

Exploring Wikipedia with Spark

 databricks training

Spark Summit West 2016

Environments

Workloads

Goal: unified engine across data **sources**,
workloads and **environments**

Data Sources

Goal: unified engine across data **sources**, workloads and environments

Environments

YARN



EC2



MESOS



databricks™

Spark

Workloads

DataFrames / SQL / Datasets APIs



Spark SQL

Spark Streaming

MLlib

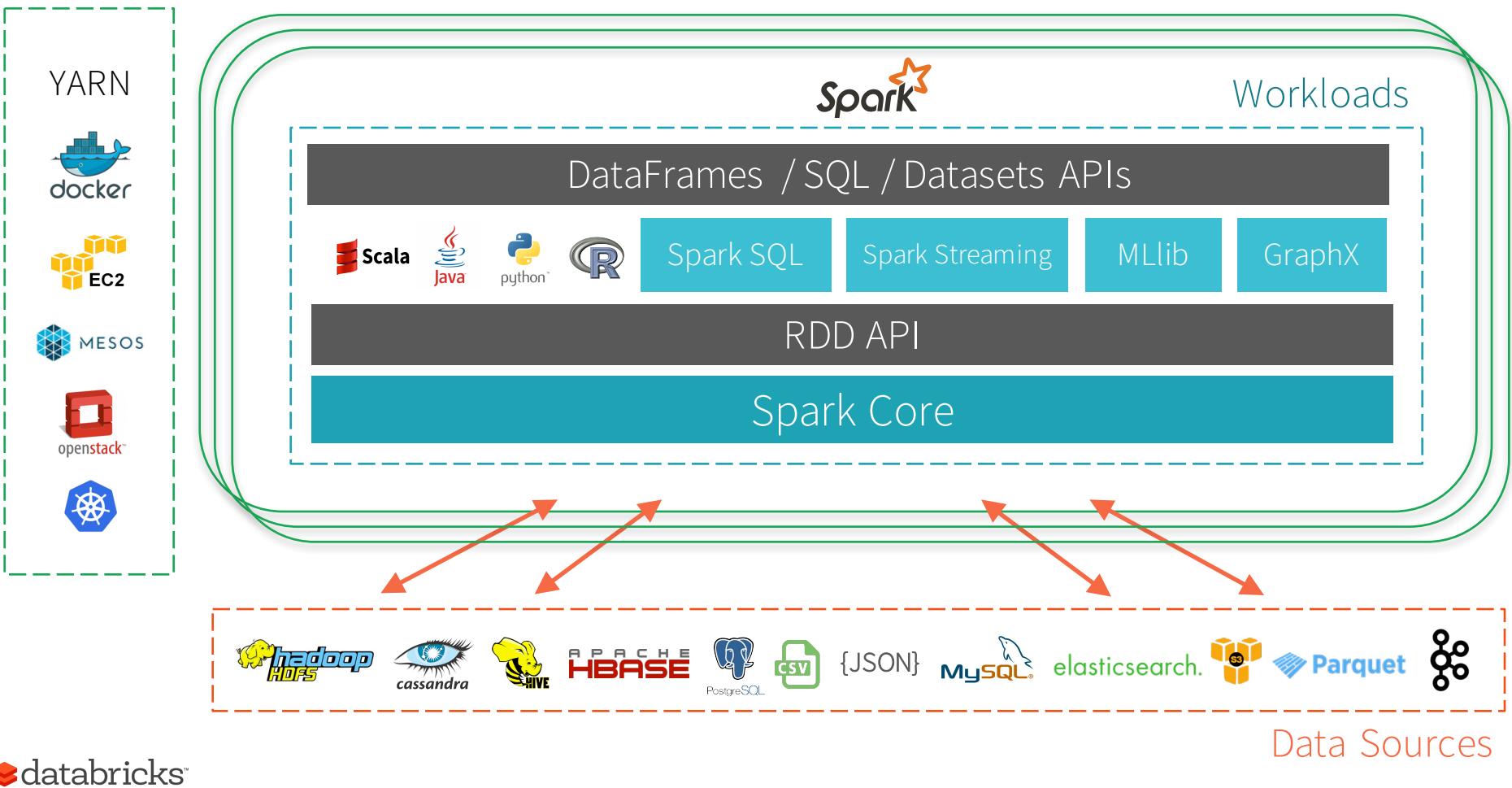
GraphX

RDD API

Spark Core



Environments



databricks™





-amplab
UC BERKELEY



Today's Objectives

- Learn just enough to build simple prototypes/POCs with Spark
- Fast paced, high level overview of all major Spark components
- Hands on with Spark's programming APIs (*DataFrame/SQL, RDD, Datasets*)
- Overview of Spark architecture: Core, Streaming, Standalone Mode, DAG
- Mix of beginner + advanced topics
- Not all slides/labs covered (*reference & homework material*)
- Lots of ideas, code & datasets to play around with after class

Timings: ~6 hours total

morning: APIs: (DataFrames, SQL), Architecture

afternoon: APIs: (Datasets, RDD), GraphX/ML, Streaming

~5pm: end!

10:30am – 11am: Break

noon – 1pm: Lunch

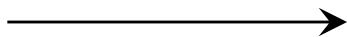
3pm – 3:30pm: Break



Data

 Spark Analytics

Pageviews (March 2015) - 255 MB



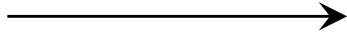
Dataframes

Clickstream (Feb 2015) - 1.2 GB



Dataframes, SQL, GraphX

Pagecounts (last hour) - ~500 MB



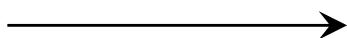
RDD, Datasets, Dataframes

English Wikipedia (Mar 5, 2016) - 54 GB



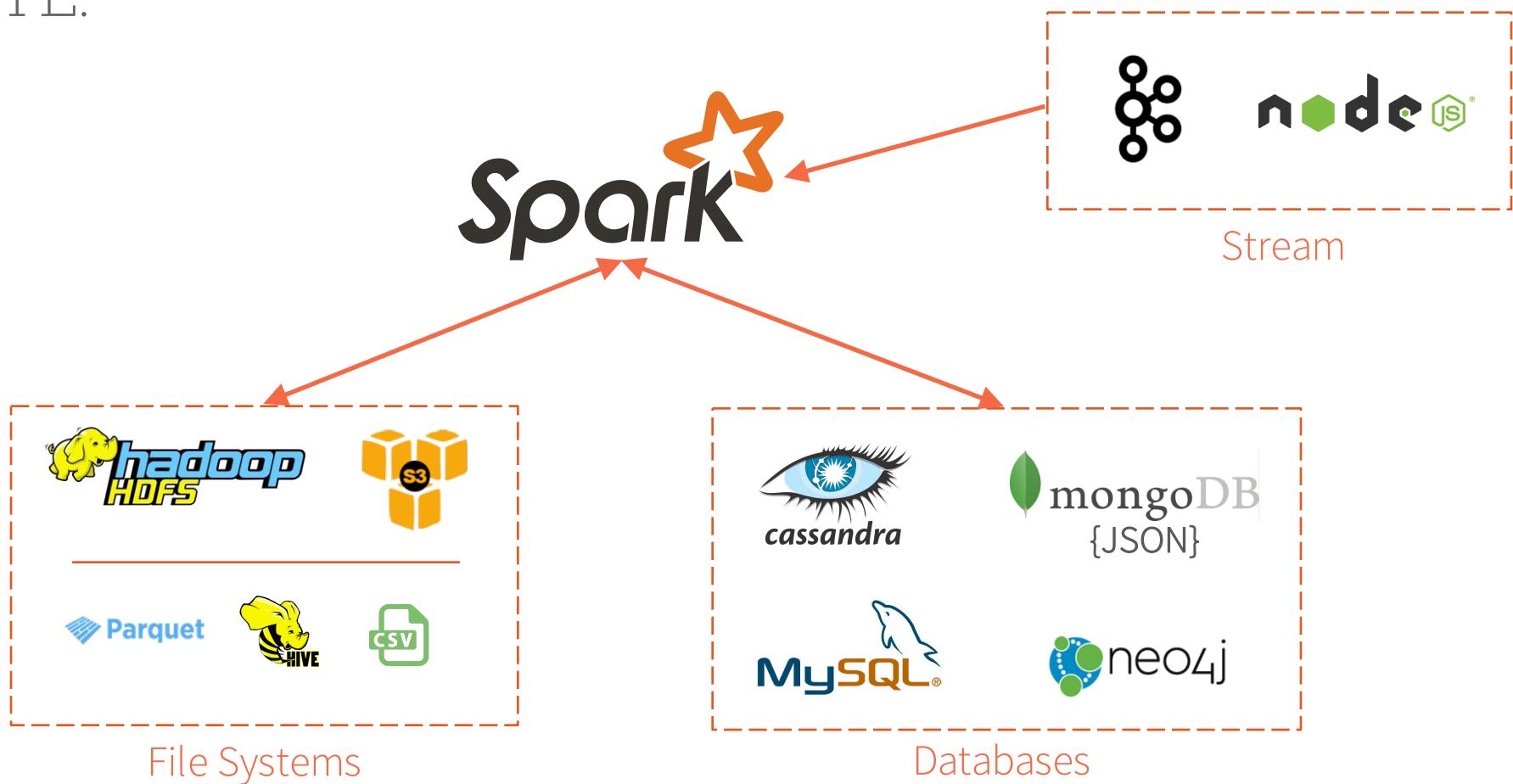
Dataframes, SQL, ML

6 Lang Wikipedias (live streams)

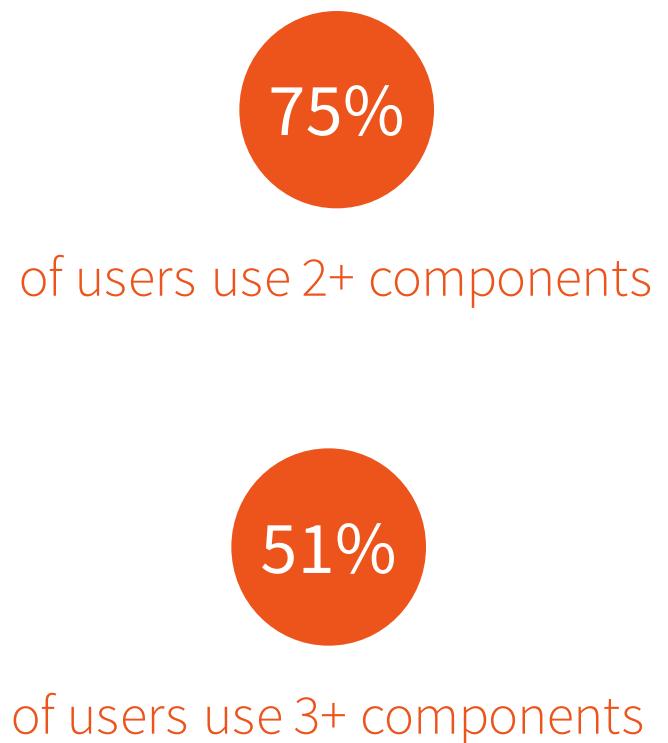
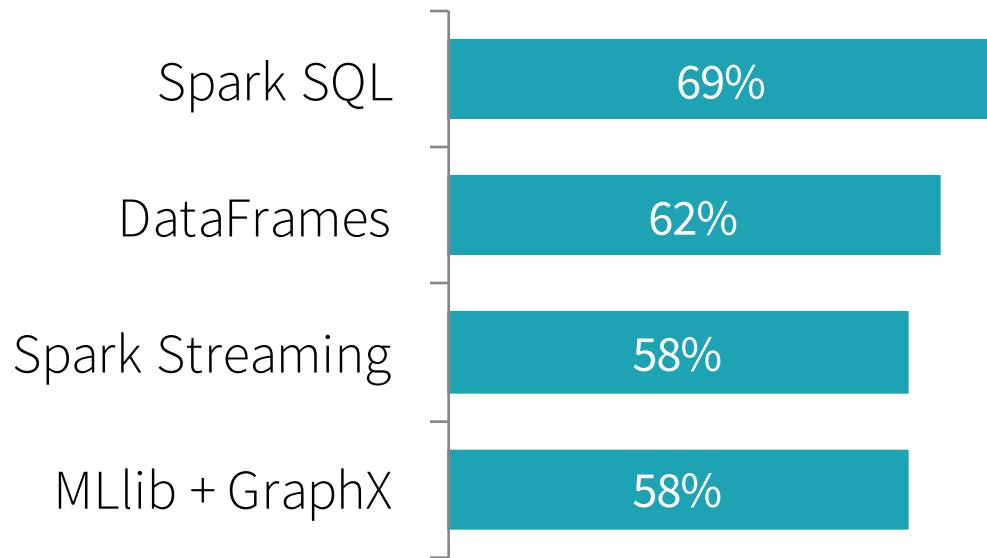


Streaming, RDD, Dataframes, SQL

ETL:



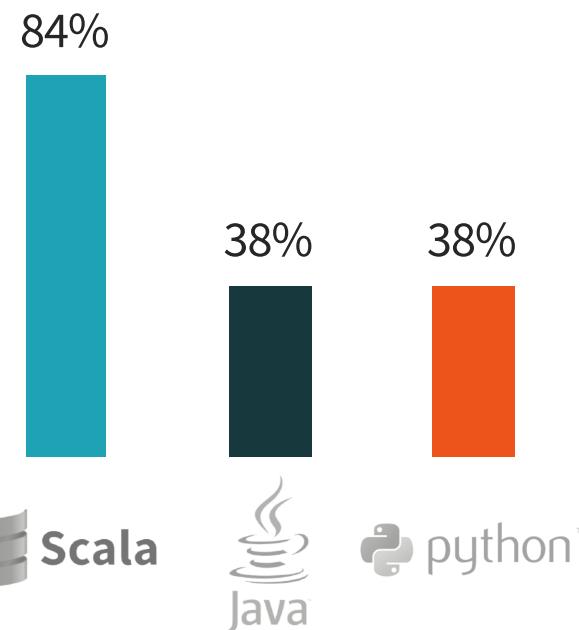
Which Libraries Do People Use?



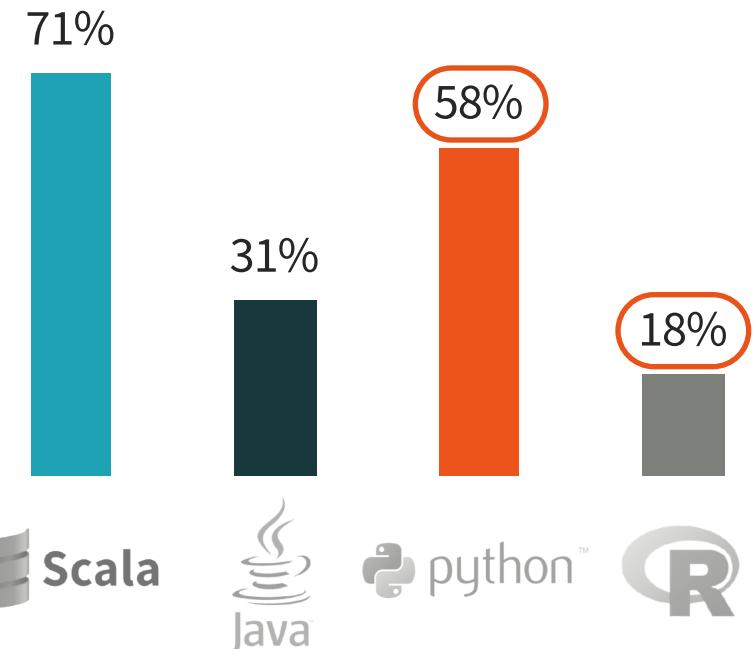
Source: <http://dbricks.co/1PtWxSu>

Which Languages Do People Use?

2014 Languages Used

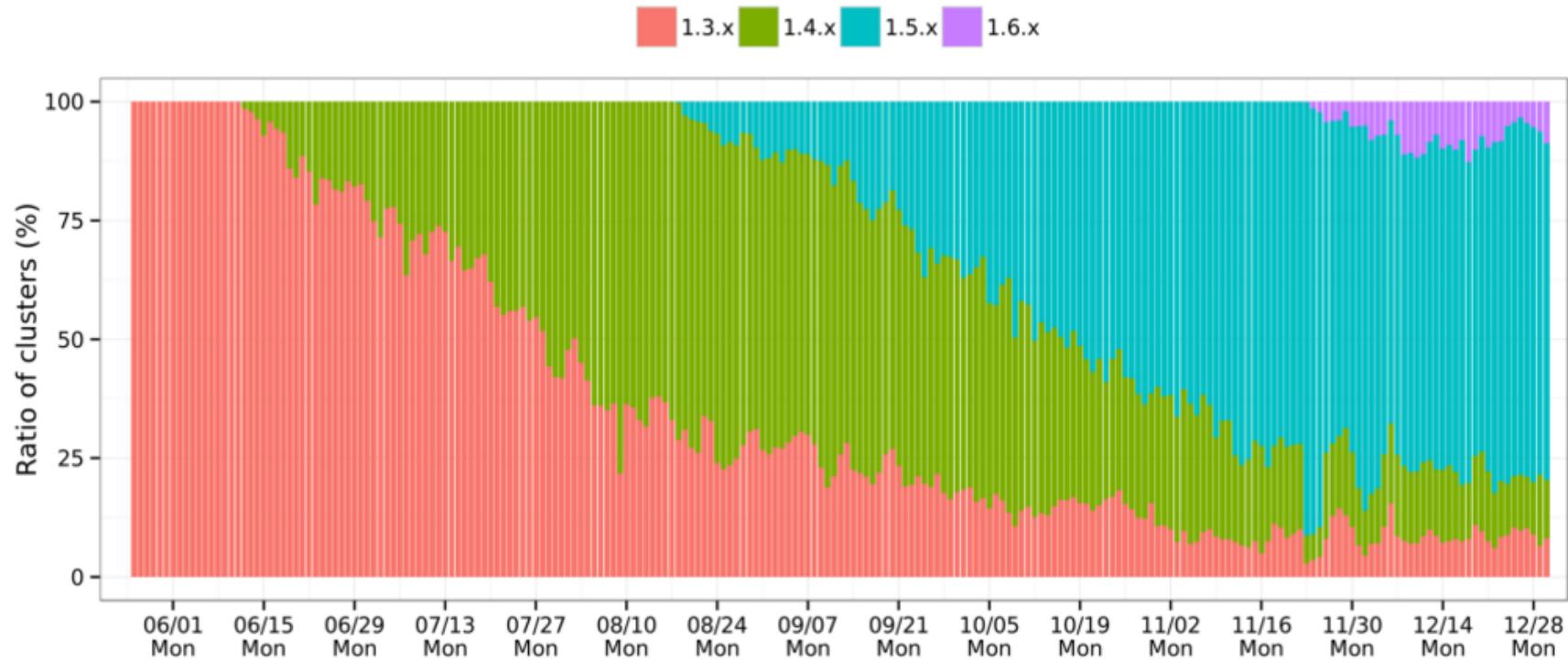


2015 Languages Used



Source: <http://dbricks.co/1PtWxSu>

Percent of clusters by Spark version in 2015



Reynold Xin @rxin · Jan 14

The most interesting thing in Spark 2015 Year in Review post is the version dist over time. databricks.com/blog/2016/01/0...



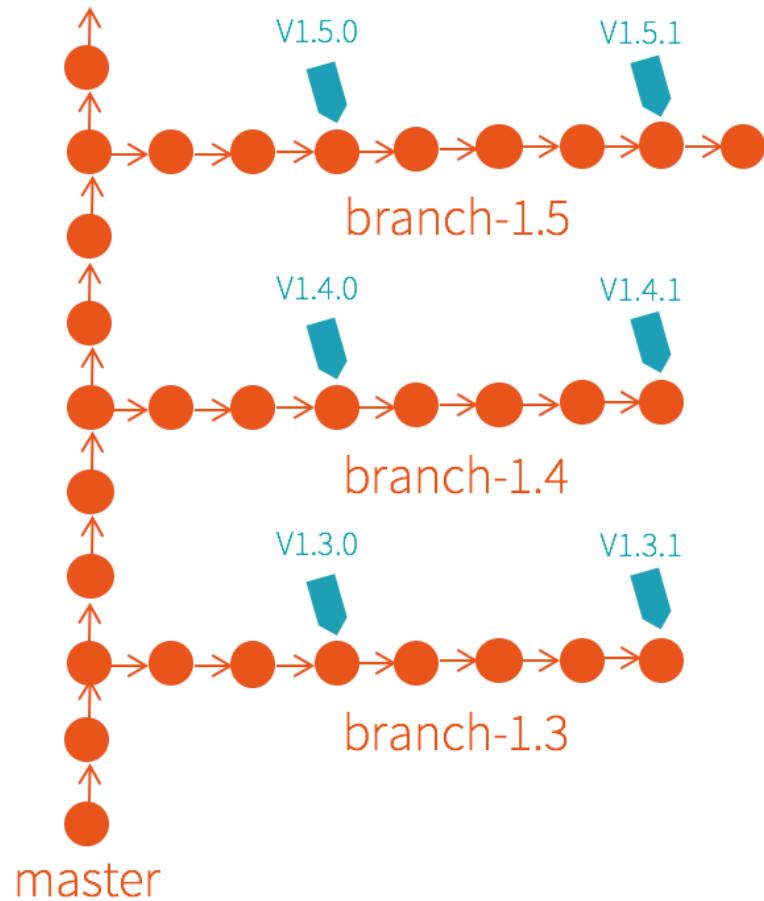
9



19

...

What do the version #s mean?



Spark project goal

Support diverse workflows:

- Batch (*jobs that run for hours*)
- Interactive / Ad Hoc (*via shell or notebooks*)
- Iterative computations (*graph and ML*)
- Streaming

Large Scale Usage:

Largest cluster: 8000 nodes 

Largest single job: 1 petabyte  

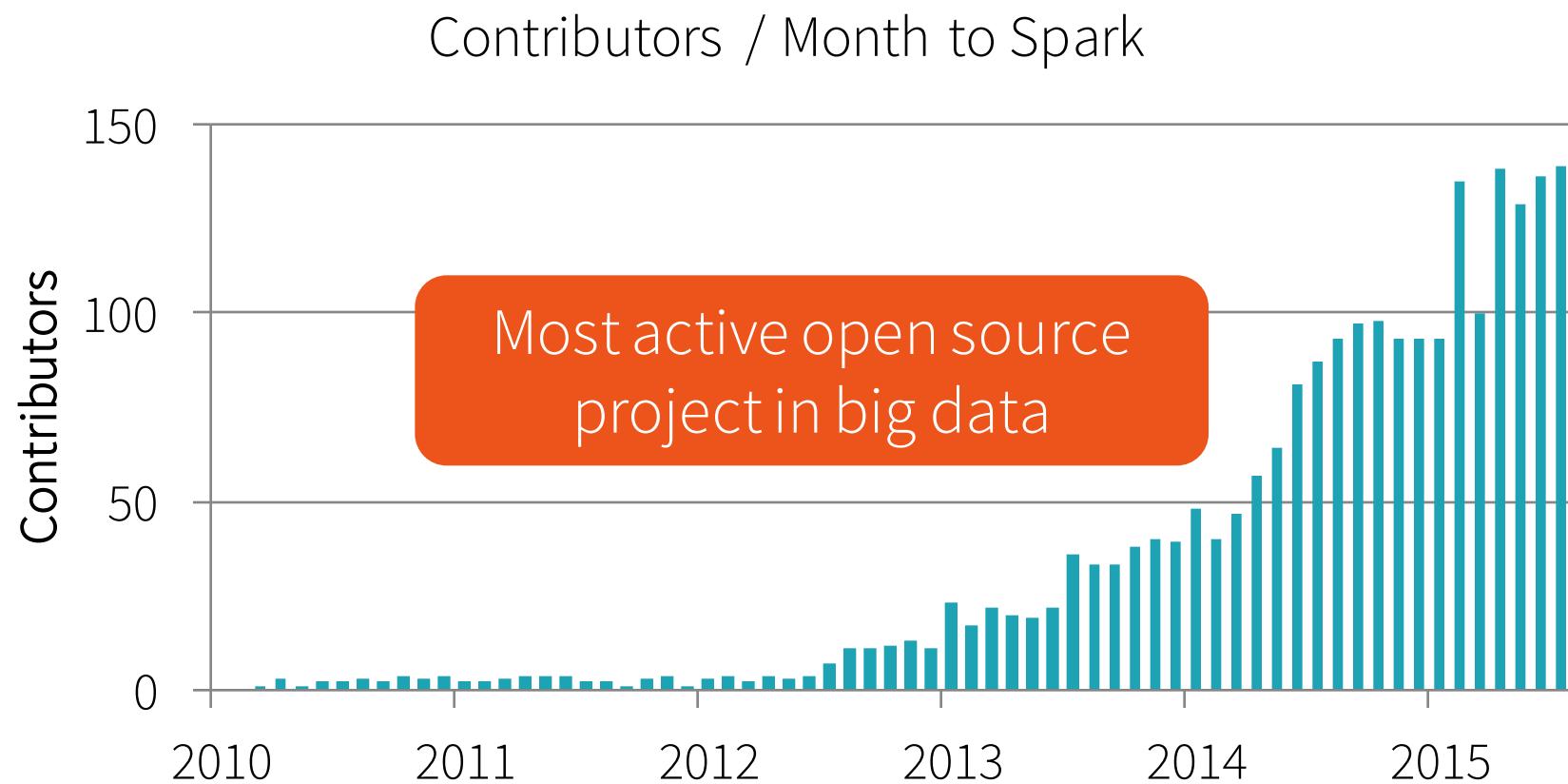
Top streaming intake: 1 TB/hour 

2014 on-disk 100 TB sort record: 23 mins / 
207 EC2 nodes

Spark Users: 1000+ production deployments



Community Growth:



Functional Programming & Scala

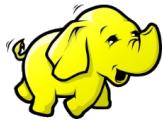
Typically used when a system's performance and integrity are both critical

2 concepts:

- **Immutable Data:** Create new data structures instead of modifying ones
that already exist
- **Stateless:** Perform every task as if for the first time, with no knowledge of earlier state

Benefits: *simplifies multithreading, removes bugs around race conditions and deadlocks, reusable functions*

Comparison to Hadoop



MapReduce

YARN

HDFS

Hadoop: combined compute,
resource management + storage



Spark

YARN

Standalone

Mesos

..

HDFS

NoSQL

SQL

..

Spark: independent of resource
management + storage layers

“Simple things should be simple, complex things should be possible”

-Alan Kay





Spark | DataFrames

I have used **DataFrames** before (in any language)...



DataFrame: pageviewsDF

Time (Str)			Site (Str)			Req (Int)			Time (Str)			Site (Str)			Req (Int)			Time (Str)			Site (Str)			Req (Int)		
ts	m	1304	ts	d	3901	ts	m	1172	ts	m	2538	ts	m	2237	ts	m	2137	ts	d	2837	ts	d	3176	ts	d	3400
ts	d	2237	ts	d	2491	ts	d	3176	ts	d	3400	ts	d	2288	ts	d	2837	ts	m	1600	ts	m	1304	ts	m	2491
ts	m	1600	ts	d	2288	ts	d	2837	ts	m	1304	ts	d	2491	ts	m	1172	ts	m	2538	ts	d	2237	ts	m	3901
ts	m	2491	ts	d	3176	ts	d	3400	ts	m	2237	ts	d	2837	ts	m	1600	ts	d	2288	ts	m	1304	ts	d	2491

Partition 1

Partition 2

Partition 3

Partition 4

`df.rdd.partitions.size = 4`

 Operations =

+



TRANSFORMATIONS



ACTIONS

DataFrame: Transformations and Actions

Transformations (<i>lazy</i>)	Actions
orderBy	show
filter	count
groupBy	take
select	collect
drop	save
join	

Transformations contribute to a query plan,
but nothing is executed until an action is called

DataF

#ABCDEFGHJKLMNOPQRSTUVWXYZ

display packages only

org.apache.spark.sql **hide focus**

- DataFrame** ←
- DataFrameHolder
- DataFrameNaFunctions
- DataFrameReader
- DataFrameStatFunctions
- DataFrameWriter

DataFrame

class **DataFrame** extends Queryable with Serializable

A distributed collection of data organized into named columns. Experimental

A **DataFrame** is equivalent to a relational table in Spark SQL. The following example creates a **DataFrame** by pointing Spark SQL to a Parquet data set.

```
val people = sqlContext.read.parquet("...") // in Scala  
DataFrame people = sqlContext.read().parquet("...") // in Java
```

Once created, it can be manipulated using the various domain-specific-language (DSL) functions defined in: [DataFrame](#) (this class), [Column](#), and [functions](#).

To select a column from the data frame, use `apply` method in Scala and `col` in Java.

```
val ageCol = people("age") // in Scala  
Column ageCol = people.col("age") // in Java
```

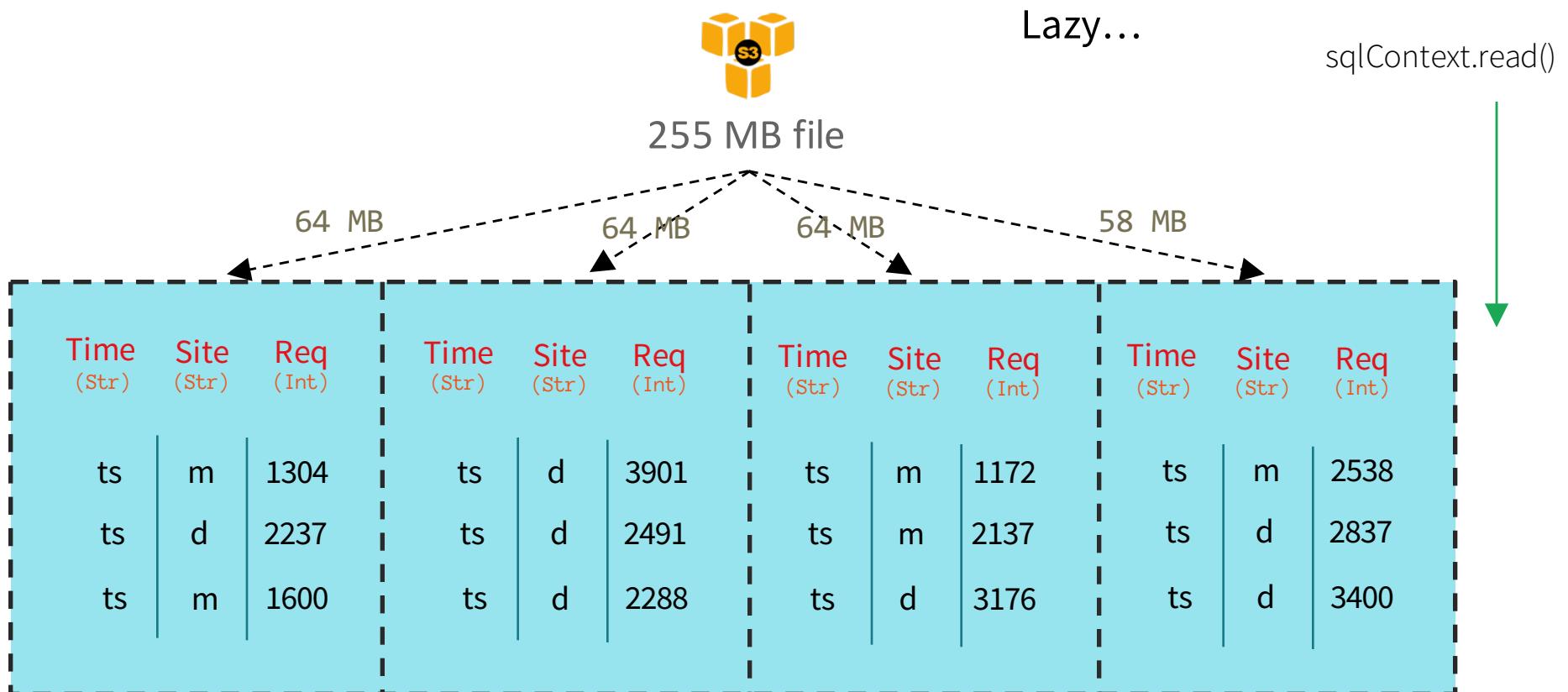
Note that the [Column](#) type can also be manipulated through its various functions.

```
// The following creates a new column that increases everybody's age by 10.  
people("age") + 10 // in Scala  
people.col("age").plus(10); // in Java
```

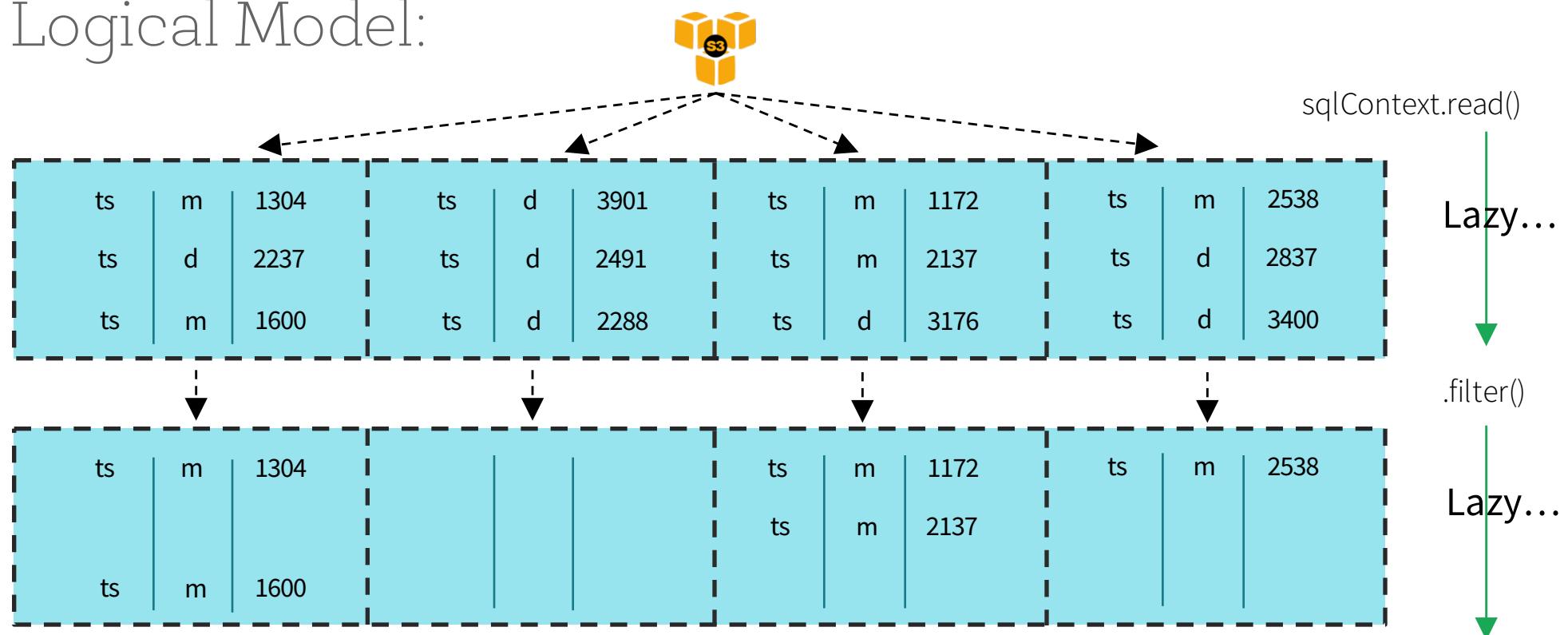
A more concrete example in Scala:

```
// To create DataFrame using SQLContext  
val people = sqlContext.read.parquet("...")  
val department = sqlContext.read.parquet("...")  
  
people.filter("age > 30")  
  .join(department, people("deptId") === department("id"))  
  .groupBy(department("name"), "gender")  
  .agg(avg(people("salary")), max(people("age")))
```

