell application UI

Spark Jobs (?)

Total Duration: 39 min
Scheduling Mode: FIFO
Active Jobs: 0
Completed Jobs: 4
Failed Jobs: 0

Active Jobs (0)

Job ID	Description	Submitted	Duration	Stages: Succeeded/Total	Tasks (for all stages): Succeeded/Total
--------	-------------	-----------	----------	-------------------------	---

Completed Jobs (4)

Job ID	Description	Submitted	Duration	Stages: Succeeded/Total	Tasks (for all stages): Succeeded/Total
3	collect at <console>:19	2014/12/01 16:18:24	38 ms	1/1 (1 skipped)	2/2 (2 skipped)
2	collect at <console>:19	2014/12/01 16:18:22	55 ms	1/1 (1 skipped)	2/2 (2 skipped)
1	collect at <console>:19	2014/12/01 16:18:07	0.2 s	2/2	4/4
0	count at <console>:15	2014/12/01 16:17:39	0.3 s	1/1	2/2

Failed Jobs (0)

Job ID	Description	Submitted	Duration	Stages: Succeeded/Total	Tasks (for all stages): Succeeded/Total
--------	-------------	-----------	----------	-------------------------	---

Spark shell - Spark Stages X

ec2-54-148-231-117.us-west-2.compute.amazonaws.com:4040/stages/

Spark Jobs Stages Storage Environment Executors Spark shell application UI

Spark Stages (for all jobs)

Total Duration: 39 min
Scheduling Mode: FIFO
Active Stages: 0
Completed Stages: 5
Failed Stages: 0

Active Stages (0)

Stage Id	Description	Submitted	Duration	Tasks: Succeeded/Total	Input	Output	Shuffle Read	Shuffle Write
----------	-------------	-----------	----------	------------------------	-------	--------	--------------	---------------

Completed Stages (5)

Stage Id	Description	Submitted	Duration	Tasks: Succeeded/Total	Input	Output	Shuffle Read	Shuffle Write
6	collect at <console>:19	+details 2014/12/01 16:18:24	28 ms	2/2	552.0 B			
4	collect at <console>:19	+details 2014/12/01 16:18:22	45 ms	2/2				
2	collect at <console>:19	+details 2014/12/01 16:18:07	69 ms	2/2				
1	map at <console>:16	+details 2014/12/01 16:18:07	76 ms	2/2	254.0 B		737.0 B	
0	count at <console>:15	+details 2014/12/01 16:17:40	0.2 s	2/2	254.0 B			



Spark shell - Storage X

ec2-54-148-231-117.us-west-2.compute.amazonaws.com:4040/storage/

Spark Jobs Stages Storage Environment Executors Spark shell application UI

Storage

RDD Name	Storage Level	Cached Partitions	Fraction Cached	Size in Memory	Size in Tachyon	Size on Disk
5	Memory Deserialized 1x Replicated	2	100%	552.0 B	0.0 B	0.0 B



The screenshot shows the Spark shell application UI for RDD Storage. The URL is <http://ec2-54-148-231-117.us-west-2.compute.amazonaws.com:4040/storage/rdd/?id=5>. The main content displays "RDD Storage Info for 5" with the following details:

- Storage Level: Memory Deserialized 1x Replicated
- Cached Partitions: 2
- Total Partitions: 2
- Memory Size: 552.0 B
- Disk Size: 0.0 B

Below this, a table shows "Data Distribution on 1 Executors":

Host	Memory Usage	Disk Usage
localhost:38329	552.0 B (265.4 MB Remaining)	0.0 B

Under "2 Partitions", there is another table:

Block Name	Storage Level	Size in Memory	Size on Disk	Executors
rdd_5_0	Memory Deserialized 1x Replicated	424.0 B	0.0 B	localhost:38329
rdd_5_1	Memory Deserialized 1x Replicated	128.0 B	0.0 B	localhost:38329

The screenshot shows the Spark shell application UI for Environment. The URL is <http://ec2-54-148-231-117.us-west-2.compute.amazonaws.com:4040/environment/>. The main content displays "Environment" with sections for "Runtime Information", "Spark Properties", and "System Properties".

Runtime Information

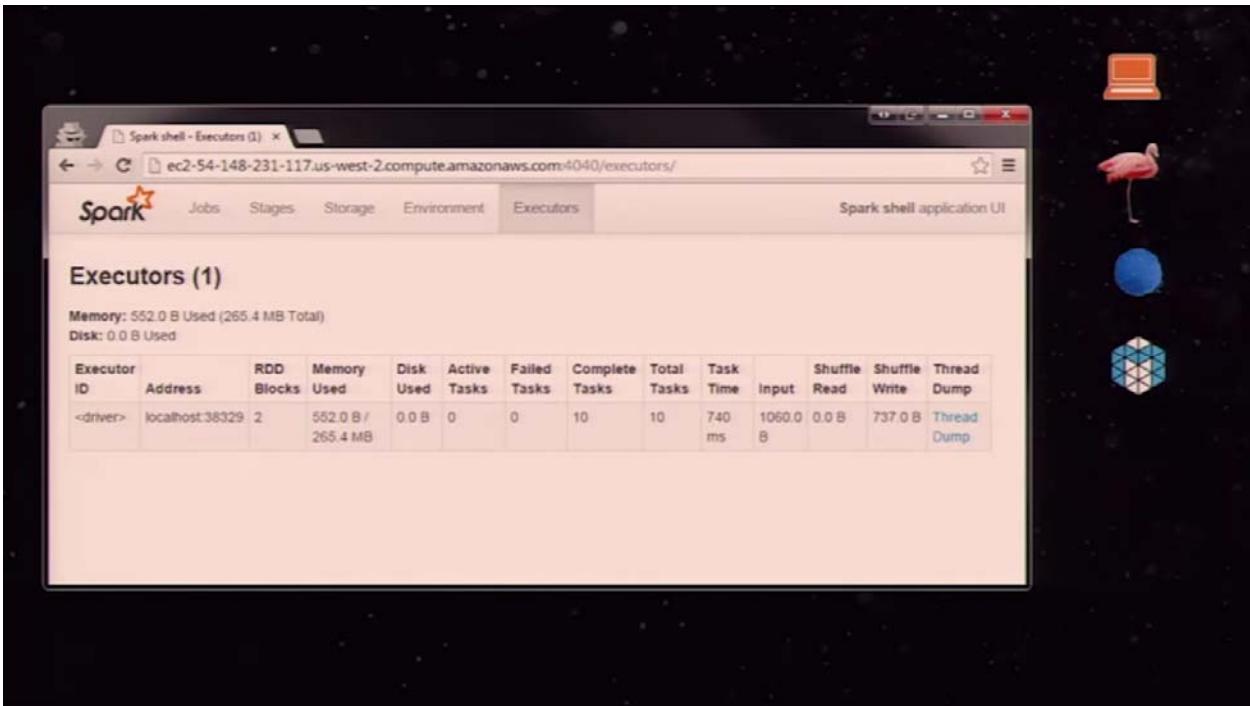
Name	Value
Java Home	/usr/java/jdk.7_0_67/jre
Java Version	1.7.0_67 (Oracle Corporation)
Scala Version	version 2.10.4

Spark Properties

Name	Value
spark.app.id	local-1417468637156
spark.app.name	Spark shell
spark.driver.host	ip-10-0-125-125.us-west-2.compute.internal
spark.driver.port	09091
spark.executor.id	driver
spark.resolver.url	http://10.0.125.125:56999
spark.jars	
spark.master	local[1]
spark.repl.class.uri	http://10.0.125.125:57870
spark.scheduler.mode	FIFO
spark.tachyonStore.tacerName	spark-a5c91951-a604-4425-babd-ate2e9146a70