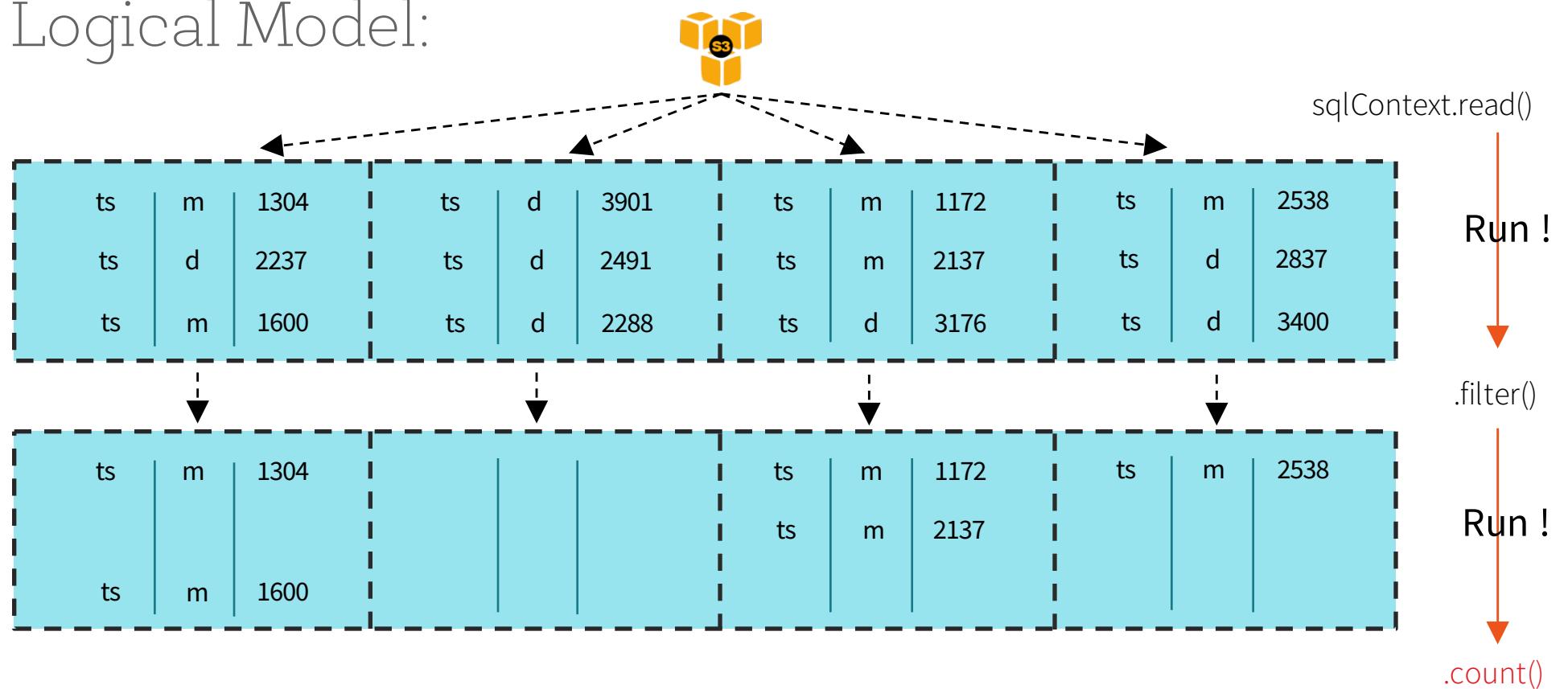
DataFrame: pageviewsDF



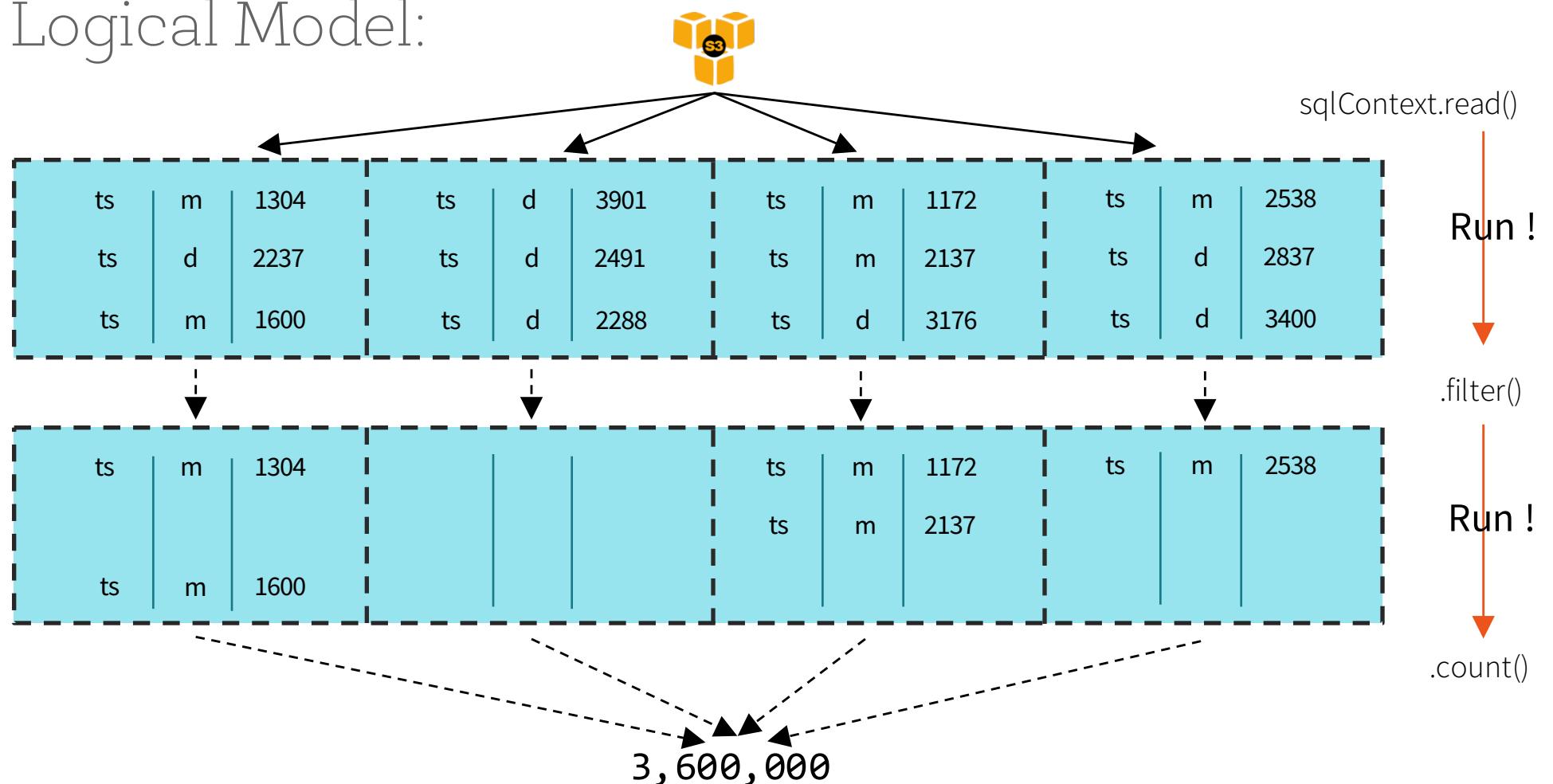
Logical Model:



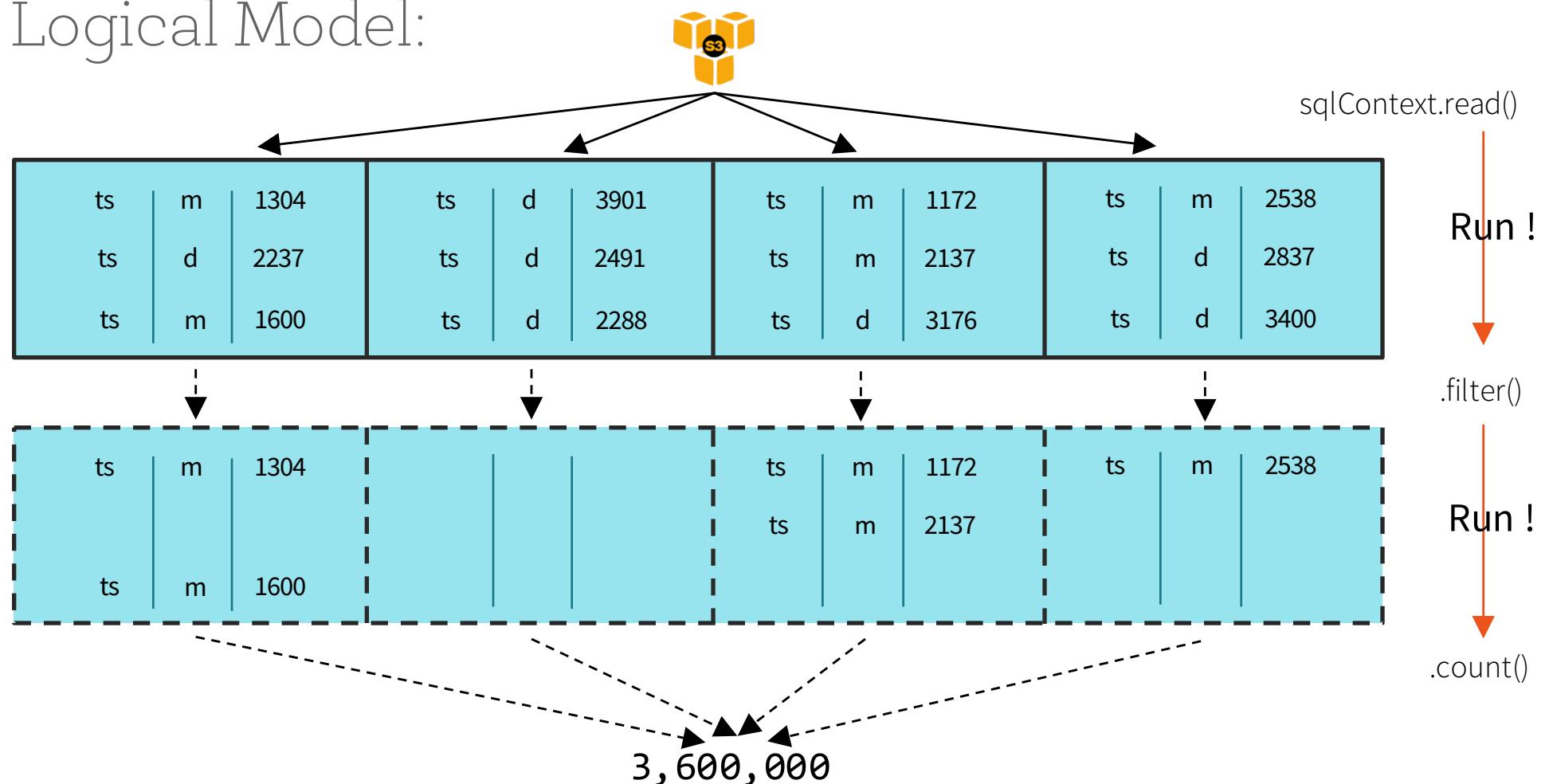
Logical Model:



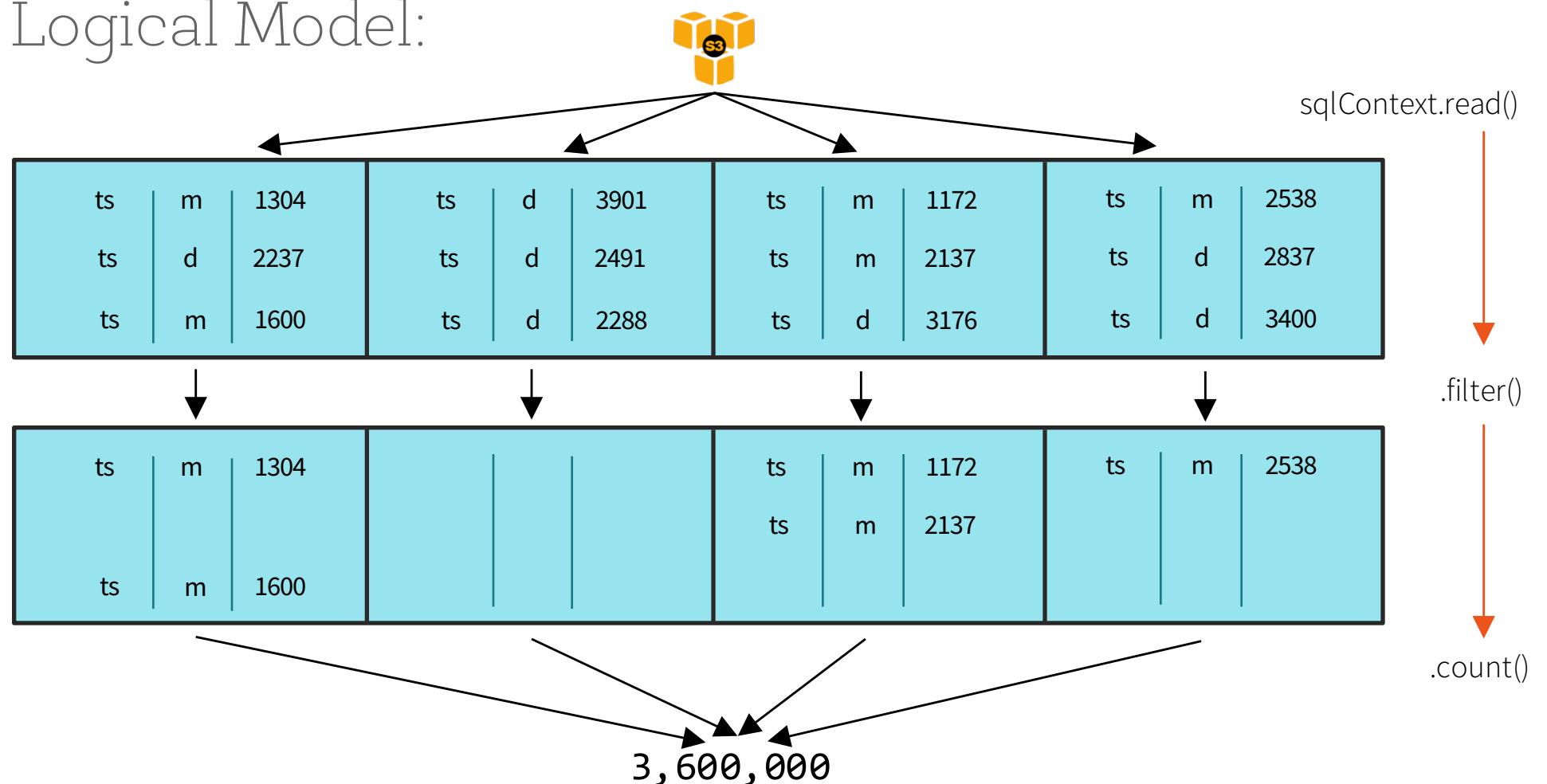
Logical Model:



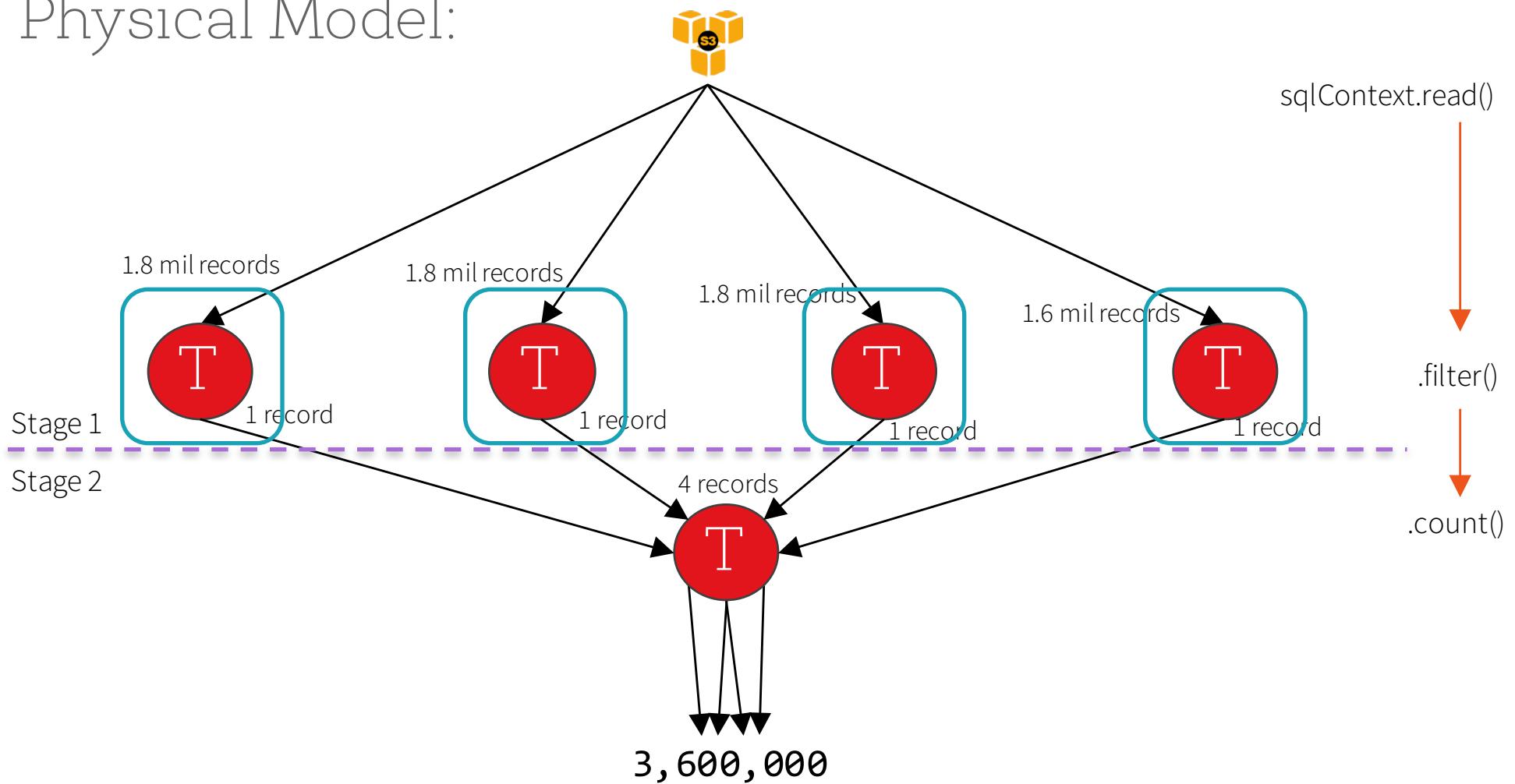
Logical Model:



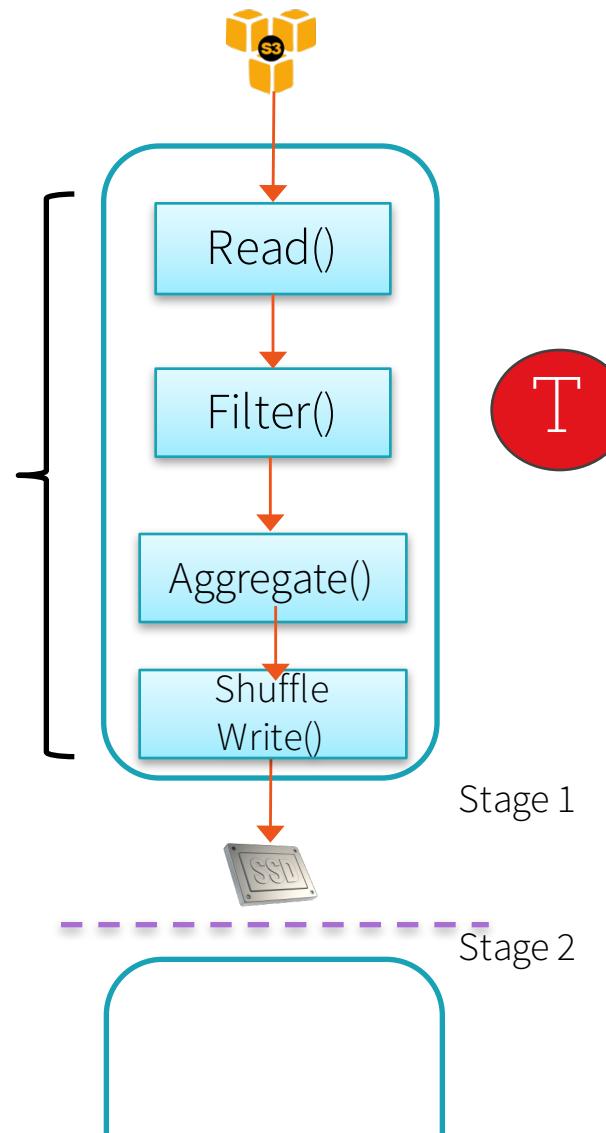
Logical Model:



Physical Model:



Pipelining



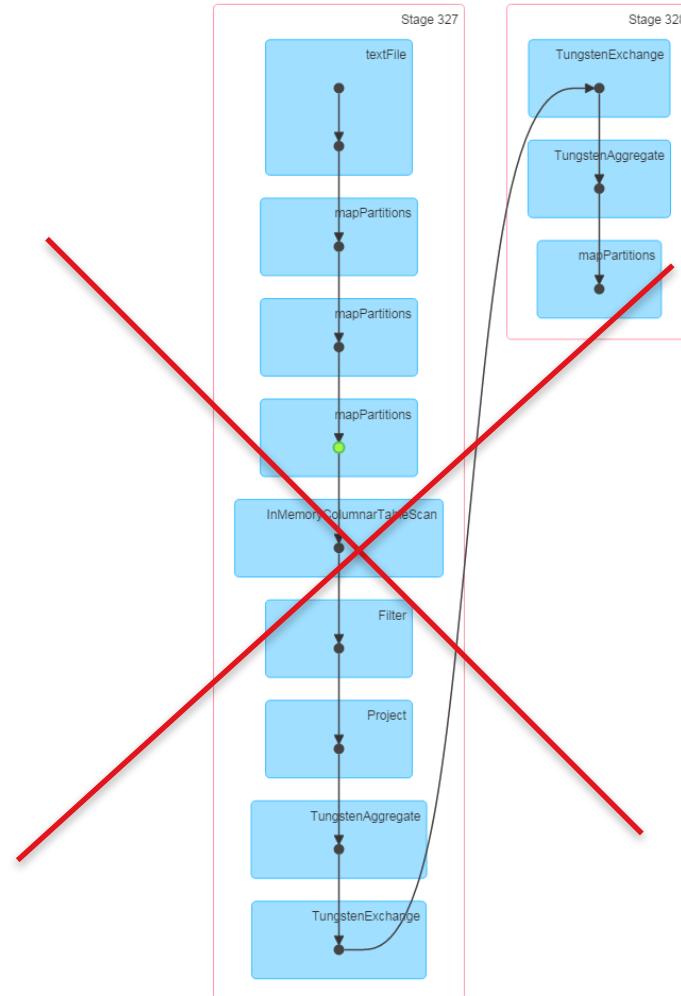
Physical Model:

Completed Stages (2)

Stage Id	Pool Name	Description		Submitted	Duration	Tasks: Succeeded/Total	Input	Output	Shuffle Read	Shuffle Write
114	1622753363940683362	pageviewsDF.filter(\$"site" === "mobile").count() count at <console>:27	+details	2015/11/30 09:39:09	7 ms	1/1			168.0 B	
113	1622753363940683362	pageviewsDF.filter(\$"site" === "mobile").count() count at <console>:27	+details	2015/11/30 09:38:56	12 s	4/4	250.0 MB		168.0 B	

Physical Model:

(DAG Visualization / Job)



Physical Model:

(DAG Visualization / SQL)

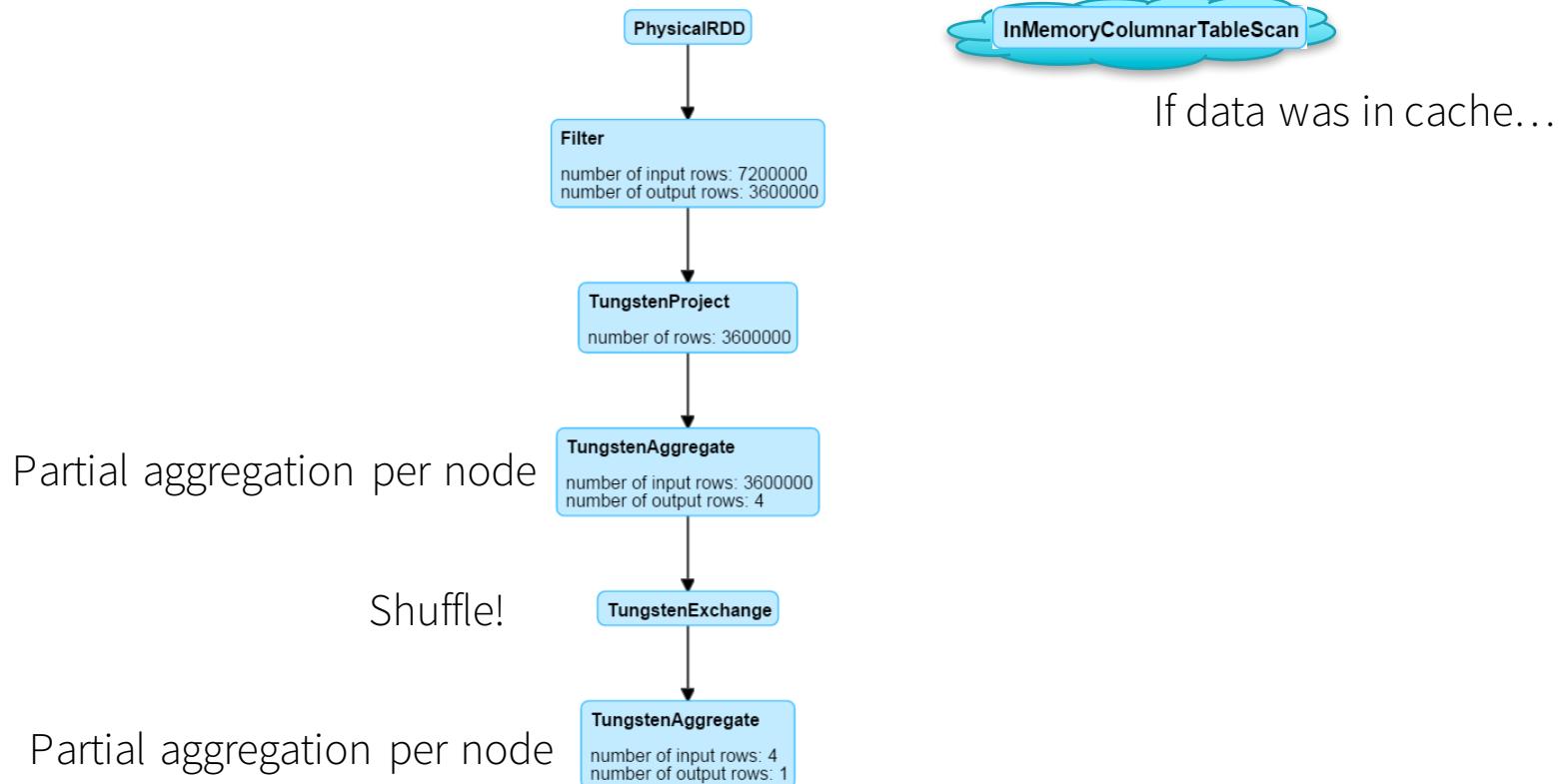
Jobs Stages Storage Environment Executors **SQL** JDBC/ODBC Server

Details for Query 5

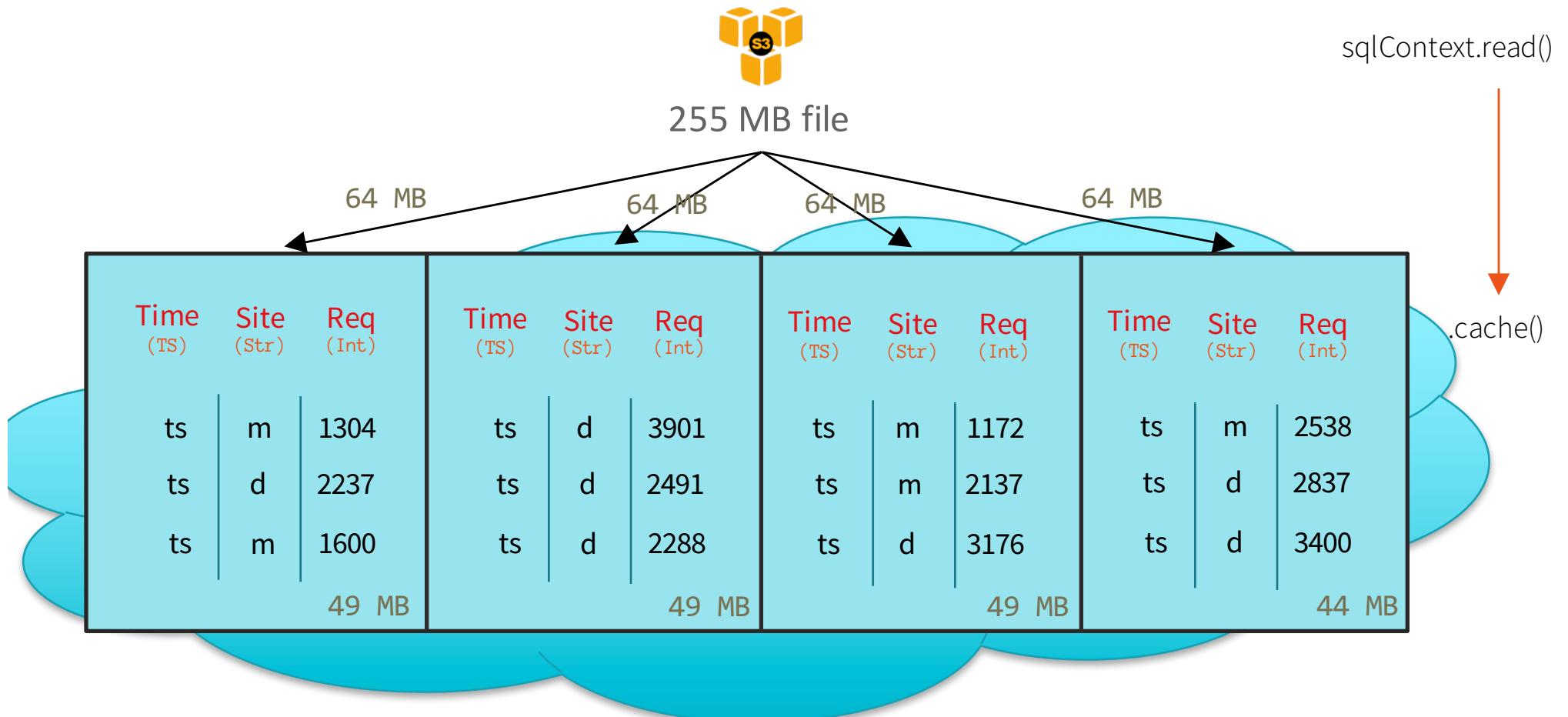
Submitted Time: 2015/11/30 09:38:56

Duration: 13 s

Succeeded Jobs: 85



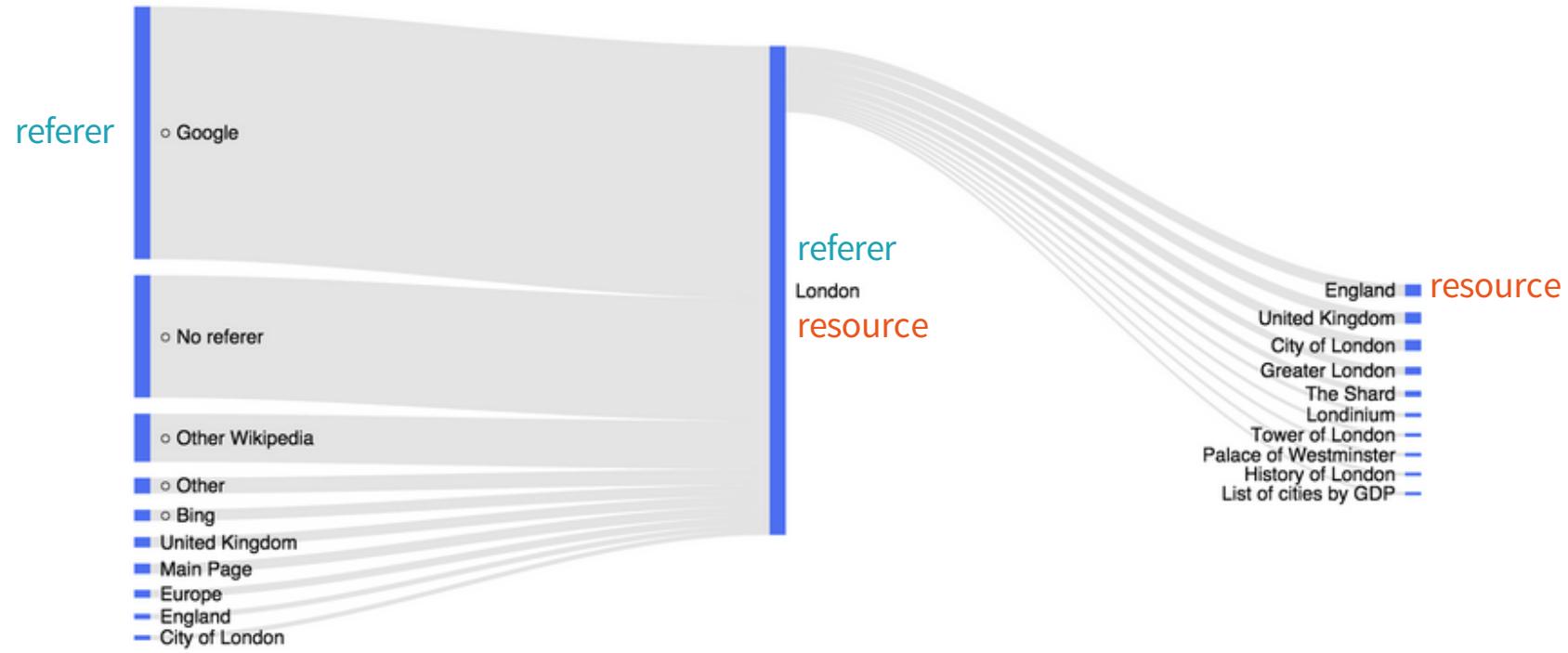
DataFrame: pageviewsDF





clickstream

(1.2 GB)



Wikipedia Clickstream

By Ellery Wulczyn
Data Scientist @ The Wikimedia Foundation

Clickstream Data: 1.2 GB tsv file

22 million

(**referer**, **resource**) pairs

out of 3.2 billion requests

in February 2015.

(only includes requests for articles in the desktop version of English Wikipedia)

```
> display(clickstreamDF)
```

(referer, resource)

prev_id	curr_id	n	prev_title	curr_title	type
	3632887	121	other-google	!!	other
	3632887	93	other-wikipedia	!!	other
	3632887	46	other-empty	!!	other
	3632887	10	other-other	!!	other
64486	3632887	11	!_(disambiguation)	!!	other
2061699	2556962	19	Louden_Up_Now	!!!_(album)	link
	2556962	25	other-empty	!!!_(album)	other
	2556962	16	other-google	!!!_(album)	other
	2556962	11	other-wikipedia	!!!_(album)	other

Showing the first 1000 rows.

February 2015

22 million records (*out of 3.2 billion requests*)

(referer)

article in the main namespace (ns) of Eng Wikipedia	→	article title
any Wikipedia page not in main Eng ns	→	other-wikipedia
an empty referer	→	other-empty
a page from any other Wikimedia project	→	other-internal
Google	→	other-google
Yahoo	→	other-yahoo
Bing	→	other-bing
Facebook	→	other-facebook
Twitter	→	other-twitter
anything else	→	other-other

prev_id	curr_id	n	prev_title	curr_title	type
	3632887	121	other-google	!!	other
	3632887	93	other-wikipedia	!!	other
	3632887	46	other-empty	!!	other

n: the number of occurrences of the (referer, resource) pair

prev_id	curr_id	n	prev_title	curr_title	type
	3632887	121	other-google	!!	other
	3632887	93	other-wikipedia	!!	other
	3632887	46	other-empty	!!	other

link: if the referer and request are both articles and the referer links to the request

redlink: links to a Wikipedia page that do not exist (*uncreated pages*)

other: if the referer + request are both articles but the referer does not link to the request.
(this can happen when clients search or spoof their refer)

The screenshot shows a web browser window with the following details:

- Title Bar:** Wikipedia Clickstream - F
- Address Bar:** datahub.io/dataset/wikipedia-clickstream/resource/be85cc68-d1e6-4134-804a-fd36b94dbb82
- Header:** datahub (The easy way to get, use and share data) with links to Datasets, Organizations, About, Blog, Help, and a Search bar.
- Breadcrumbs:** Home / Organizations / Wikimedia / Wikipedia Clickstream / February 2015 English ...
- Section Header:** February 2015 English Wikipedia Clickstream
- Text:** URL: http://files.figshare.com/1927919/2015_02_clickstream.tsv.gz
- Section:** Data Preparation
- Description:** The data contains counts of (referrer, resource) pairs extracted from the request logs of English Wikipedia. When a client requests a resource by following a link or performing a search, the URI of the webpage that linked to the resource is included with the request in an HTTP header called the "referer". This data captures 22 million (referrer, resource) pairs from a total of 3.2 billion requests collected during the month of February 2015.
- Note:** The dataset only includes requests for articles in the main namespace of the desktop version of English Wikipedia.
- Referrer Mapping:** Refers to a fixed set of values corresponding to internal traffic or external traffic from one of the top 5 global traffic sources to English Wikipedia, based on this scheme:
 - an article in the main namespace of English Wikipedia -> the article title
 - any Wikipedia page that is not in the main namespace of English Wikipedia -> 'other-wikipedia'
 - an empty referer -> 'other-empty'
 - a page from any other Wikimedia project -> 'other-internal'
 - Google -> 'other-google'
 - Yahoo -> 'other-yahoo'
 - Bing -> 'other-bing'
 - Facebook -> 'other-facebook'
 - Twitter -> 'other-twitter'
 - anything else -> 'other-other'
- Buttons:** Go to resource



What does the traffic inflow vs outflow look like for the **San Diego** article?

Not logged in | Talk | Contributions | Create account | Log in

Article | Talk | Read | View source | View history | Search |

From Wikipedia, the free encyclopedia

San Diego

City of San Diego

San Diego is a city in California. For the county in California, see San Diego County. For other uses, see San Diego (disambiguation).
"San Diego" redirects here. For the historical town, see San Diego (town).

San Diego (seen it: *español* (Spanish for "Saint Didacus")) is a major city in California, on the coast of the Pacific Ocean in Southern California, approximately 123 miles (190 km) south of Los Angeles and immediately adjacent to the border with Mexico.

With an estimated population of 1,381,069 as of July 1, 2014,^[1] San Diego is the second largest city in the state of California and second largest in California after San Diego is the "Empire of California"^[2] and is known for its mild year-round climate, natural deepwater harbor, extensive beaches, long association with the U.S. Navy, and recent emergence as a healthcare and biotechnology center.

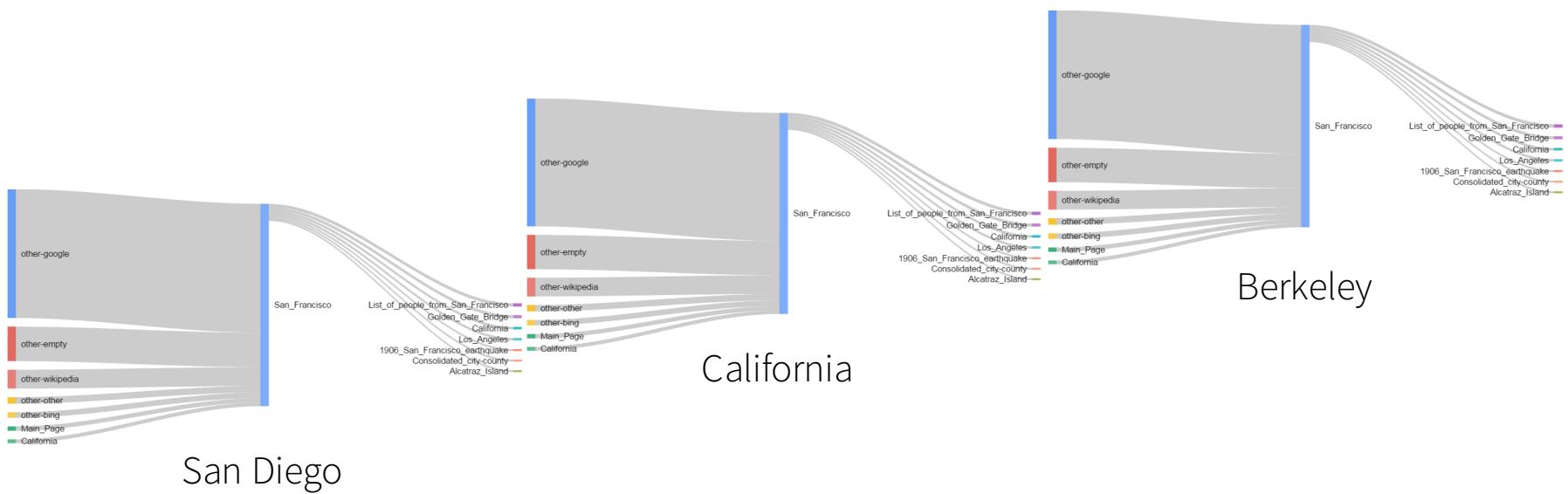
Historically home to the Kumeyaay people, San Diego was the first site visited by Europeans on what is now the West Coast of the United States. Upon landing in San Diego Bay in 1542, Juan Rodriguez Cabrillo claimed the entire area for Spain, forming the basis for the settlement of Alta California 200 years later. The Port of San Diego was established in 1769, and the first non-native European settlement in what is now California. In 1821, San Diego became part of newly independent Mexico, and in 1850, became part of the United States following the Mexican-American War and the admission of California to the union.