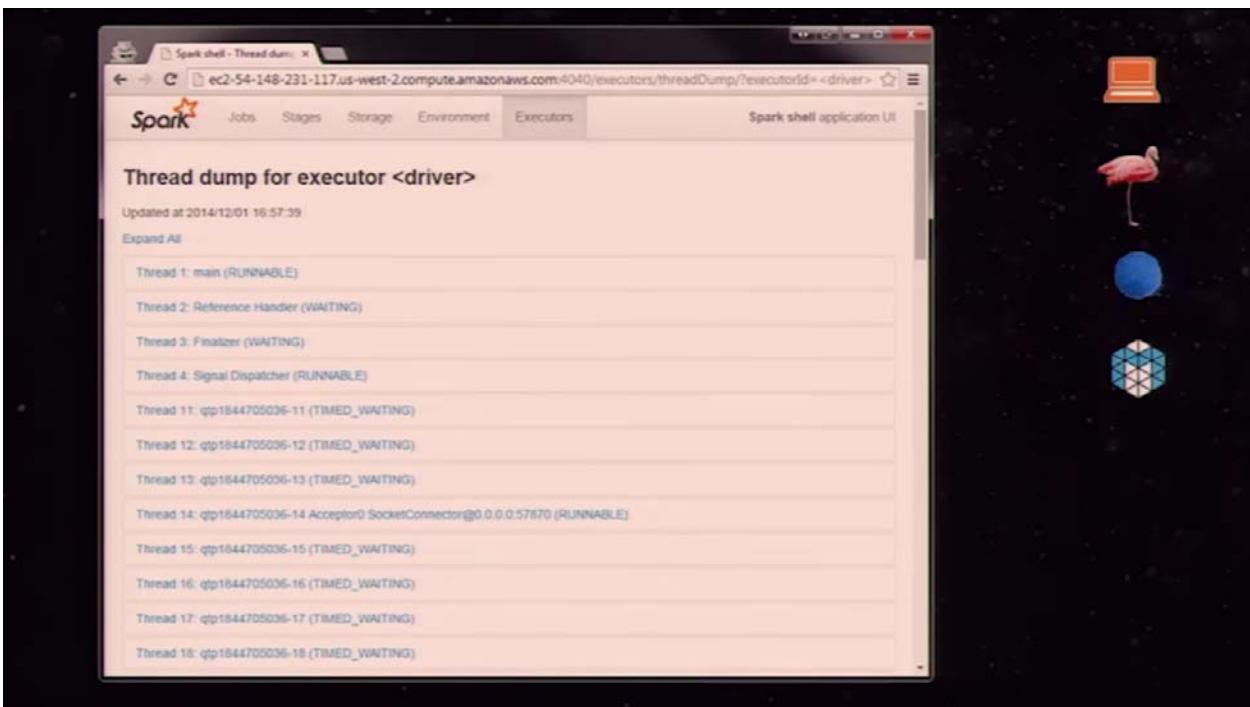
System Properties

Name	Value



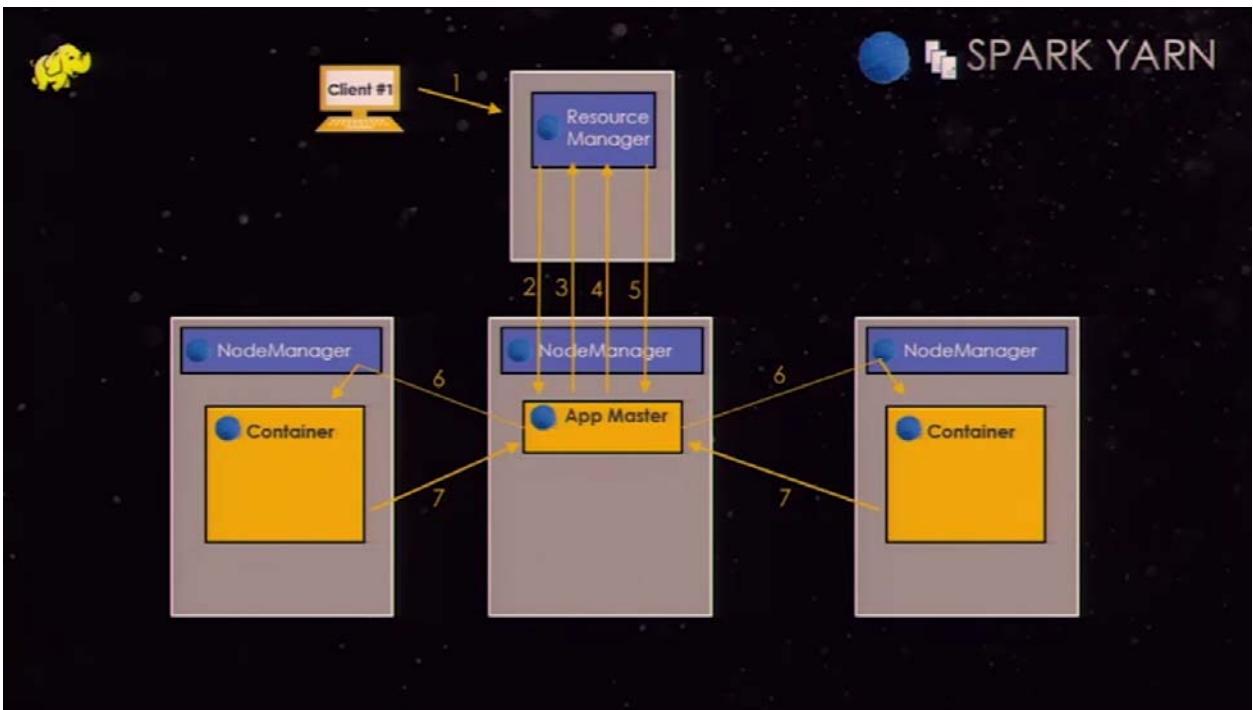
A screenshot of the Spark Executors UI. The title bar says "Spark shell - Executors (1)" and the URL is "ec2-54-148-231-117.us-west-2.compute.amazonaws.com:4040/executors/". The UI has tabs for Jobs, Stages, Storage, Environment, and Executors, with Executors selected. The main area shows the title "Executors (1)". Below it, memory usage is listed: "Memory: 552.0 B Used (265.4 MB Total)" and "Disk: 0.0 B Used". A table provides detailed executor information:

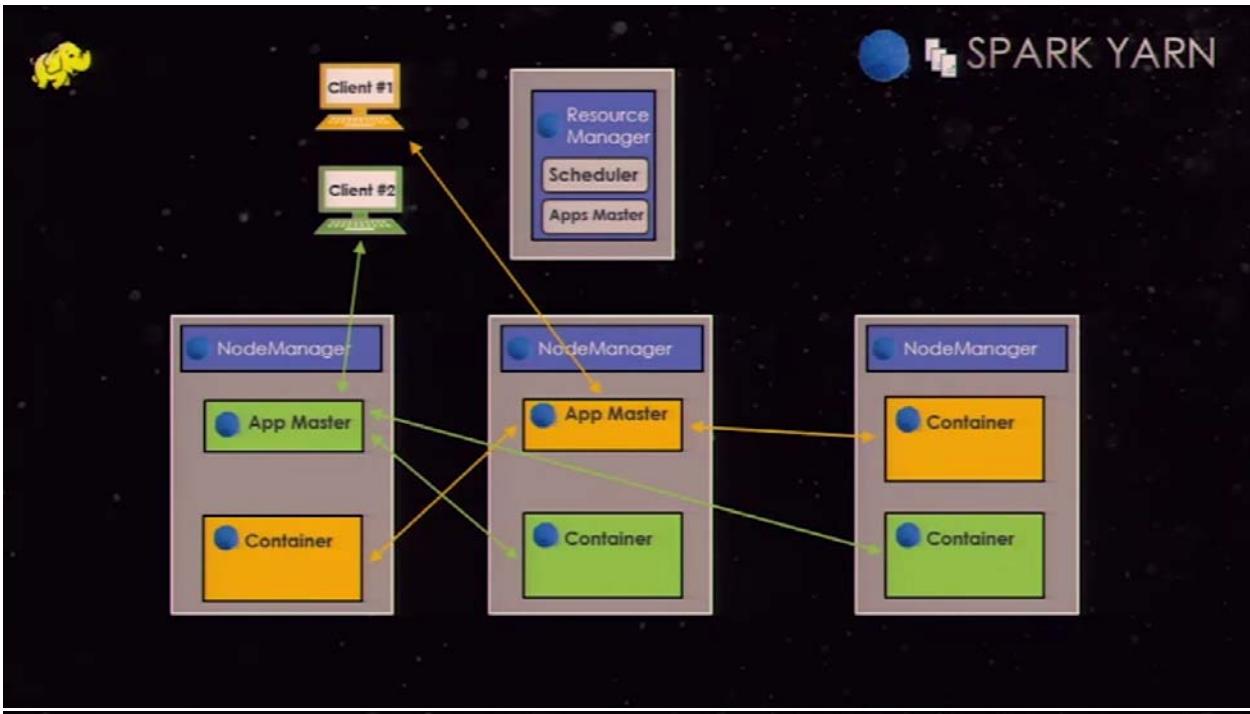
Executor ID	Address	RDD Blocks	Memory Used	Disk Used	Active Tasks	Failed Tasks	Complete Tasks	Total Tasks	Task Time	Shuffle Input	Shuffle Read	Shuffle Write	Thread Dump
<driver>	localhost:38329	2	552.0 B / 265.4 MB	0 B	0	0	10	10	740 ms	1060.0 B	0.0 B	737.0 B	Thread Dump