Images from top, left to right: San Diego skyline; San Diego bridge; San Diego cathedral; San Diego park.

Main page | Contents | Featured content | Current events | Recent changes | Random page | Donate to Wikipedia | Wikipedia store | Help | About Wikipedia | Community portal | Recent changes | Contact page | Tools | What links here | Related changes | Upload files | Special pages | Permanent link | Page information | Improve item | Cite this page | Print/export | Download as PDF | Portable version | Other projects | Wikipedia Commons

Riding the Clickstream:

Apache Spark?

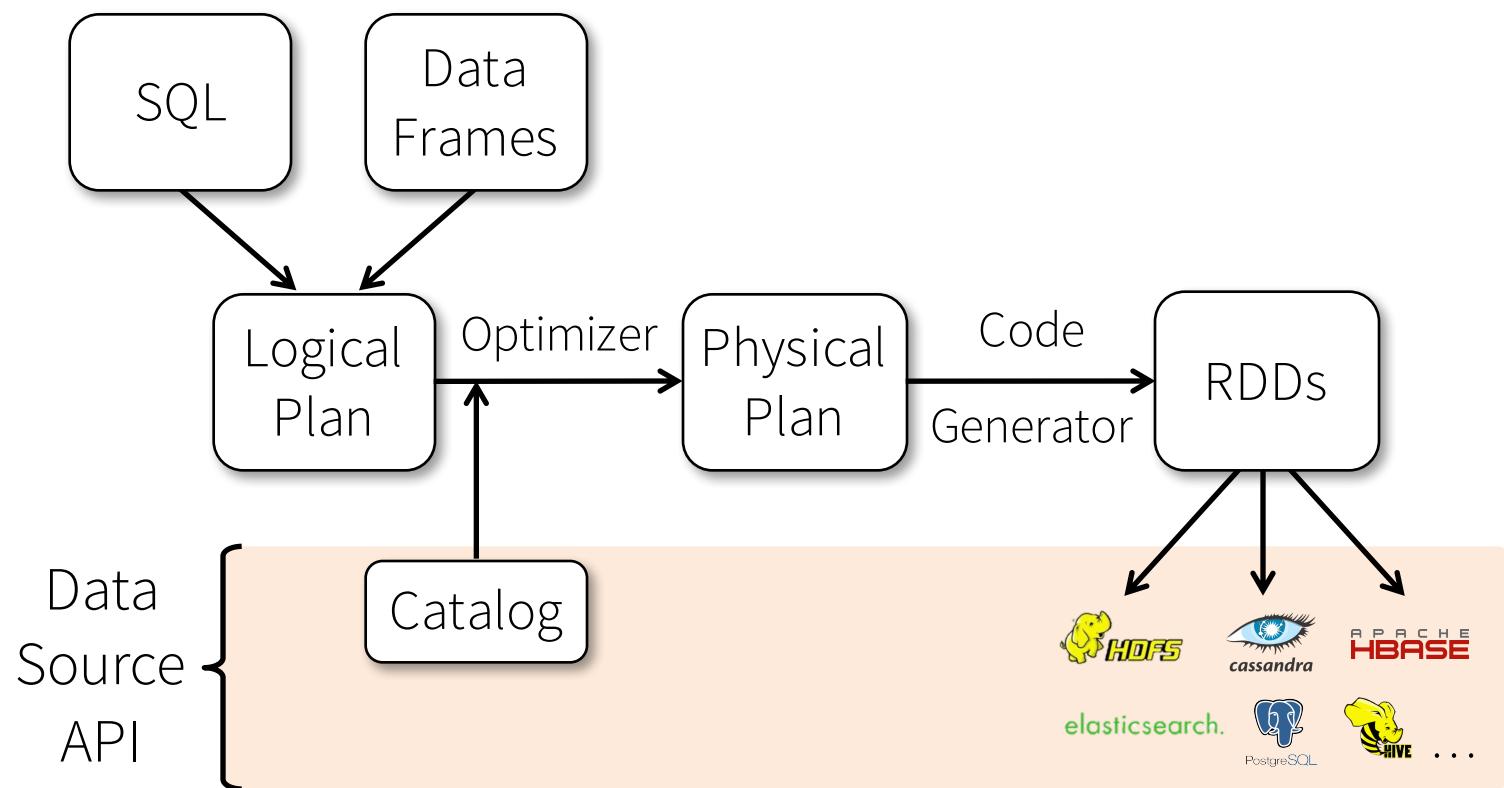




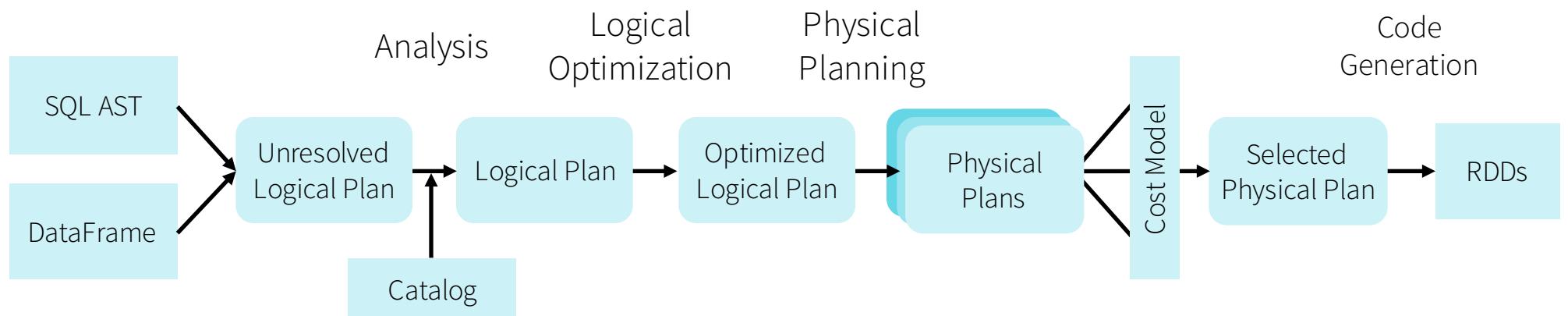
Spark Catalyst

 databricks training

Spark SQL Architecture

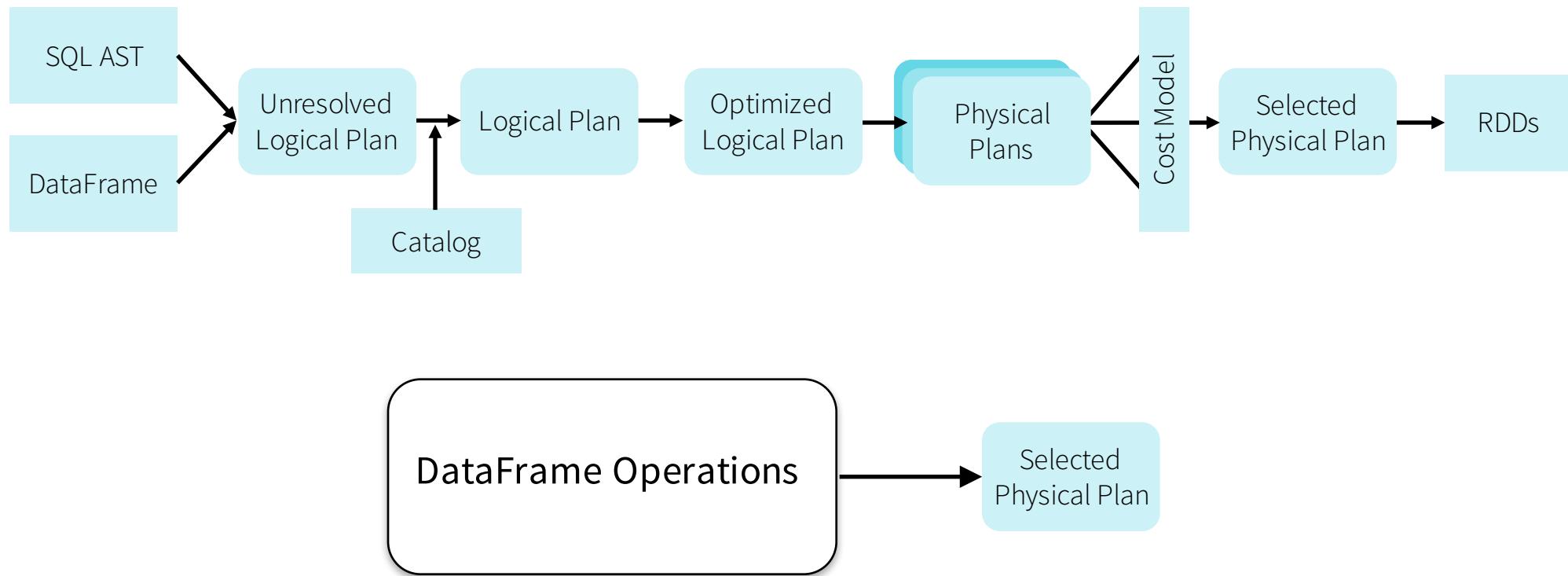


DataFrames + SQL: Under the hood



DataFrames and SQL share the same optimization/execution pipeline

DataFrames: Under the hood



▼ Details

```
== Parsed Logical Plan ==
Aggregate [count(1) AS count#25L]
  Filter (site#11 = mobile)
  Subquery pageviews_by_second
    Relation[timestamp#10,site#11,requests#12] CsvRelation(/FileStore/tables/rv3hb79r1449478143619,true, ,null,#,PERMISSIVE,COMMONS,false,false,StructType(StructField(timestamp,StringType,true), StructField(site,StringType,true), StructField(requests,IntegerType,true)),UTF-8,false)

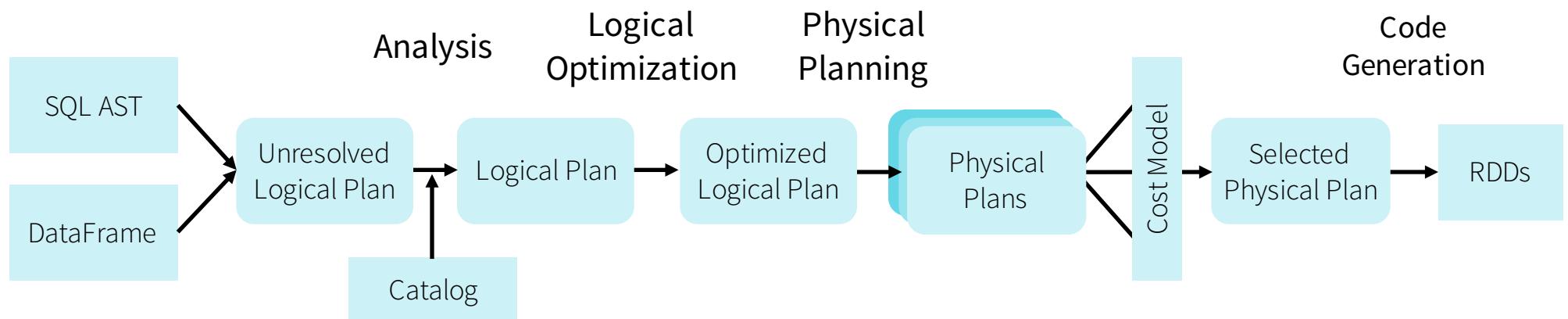
== Analyzed Logical Plan ==
count: bigint
Aggregate [count(1) AS count#25L]
  Filter (site#11 = mobile)
  Subquery pageviews_by_second
    Relation[timestamp#10,site#11,requests#12] CsvRelation(/FileStore/tables/rv3hb79r1449478143619,true, ,null,#,PERMISSIVE,COMMONS,false,false,StructType(StructField(timestamp,StringType,true), StructField(site,StringType,true), StructField(requests,IntegerType,true)),UTF-8,false)

== Optimized Logical Plan ==
Aggregate [count(1) AS count#25L]
  Project
    Filter (site#11 = mobile)
    Relation[timestamp#10,site#11,requests#12] CsvRelation(/FileStore/tables/rv3hb79r1449478143619,true, ,null,#,PERMISSIVE,COMMONS,false,false,StructType(StructField(timestamp,StringType,true), StructField(site,StringType,true), StructField(requests,IntegerType,true)),UTF-8,false)

== Physical Plan ==
TungstenAggregate(key=[], functions=[(count(1),mode=Final,isDistinct=false)], output=[count#25L])
  TungstenExchange SinglePartition
    TungstenAggregate(key=[], functions=[(count(1),mode=Partial,isDistinct=false)], output=[currentCount#28L])
      TungstenProject
        Filter (site#11 = mobile)
        Scan CsvRelation(/FileStore/tables/rv3hb79r1449478143619,true, ,null,#,PERMISSIVE,COMMONS,false,false,StructType(StructField(timestamp,StringType,true), StructField(site,StringType,true), StructField(requests,IntegerType,true)),UTF-8,false)[timestamp#10,site#11,requests#12]

Code Generation: true
```

Using Catalyst in Spark SQL



Analysis: analyzing a logical plan to resolve references

Logical Optimization: logical plan optimization

Physical Planning: Physical planning

Code Generation: Compile parts of the query to Java bytecode

Catalyst Optimizations

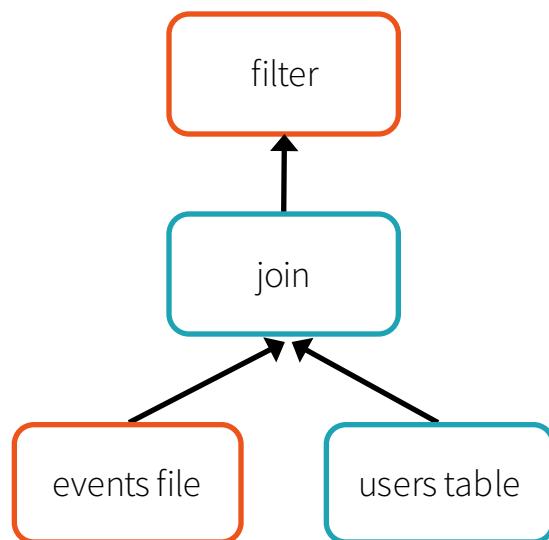
Logical Optimizations

- Push filter predicates down to data source, so irrelevant data can be skipped
- **Parquet:** skip entire blocks, turn comparisons on strings into cheaper integer comparisons via dictionary encoding
- **RDBMS:** reduce amount of data traffic by pushing predicates down

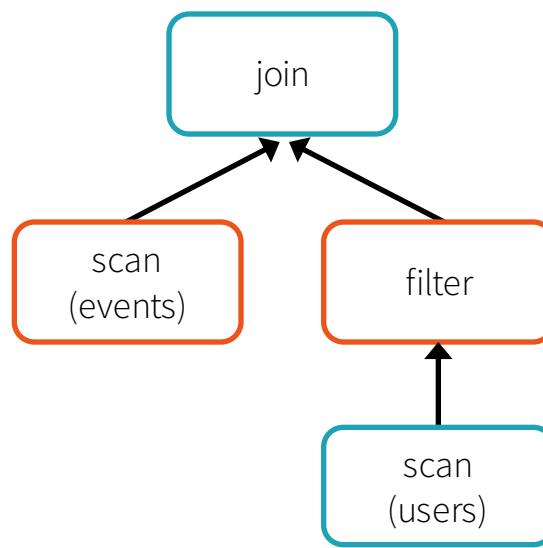
Create Physical Plan & generate JVM bytecode

- Catalyst compiles operations into physical plans for execution and generates JVM bytecode
- Intelligently choose between broadcast joins and shuffle joins to reduce network traffic
- **Lower level optimizations:** eliminate expensive object allocations and reduce virtual function calls

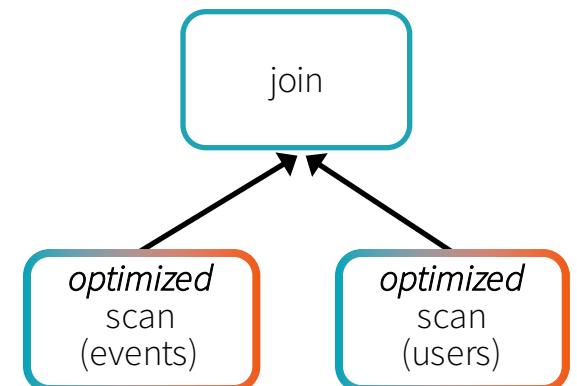
Logical Plan



Physical Plan



Optimized Physical Plan
with Predicate Pushdown
and Column Pruning



{JSON}

{JSON}

 Parquet



PRODUCT SPARK SOLUTIONS CUSTOMERS COMPANY BLOG RESOURCES

Trai

COMPANY

All Posts

Partners

Events

Press Releases

DEVELOPER

All Posts

Spark

Spark SQL

Spark Streaming

MLlib

Spark Summit

Search Blog



Subscribe



Deep Dive into Spark SQL's Catalyst Optimizer

April 13, 2015 | by Michael Armbrust, Yin Huai, Cheng Liang, Reynold Xin and Matei Zaharia



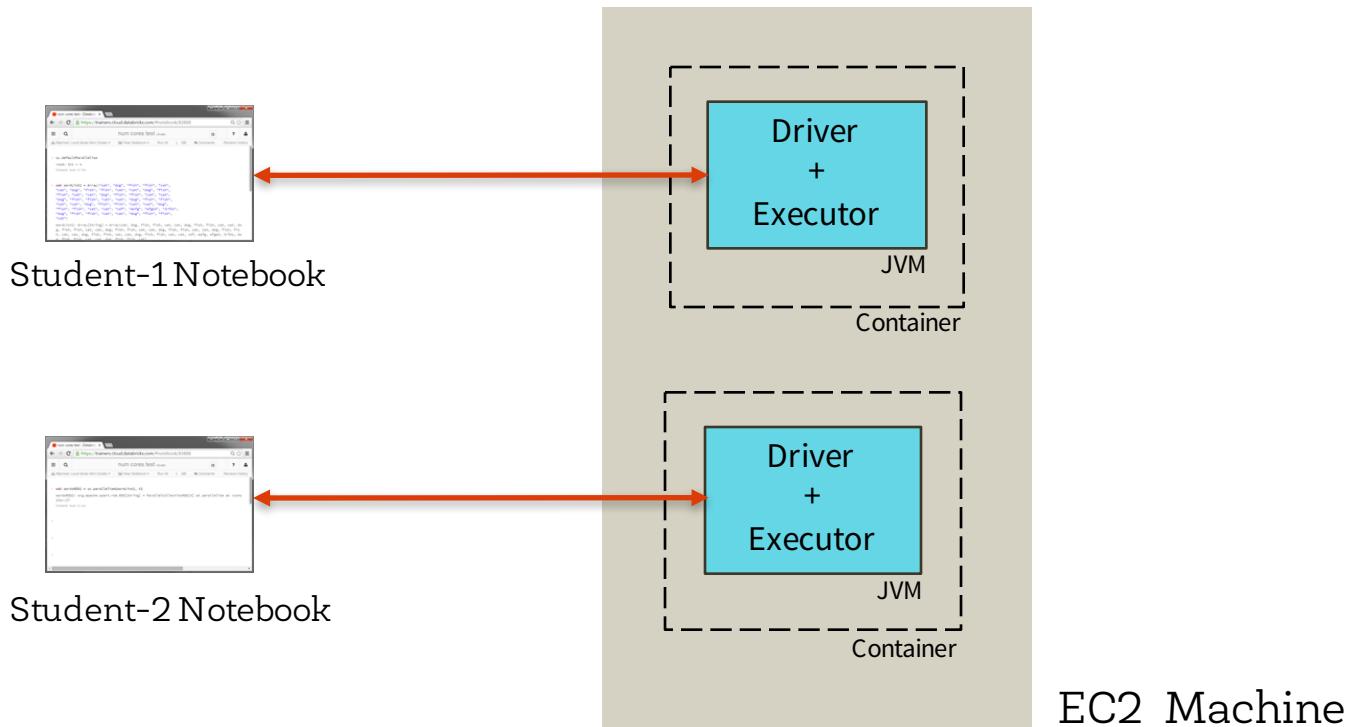
Spark SQL is one of the newest and most technically involved components of Spark. It powers both SQL queries and the new [DataFrame API](#). At the core of Spark SQL is the Catalyst optimizer, which leverages advanced programming language features (e.g. Scala's [pattern matching](#) and [quasiquotes](#)) in a novel way to build an extensible query optimizer.

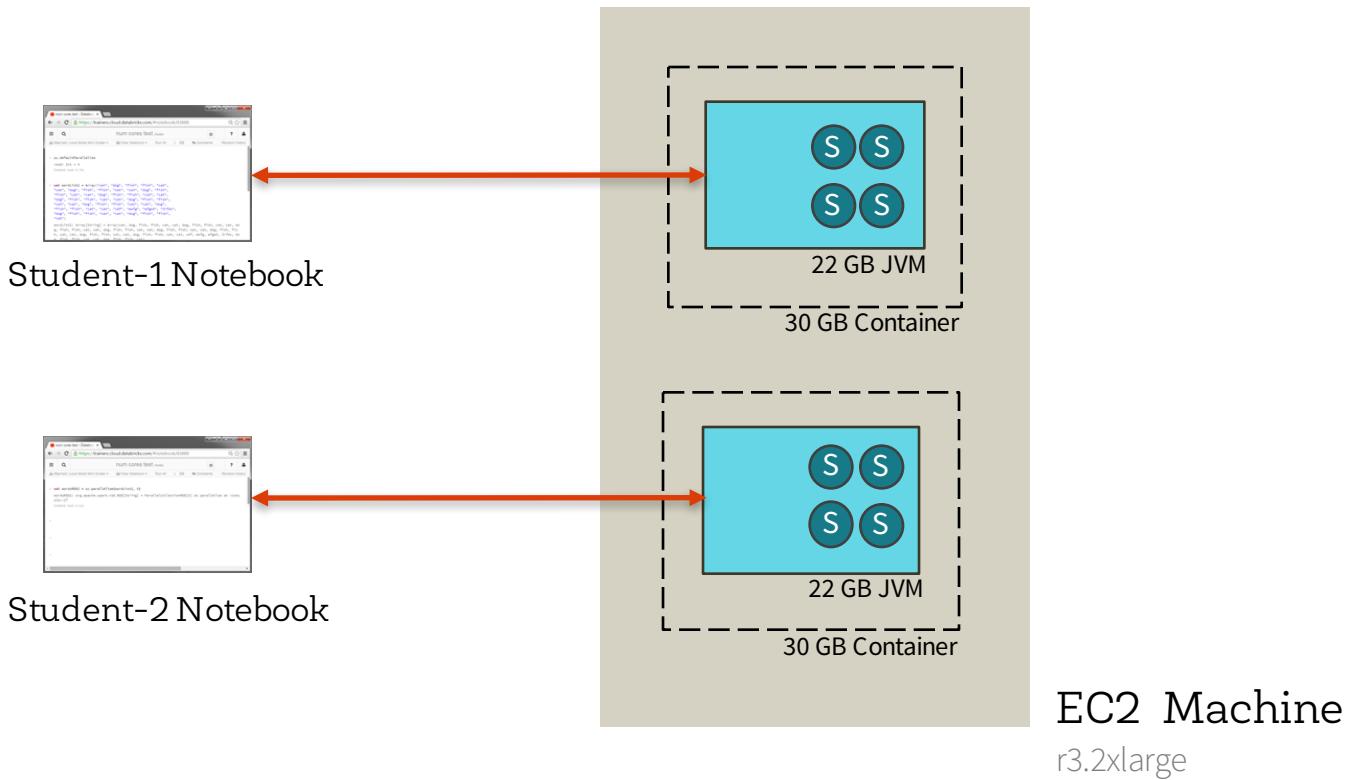
We recently published a [paper](#) on Spark SQL that will appear in [SIGMOD 2015](#) (co-authored with Davies Liu, Joseph K. Bradley, Xiangrui Meng, Tomer Kaftan, Michael J. Franklin, and Ali Ghodsi). In this blog post we are republishing a section in the paper that explains the internals of the Catalyst optimizer for broader consumption.

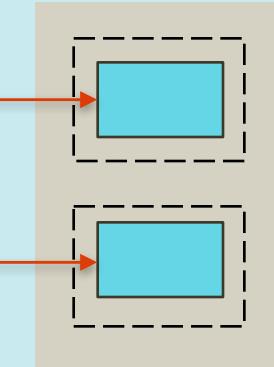
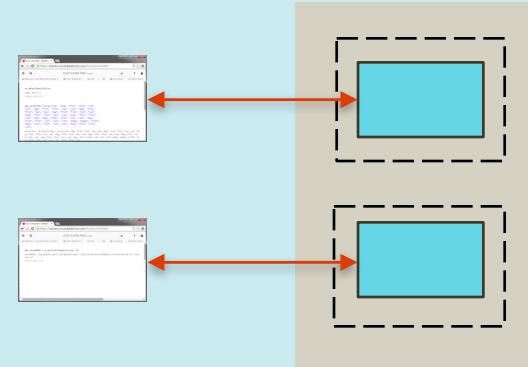
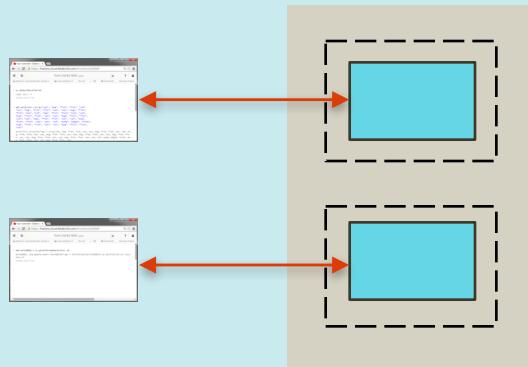
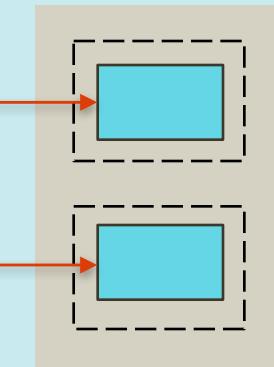
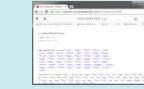
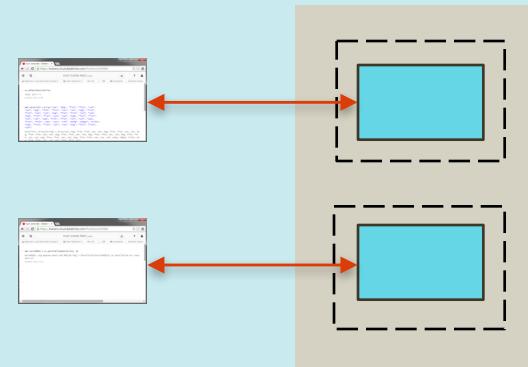
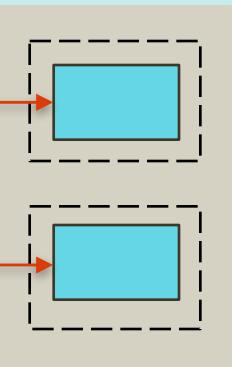
To implement Spark SQL, we designed a new extensible optimizer, Catalyst, based on functional programming constructs in Scala. Catalyst's extensible design had two purposes. First, we wanted to make it easy to add new optimization techniques and features to Spark SQL, especially for the purpose of tackling various problems we were seeing with big data (e.g., semistructured data and advanced analytics). Second, we wanted to enable external developers to extend the optimizer — for example, by adding data source specific rules that can push filtering or aggregation into external storage systems, or support for new data types. Catalyst supports both rule-based and cost-based optimization.

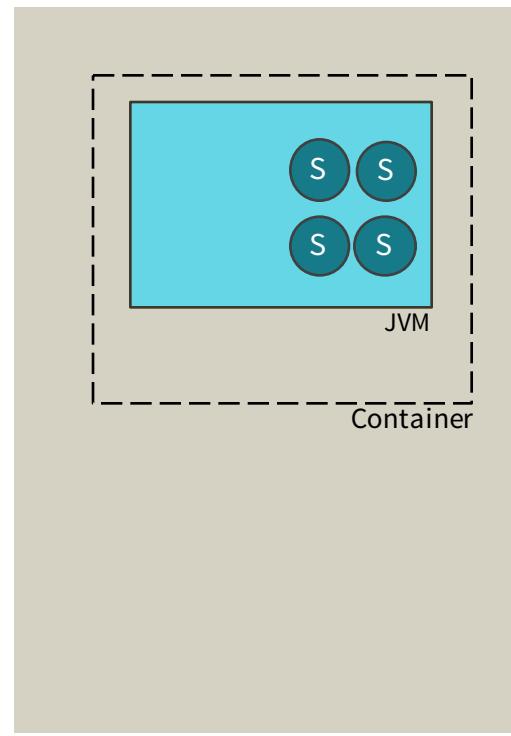


Architecture (local mode)

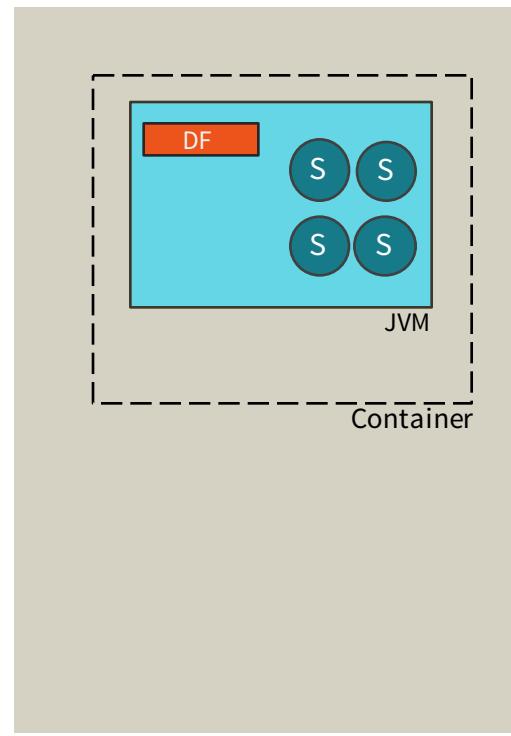




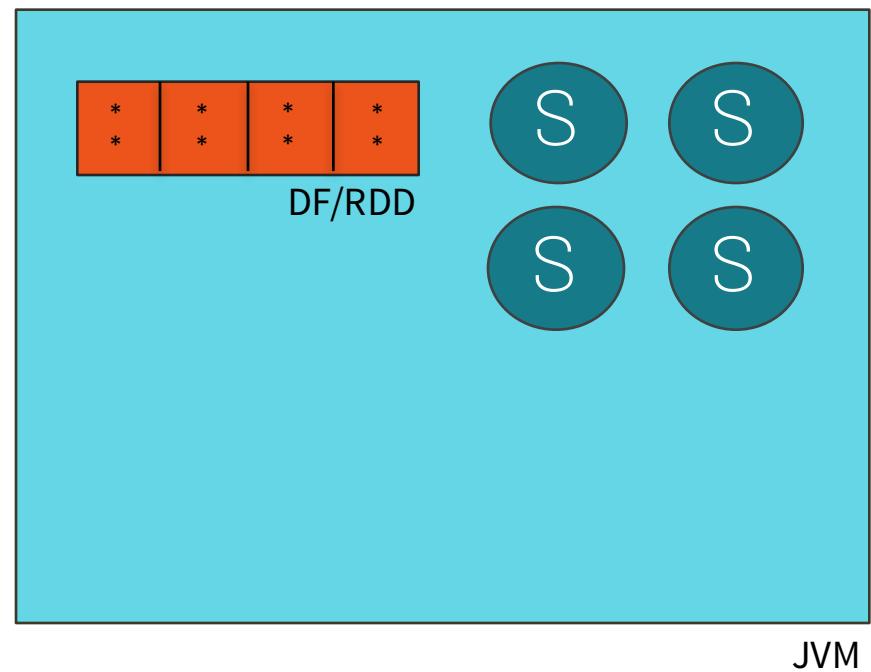


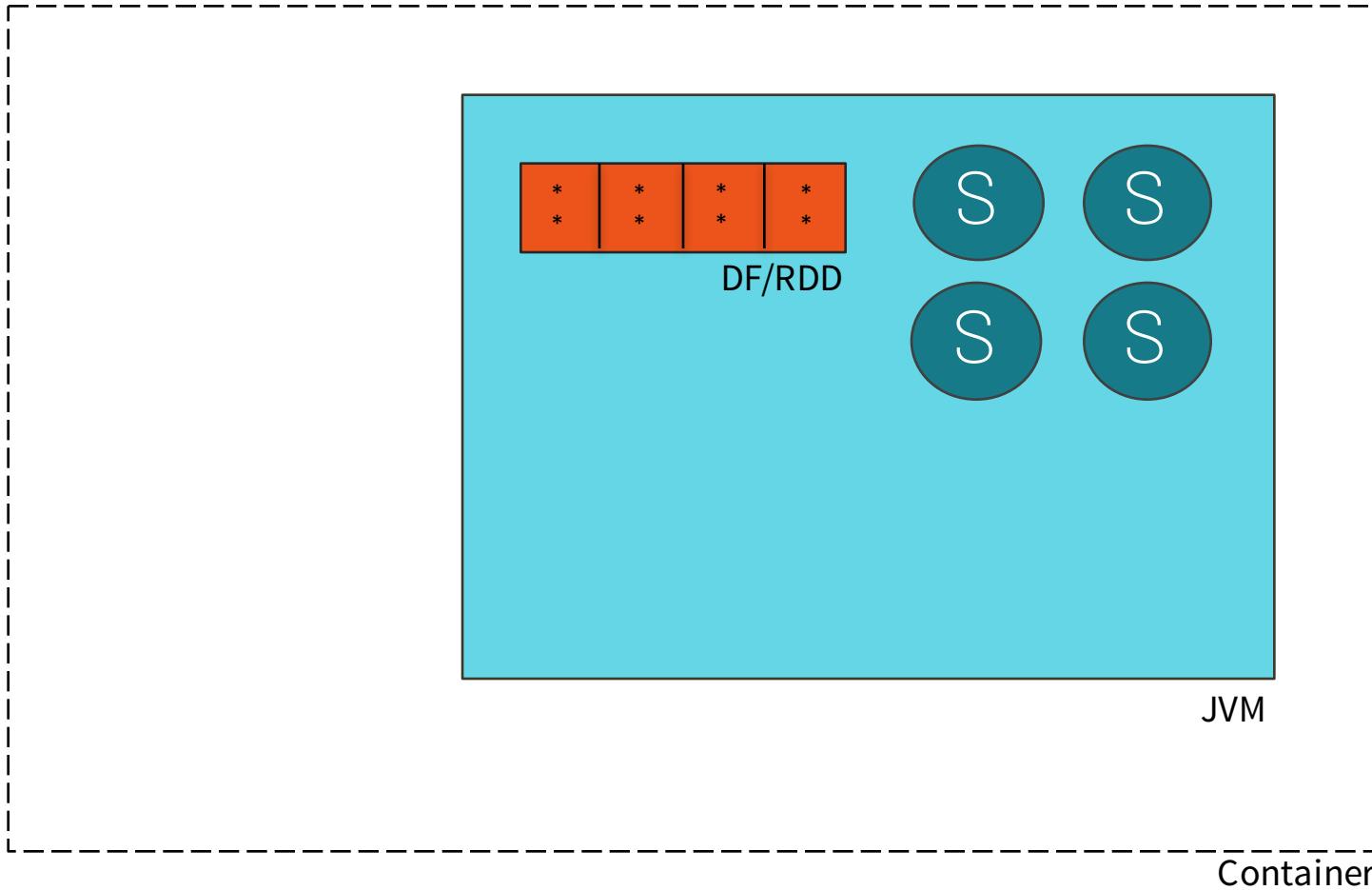


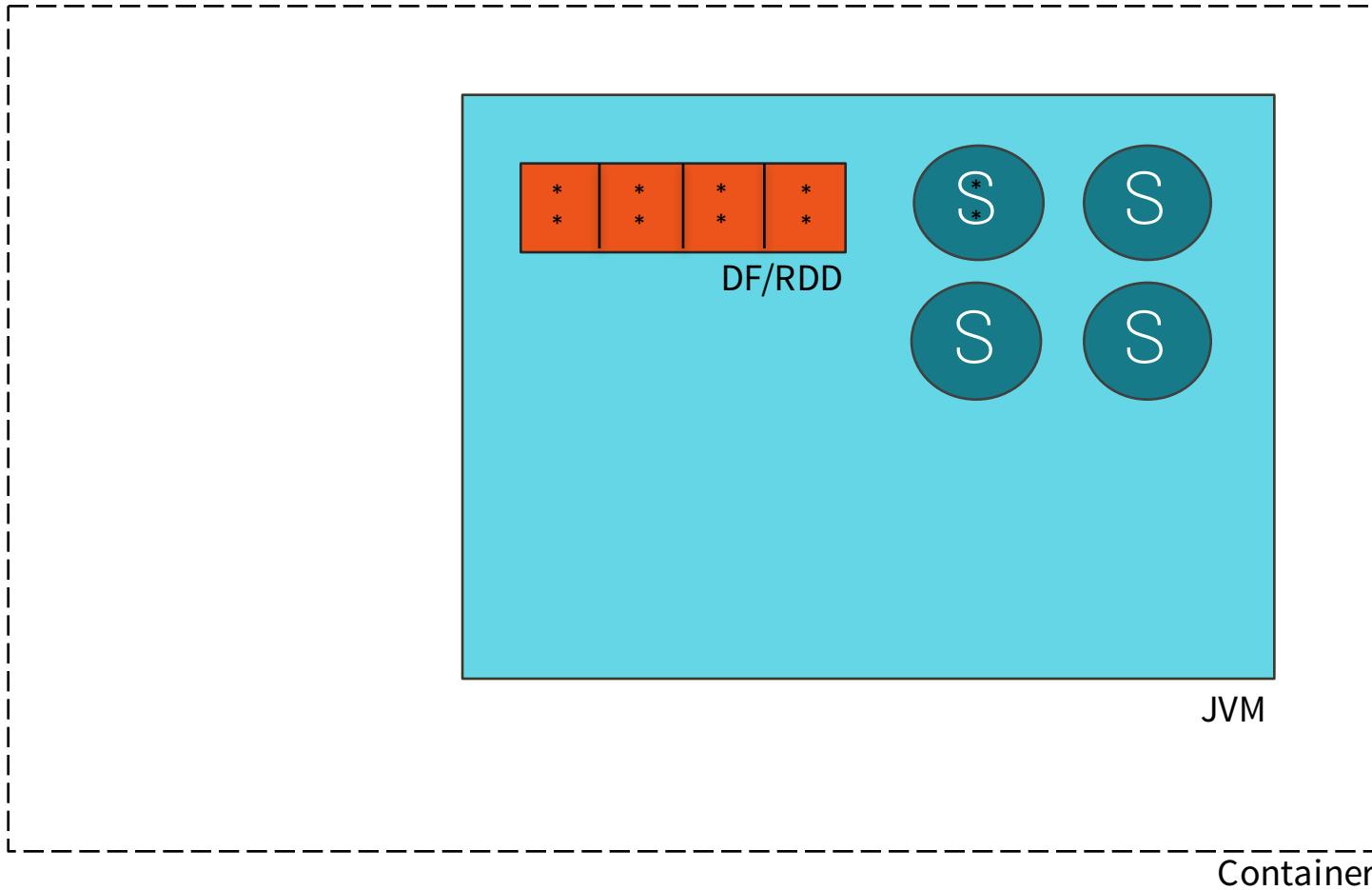
EC2 Machine

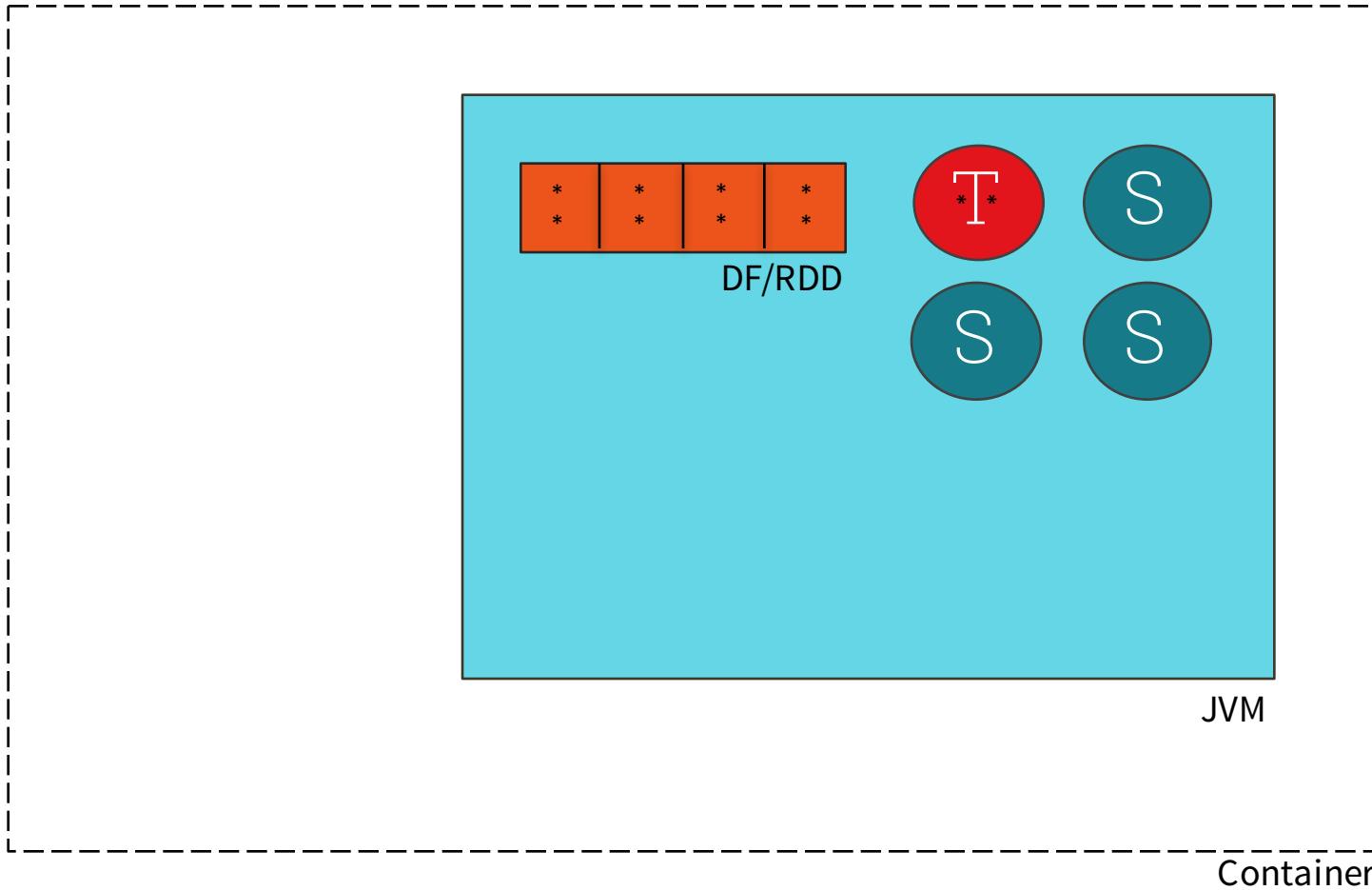


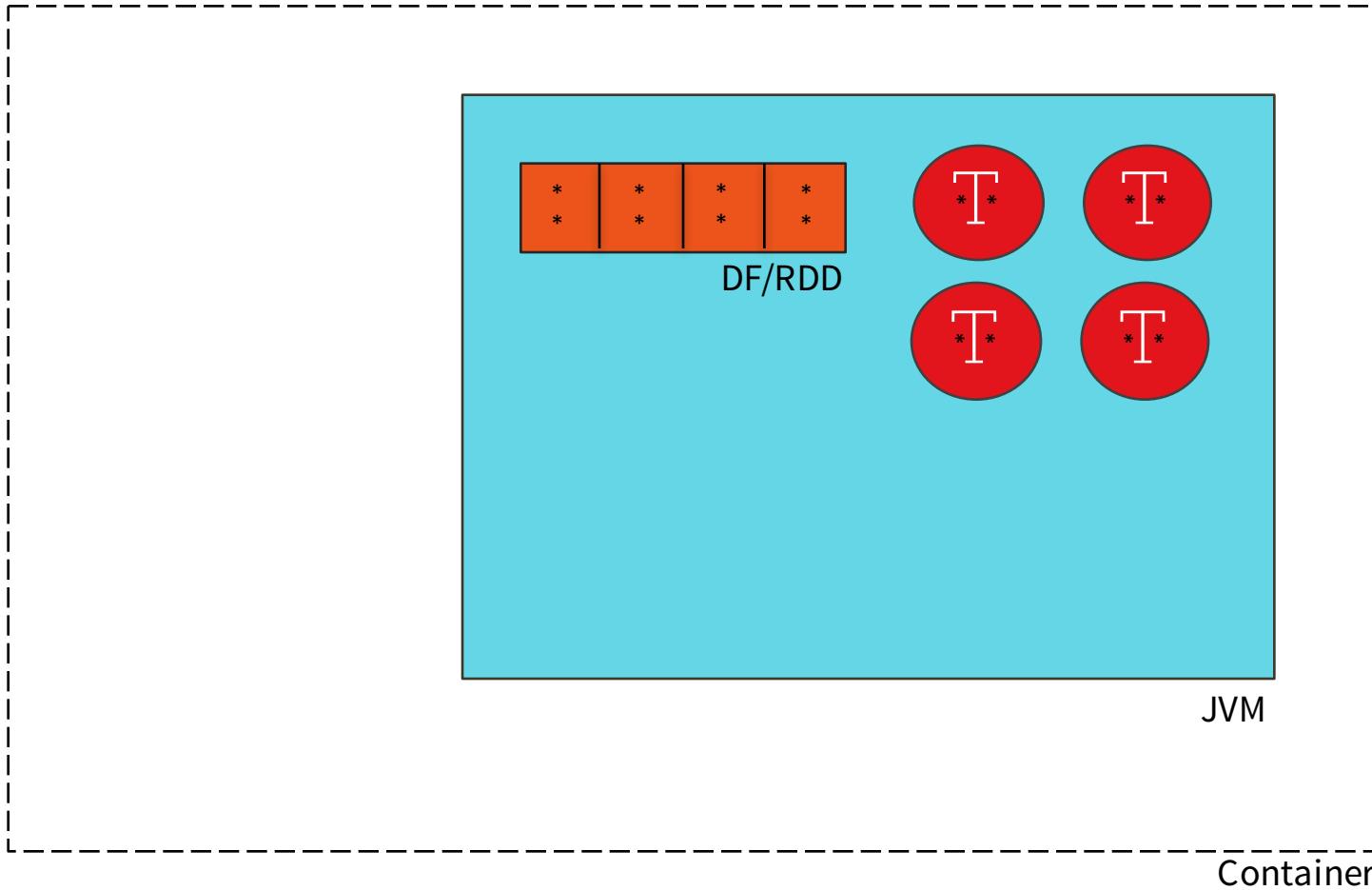
EC2 Machine

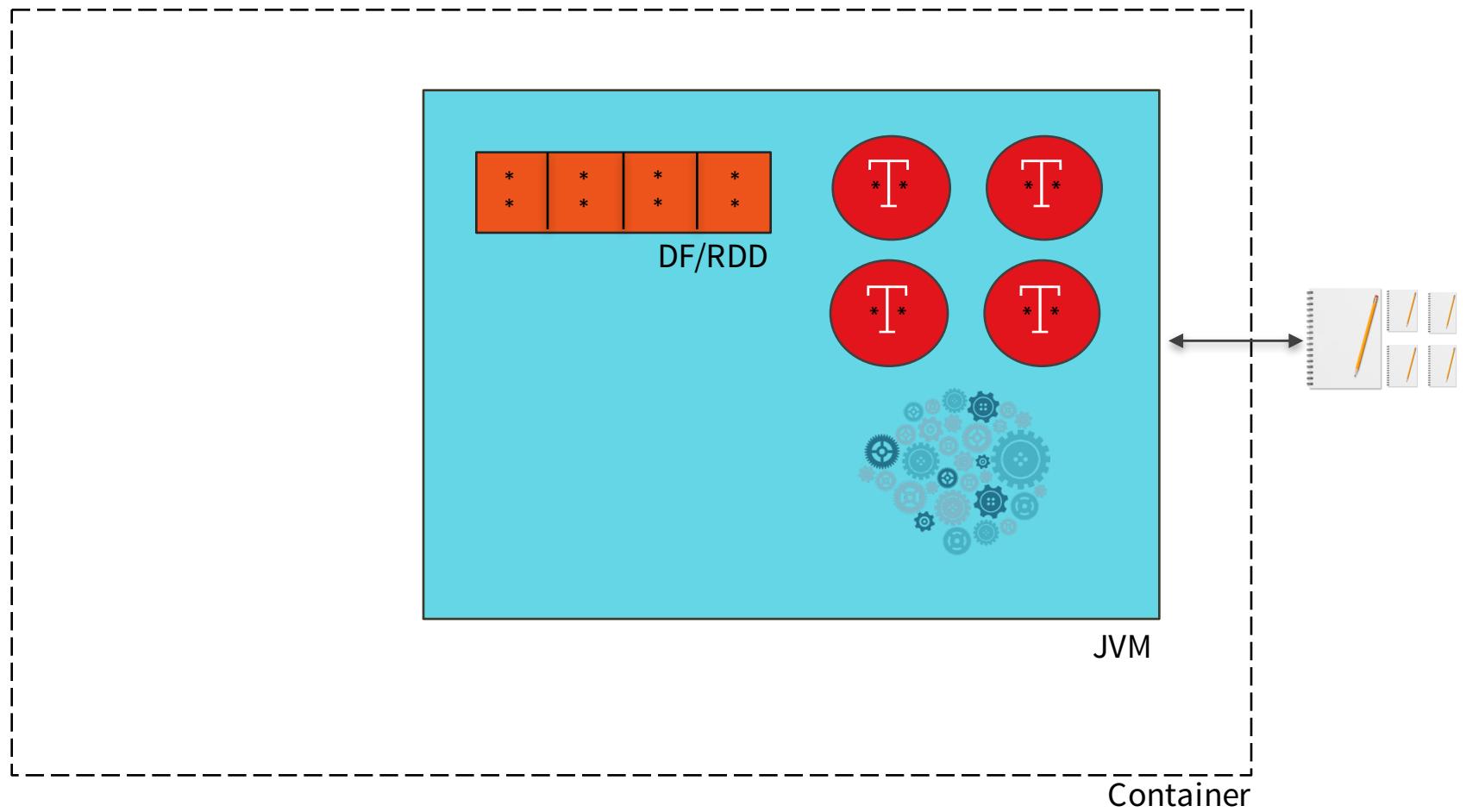


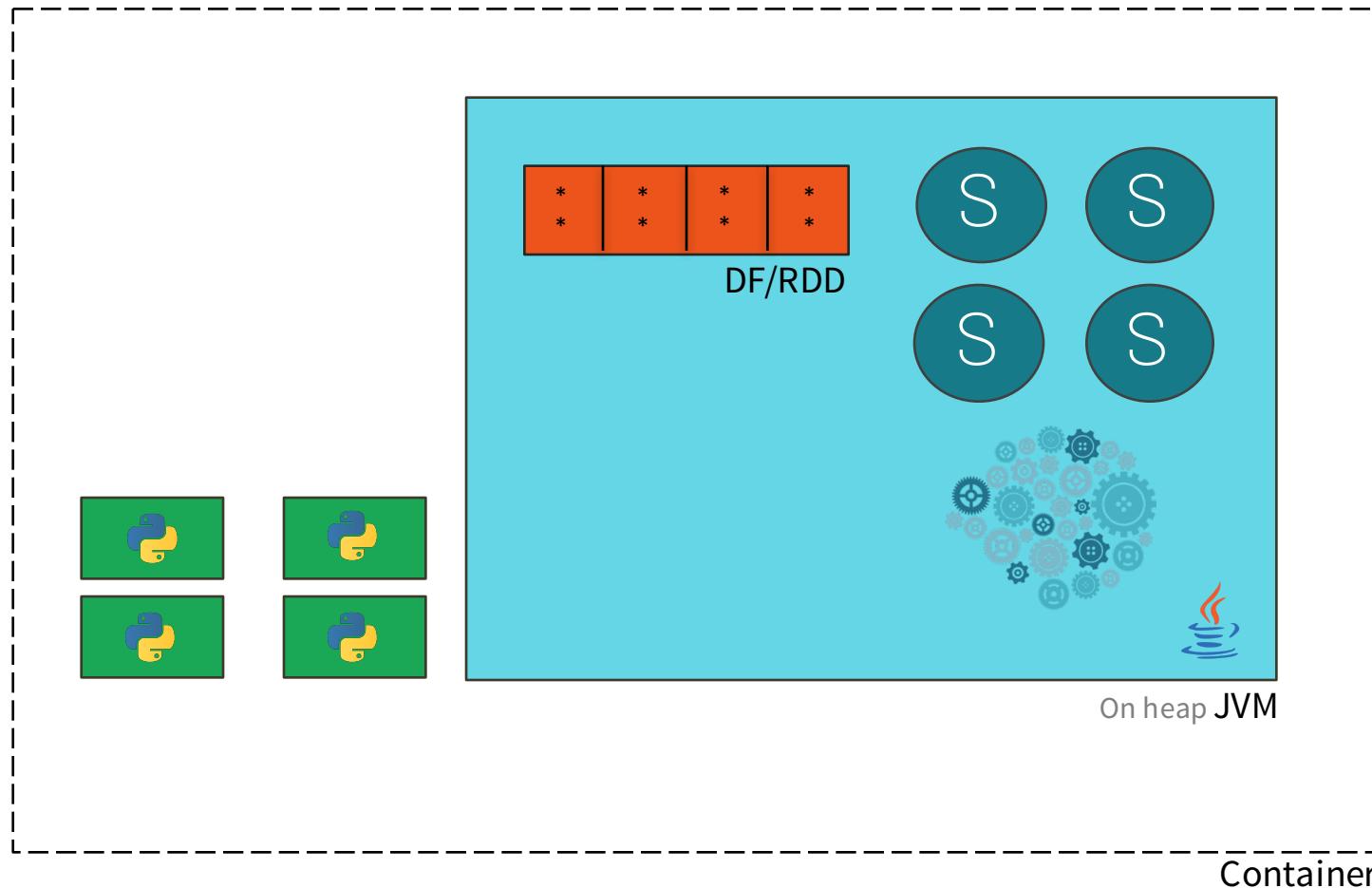


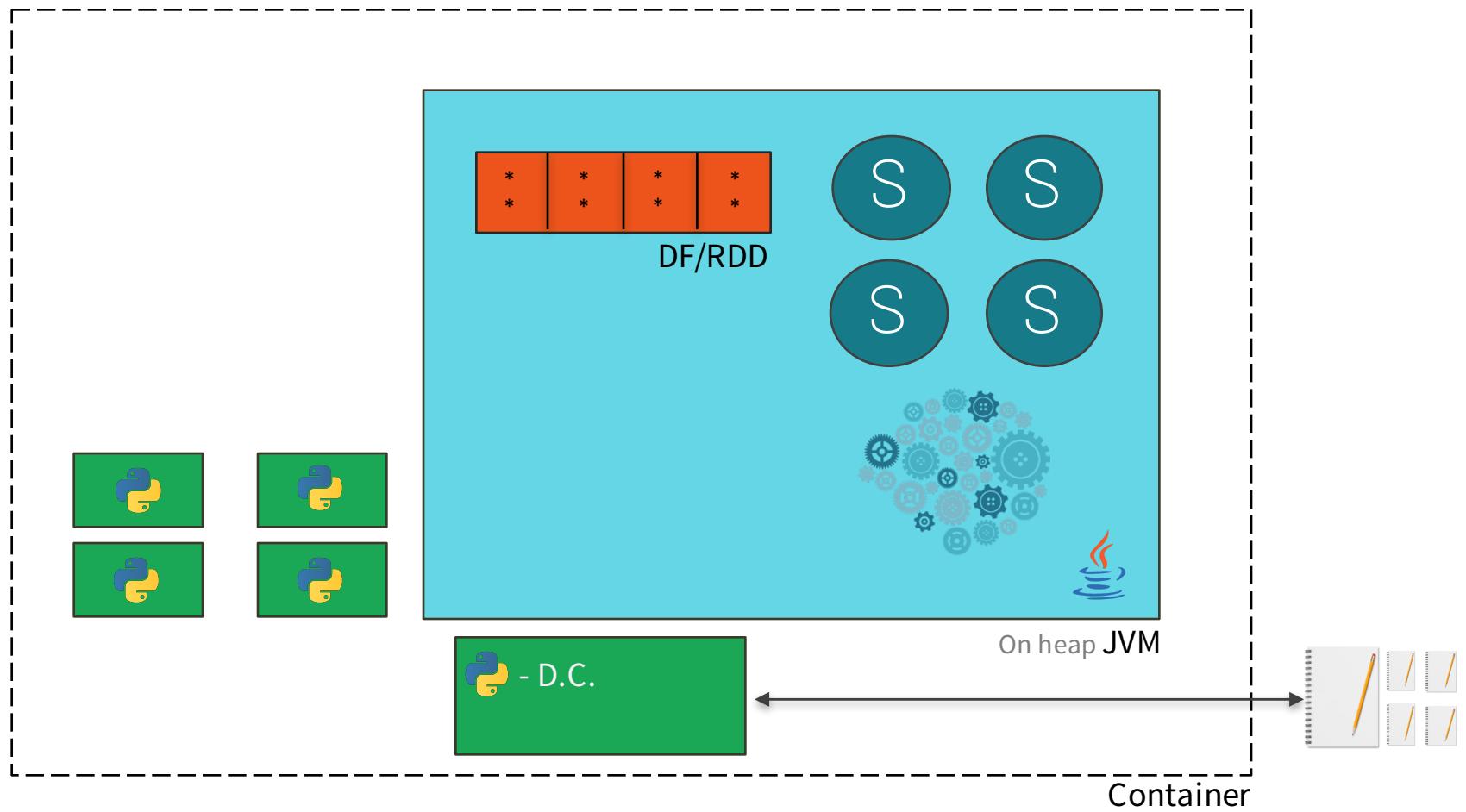


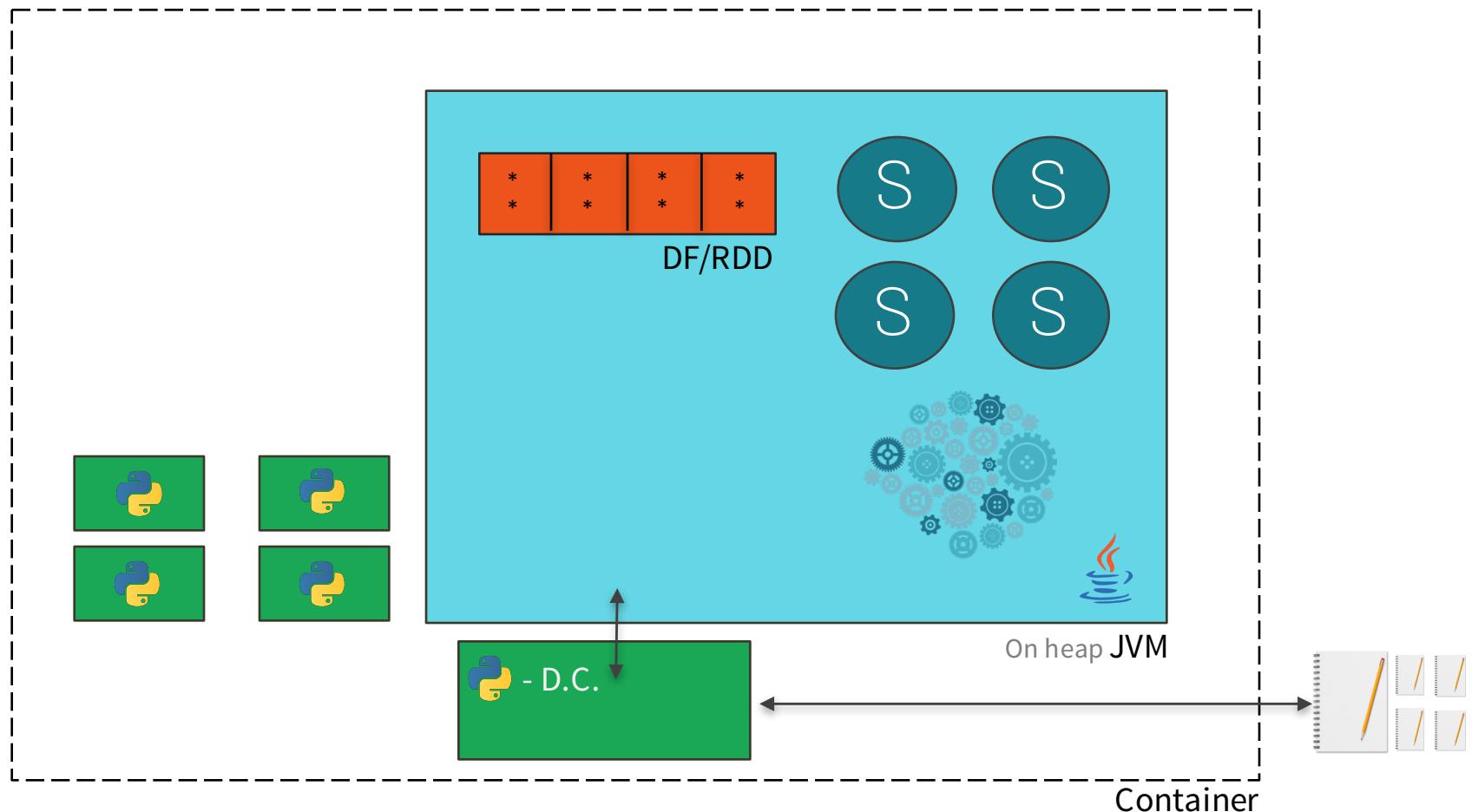




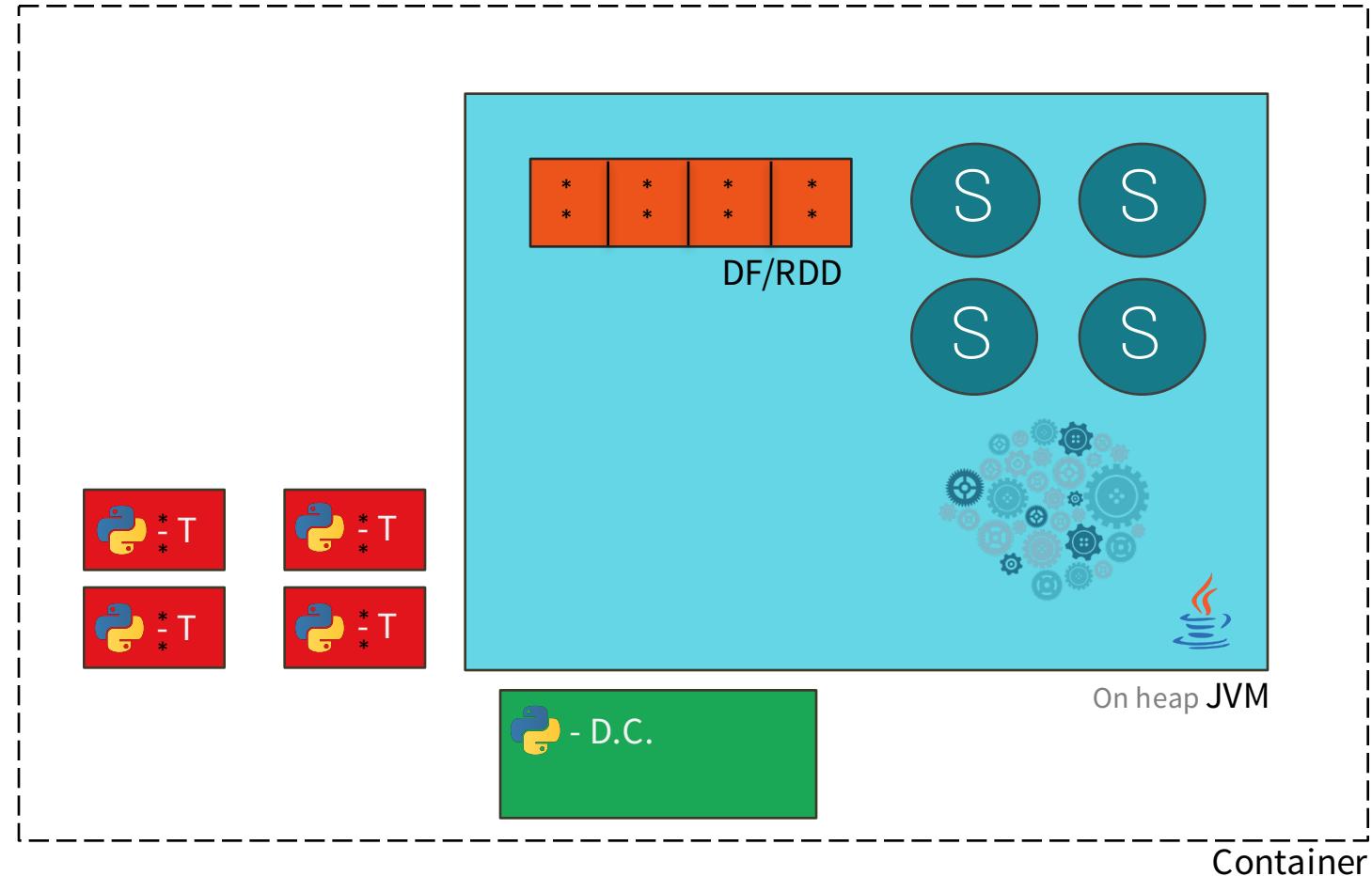




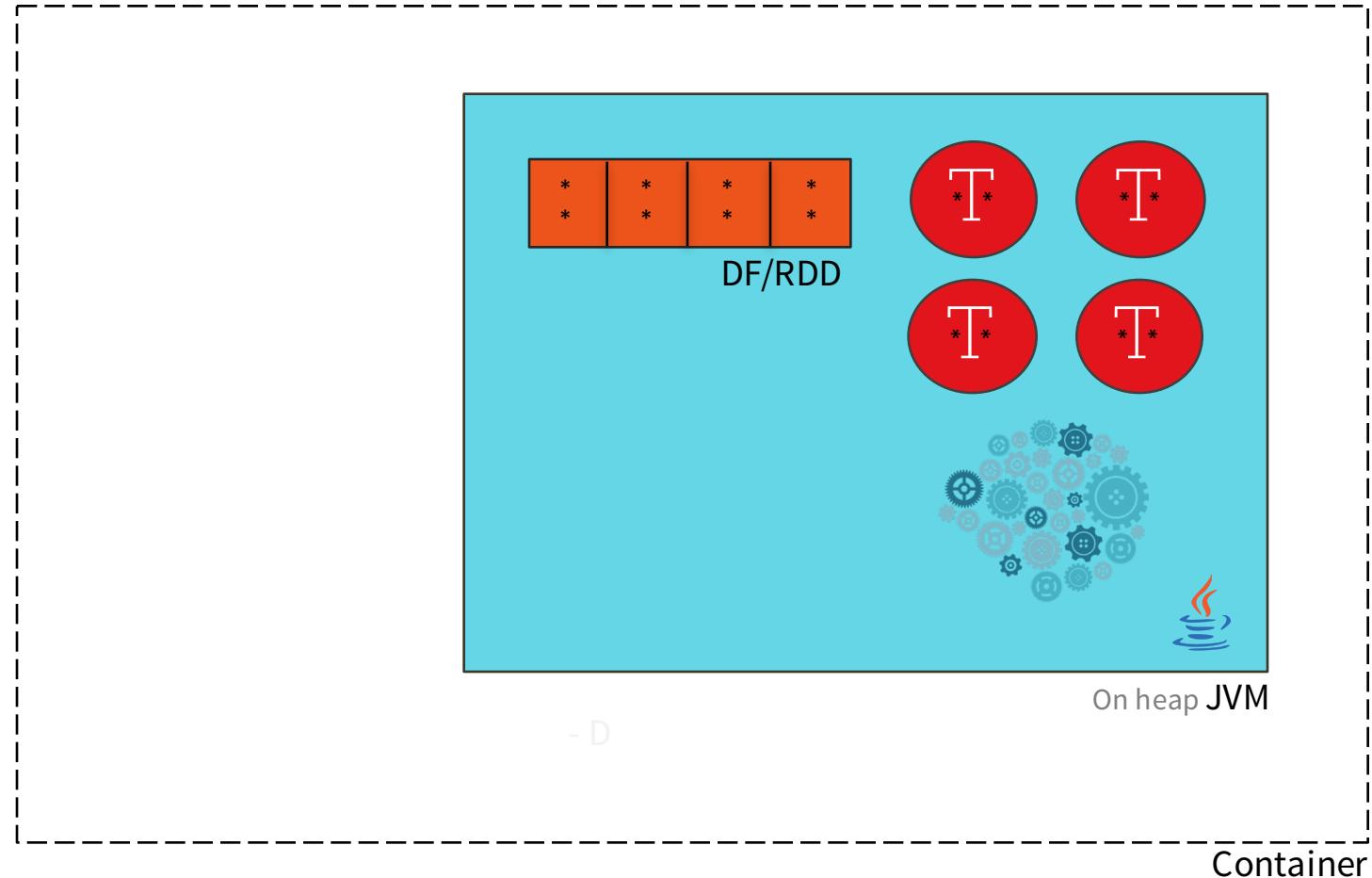




Processing
Python UDFs
(slower)



Processing
Scala / Java



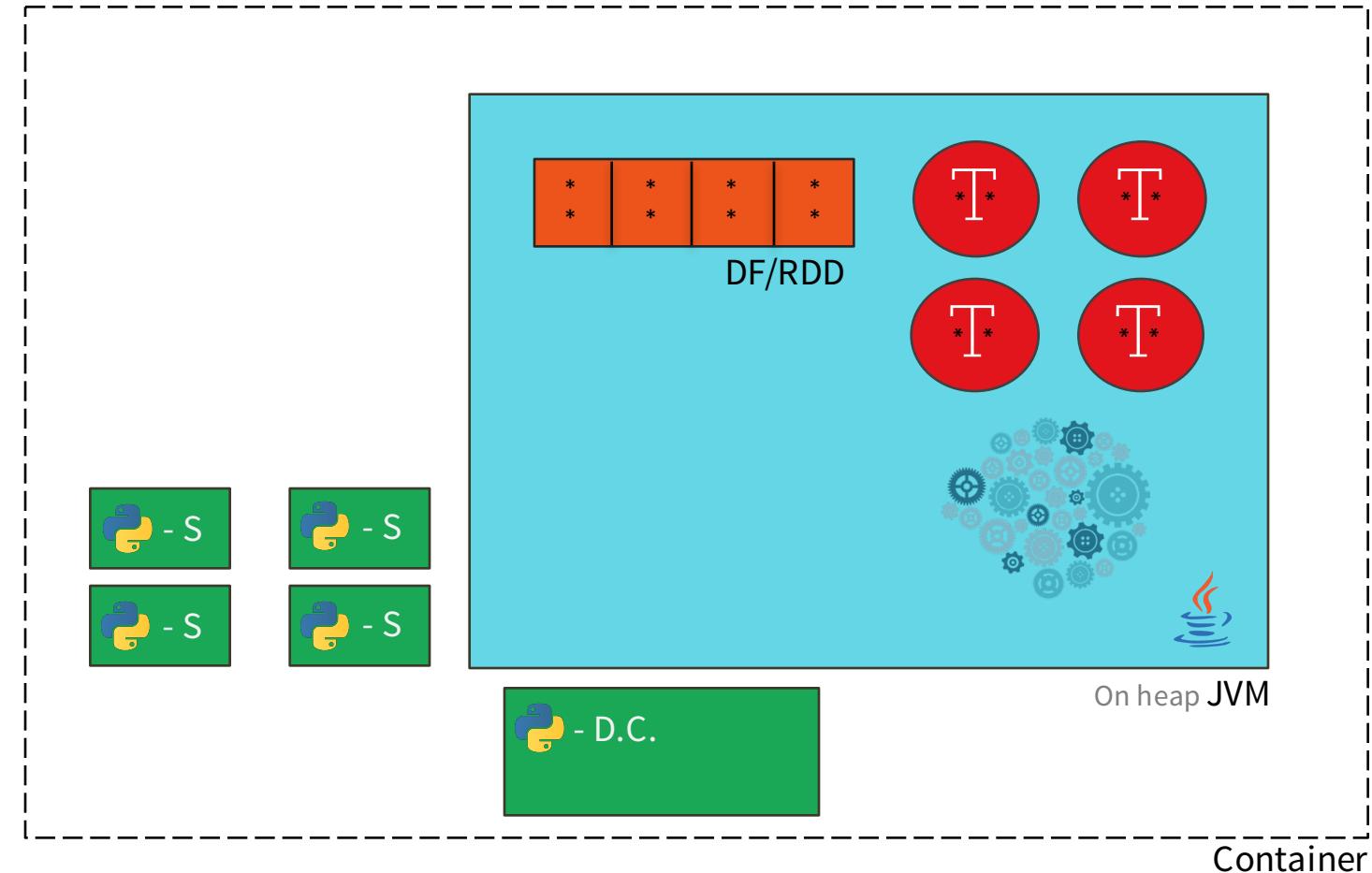
Processing DataFrames

 Scala

 Java

 python™

 R



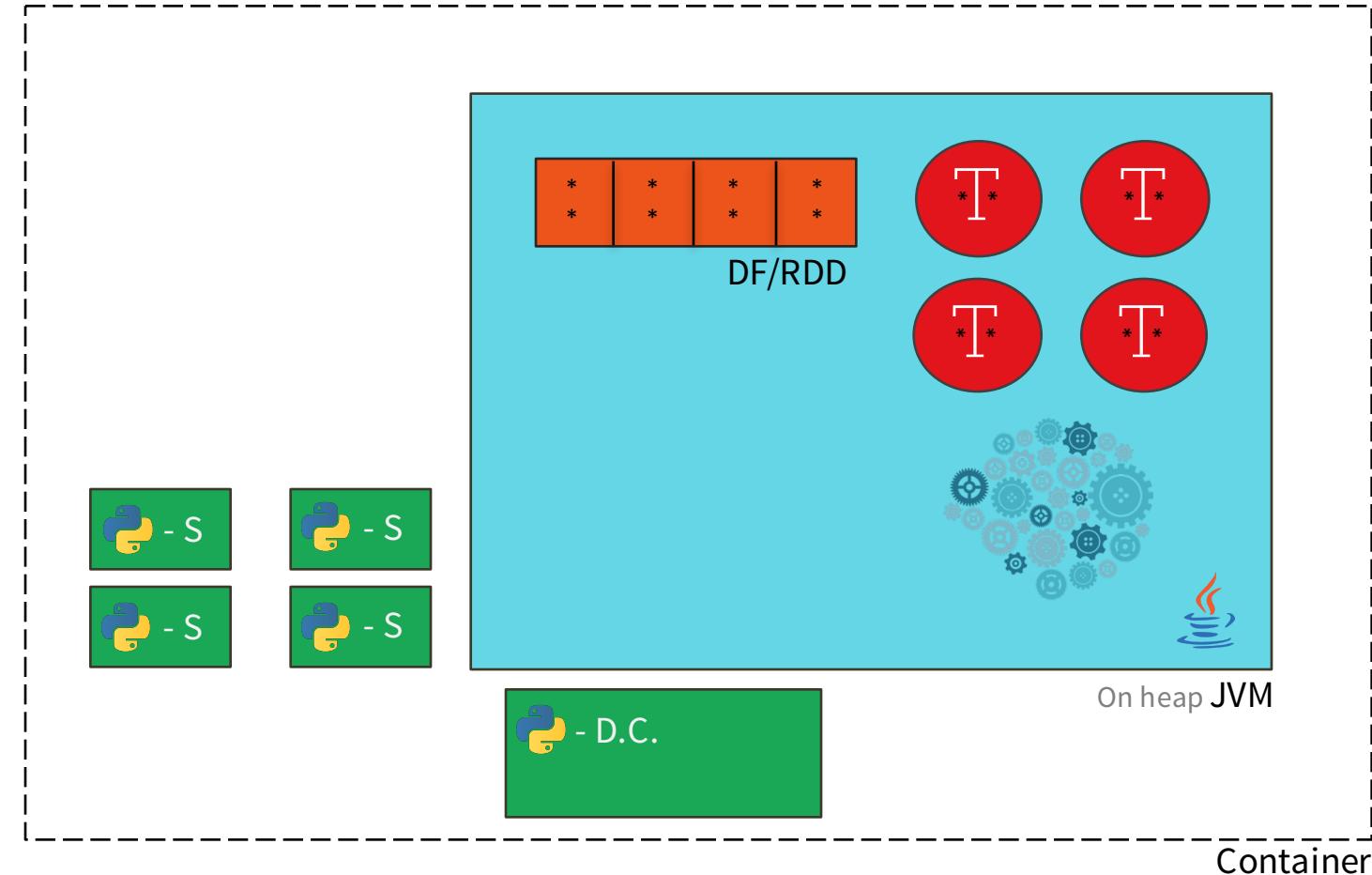
Processing MLlib

 Scala

 Java

 python™

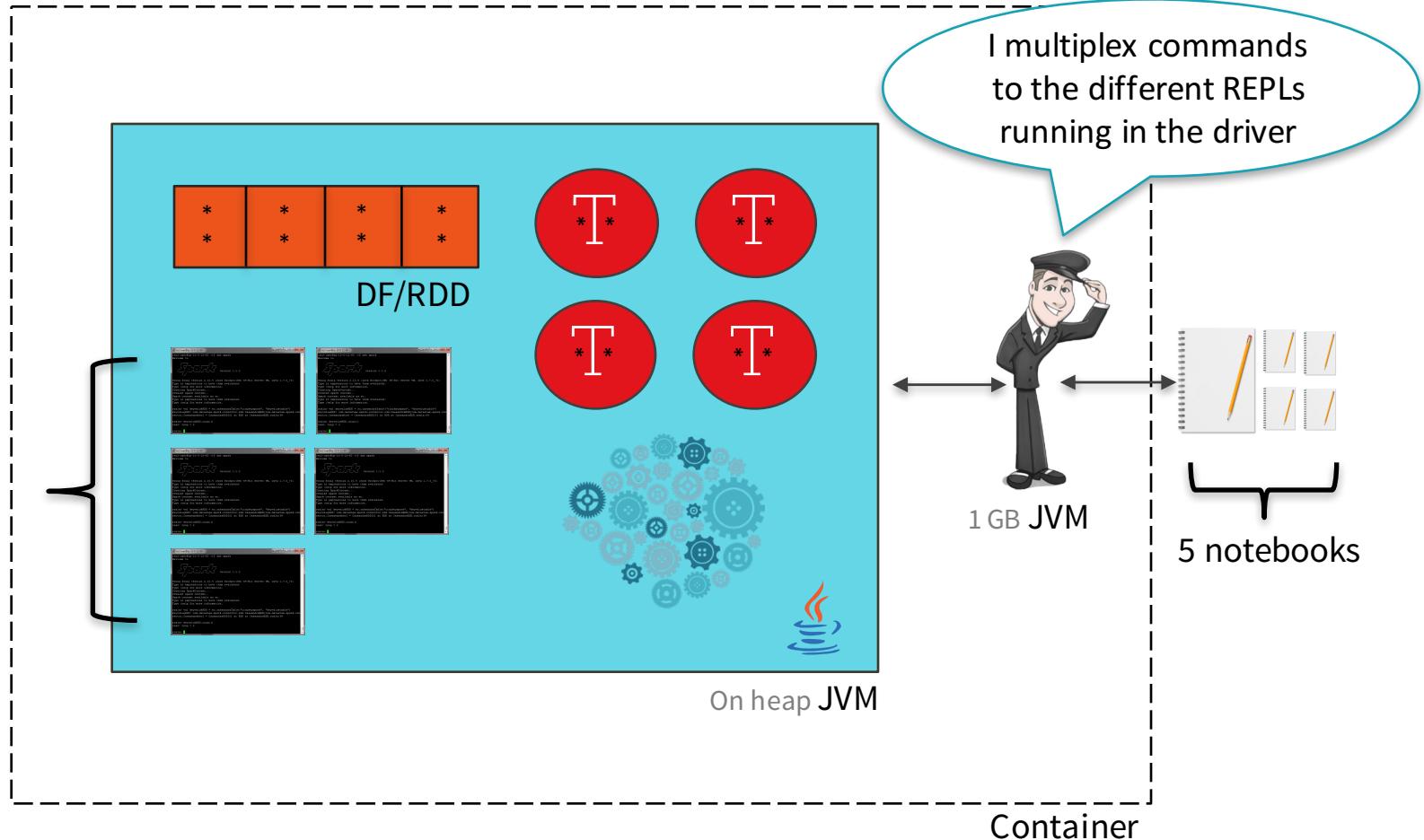
 R



Chauffeur
“one last trick”