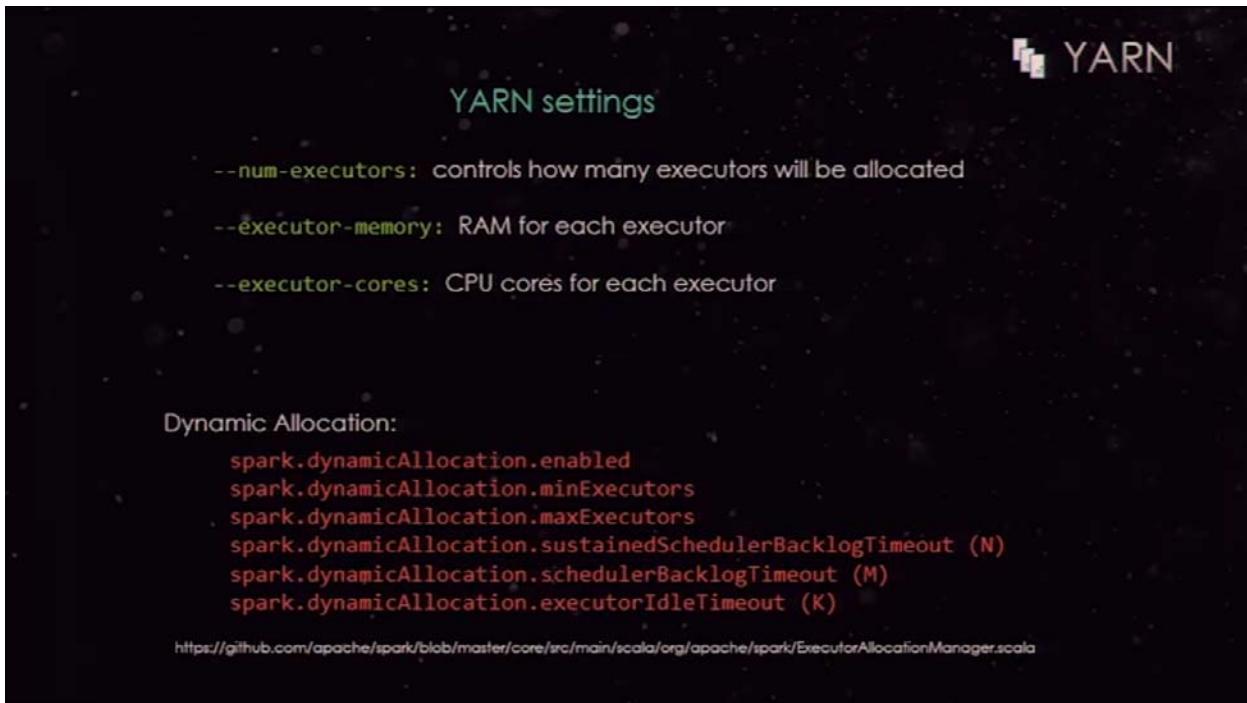
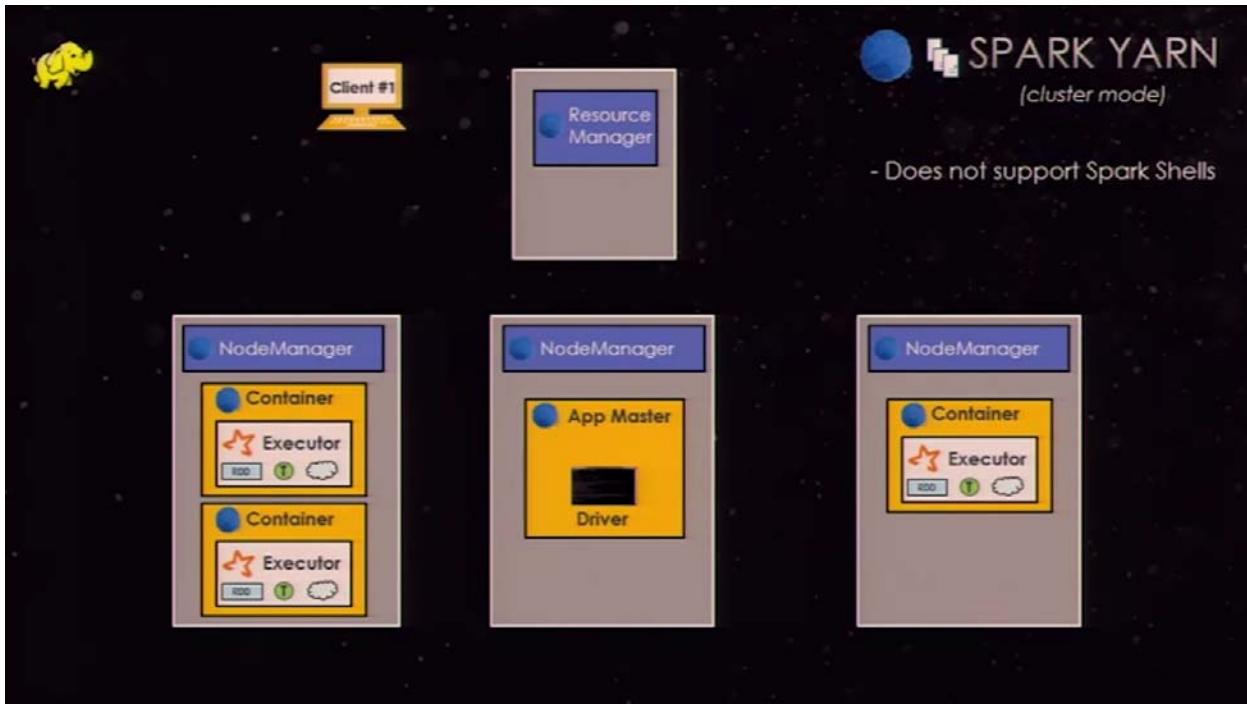


A screenshot of the Spark Thread dump UI. The title bar says "Spark shell - Thread dump X" and the URL is "ec2-54-148-231-117.us-west-2.compute.amazonaws.com:4040/executors/threadDump/?executorId=<driver>". The UI has tabs for Jobs, Stages, Storage, Environment, and Executors, with Executors selected. The main area shows the title "Thread dump for executor <driver>". It indicates the dump was updated at "2014/12/01 16:57:39". A "Expand All" button is present. A list of threads is shown:

- Thread 1: main (RUNNABLE)
- Thread 2: Reference Handler (WAITING)
- Thread 3: Finalizer (WAITING)
- Thread 4: Signal Dispatcher (RUNNABLE)
- Thread 11: qtp1844705036-11 (TIMED_WAITING)
- Thread 12: qtp1844705036-12 (TIMED_WAITING)
- Thread 13: qtp1844705036-13 (TIMED_WAITING)
- Thread 14: qtp1844705036-14 Acceptor@0.0.0.57870 (RUNNABLE)
- Thread 15: qtp1844705036-15 (TIMED_WAITING)
- Thread 16: qtp1844705036-16 (TIMED_WAITING)
- Thread 17: qtp1844705036-17 (TIMED_WAITING)
- Thread 18: qtp1844705036-18 (TIMED_WAITING)







YARN resource manager UI: <http://<ip address>:8088>

(No apps running)

Cluster Metrics

Apps Submitted	Apps Pending	Apps Running	Apps Completed	Containers Running	Memory Used	Memory Total	Memory Reserved	VCores Used	VCores Total	VCores Reserved	Active Nodes	De
0	0	0	0	0	0 B	2.71 GB	0 B	0	4	0	1	0

User Metrics for dr who

Apps Submitted	Apps Pending	Apps Running	Apps Completed	Containers Running	Containers Pending	Containers Reserved	Memory Used	Memory Pending	Memory Reserved
0	0	0	0	0	0	0	0 B	0 B	0 B

Show: 20 entries

ID	User	Name	Application Type	Queue	StartTime	FinishTime	State	FinalStatus
No data available in table								

Showing 0 to 0 of 0 entries

```
[ec2-user@ip-10-0-72-36 ~]$ spark-submit --class org.apache.spark.examples.SparkPi --deploy-mode client --master yarn /opt/cloudera/parcels/CDH-5.2.1-1.cdh5.2.1.p0.12/jars/spark-examples-1.1.0-cdh5.2.1-hadoop2.5.0-cdh5.2.1.jar 10
```

App running in **client** mode

All Applications

ID	User	Name	Application Type	Queue	StartTime	FinishTime	State	FinalStatus	Progress	Tracking URL
application_1417641624005_0001	ec2-user	Spark Pi	SPARK	root,ec2-user	Thu, 04 Dec 2014 15:30:43 GMT	Thu, 04 Dec 2014 15:31:14 GMT	FINISHED	SUCCEEDED	<div style="width: 100%;"></div>	History
application_1417641624005_0002	ec2-user	Spark Pi	SPARK	root,ec2-user	Thu, 04 Dec 2014 15:25:48 GMT	Thu, 04 Dec 2014 15:26:19 GMT	FINISHED	SUCCEEDED	<div style="width: 100%;"></div>	History
application_1417641624005_0001	ec2-user	Spark Pi	SPARK	root,ec2-user	Thu, 04 Dec 2014 15:25:18 GMT	Thu, 04 Dec 2014 15:25:35 GMT	FINISHED	SUCCEEDED	<div style="width: 100%;"></div>	History

```
[ec2-user@ip-10-0-72-36 ~]$ spark-submit --class
org.apache.spark.examples.SparkPi --deploy-mode cluster --master
yarn /opt/cloudera/parcels/CDH-5.2.1-1.cdh5.2.1.p0.12/jars/spark-
examples-1.1.0-cdh5.2.1-hadoop2.5.0-cdh5.2.1.jar 10
```