5 REPLs



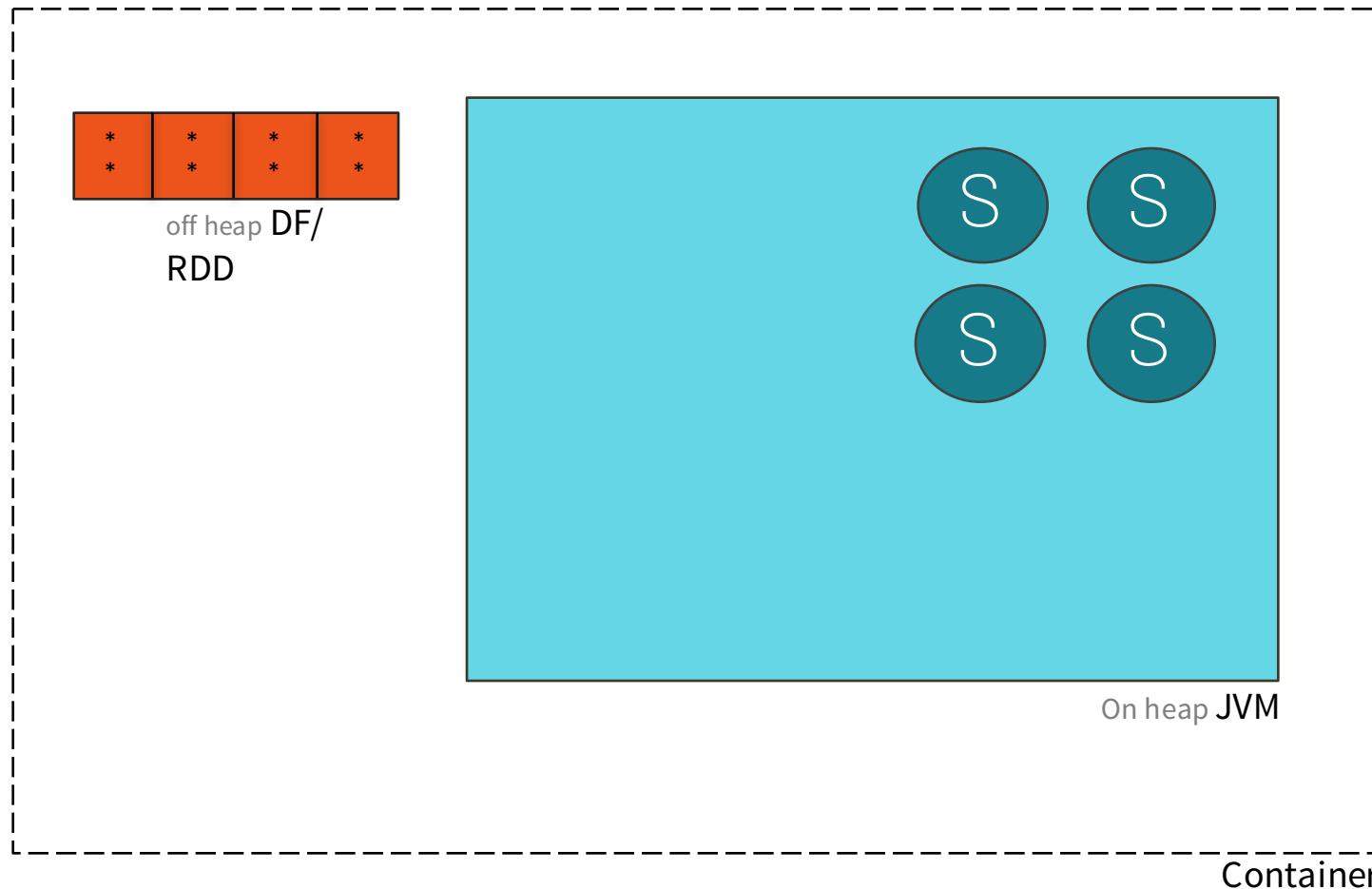
Demo!

- Local mode JVMs
- Executors tab
- %sh free -mh
- %sh jps
- %sh jps -v
- %sh ps -fe
- Events UI for tasks



Spark >= 1.7

Project
Tungsten



Local Mode UI:

Hostname: ec2-52-32-199-8.us-west-2.compute.amazonaws.com

Jobs Stages Storage Environment **Executors** SQL JDBC/ODBC Server

Executors (1)

Memory: 0.0 B Used (9.6 GB Total)
Disk: 0.0 B Used

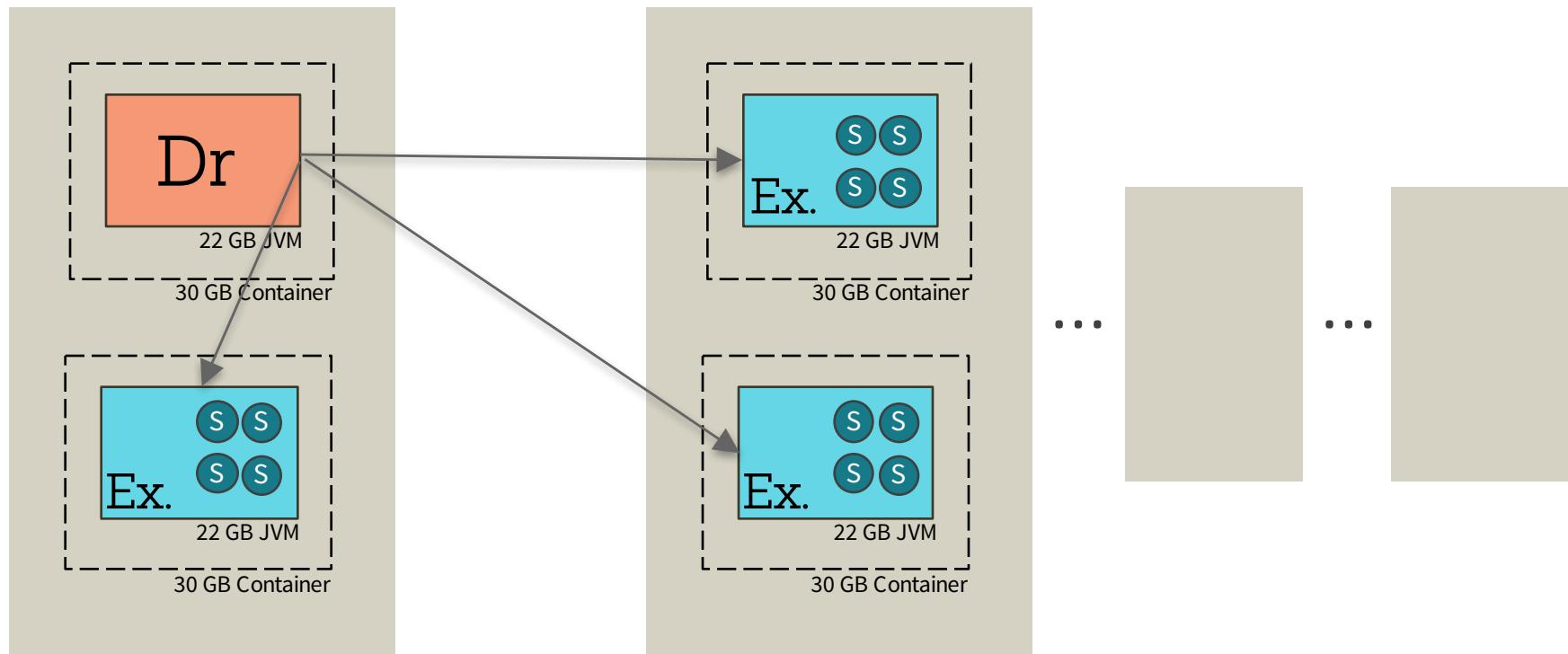
Executor ID	Address	RDD Blocks	Storage Memory	Disk Used	Active Tasks	Failed Tasks	Complete Tasks	Total Tasks	Task Time	Input	Shuffle Read	Shuffle Write	Thread Dump
driver	10.55.241.199:56517	0	0.0 B / 9.6 GB	0.0 B	0	0	0	0	0 ms	0.0 B	0.0 B	0.0 B	Thread Dump



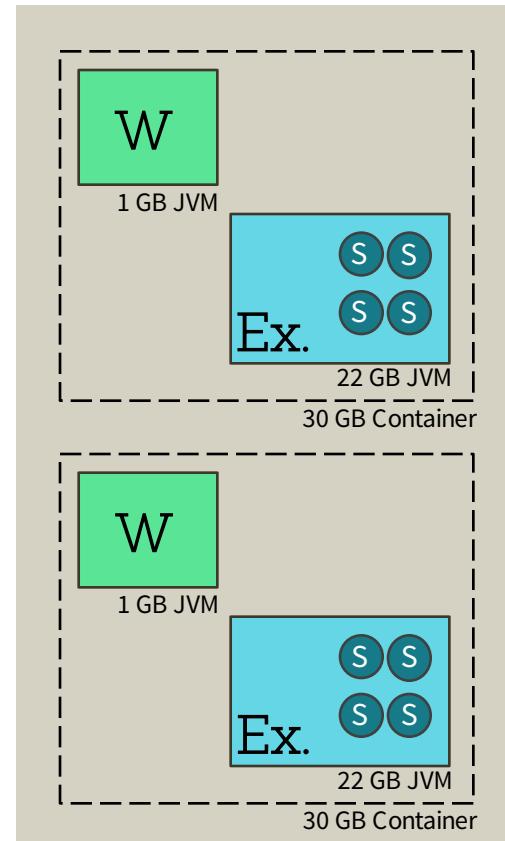
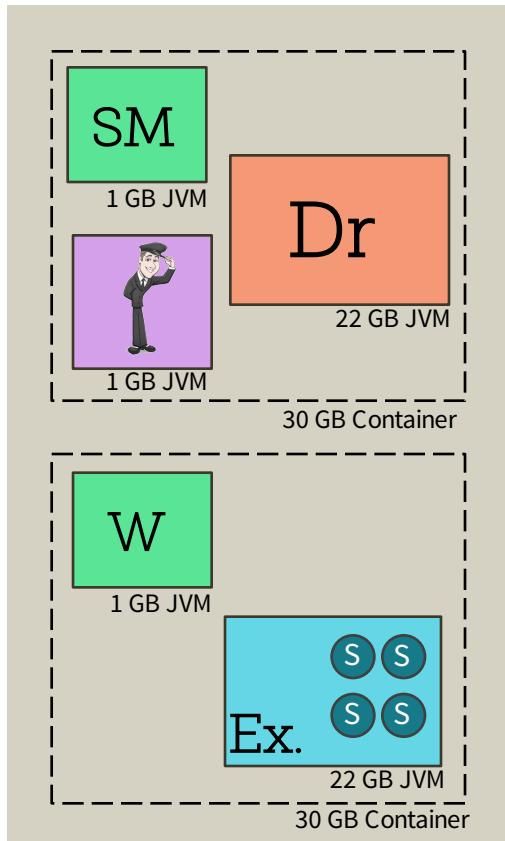
Architecture [revisited]

(standalone mode)

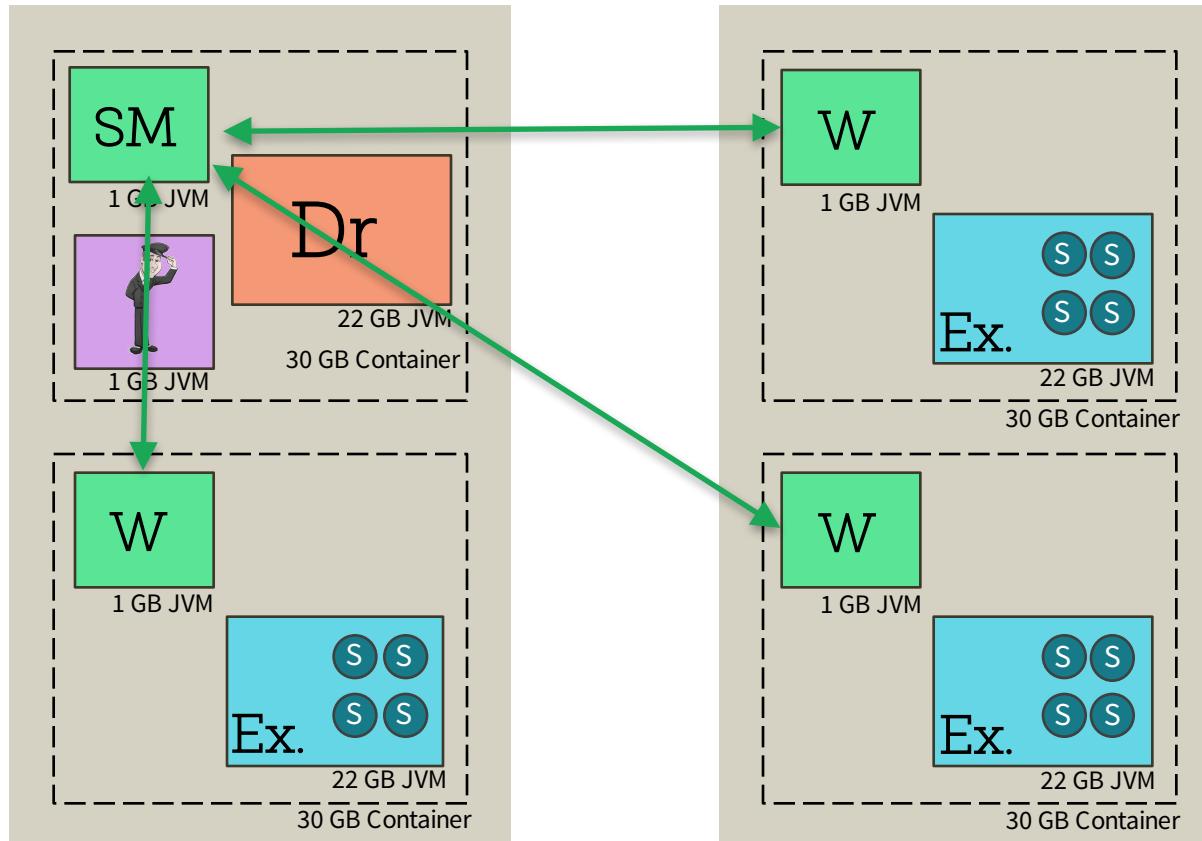
Standalone Mode:



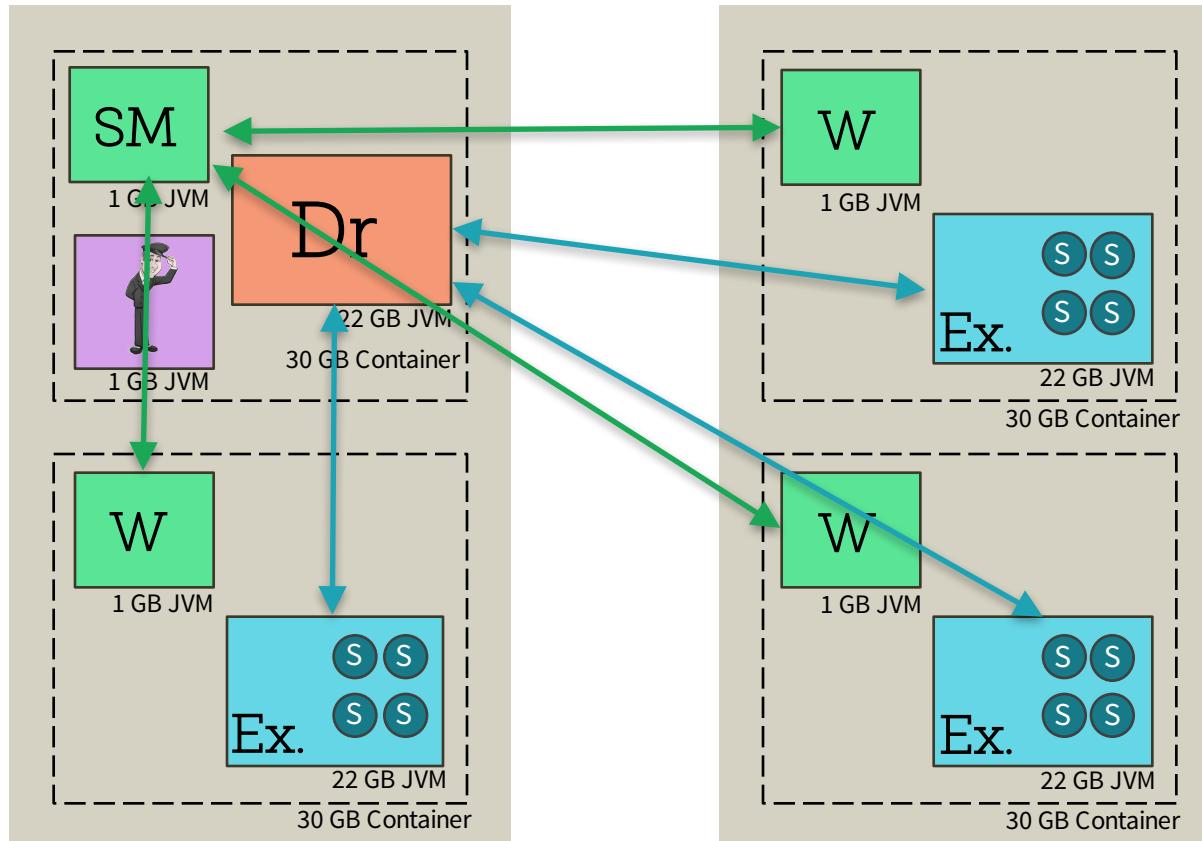
Standalone Mode:



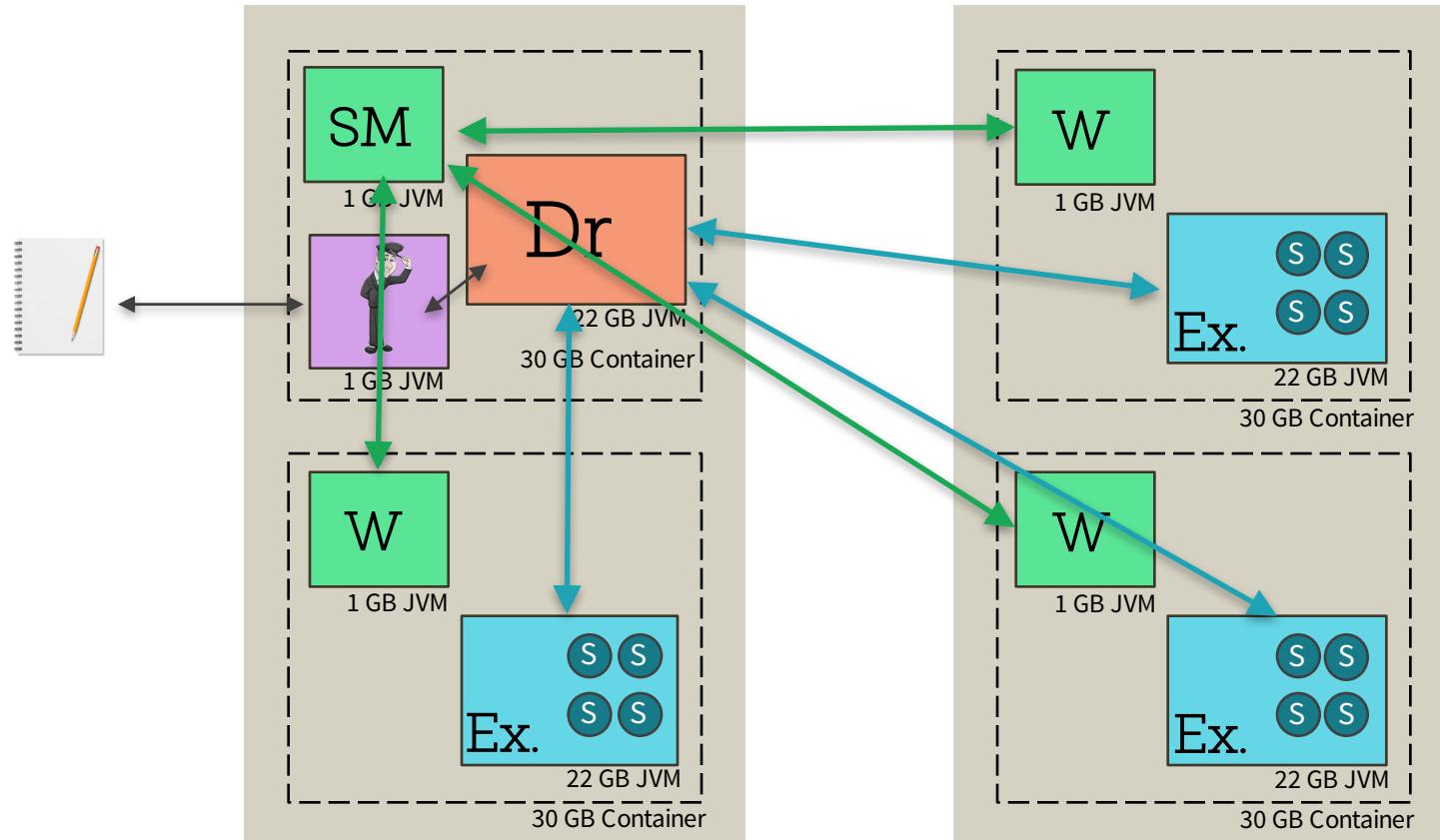
Standalone Mode:

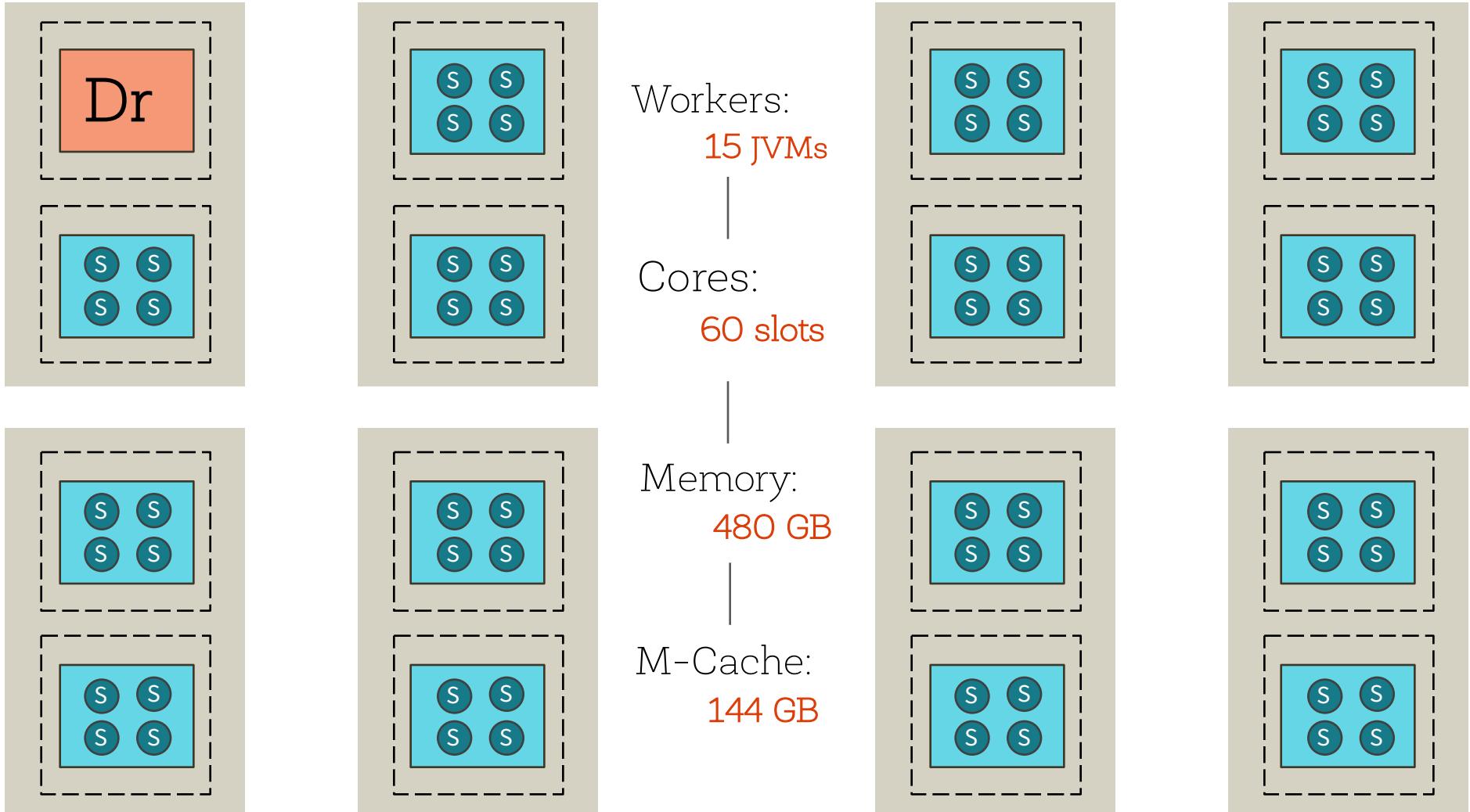


Standalone Mode:



Standalone Mode:





Demo!

- Standalone mode JVMs
- Executors tab
- Spark Master UI
- Spark Worker UI
- Print to stout in Driver & Executor
- conf() settings

Standalone Mode UI:

Spark UI

Hostname: ec2-52-32-199-8.us-west-2.compute.amazonaws.com

Jobs Stages Storage Environment **Executors** SQL JDBC/ODBC Server

Executors (16) 15

Memory: 0.0 B Used (153.5 GB Total)
Disk: 0.0 B Used

Executor ID	Address	RDD Blocks	Storage Memory	Disk Used	Active Tasks	Failed Tasks	Complete Tasks	Total Tasks	Task Time	Input	Shuffle Read	Shuffle Write	Logs	Thread Dump
0	10.55.239.54:45718	0	0.0 B / 9.6 GB	0.0 B	0	0	0	0	0 ms	0.0 B	0.0 B	0.0 B	stdout stderr	Thread Dump
1	10.55.232.95:55638	0	0.0 B / 9.6 GB	0.0 B	0	0	0	0	0 ms	0.0 B	0.0 B	0.0 B	stdout stderr	Thread Dump
10	10.55.252.176:36039	0	0.0 B / 9.6 GB	0.0 B	0	0	0	0	0 ms	0.0 B	0.0 B	0.0 B	stdout stderr	Thread Dump
11	10.55.232.94:44122	0	0.0 B / 9.6 GB	0.0 B	0	0	0	0	0 ms	0.0 B	0.0 B	0.0 B	stdout stderr	Thread Dump
• • •														
9	10.55.226.83:38546	0	0.0 B / 9.6 GB	0.0 B	0	0	0	0	0 ms	0.0 B	0.0 B	0.0 B	stdout stderr	Thread Dump
driver	10.55.241.196:38254	0	0.0 B / 9.6 GB	0.0 B	0	0	0	0	0 ms	0.0 B	0.0 B	0.0 B	Thread Dump	



Spark : RDDs



= easy



= medium

Essential Core & Intermediate Spark Operations

General

- map
- filter
- flatMap
- mapPartitions
- mapPartitionsWithIndex
- groupBy
- sortBy

Math / Statistical

- sample
- randomSplit

Set Theory / Relational

- union
- intersection
- subtract
- distinct
- cartesian
- zip

Data Structure / I/O

- keyBy
- zipWithIndex
- zipWithUniqueID
- zipPartitions
- coalesce
- repartition
- repartitionAndSortWithinPartitions
- pipe

- reduce
- collect
- aggregate
- fold
- first
- take
- foreach
- top
- treeAggregate
- treeReduce
- foreachPartition
- collectAsMap

- count
- takeSample
- max
- min
- sum
- histogram
- mean
- variance
- stdev
- sampleVariance
- countApprox
- countApproxDistinct

- takeOrdered

- saveAsTextFile
- saveAsSequenceFile
- saveAsObjectFile
- saveAsHadoopDataset
- saveAsHadoopFile
- saveAsNewAPIHadoopDataset
- saveAsNewAPIHadoopFile



= easy



= medium

Essential Core & Intermediate Spark Operations

- General**
- flatMapValues
 - groupByKey
 - reduceByKey
 - reduceByKeyLocally
 - foldByKey
 - aggregateByKey
 - sortByKey
 - combineByKey
 - keys
 - values

Math / Statistical

- sampleByKey

Set Theory / Relational

- cogroup (=groupWith)
- join
- subtractByKey
- fullOuterJoin
- leftOuterJoin
- rightOuterJoin

Data Structure / I/O

- partitionBy

-
- countByKey
 - countByValue
 - countByValueApprox
 - countApproxDistinctByKey
 - countApproxDistinctByKey
 - countByKeyApprox
 - sampleByKeyExact

Types of RDDs:

HadoopRDD: core functionality for reading data in HDFS using older MR API

MapPartitionsRDD: a result of calling actions like `map`, `flatMap`, `filter`, `mapPartitions`, etc

PairRDD: holds key/value pairs, a result of `groupByKey` and `join` operations

PipedRDD: result of piping elements to a forked external process

ShuffledRDD: a result after shuffling (`repartition` or `coalesce`)



TRANSFORMATIONS AND ACTIONS

<http://training.databricks.com/visualapi.pdf>



A Visual Guide of the API

- Most transformations operate on a **per item basis**, applying the user function to each item (`map`, `filter`, etc)
- A few transformations operate on a **per partition basis**, applying the user function (*of type Iterator*) separately on each partition (`mapPartitions`, `mapPartitionsWithIndex`, etc)

2 kinds of Actions



VS

distributed

occurs across the cluster

`saveAsTextFile`, (`HDFS`, `S3`, `SQL`, `NoSQL`, etc.)

driver

result must fit in driver JVM

collect, count, reduce, take, show..

Pagecounts:



The screenshot shows a web browser window with the address bar containing the URL <https://dumps.wikimedia.org/other/pagecounts-raw/2015/2015-11/>. The page title is "Index of page view statistics for 2015-11".

Index of page view statistics for 2015-11

Pagecount files for 2015-11

Check the [hashes](#) after your download, to make sure your files arrived intact.

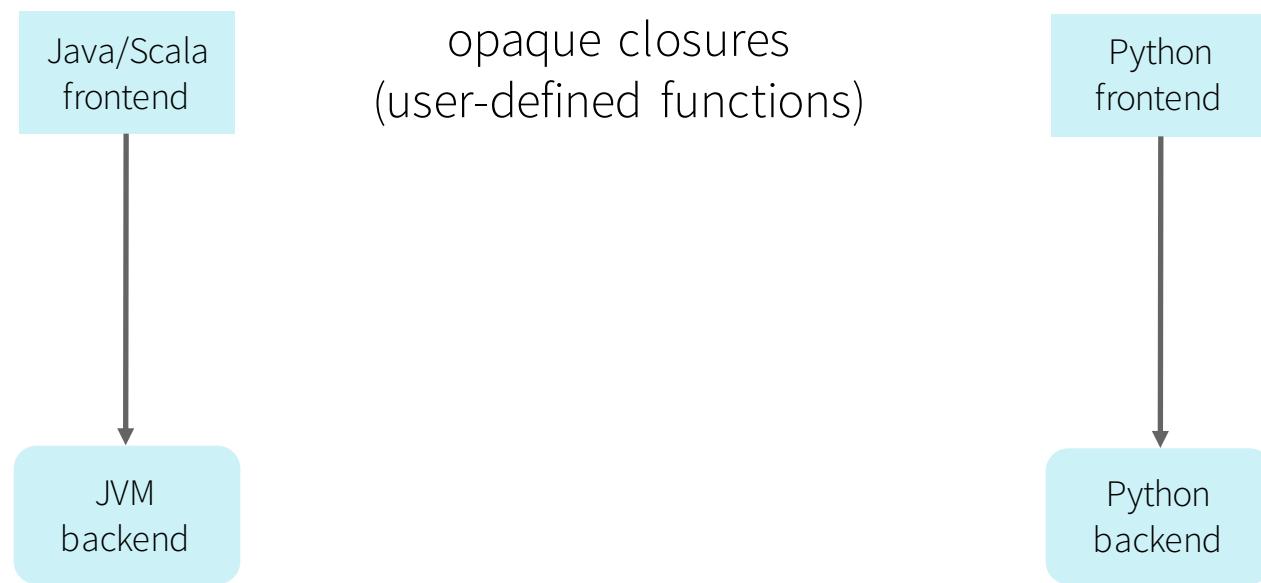
- [pagecounts-20151101-000000.gz](#), size 77M
- [pagecounts-20151101-010000.gz](#), size 77M
- [pagecounts-20151101-020000.gz](#), size 79M
- [pagecounts-20151101-030000.gz](#), size 76M
- [pagecounts-20151101-040000.gz](#), size 76M
- [pagecounts-20151101-050000.gz](#), size 80M
- [pagecounts-20151101-060000.gz](#), size 78M
- [pagecounts-20151101-070000.gz](#), size 77M
- [pagecounts-20151101-080000.gz](#), size 79M
- [pagecounts-20151101-090000.gz](#), size 83M
- [pagecounts-20151101-100000.gz](#), size 88M
- [pagecounts-20151101-110000.gz](#), size 88M

Pagecounts:

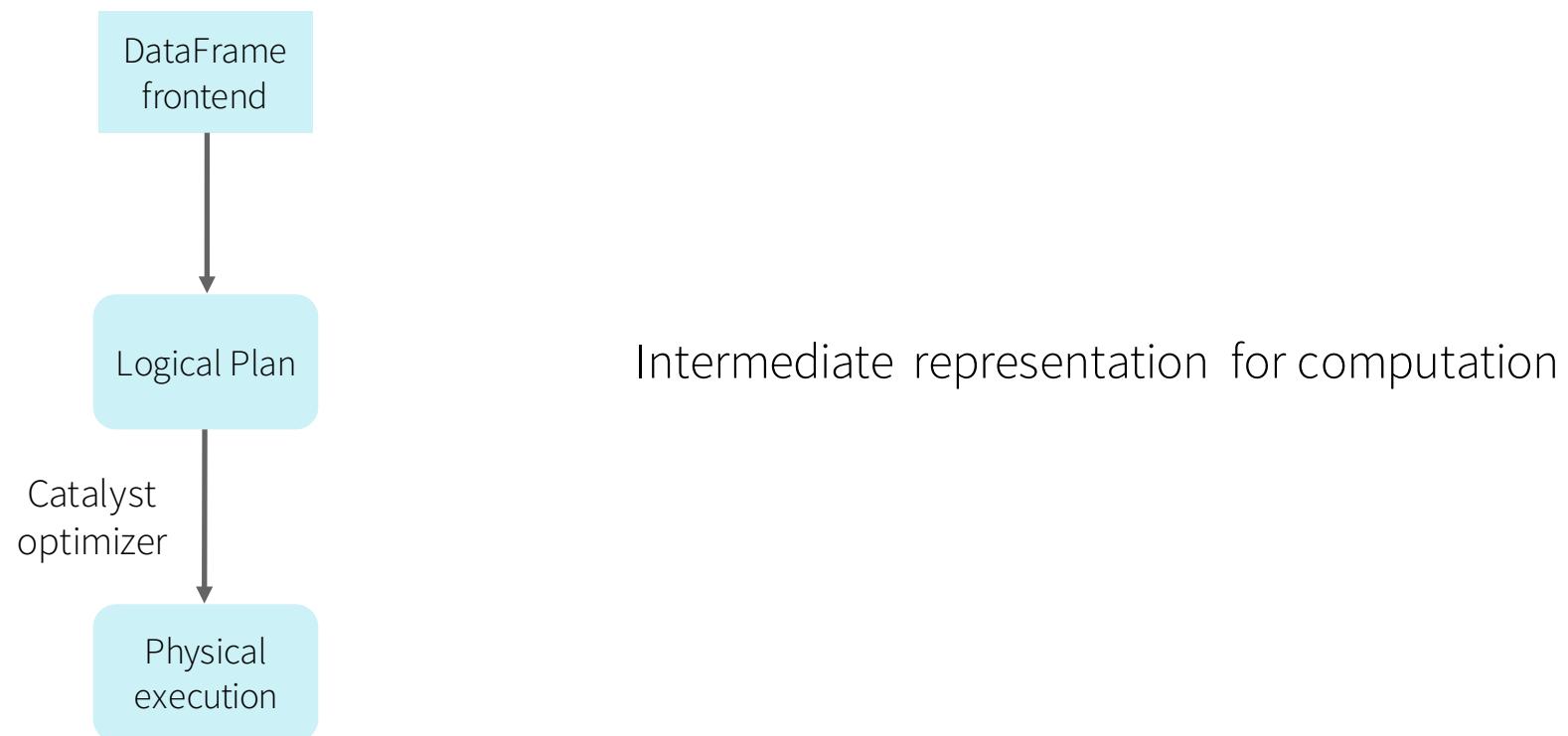
- Hourly Statistics (not unique visits)
- ~500 MB each file, each hour (uncompressed)
- Partitioned Parquet files

Project name	Page title	# of requests	Total size of content returned
en	Main_Page	245839	4737756101
en	Apache_Spark	1370	29484844
en.mw	Apache_Spark		
en.d	discombobulate	200	284834
fr.b	Special:Recherche/Acteurs_et_actrices_N	1	739

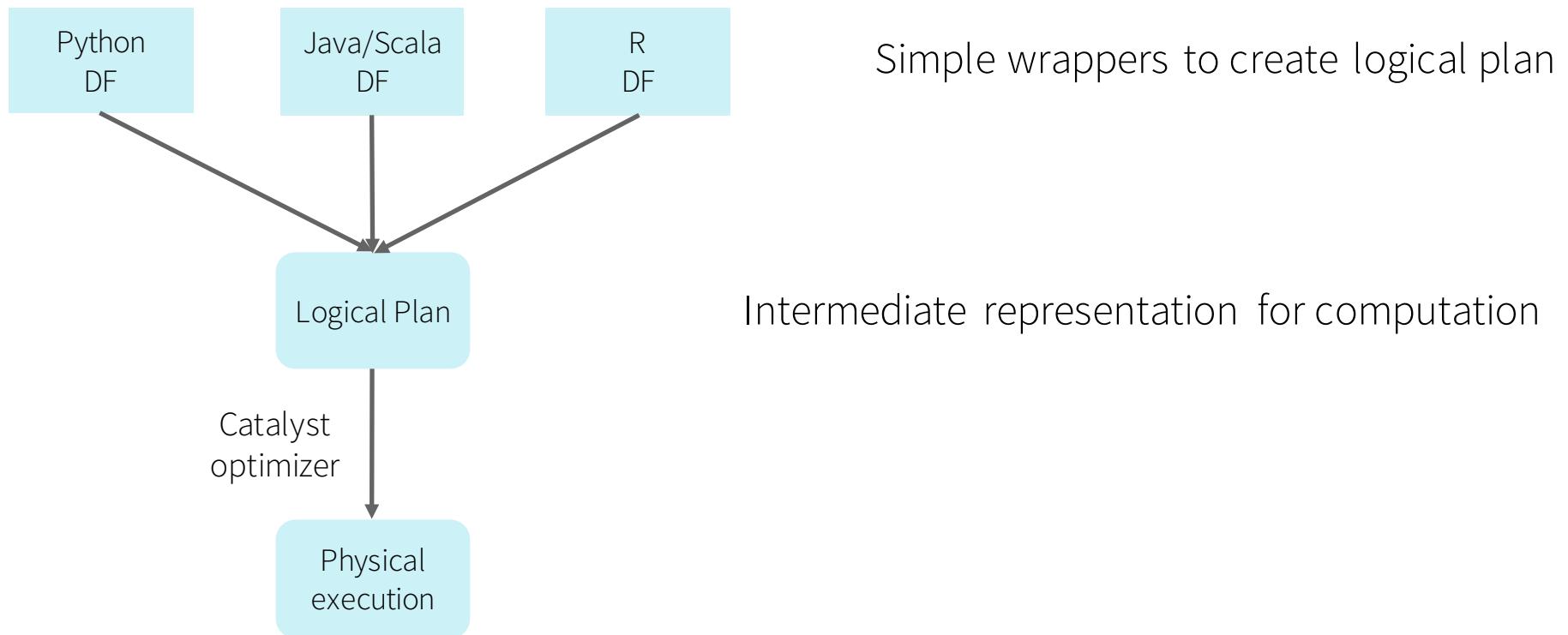
Spark RDD Execution



Spark DataFrame Execution



Spark DataFrame Execution



RDDs

Collections of Native JVM Objects

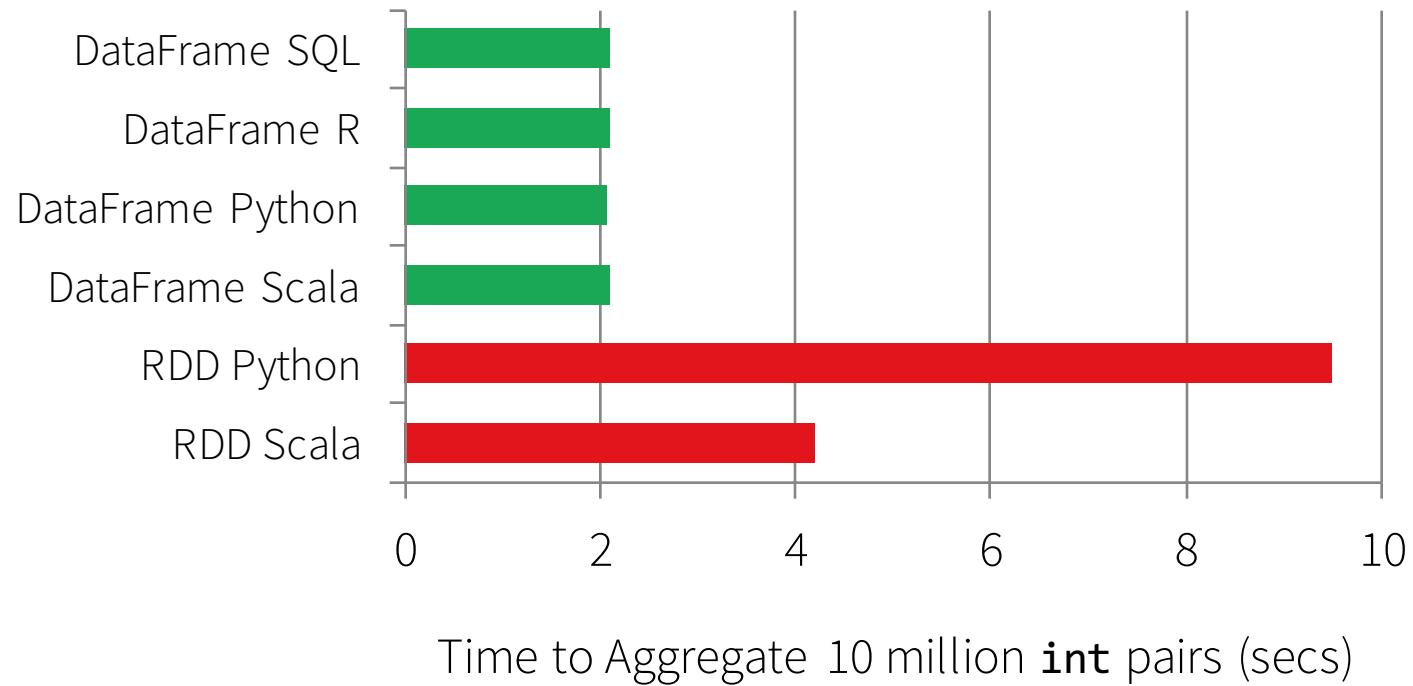
- Compile-time type-safety
- Easy to express certain types of logic
- Lots of existing code + users
- Lower level control of Spark
- Imperative

DataFrames / SQL

Structured Binary Data (Tungsten)

- Lower memory pressure (gc & space)
- Memory accounting (avoids OOMs)
- Faster sorting / hashing / serialization
- More opportunities for automatic optimization
- Declarative

Not Just Less Code: Faster Implementations



<https://gist.github.com/rxin/c1592c133e4bccf515dd>

Dataset API

- DataFrames are useful, but their API is dynamically typed: can we get their performance advantages while keeping static typing?
- New API arrived in Spark 1.6

Datasets

RDDs

- Functional Programming
- Type-safe

Dataframes

- Relational
- Catalyst query optimization
- Tungsten direct/packed RAM
- JIT code generation
- Sorting/suffling without deserializing

Spark 

Source: michaelmalak

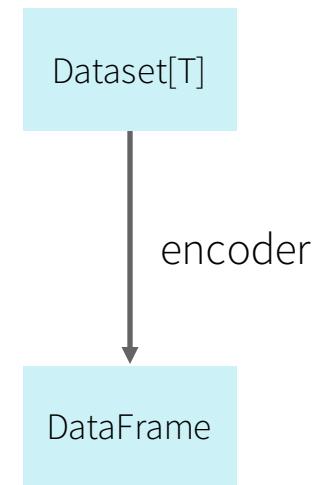
Dataset API in Spark 1.6

- Typed interface over DataFrames / Tungsten

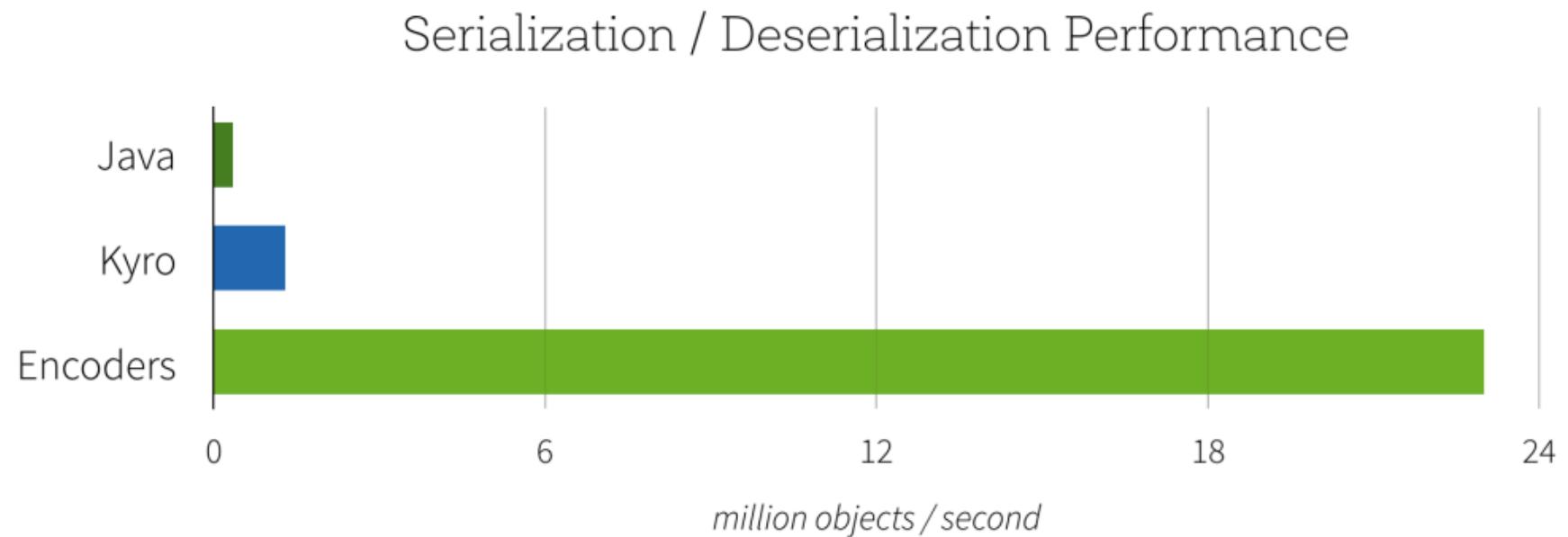
- `case class Person(name: String, age: Int)`
- `val dataframe = read.json("people.json")`
- `val ds: Dataset[Person] = dataframe.as[Person]`
- `ds.filter(p => p.name.startsWith("M"))`
`.groupBy("name")`
`.avg("age")`

Dataset

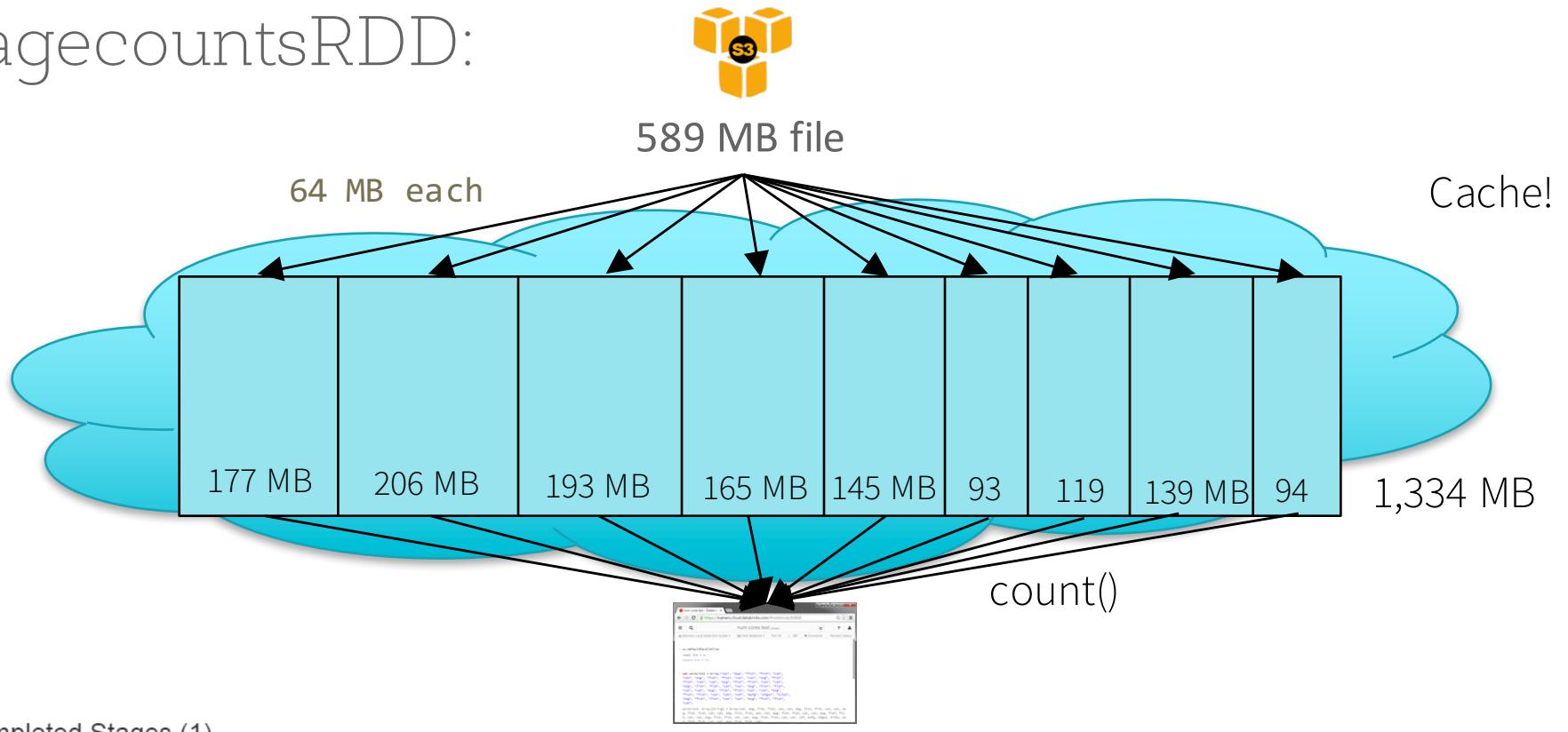
- “Encoder” to specify type information so Spark can translate it into DataFrame and generate optimized memory layouts
- Checkout SPARK-9999



Datasets: Lightning-fast Serialization with Encoders



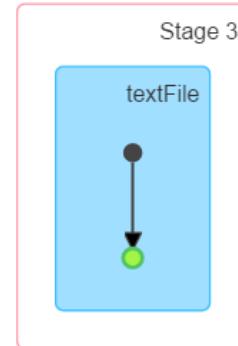
pagecountsRDD:



Completed Stages (1)

Stage Id	Pool Name	Description	Submitted	Duration	Tasks: Succeeded/Total	Input	Output	Shuffle Read	Shuffle Write
3	7463896219951833835	pagecountsRDD.setName("pagecountsRDD").cache.count count at <console>:29	+details	2015/11/30 10:26:45	5 s	9/9	562.4 MB		

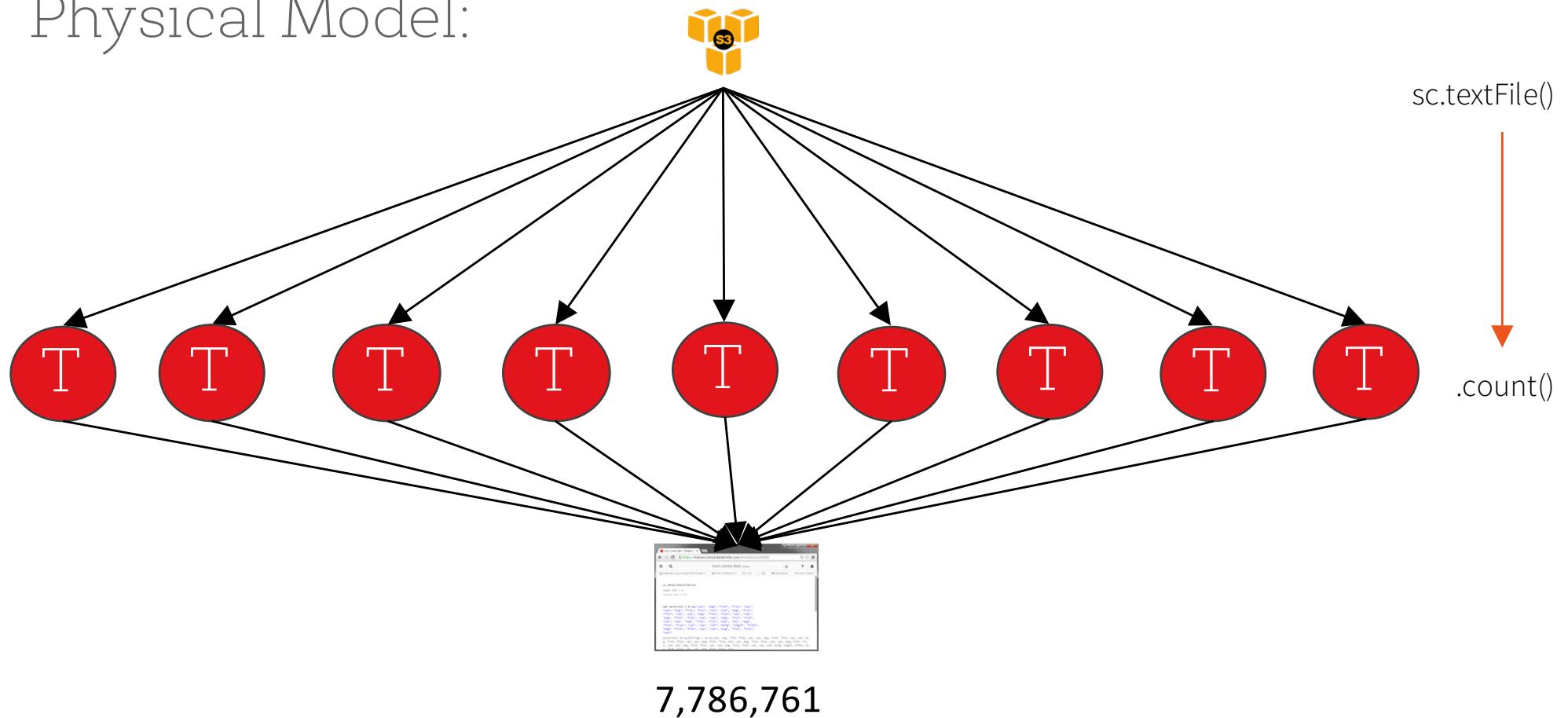
pagecountsRDD:



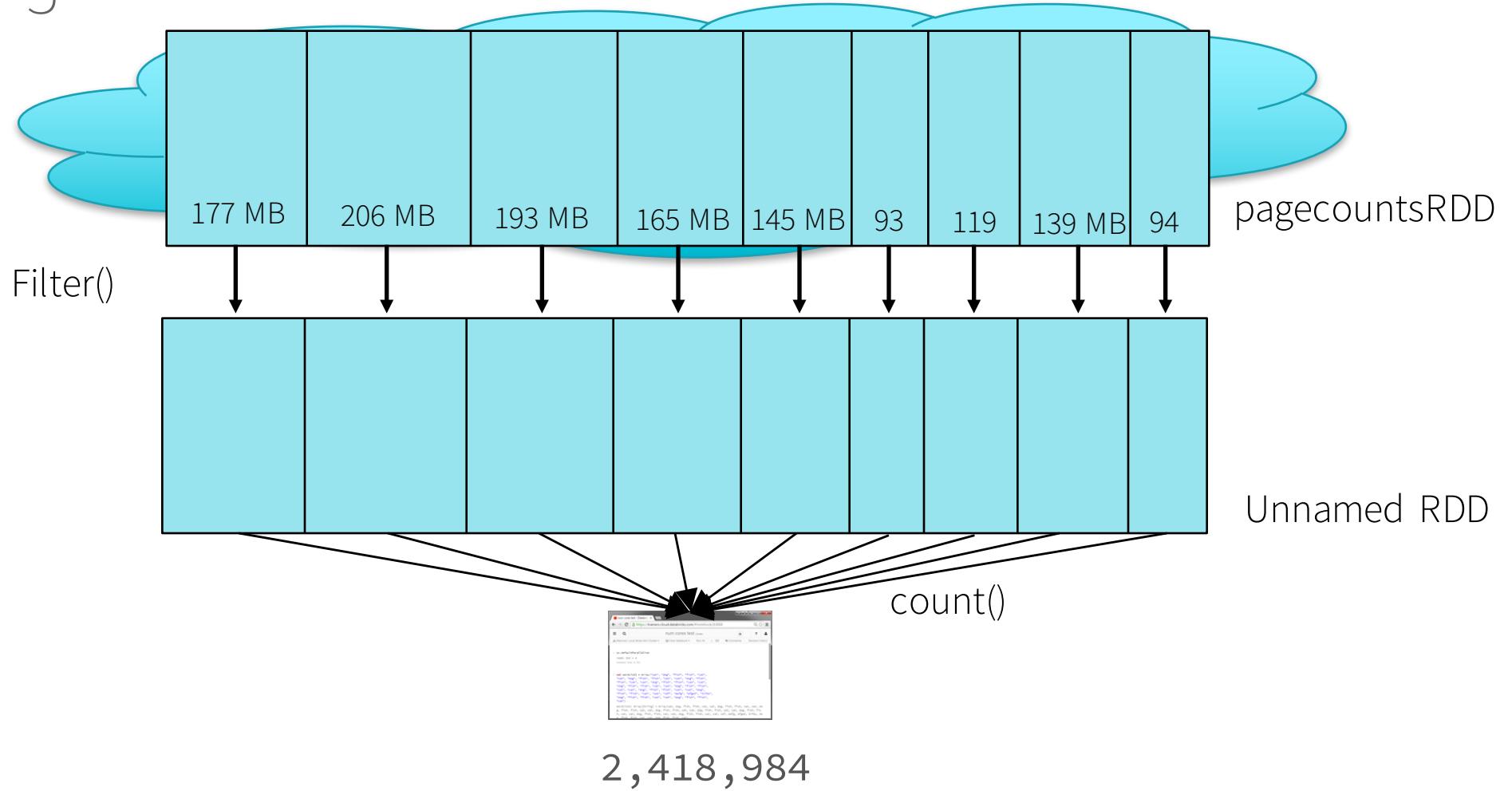
Tasks

Index ?	ID	Attempt	Status	Locality Level	Executor ID / Host	Launch Time	Duration	GC Time	Input Size / Records	Errors
0	13	0	SUCCESS	PROCESS_LOCAL	driver / localhost	2015/11/30 10:26:45	2 s		64.0 MB (hadoop) / 1398177	
1	14	0	SUCCESS	PROCESS_LOCAL	driver / localhost	2015/11/30 10:26:45	2 s		64.0 MB (hadoop) / 1906223	
2	15	0	SUCCESS	PROCESS_LOCAL	driver / localhost	2015/11/30 10:26:45	2 s		64.0 MB (hadoop) / 1709442	
3	16	0	SUCCESS	PROCESS_LOCAL	driver / localhost	2015/11/30 10:26:45	2 s		64.0 MB (hadoop) / 1226943	
4	17	0	SUCCESS	PROCESS_LOCAL	driver / localhost	2015/11/30 10:26:46	2 s	0.4 s	64.0 MB (hadoop) / 778096	
5	18	0	SUCCESS	PROCESS_LOCAL	driver / localhost	2015/11/30 10:26:46	2 s	0.4 s	64.0 MB (hadoop) / 306501	
6	19	0	SUCCESS	PROCESS_LOCAL	driver / localhost	2015/11/30 10:26:46	2 s	0.4 s	64.0 MB (hadoop) / 90926	
7	20	0	SUCCESS	PROCESS_LOCAL	driver / localhost	2015/11/30 10:26:46	2 s	0.4 s	64.0 MB (hadoop) / 161412	
8	21	0	SUCCESS	PROCESS_LOCAL	driver / localhost	2015/11/30 10:26:48	1 s		50.4 MB (hadoop) / 209041	

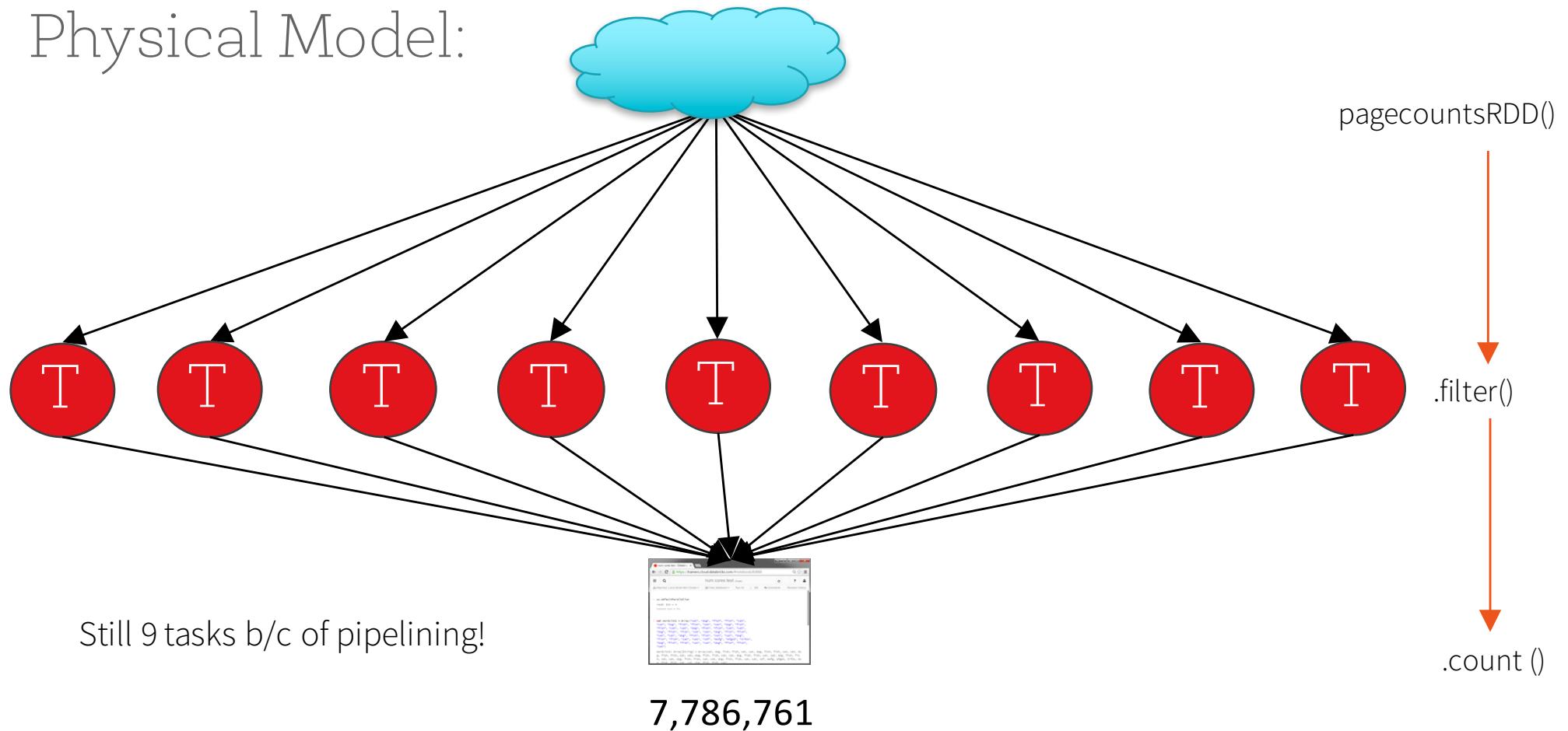
Physical Model:



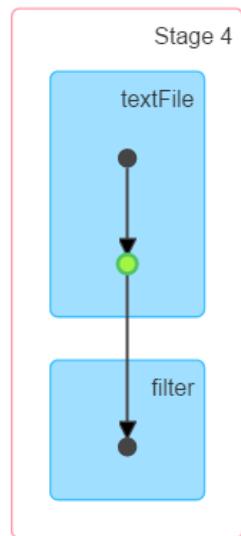
pagecountsRDD:



Physical Model:



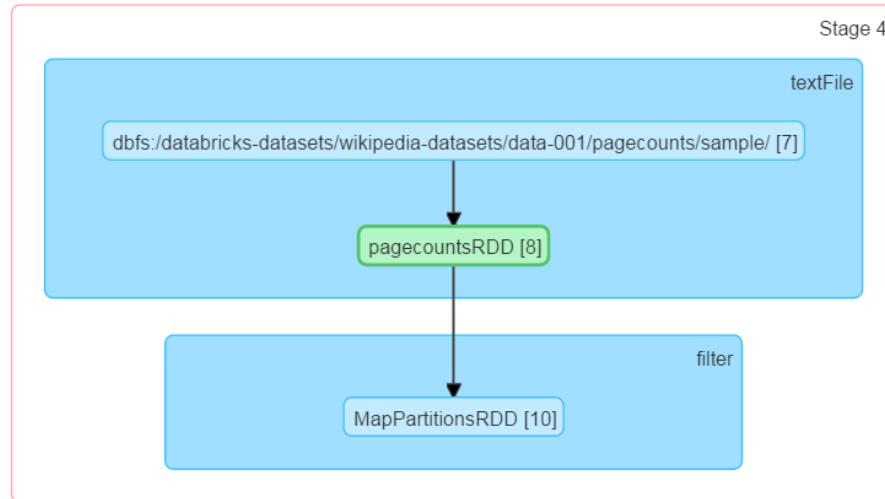
pagecountsRDD:



Completed Stages (1)

Stage Id	Pool Name	Description	Submitted	Duration	Tasks: Succeeded/Total	Input	Output	Shuffle Read	Shuffle Write
4	7463896219951833835	enPagecountsRDD.count() count at <console>:31	+details	2015/11/30 10:46:55	0.2 s	9/9	1334.3 MB		

pagecountsRDD:



Tasks

Index ?	ID	Attempt	Status	Locality Level	Executor ID / Host	Launch Time	Duration	GC Time	Input Size / Records	Errors
0	22	0	SUCCESS	PROCESS_LOCAL	driver / localhost	2015/11/30 10:46:55	62 ms		177.2 MB (memory) / 1398177	
1	23	0	SUCCESS	PROCESS_LOCAL	driver / localhost	2015/11/30 10:46:55	94 ms		206.1 MB (memory) / 1906223	
2	24	0	SUCCESS	PROCESS_LOCAL	driver / localhost	2015/11/30 10:46:55	99 ms		193.8 MB (memory) / 1709442	
3	25	0	SUCCESS	PROCESS_LOCAL	driver / localhost	2015/11/30 10:46:55	82 ms		165.2 MB (memory) / 1226943	
4	26	0	SUCCESS	PROCESS_LOCAL	driver / localhost	2015/11/30 10:46:55	55 ms		145.2 MB (memory) / 778096	
5	27	0	SUCCESS	PROCESS_LOCAL	driver / localhost	2015/11/30 10:46:55	21 ms		93.0 MB (memory) / 306501	
6	28	0	SUCCESS	PROCESS_LOCAL	driver / localhost	2015/11/30 10:46:55	9 ms		119.8 MB (memory) / 90926	
7	29	0	SUCCESS	PROCESS_LOCAL	driver / localhost	2015/11/30 10:46:55	14 ms		139.6 MB (memory) / 161412	
8	30	0	SUCCESS	PROCESS_LOCAL	driver / localhost	2015/11/30 10:46:55	16 ms		94.4 MB (memory) / 209041	

pagecountsRDD:



Summary Metrics for 9 Completed Tasks

Metric	Min	25th percentile	Median	75th percentile	Max
Duration	9 ms	16 ms	55 ms	82 ms	99 ms
GC Time	0 ms	0 ms	0 ms	0 ms	0 ms
Input Size / Records	93.0 MB / 90926	119.8 MB / 209041	145.2 MB / 778096	177.2 MB / 1398177	206.1 MB / 1906223

Lab!

- Analyzing Pagecounts with RDDs

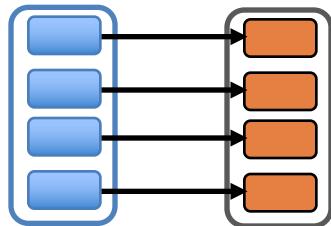


VS



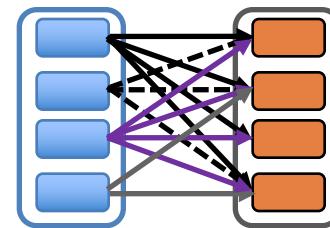
narrow

each partition of the parent RDD is used by at most one partition of the child RDD



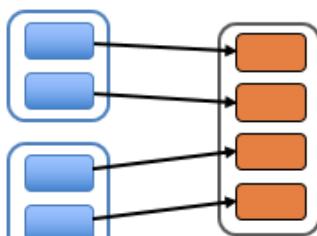
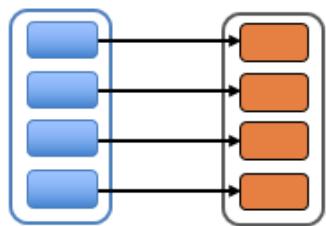
wide

multiple child RDD partitions may depend on a single parent RDD partition

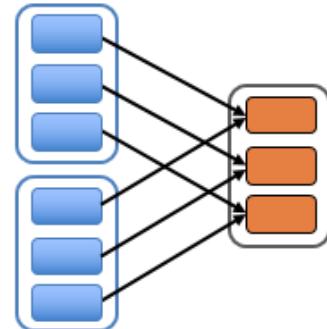


narrow

each partition of the parent RDD is used by at most one partition of the child RDD

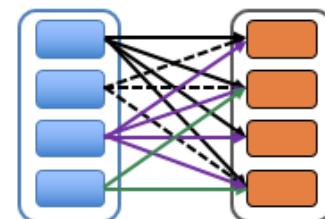


join w/ inputs co-partitioned

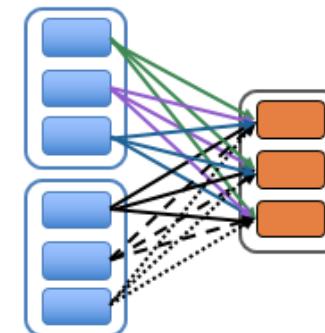


wide

multiple child RDD partitions may depend on a single parent RDD partition



groupByKey



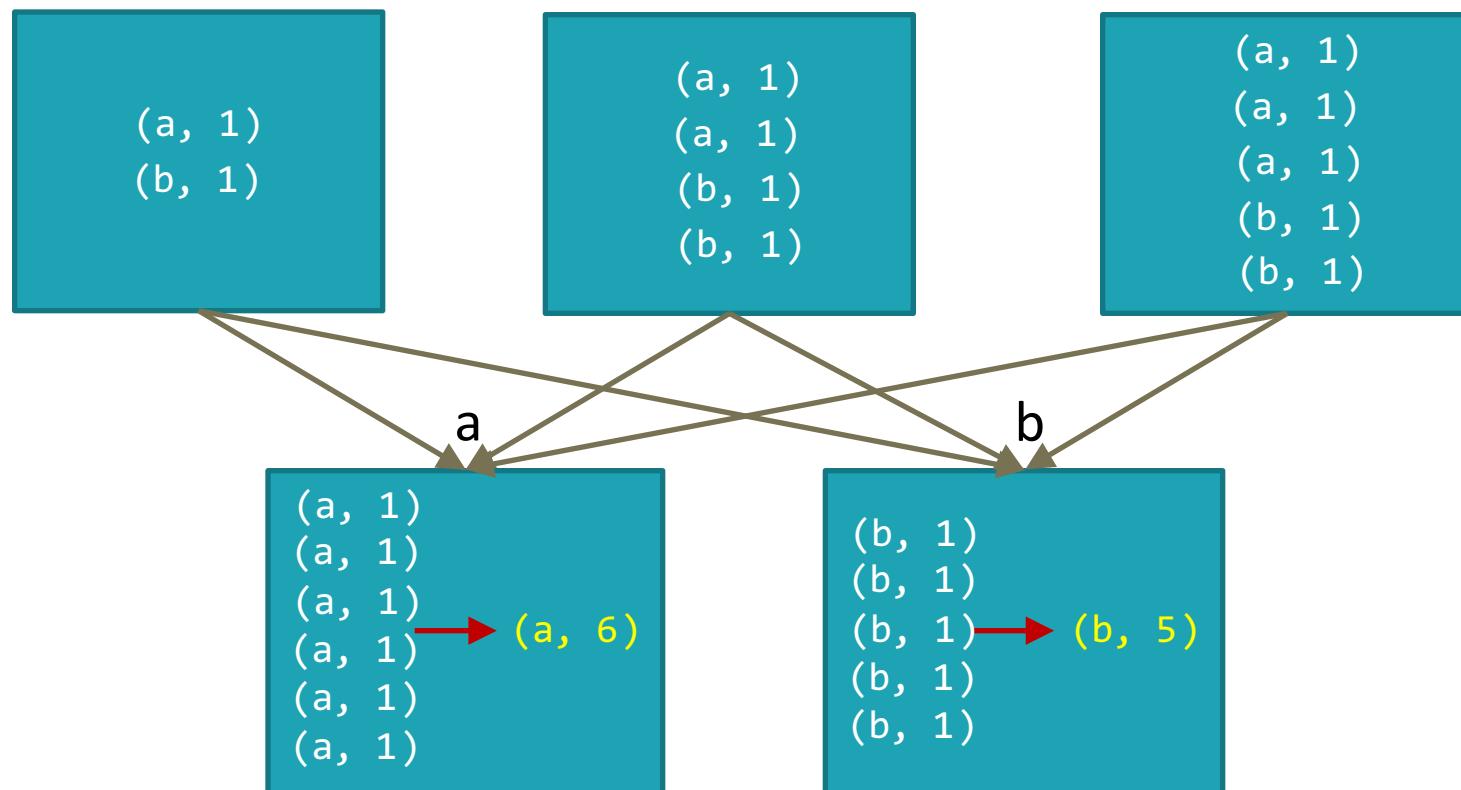
join w/ inputs not co-partitioned

Which API Call Causes Most Tickets?