Cluster 1 - Spark - CloudBees Jenkins CI

History Server

ec2-54-149-62-154.us-west-2.compute.amazonaws.com:18088

Spark History Server

Event Log Location: hdfs://ip-10-0-72-36.us-west-2.compute.internal:8020/user/spark/applicationHistory

Showing 1-2 of 2

App Name	Started	Completed	Duration	Spark User	Last Updated
Spark shell	2014/12/04 09:14:01	2014/12/04 09:21:19	7.3 min	ec2-user	2014/12/04 09:21:20
Spark shell	2014/12/04 09:07:36	2014/12/04 09:13:47	6.2 min	ec2-user	2014/12/04 09:13:48

PLUGGABLE RESOURCE MANAGEMENT

	Spark Central Master	Who starts Executors?	Tasks run in
Local	[none]	Human being	Executor
Standalone	Standalone Master	Worker JVM	Executor
YARN	YARN App Master	Node Manager	Executor
Mesos	Mesos Master	Mesos Slave	Executor

DEPLOYING AN APP TO THE CLUSTER

`spark-submit` provides a uniform interface for submitting jobs across all cluster managers

```
bin/spark-submit --master spark://host:7077  
--executor-memory 10g  
my_script.py
```

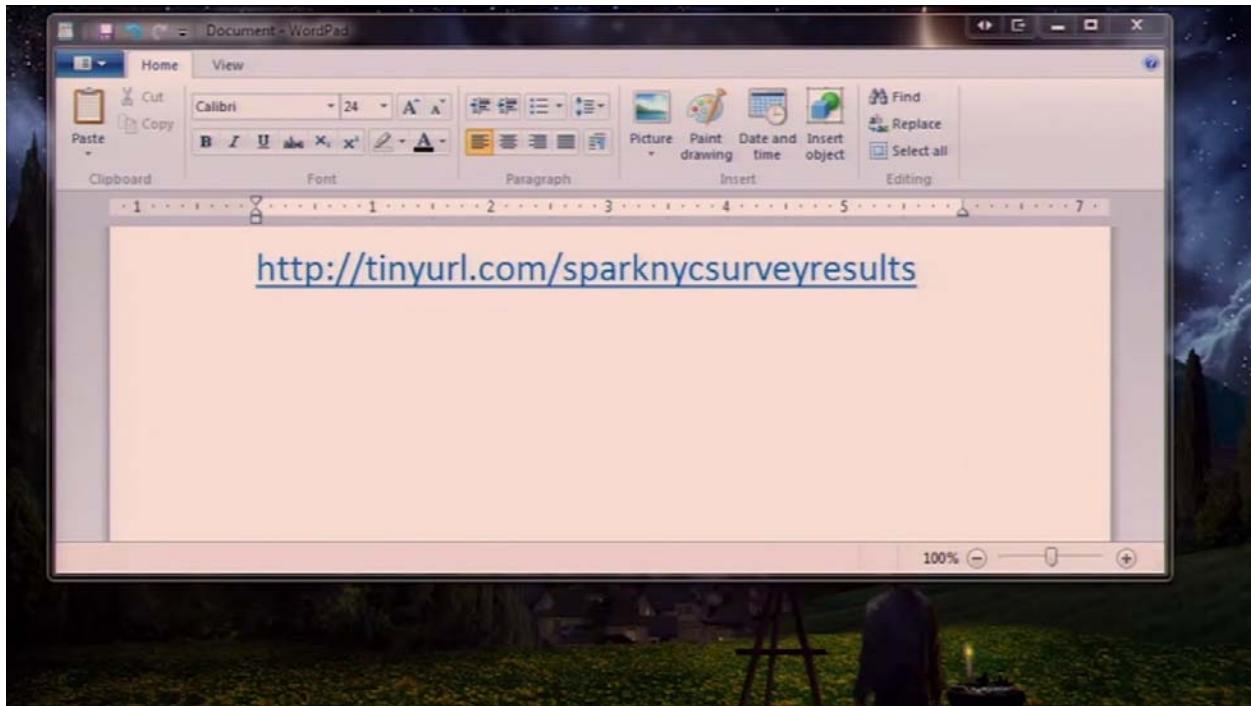
Table 7-2. Possible values for the `--master` flag in `spark-submit`



Value	Explanation
spark://host:port	Connect to a Spark Standalone master at the specified port. By default Spark Standalone master's listen on port 7077 for submitted jobs.
mesos://host:port	Connect to a Mesos cluster master at the specified port. By default Mesos masters listen on port 5050 for submitted jobs.
yarn	Indicates submission to YARN cluster. When running on YARN you'll need to export HADOOP_CONF_DIR to point the location of your Hadoop configuration directory.
local	Run in local mode with a single core.
local[N]	Run in local mode with N cores.
local[*]	Run in local mode and use as many cores as the machine has.

Source: Learning Spark

HDFS	121	69.9%
MapReduce	102	59%
YARN	66	38.2%
Mesos	7	4%
Cascading	7	4%
Kafka	43	24.9%
Storm	20	11.6%
Flume	23	13.3%
HBase	35	20.2%
Cassandra	42	24.3%
Hive	69	39.9%
Impala	23	13.3%
Pig	38	22%
Parquet	33	19.1%
ZooKeeper	48	27.7%
MongoDB	41	23.7%
Couchbase	2	1.2%
Neo4j	14	8.1%
Titan	3	1.7%



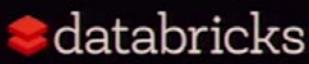
AGENDA

Before Lunch

- History of Spark
- RDD fundamentals
- Spark Runtime Architecture
Integration with Resource Managers
(Standalone, YARN)
- GUIs
- Lab: DevOps 101 

After Lunch

- Memory and Persistence
- Jobs -> Stages -> Tasks
- Broadcast Variables and Accumulators
- PySpark
- DevOps 102 
- Shuffle
- Spark Streaming



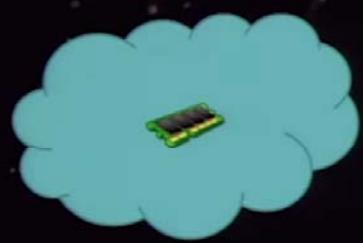


Recommended to use at most only 75% of a machine's memory for Spark

Minimum Executor heap size should be 8 GB

Max Executor heap size depends... maybe 40 GB (watch GC)

Memory usage is greatly affected by storage level and serialization format



Vs.





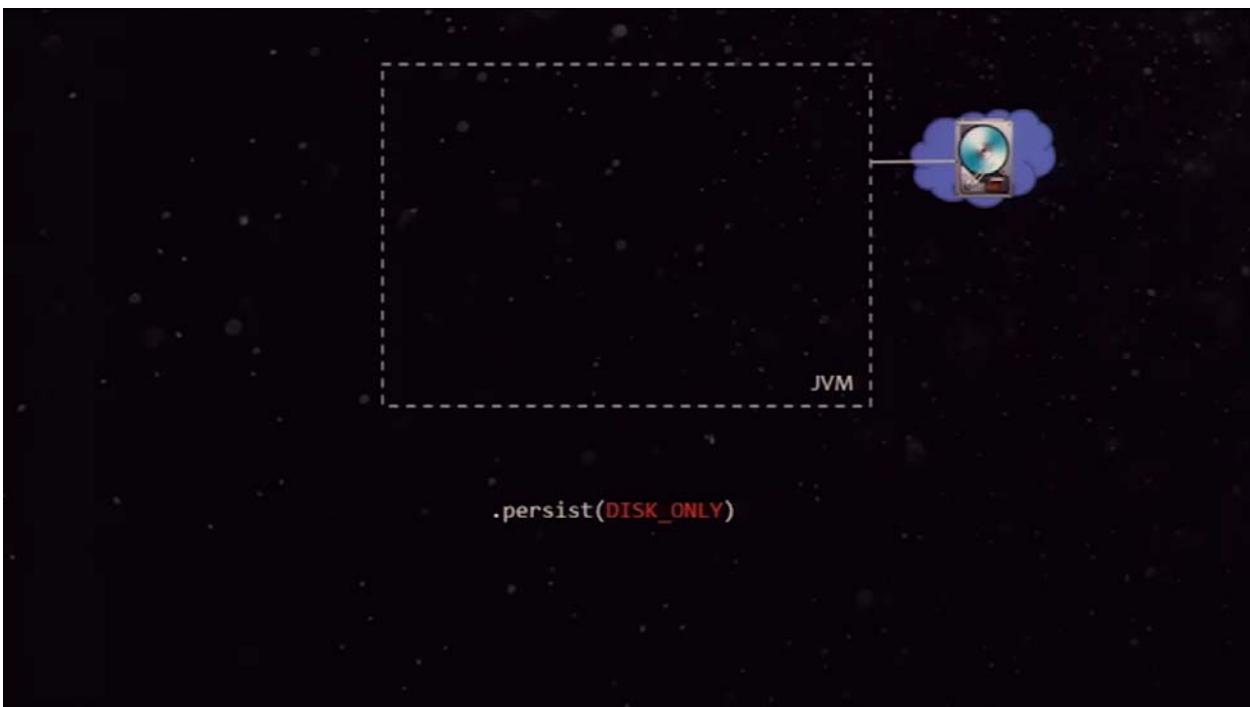
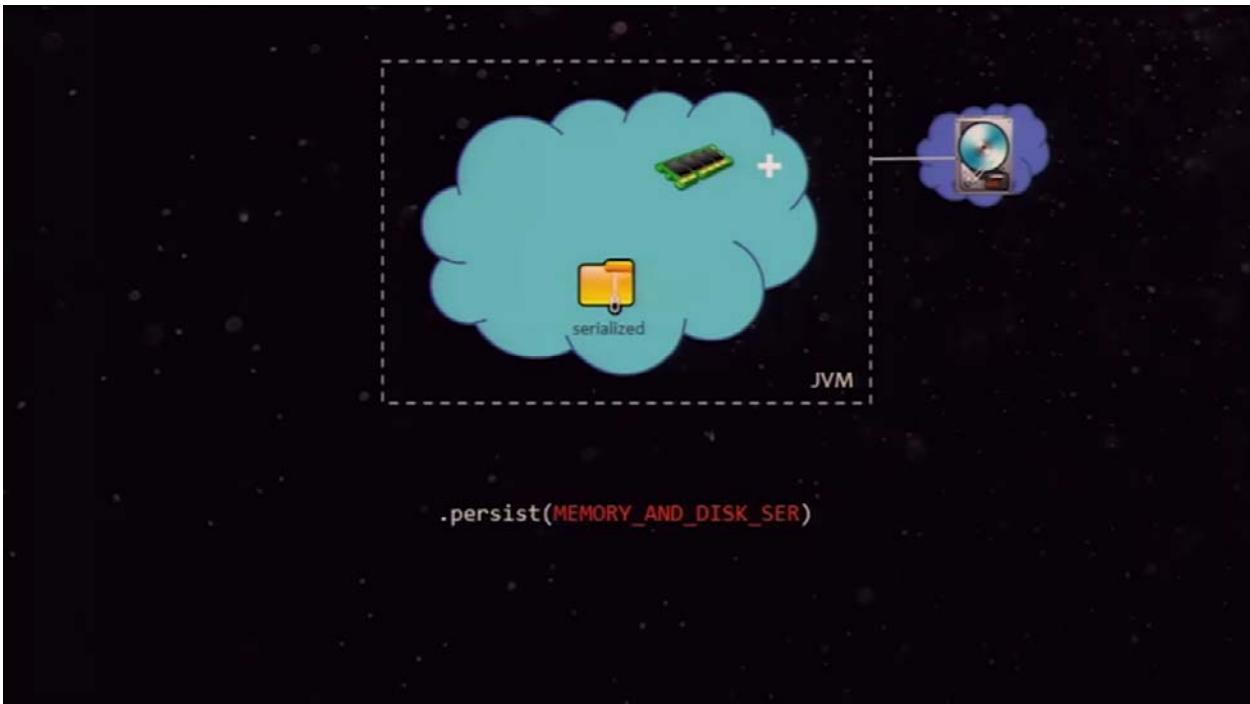
```
RDD.cache() == RDD.persist(MEMORY_ONLY)
```