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

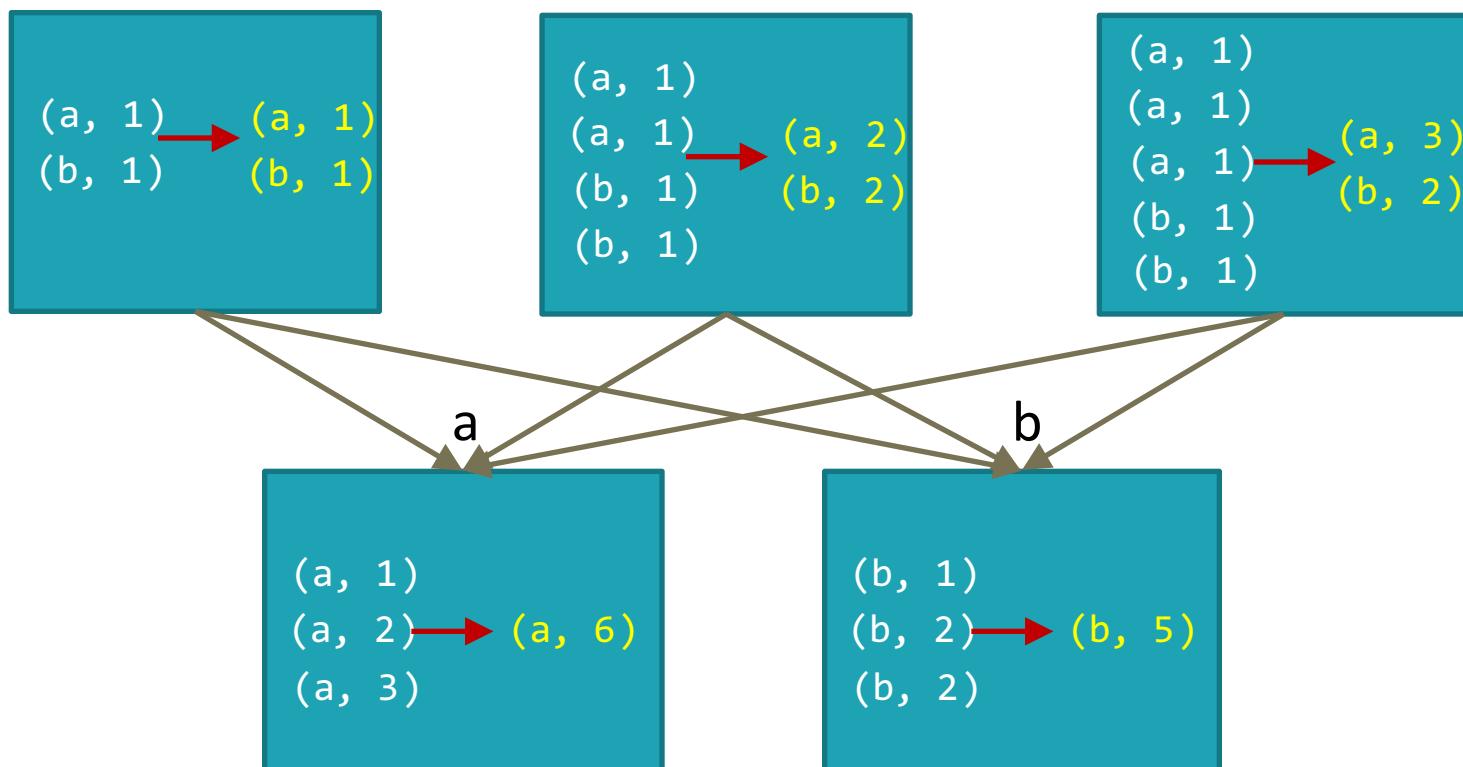
GroupByKey

Shuffles all the data



ReduceByKey

Shuffles only the results of sub-aggregations
in each partition of the data





Spark | Memory

Memory Storage Levels:

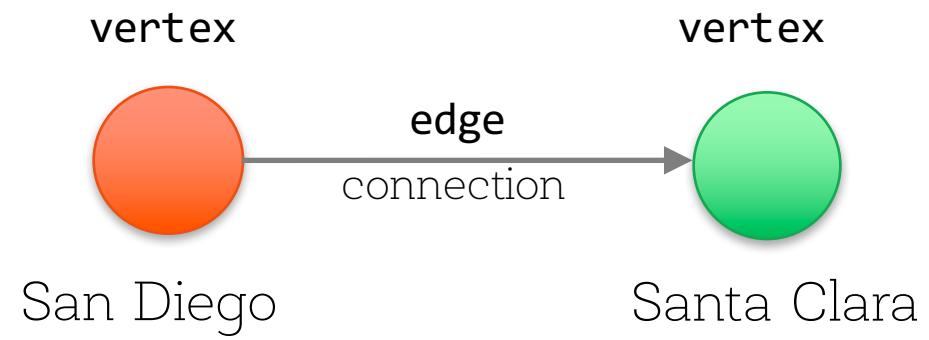
- MEMORY_ONLY
- MEMORY_AND_DISK
- MEMORY_ONLY_SER
- MEMORY_AND_DISK_SER
- DISK_ONLY
- MEMORY_ONLY_2, MEMORY_AND_DISK_2, etc
- OFF_HEAP

Special Case:

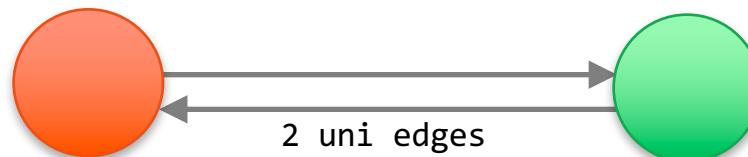
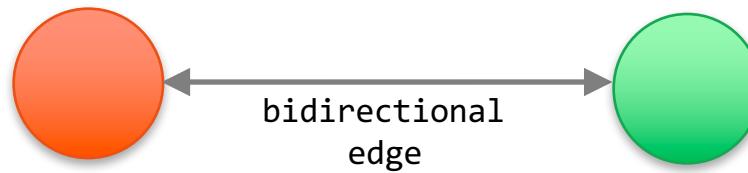
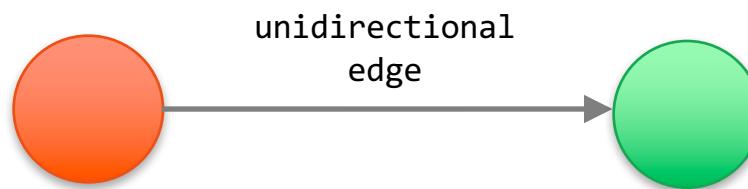




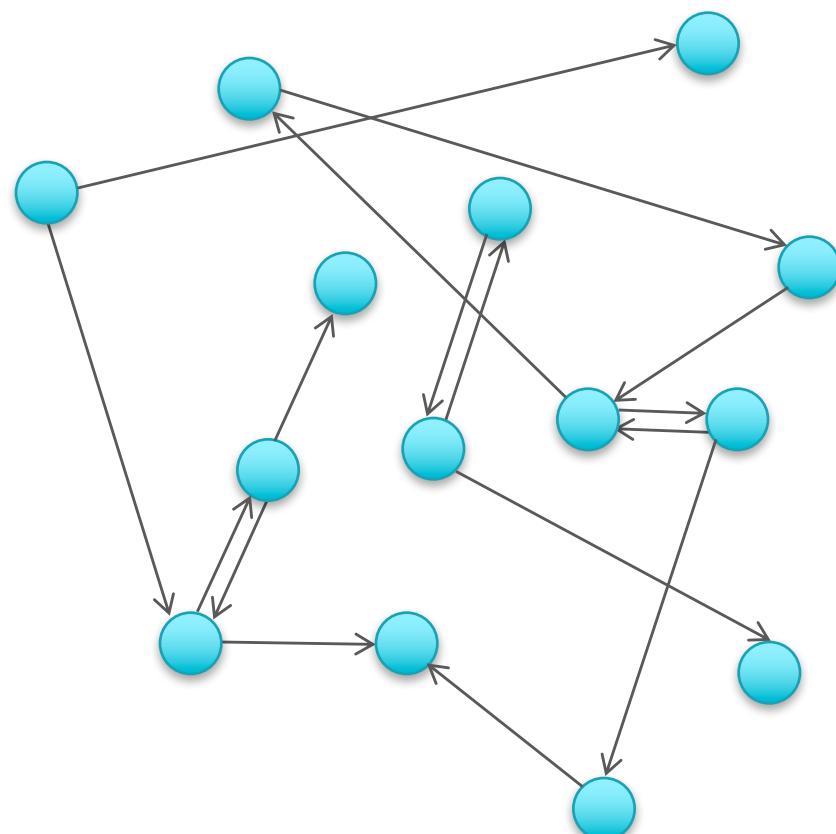
Spark | GraphX

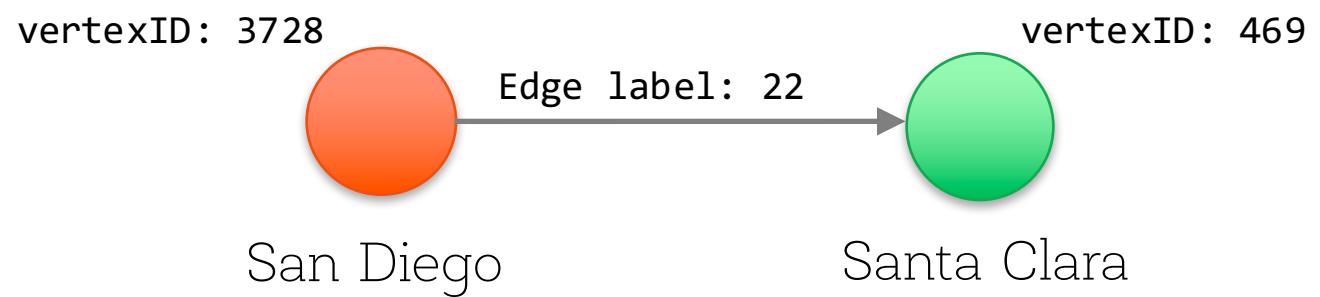


Edge types:

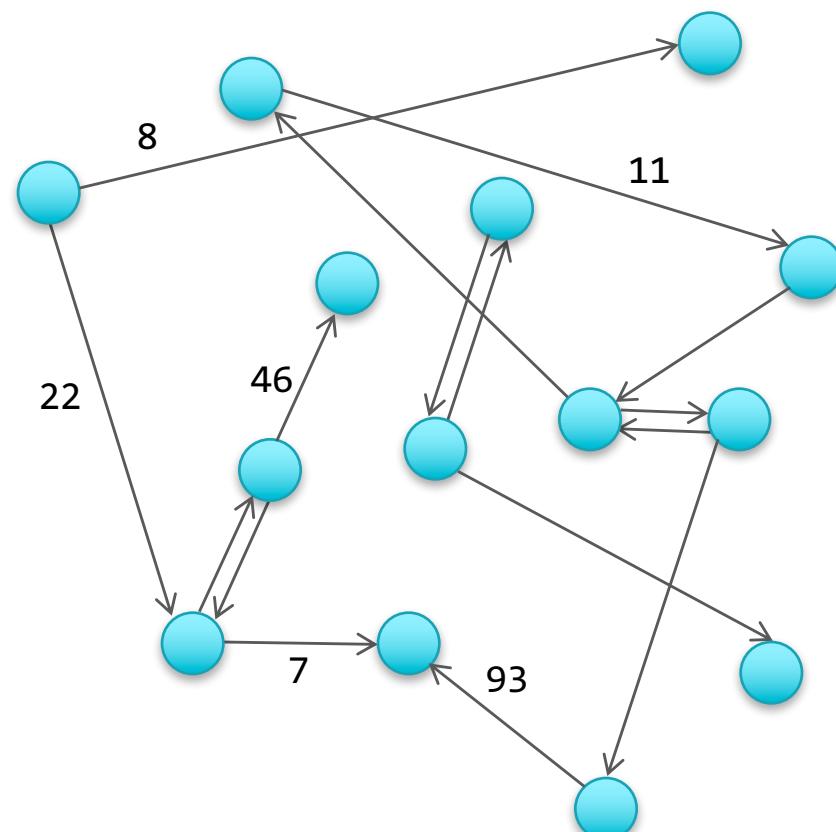


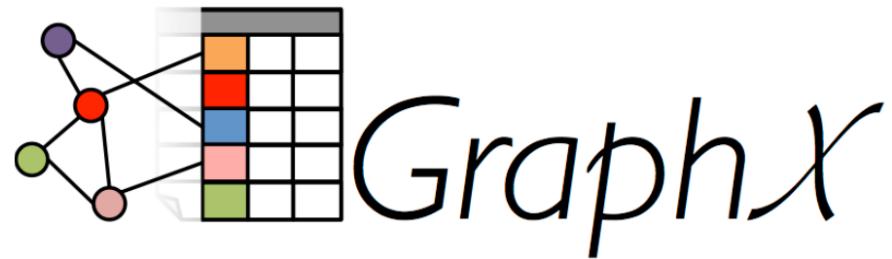
GraphX: directed multigraph





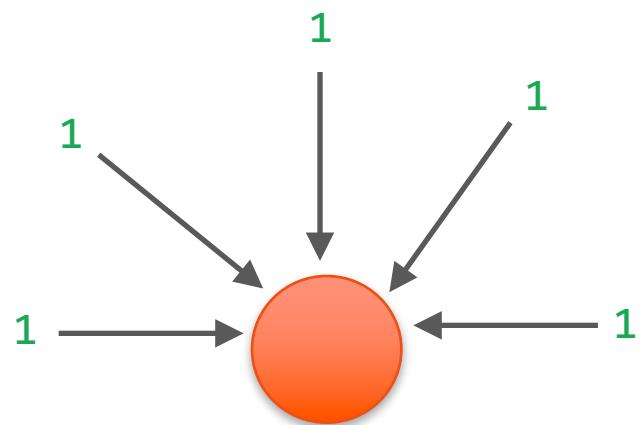
GraphX: directed multigraph w/ labels





Algorithms Tour

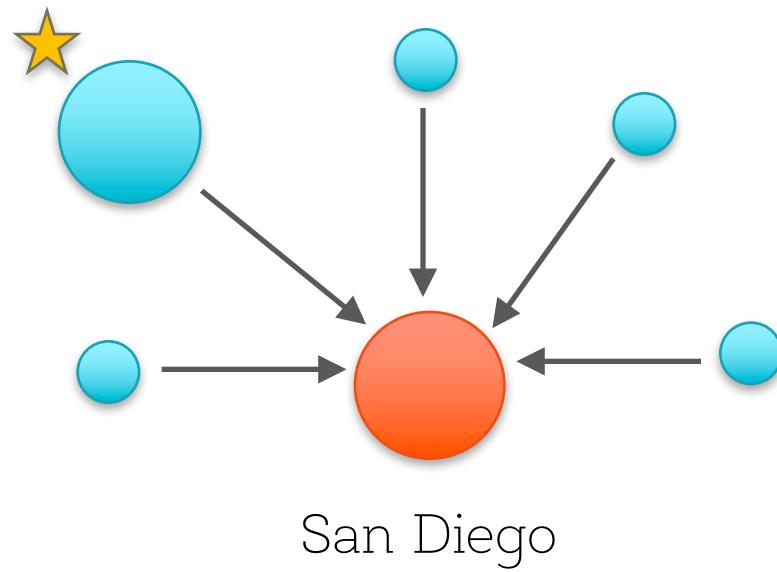
In-degree:



San Diego

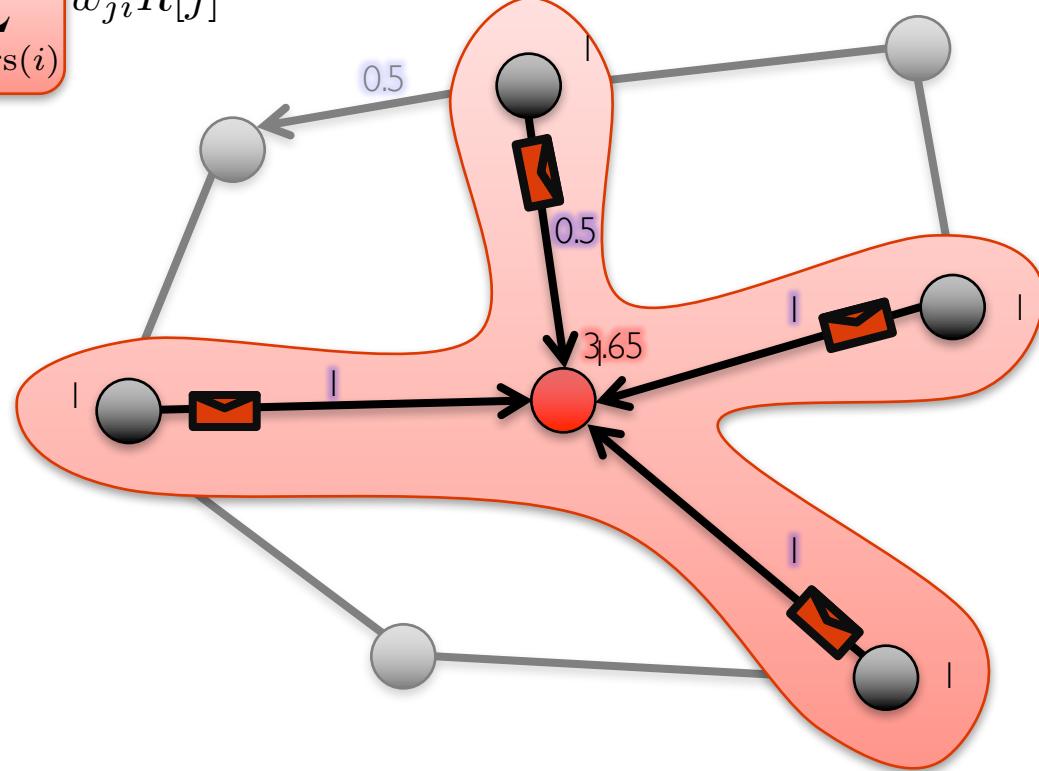
inDegree = 5

PageRank:

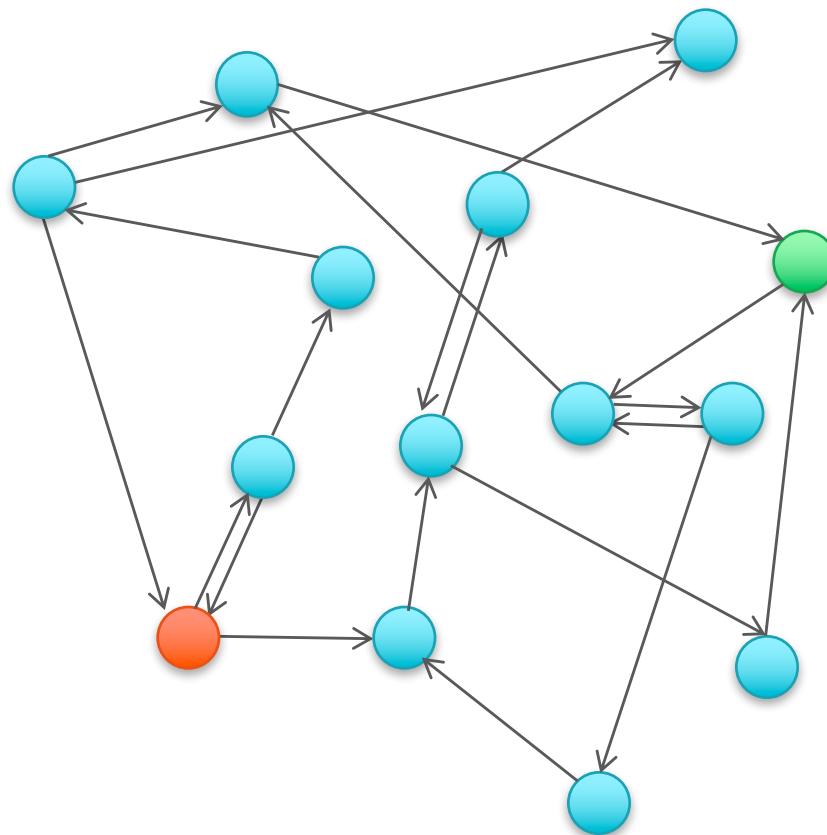


PageRank

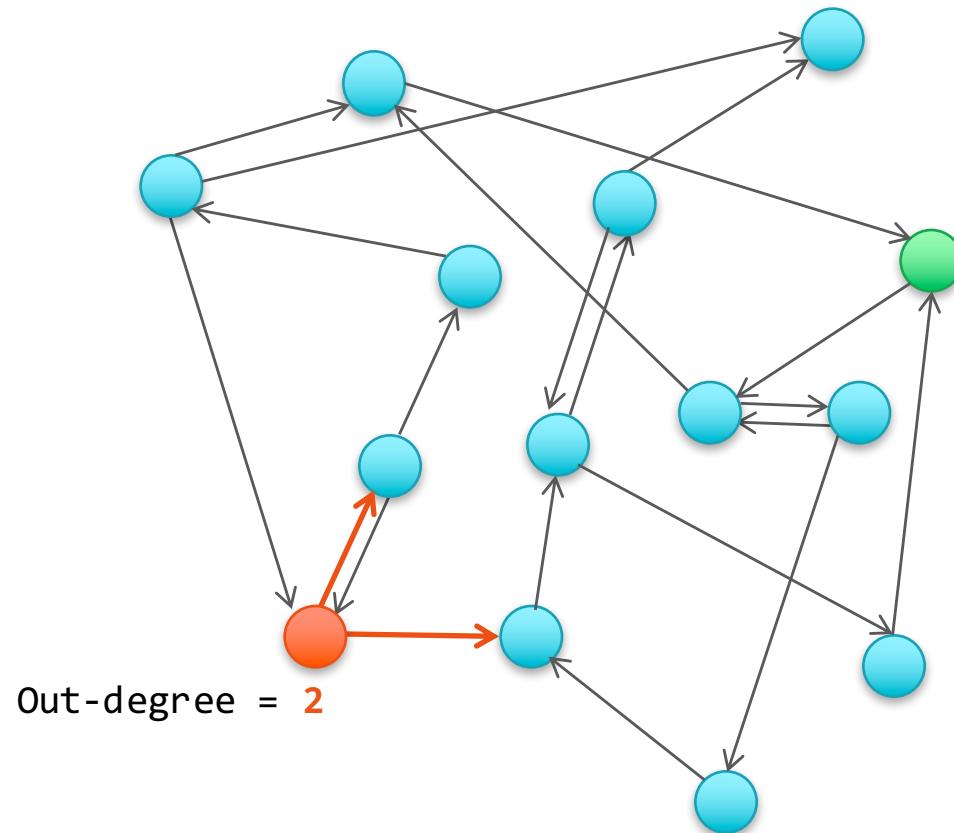
$$R[i] = 0.15 + \sum_{j \in \text{Nbrs}(i)} w_{ji} R[j]$$



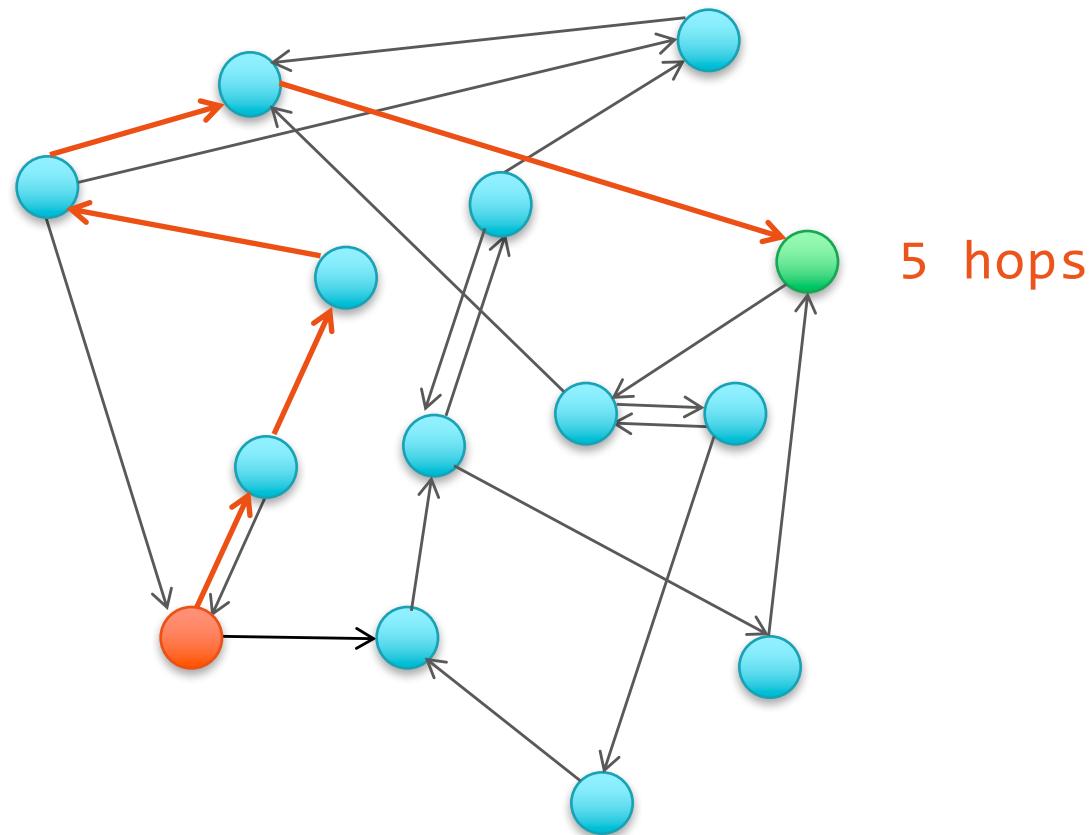
Shortest Path:



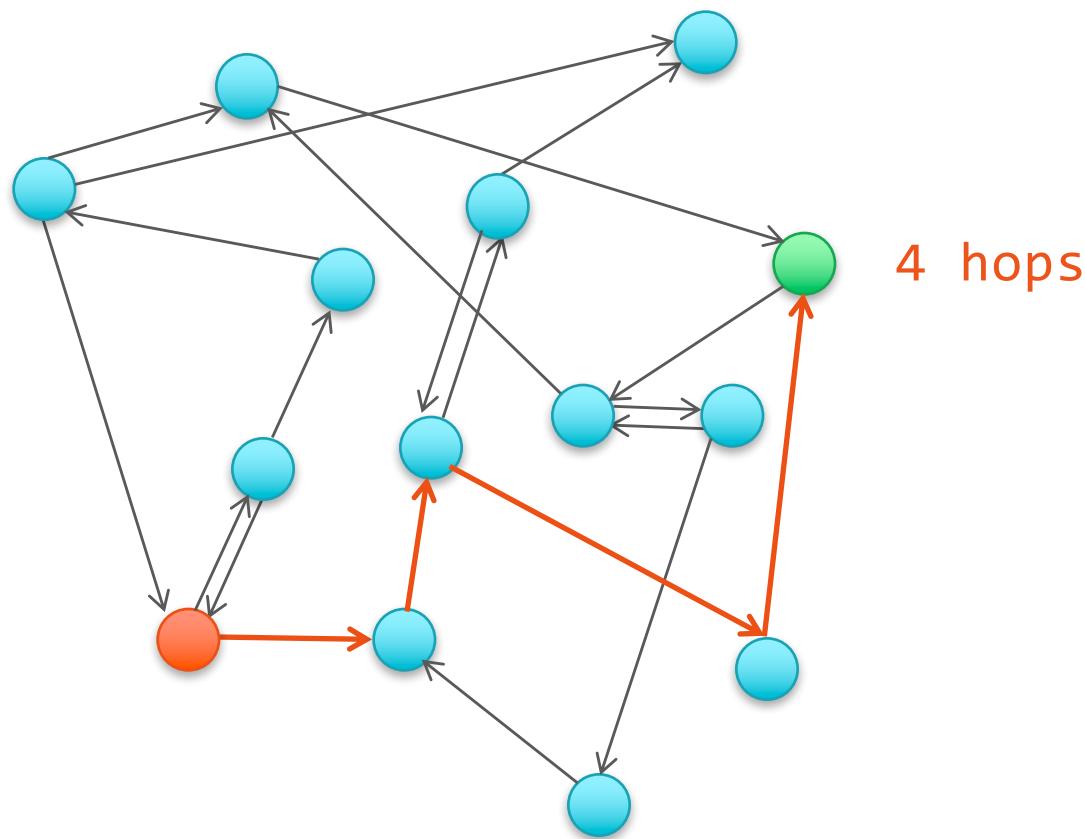
Shortest Path:



Shortest Path:

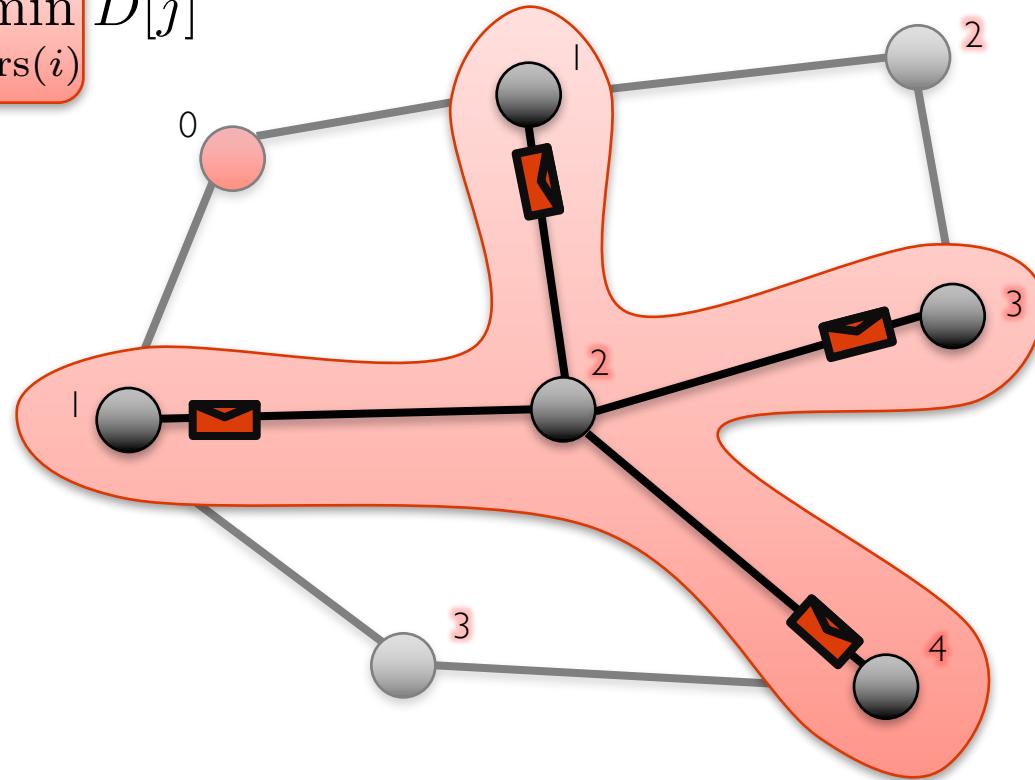


Shortest Path:

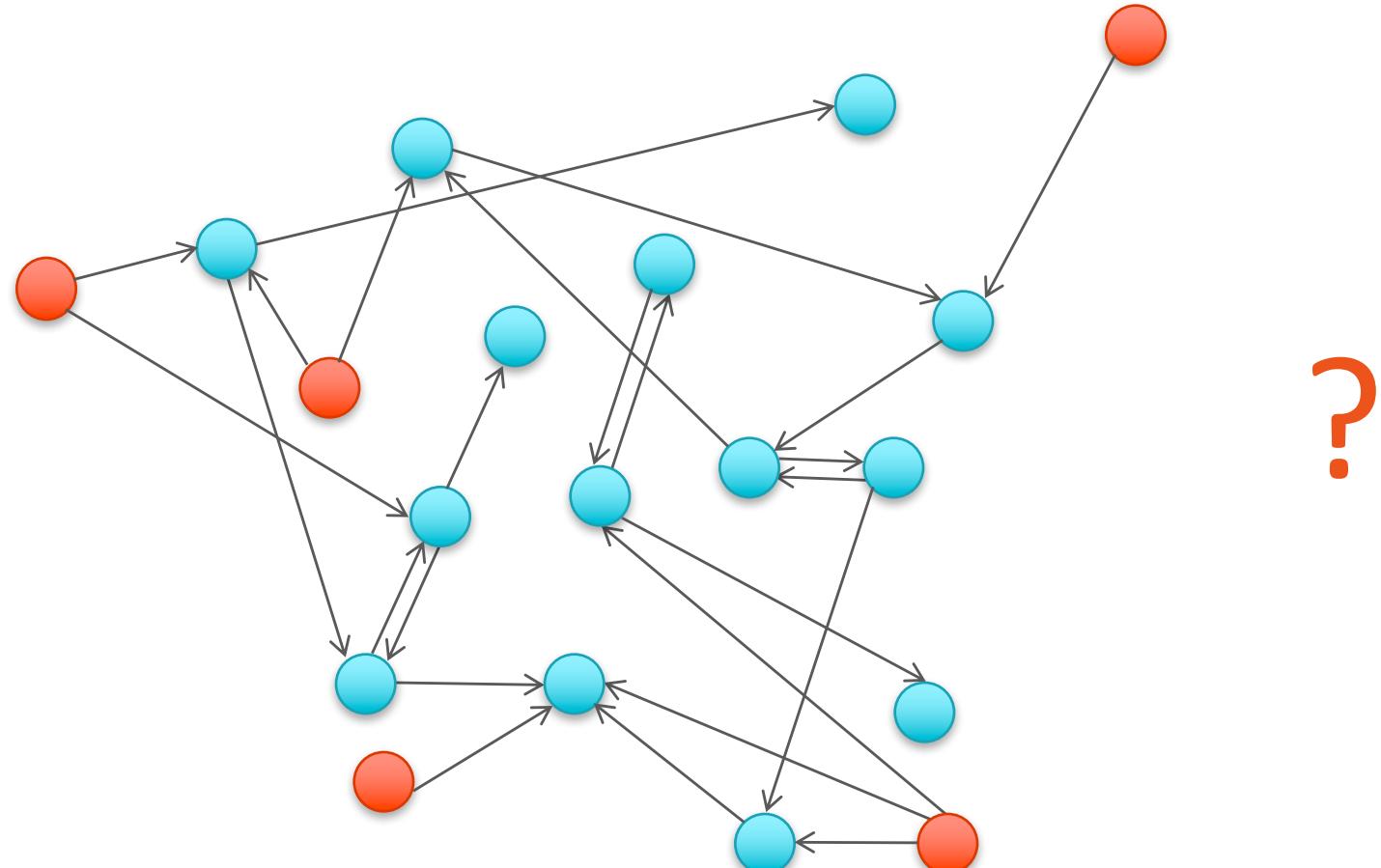


Single-Source Shortest Path

$$D[i] = 1 + \arg \min_{j \in \text{Nbrs}(i)} D[j]$$



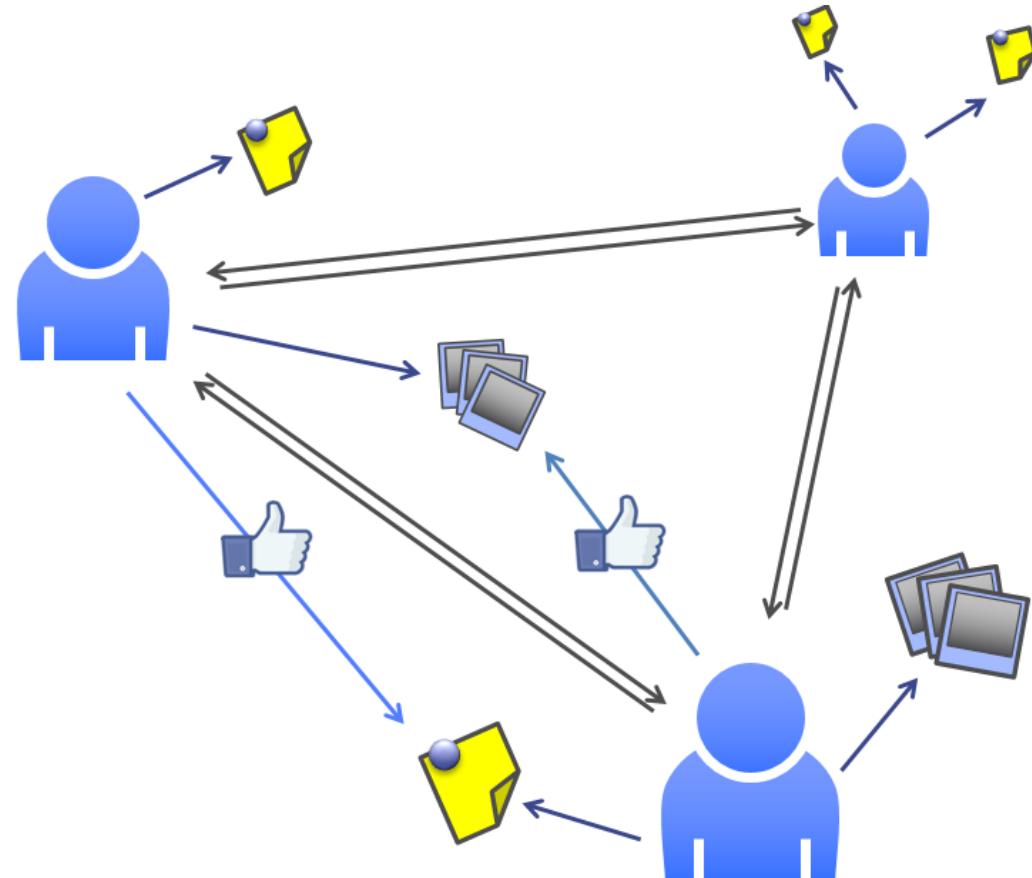
GraphX: 2 types of vertices



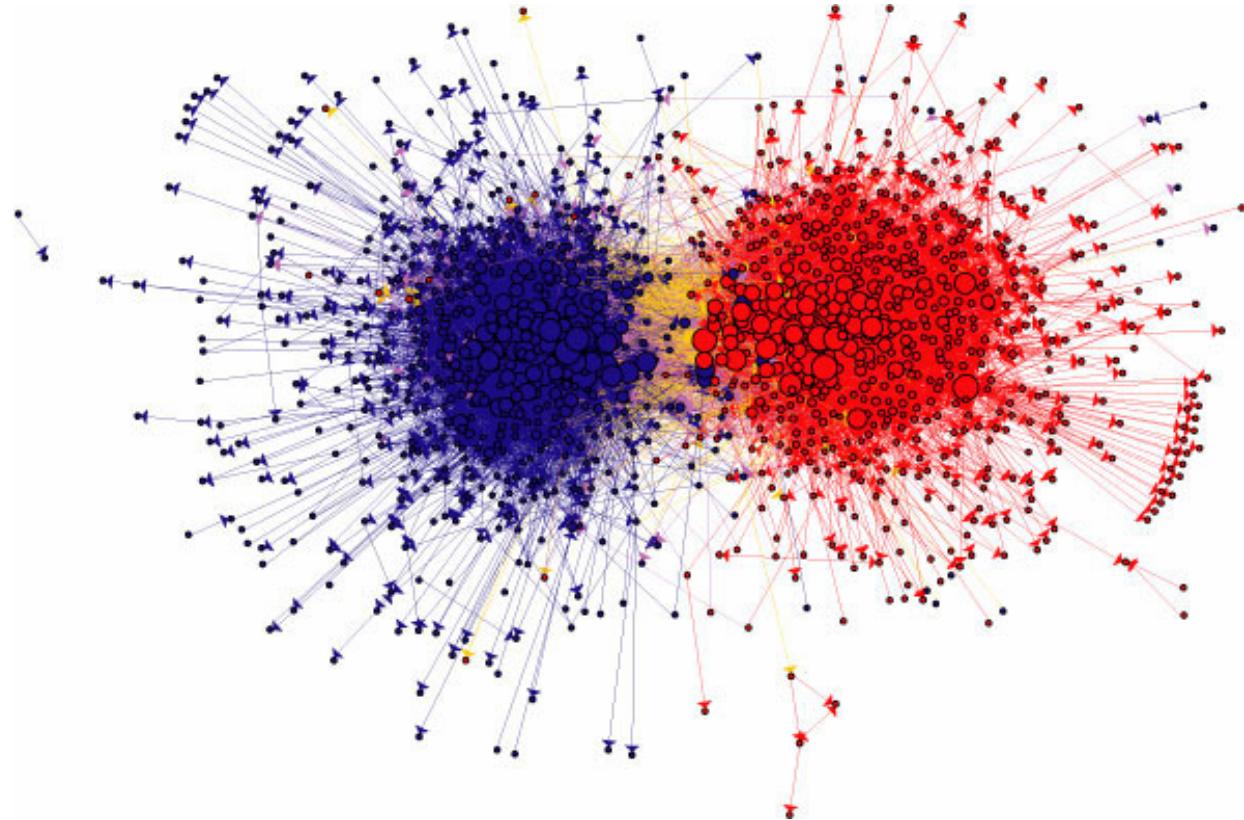
Graphs are everywhere...