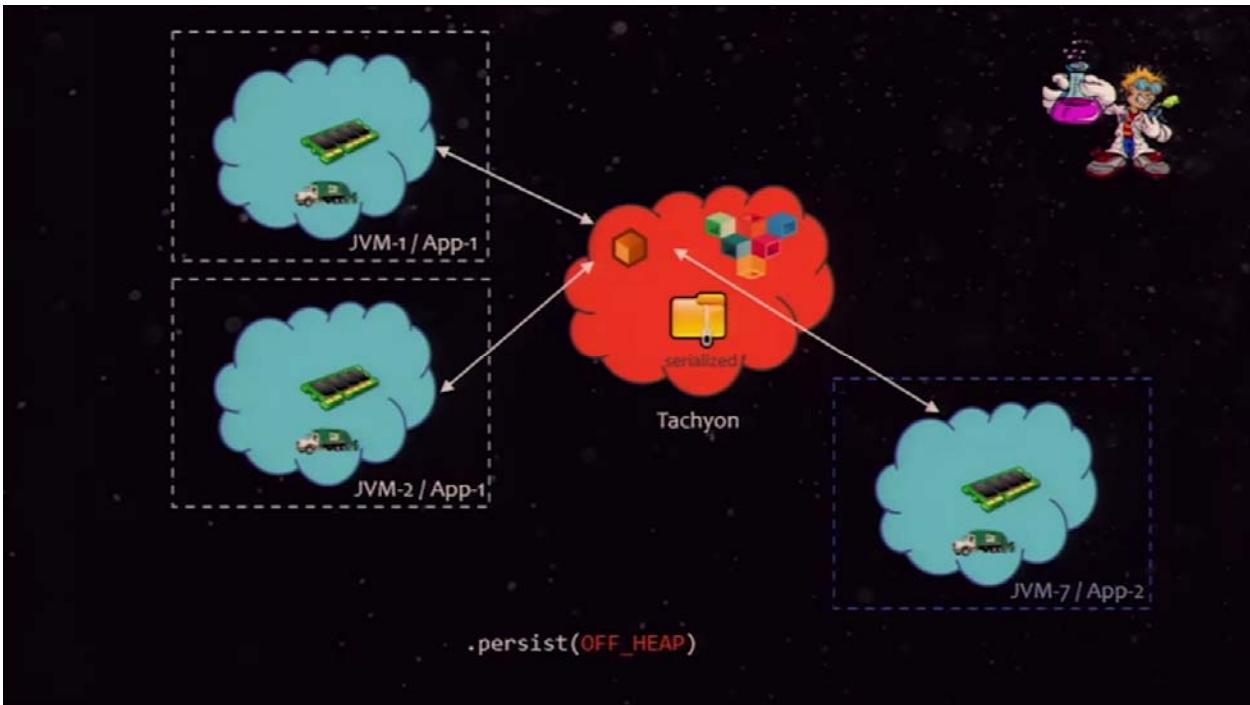
most CPU-efficient option

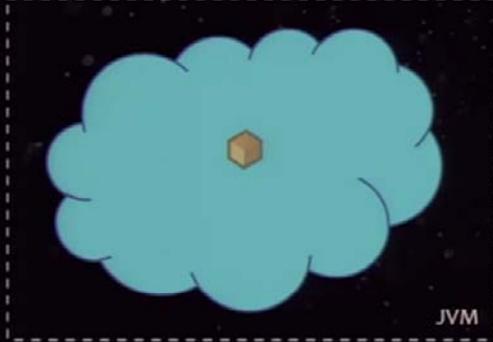


```
RDD.persist(MEMORY_ONLY_SER)
```









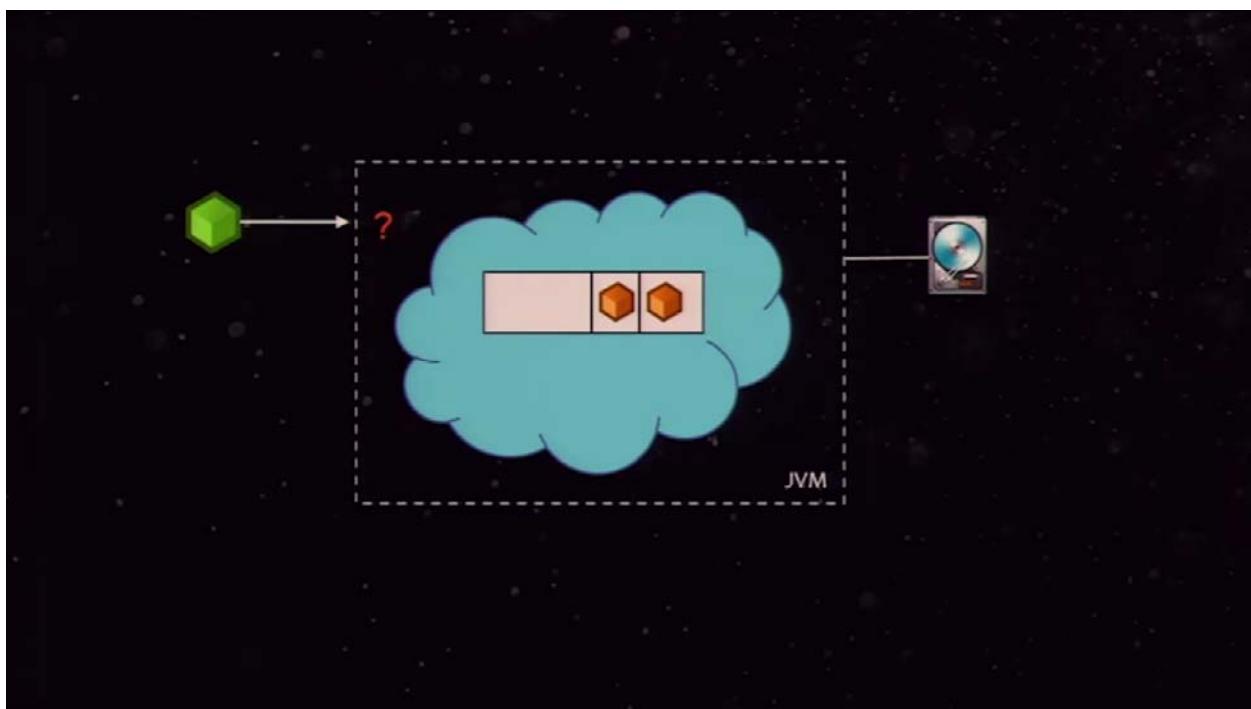
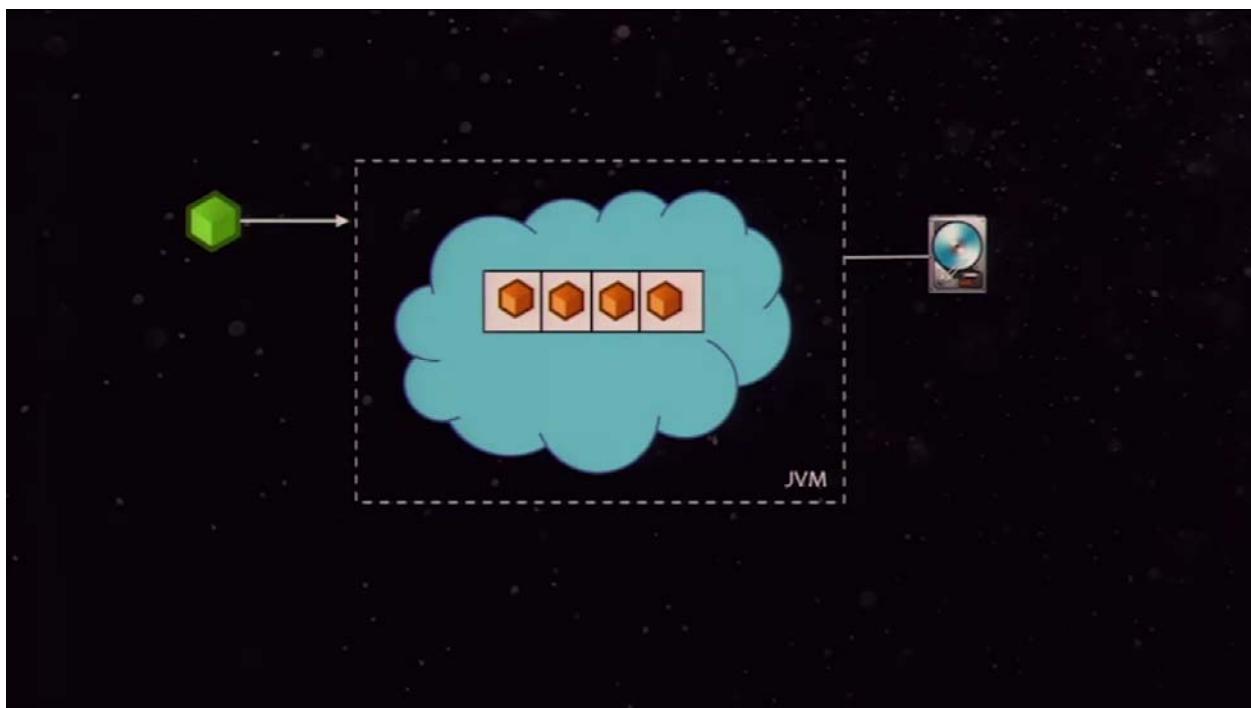
JVM

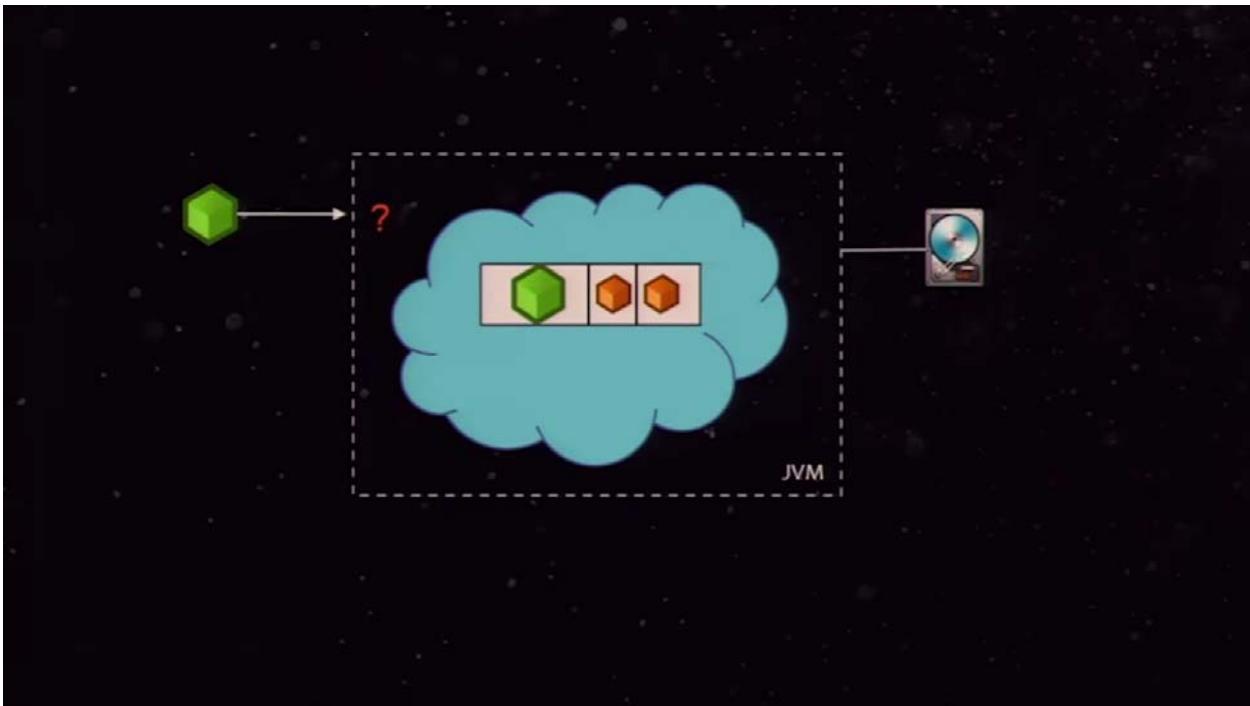
```
.unpersist()
```



JVM

```
.unpersist()
```





?

- If RDD fits in memory, choose `MEMORY_ONLY`
- If not, use `MEMORY_ONLY_SER` w/ fast serialization library
- Don't spill to disk unless functions that computed the datasets are very expensive or they filter a large amount of data.
(recomputing may be as fast as reading from disk)
- Use replicated storage levels sparingly and only if you want fast fault recovery (maybe to serve requests from a web app)



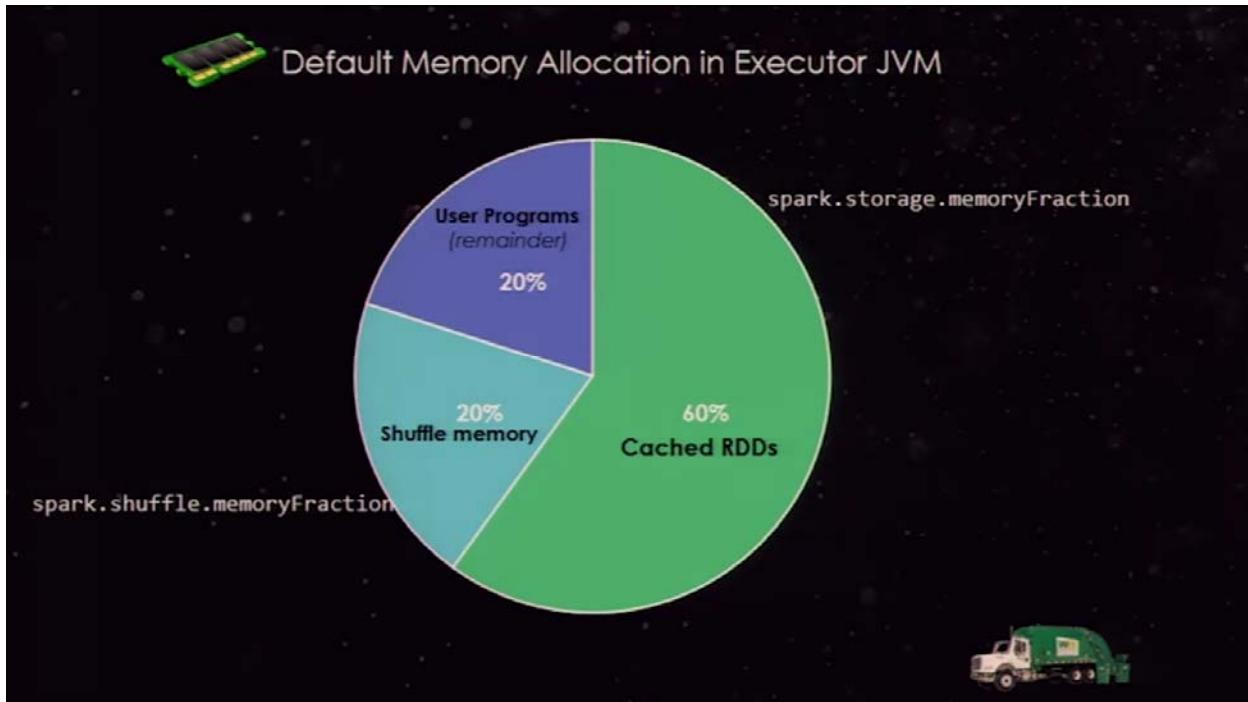
Remember!



Intermediate data is automatically persisted during shuffle operations



PySpark: stored objects will always be serialized with Pickle library, so it does not matter whether you choose a serialized level.



MEMORY

Spark uses memory for:

RDD Storage: when you call `.persist()` or `.cache()`. Spark will limit the amount of memory used when caching to a certain fraction of the JVM's overall heap, set by `spark.storage.memoryFraction`

Shuffle and aggregation buffers: When performing shuffle operations, Spark will create intermediate buffers for storing shuffle output data. These buffers are used to store intermediate results of aggregations in addition to buffering data that is going to be directly output as part of the shuffle.

User code: Spark executes arbitrary user code, so user functions can themselves require substantial memory. For instance, if a user application allocates large arrays or other objects, these will contribute to overall memory usage. User code has access to everything "left" in the JVM heap after the space for RDD storage and shuffle storage are allocated.

DETERMINING MEMORY CONSUMPTION

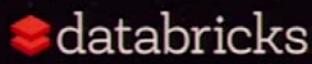
1. Create an RDD
2. Put it into cache
3. Look at SparkContext logs
on the driver program or
Spark UI

logs will tell you how much memory each partition is consuming, which you can aggregate to get the total size of the RDD

```
INFO BlockManagerMasterActor: Added rdd_0_1 in memory on mbk.local:50311 (size: 717.5 KB, free: 332.3 MB)
```

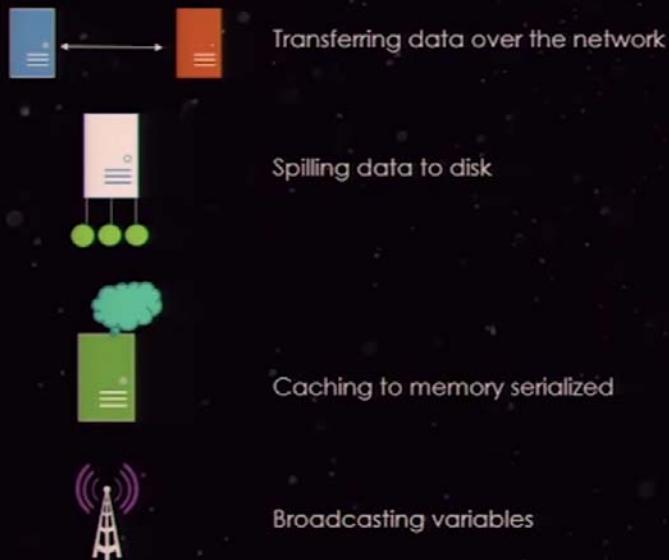


DATA SERIALIZATION



SERIALIZATION

Serialization is used when:



Java serialization



Kryo serialization

vs.

- Uses Java's `ObjectOutputStream` framework
 - Works with any class you create that implements `java.io.Serializable`
 - You can control the performance of serialization more closely by extending `java.io.Externalizable`
 - Flexible, but quite slow
 - Leads to large serialized formats for many classes
- Recommended serialization for production apps
 - Use Kryo version 2 for speedy serialization (10x) and more compactness
 - Does not support all `Serializable` types
 - Requires you to register the classes you'll use in advance
 - If set, will be used for serializing shuffle data between nodes and also serializing RDDs to disk



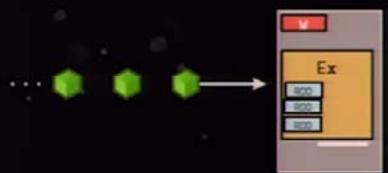
To register your own custom classes with Kryo, use the `registerKryoClasses` method:

```
val conf = new SparkConf().setMaster(...).setAppName(...)  
conf.registerKryoClasses(Seq(classOf[MyClass1], classOf[MyClass2]))  
val sc = new SparkContext(conf)
```

- If your objects are large, you may need to increase `spark.kryoserializer.buffer.mb` config property
 - The default is 2, but this value needs to be large enough to hold the *largest* object you will serialize.



TUNING FOR



High churn



Low churn



TUNING FOR Spark

Cost of GC is proportional to the # of Java objects
 (so use an array of Ints instead of a LinkedList)



High chum

To measure GC impact:

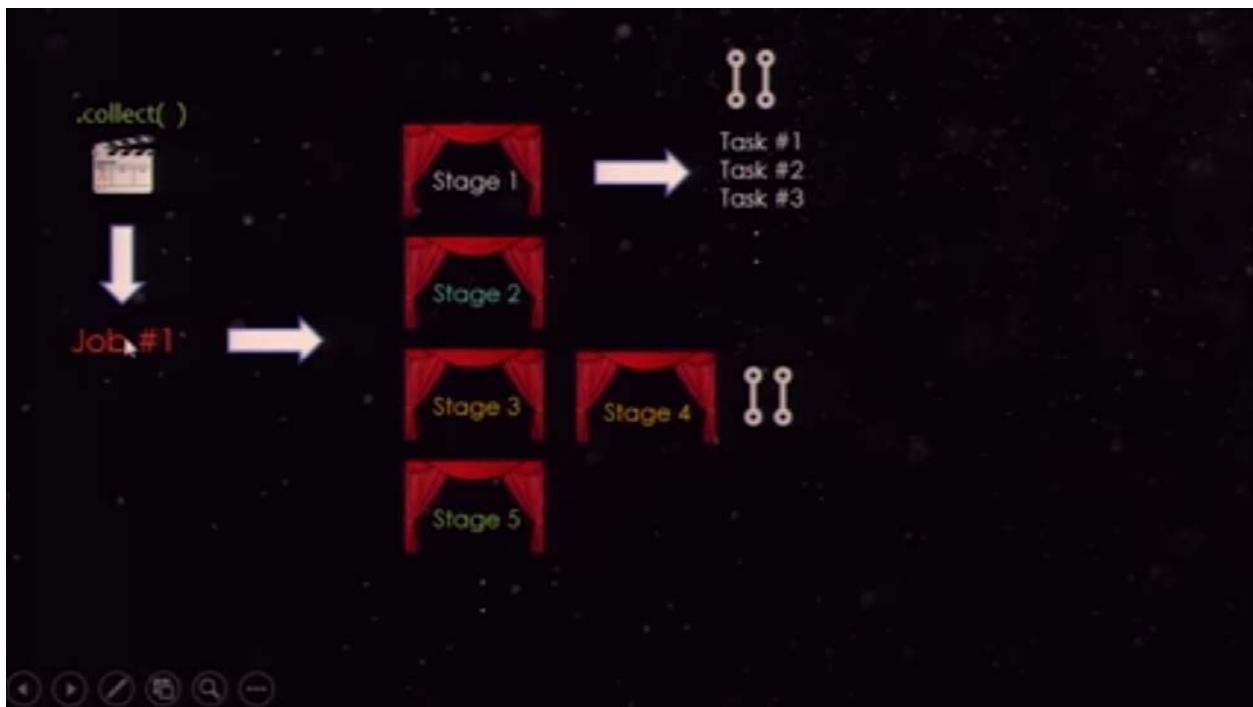
```
-verbose:gc -XX:+PrintGCDetails -XX:+PrintGCTimeStamps
```

		TUNING
Parallel GC	Parallel Old GC	CMS GC
<pre>-XX:+UseParallelGC -XX:ParallelGCThreads=<n></pre> <ul style="list-style-type: none"> Uses multiple threads to do young gen GC Will default to Serial on single core machines Aka "throughput collector" Good for when a lot of work is needed and long pauses are acceptable Use cases: batch processing 	<pre>-XX:+UseParallelOldGC</pre> <ul style="list-style-type: none"> Uses multiple threads to do both young gen and old gen GC Also a multithreading compacting collector HotSpot does compaction only in old gen 	<pre>-XX:+UseConcMarkSweepGC -XX:ParallelCMSThreads=<n></pre> <ul style="list-style-type: none"> Concurrent Mark Sweep aka "Concurrent low pause collector" Tries to minimize pauses due to GC by doing most of the work concurrently with application threads Uses same algorithm on young gen as parallel collector Use cases: 
		G1 GC
		<pre>-XX:+UseG1GC</pre> <ul style="list-style-type: none"> Garbage First is available starting Java 7 Designed to be long term replacement for CMS Is a parallel, concurrent and incrementally compacting low-pause GC

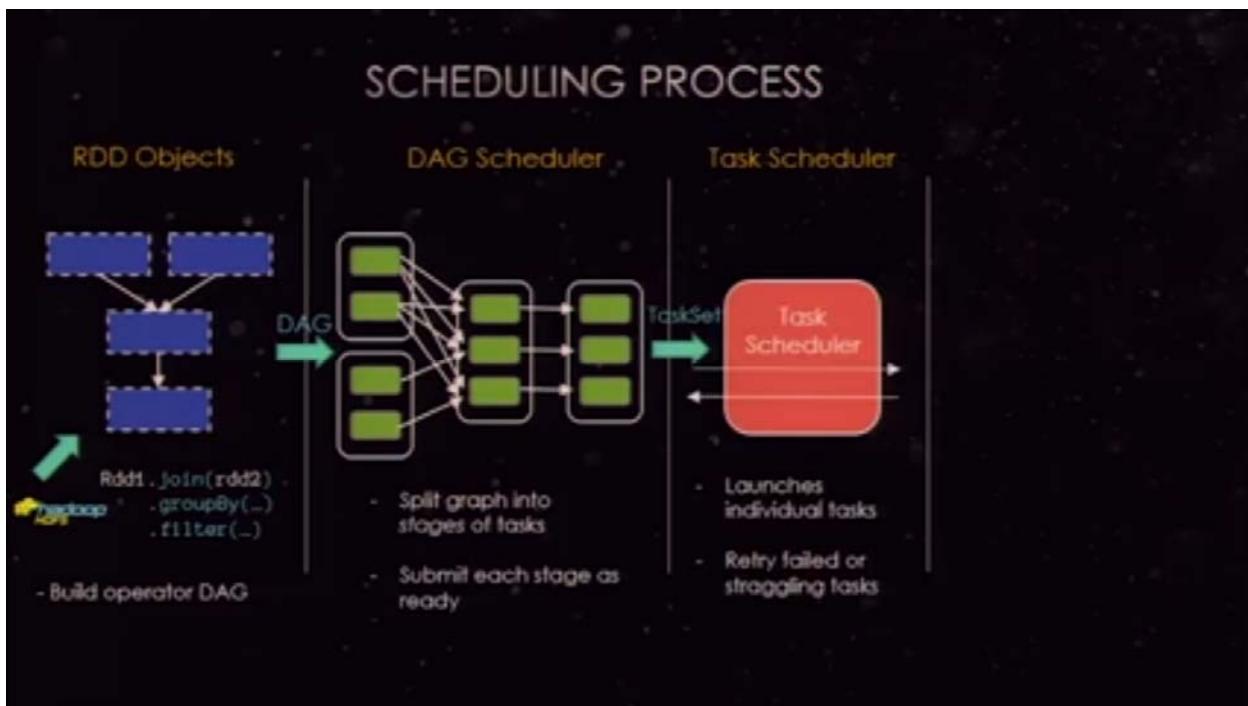
?

JOB → STAGES → TASKS

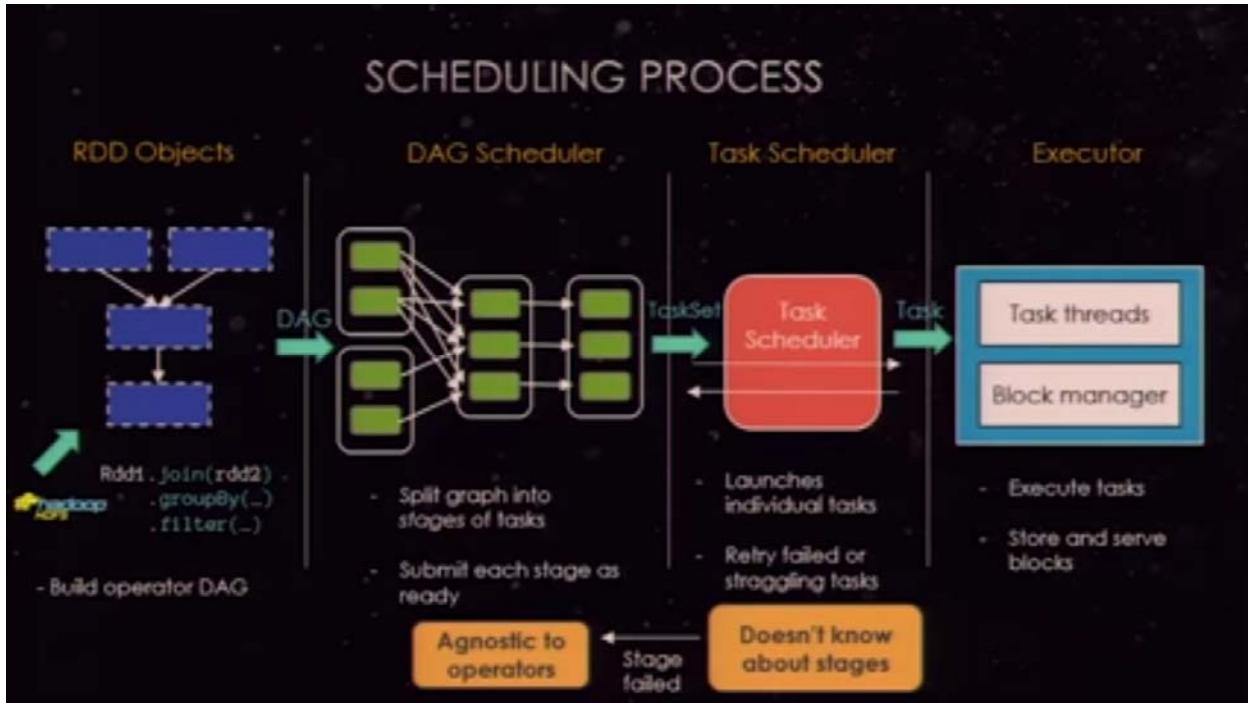
databricks



SCHEDULING PROCESS



SCHEDULING PROCESS



LINEAGE

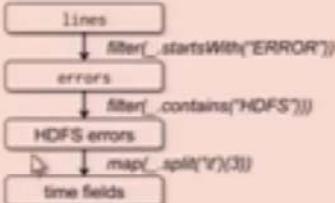


Figure 1: Lineage graph for the third query in our example.
Boxes represent RDDs and arrows represent transformations.

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



LINEAGE

"One of the challenges in providing RDDs as an abstraction is choosing a representation for them that can track lineage across a wide range of transformations."



"The most interesting question in designing this interface is how to represent dependencies between RDDs."

"We found it both sufficient and useful to classify dependencies into two types:

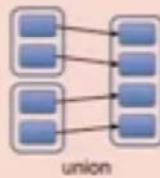
- narrow dependencies, where each partition of the parent RDD is used by at most one partition of the child RDD
- wide dependencies, where multiple child partitions may depend on it."

LINEAGE DEPENDENCIES

Narrow Dependencies:



map, filter

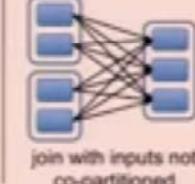


union

Wide Dependencies:



groupByKey



join with inputs not
co-partitioned

Requires
shuffle

Examples of narrow and wide dependencies.

Each box is an RDD, with partitions shown as shaded rectangles.