Vertices: Users, Images, Posts

Edges: Social Relationships, Likes



Vertices: Web Pages
Edges: Links

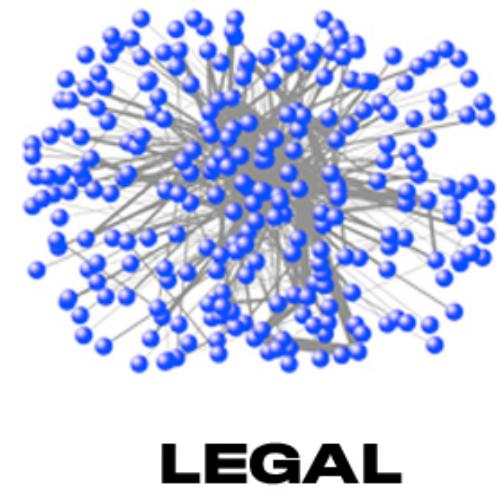
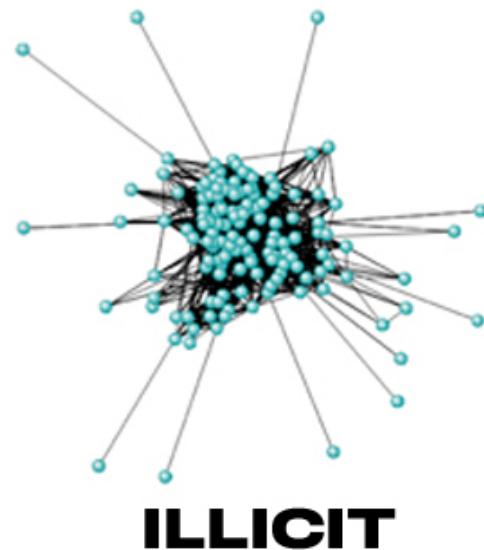


“The political blogosphere and the 2004 U.S. election: divided they blog.”

Source: Adamic and Glance, 2005.

Vertices: Users

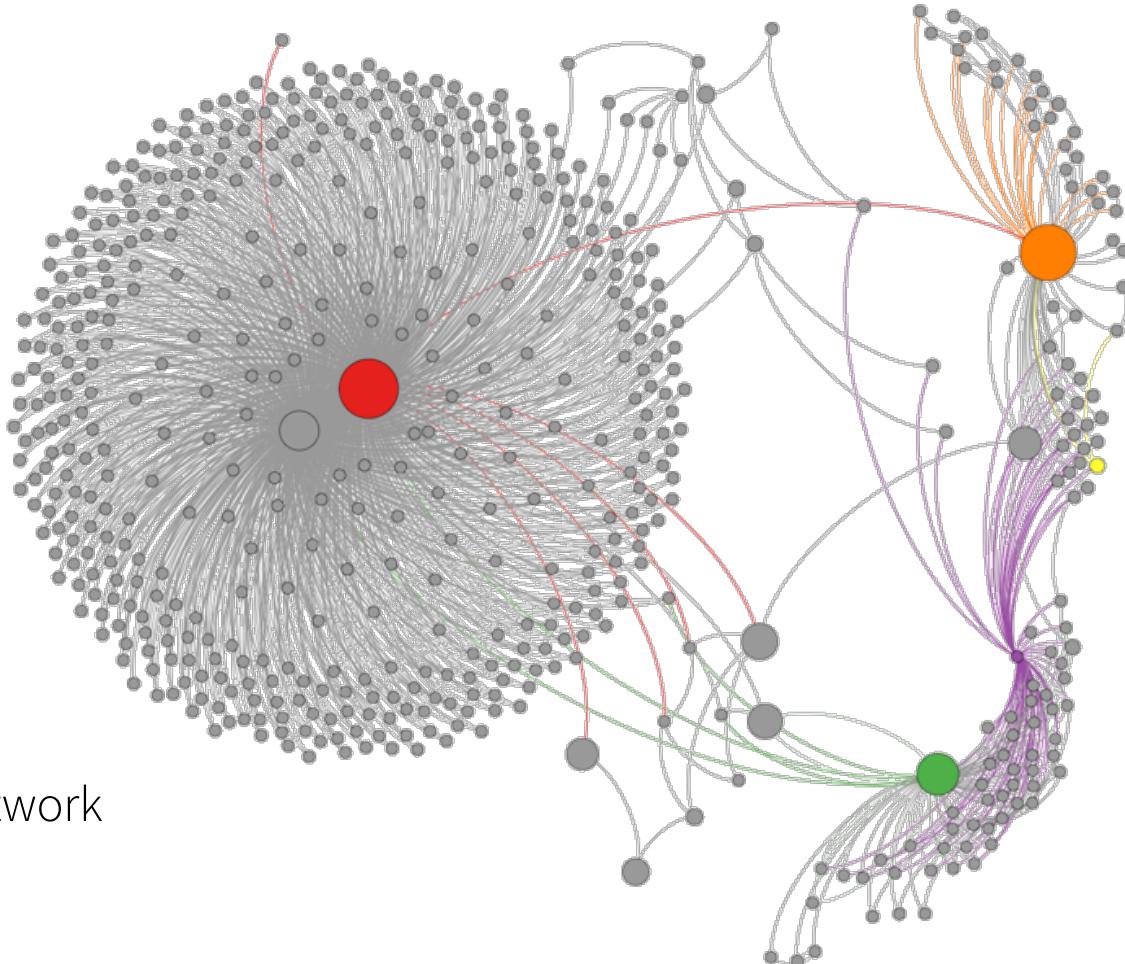
Directed Edges: Email from → to



Enron Email Graph

Source: <https://www.sciencenews.org/article/information-flow-can-reveal-dirty-deeds>

Vertices: People, Orgs
Edges: Exchange of currency

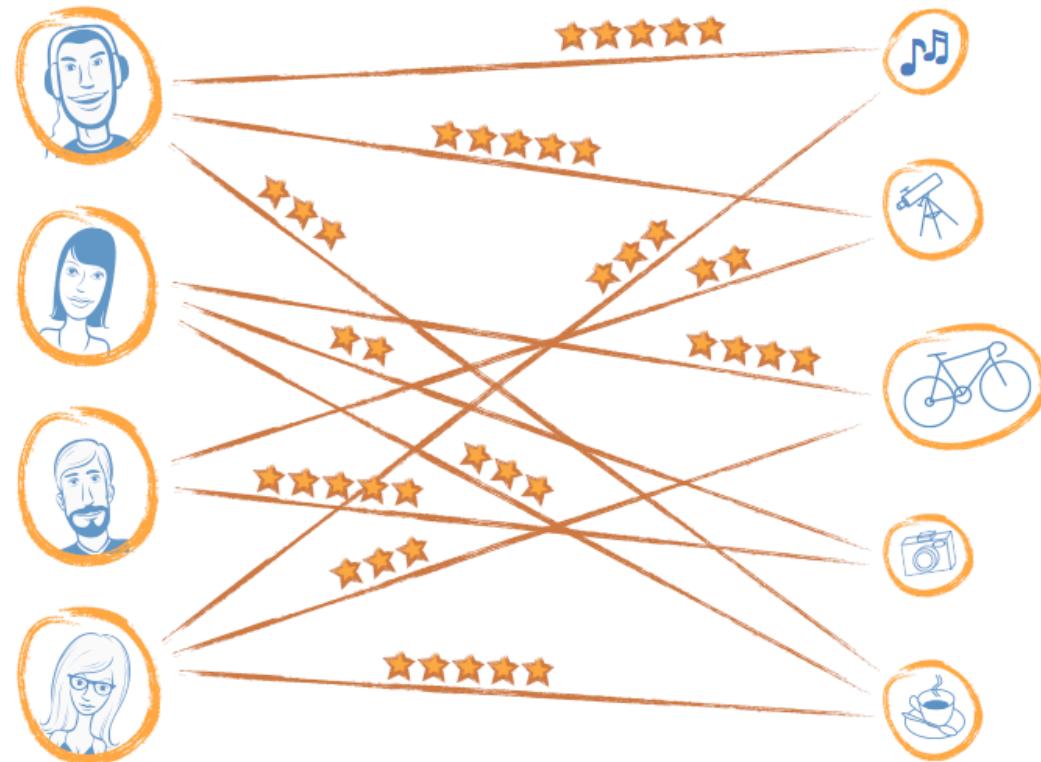


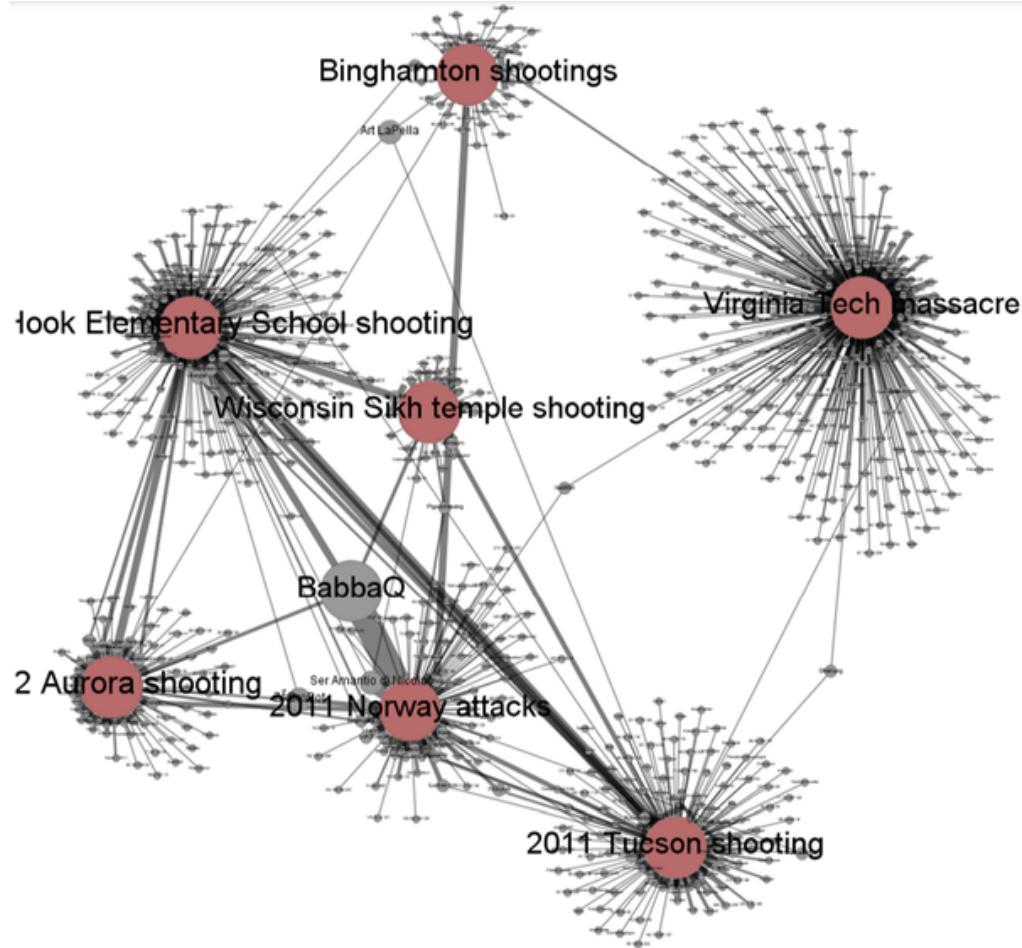
Case Study: Bitcoin Theft Transaction Network

Source: <http://anonymity-in-bitcoin.blogspot.com/2011/07/bitcoin-is-not-anonymous.html>

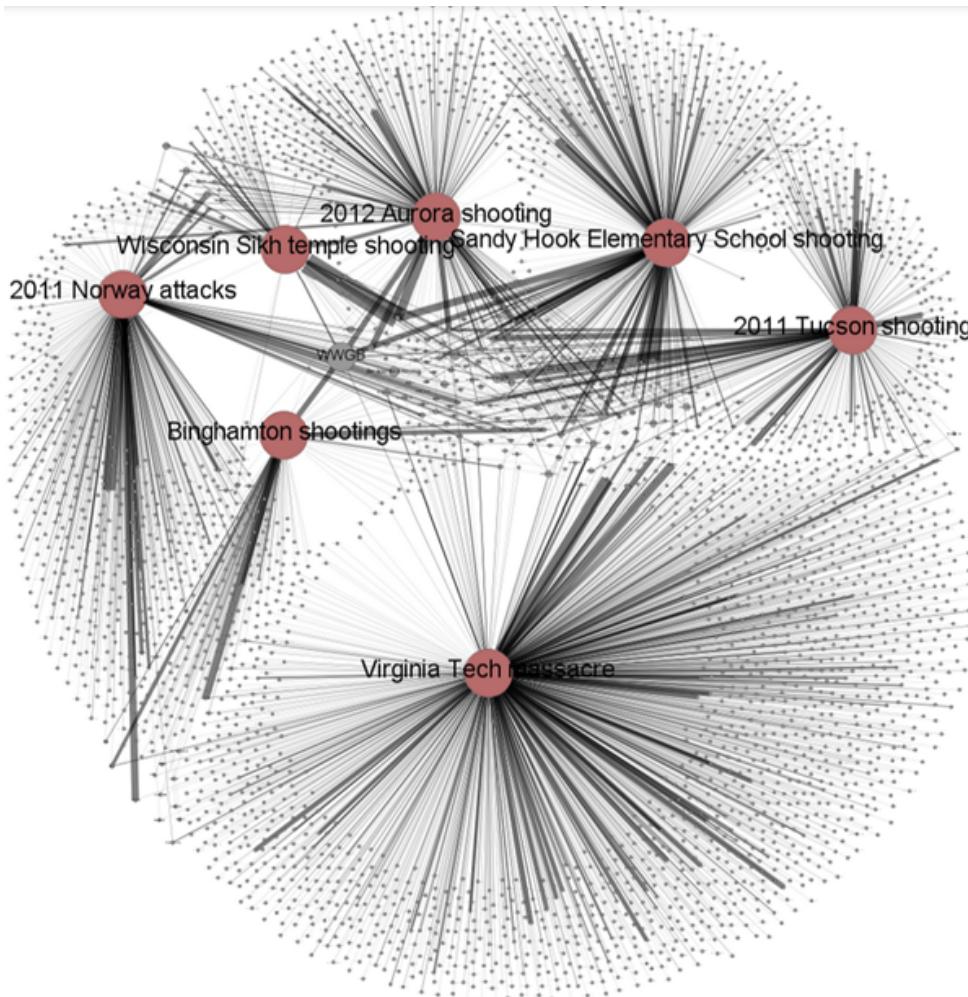
Vertices: Users, Items

Edges: Ratings

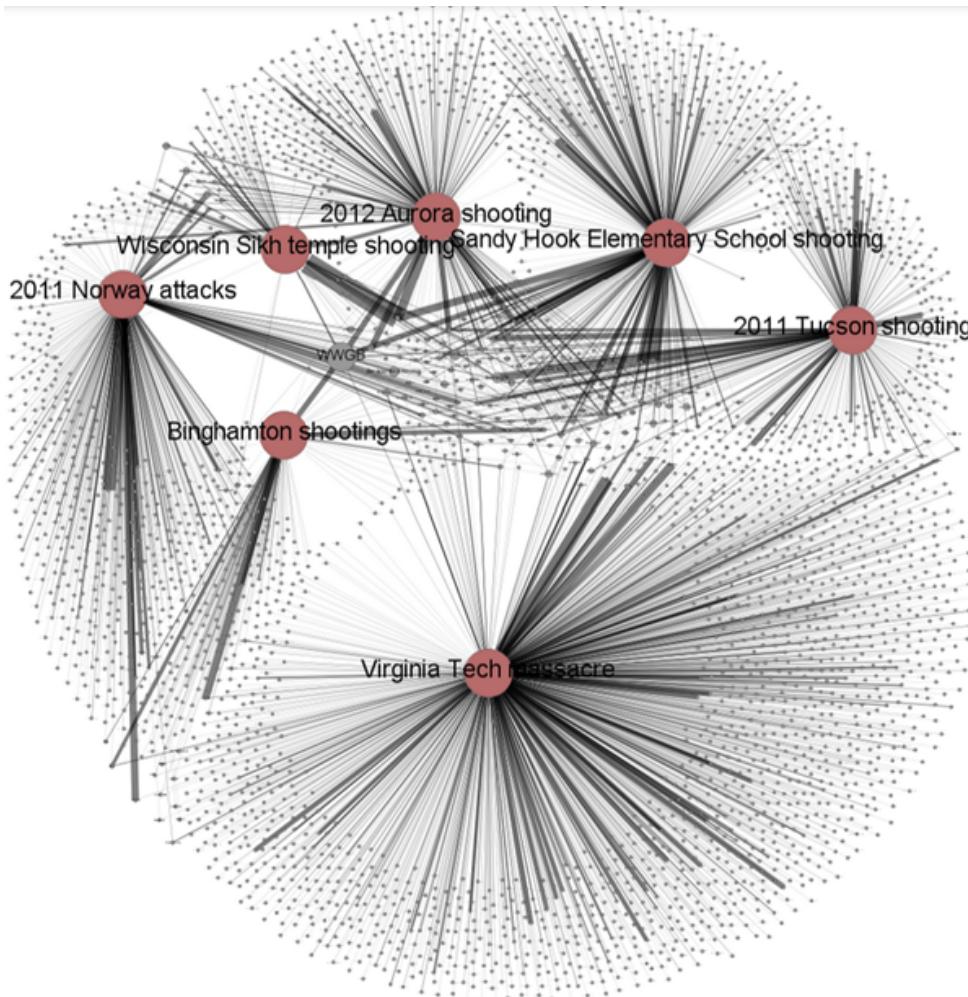




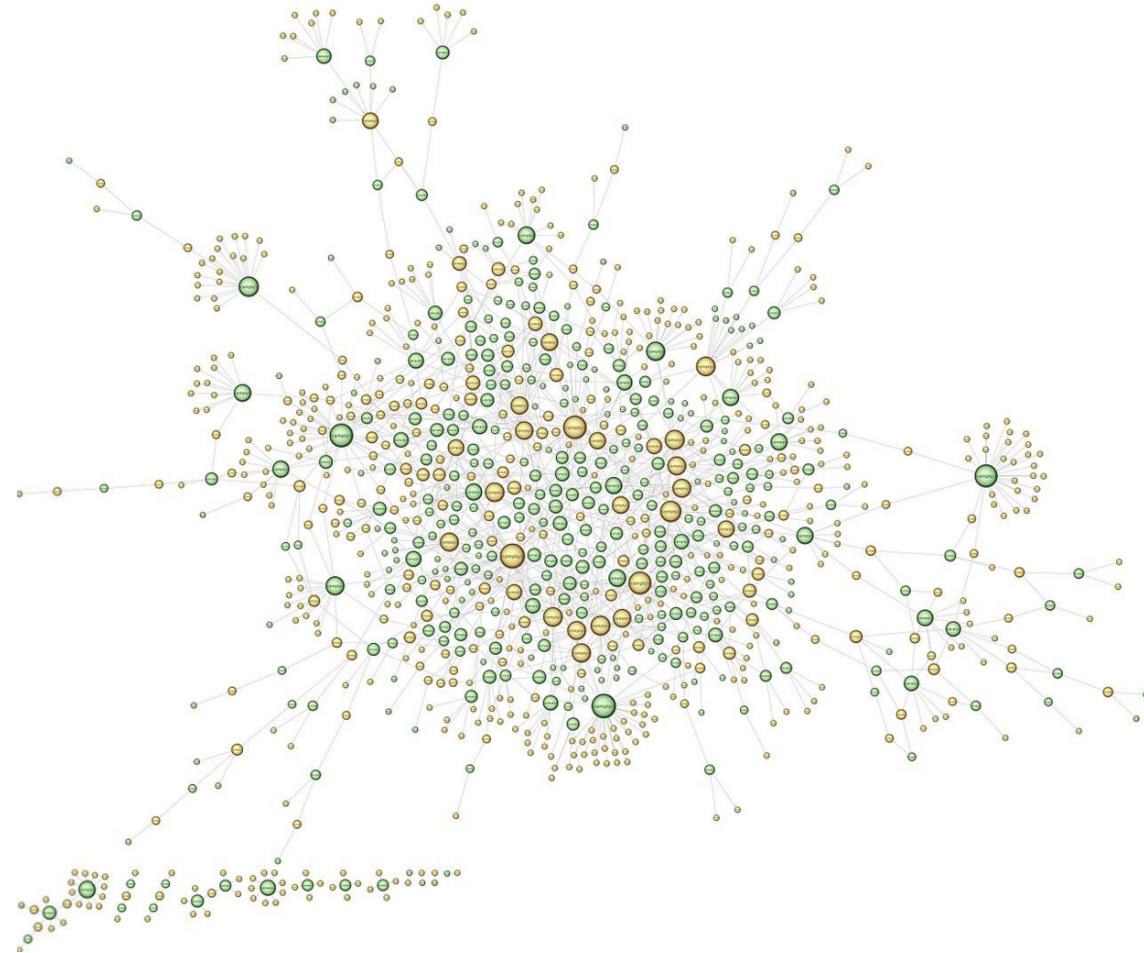
Source: Nieman Journalism Lab



Source: Nieman Journalism Lab

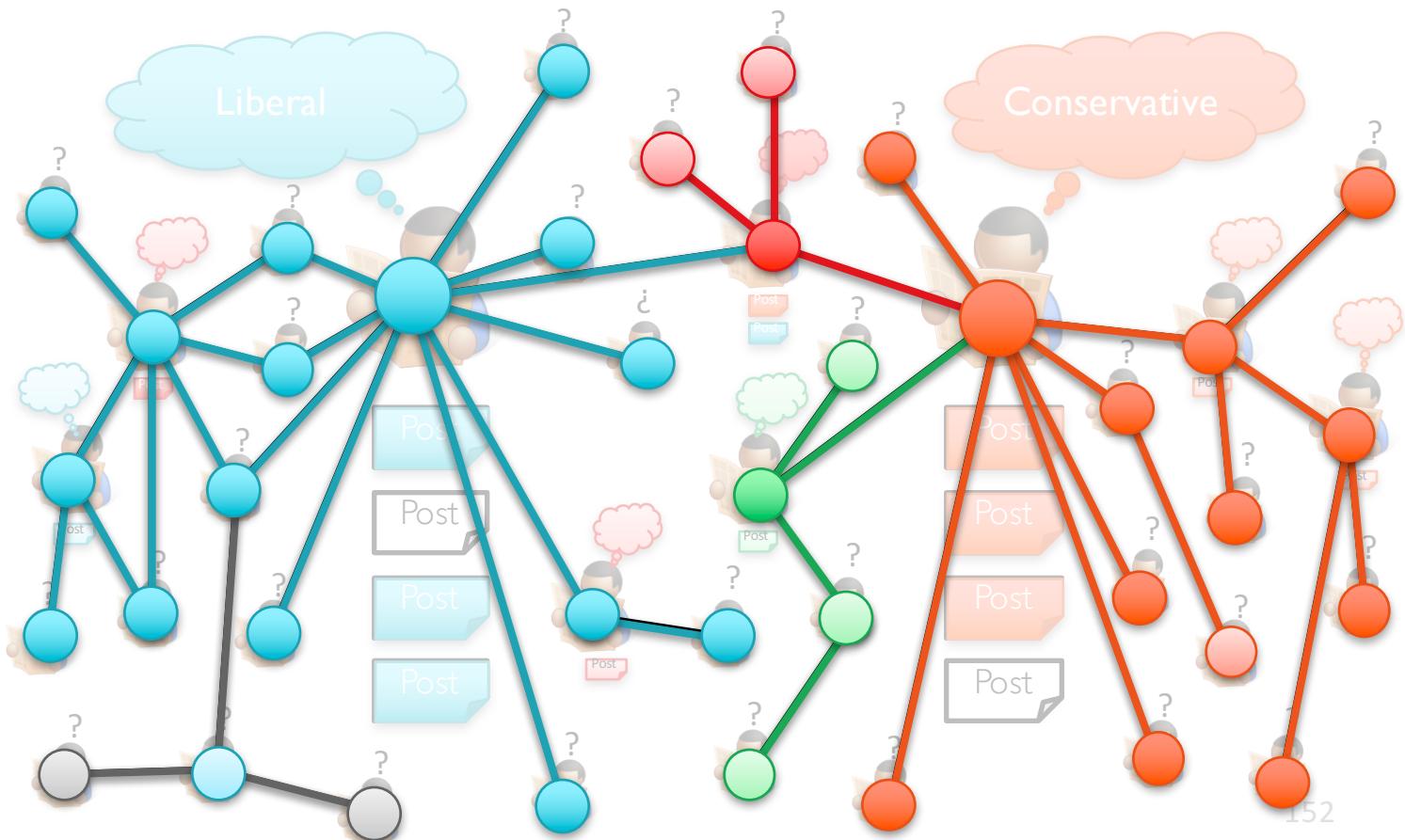


Source: Nieman Journalism Lab

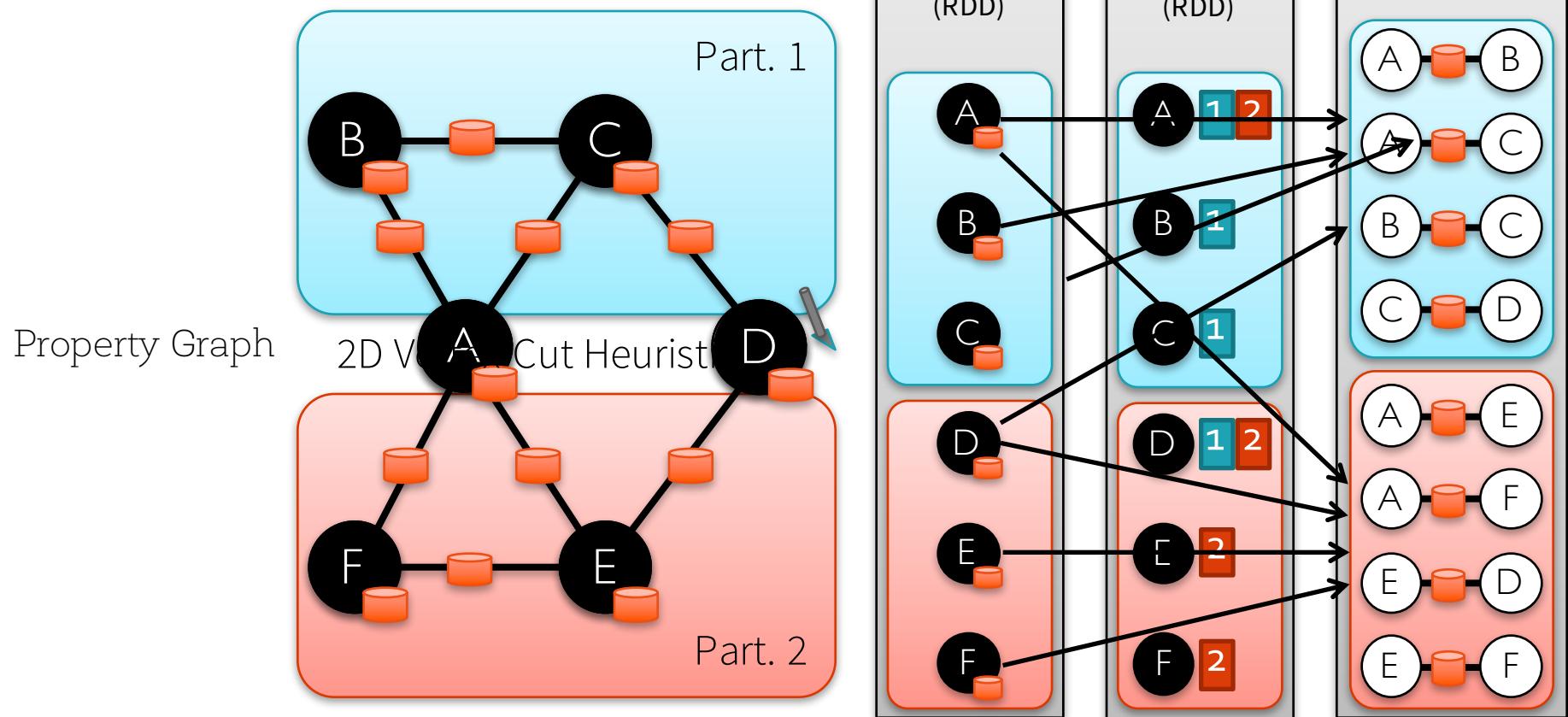


Source: Wikimedia Foundation

Belief Propagation: Predicting User Behavior



Distributed Graphs as Tables (RDDs)

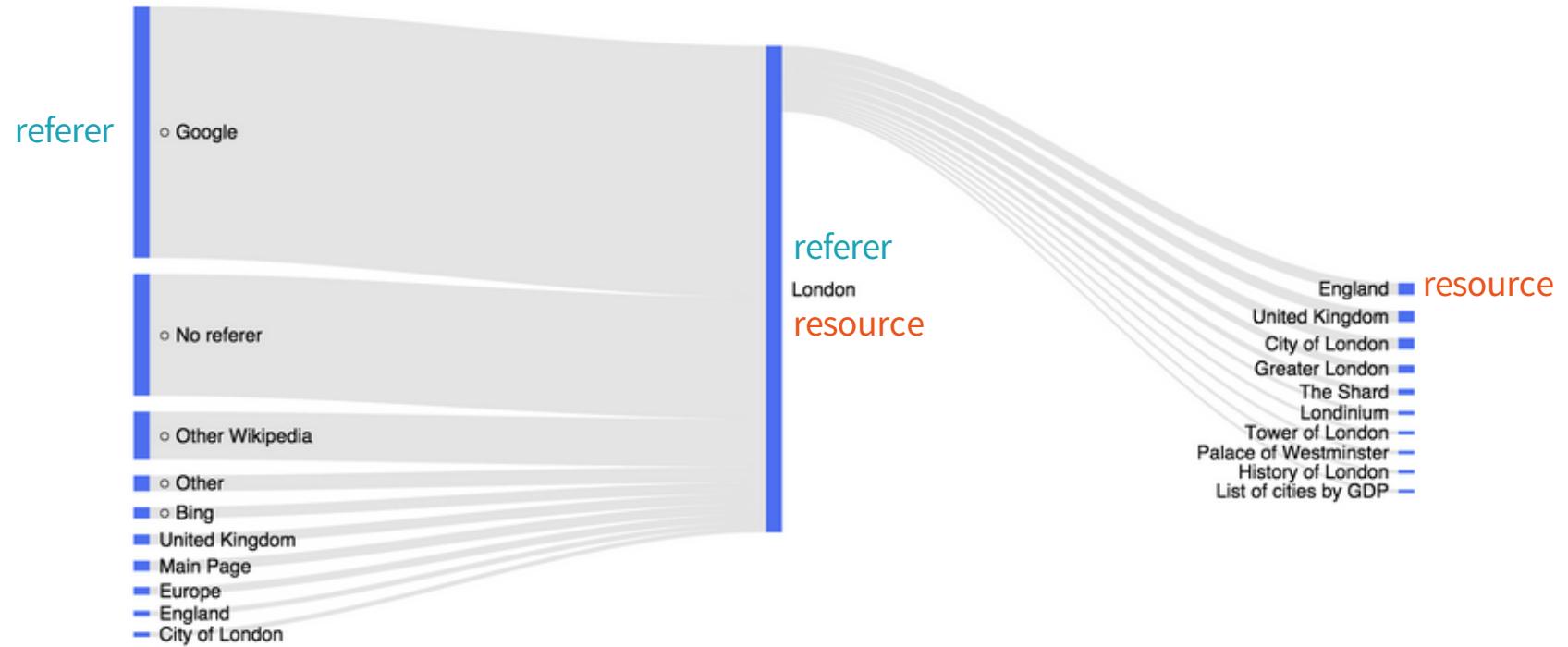




[revisited]

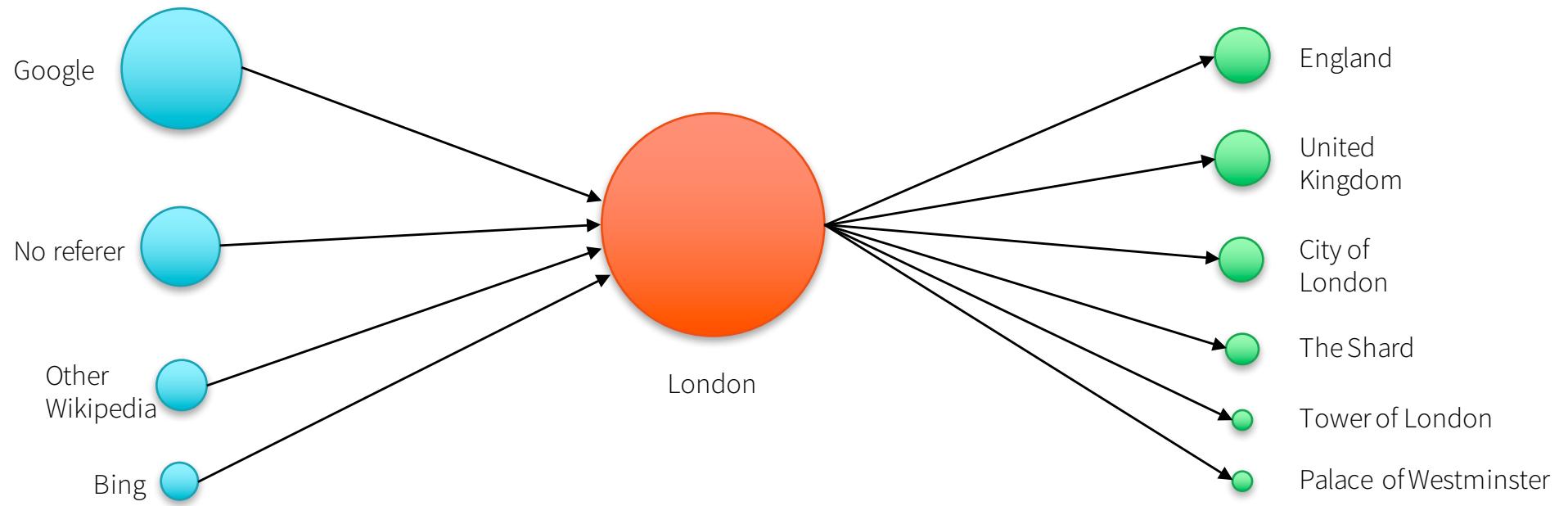
clickstream

(1.2 GB)



Wikipedia Clickstream

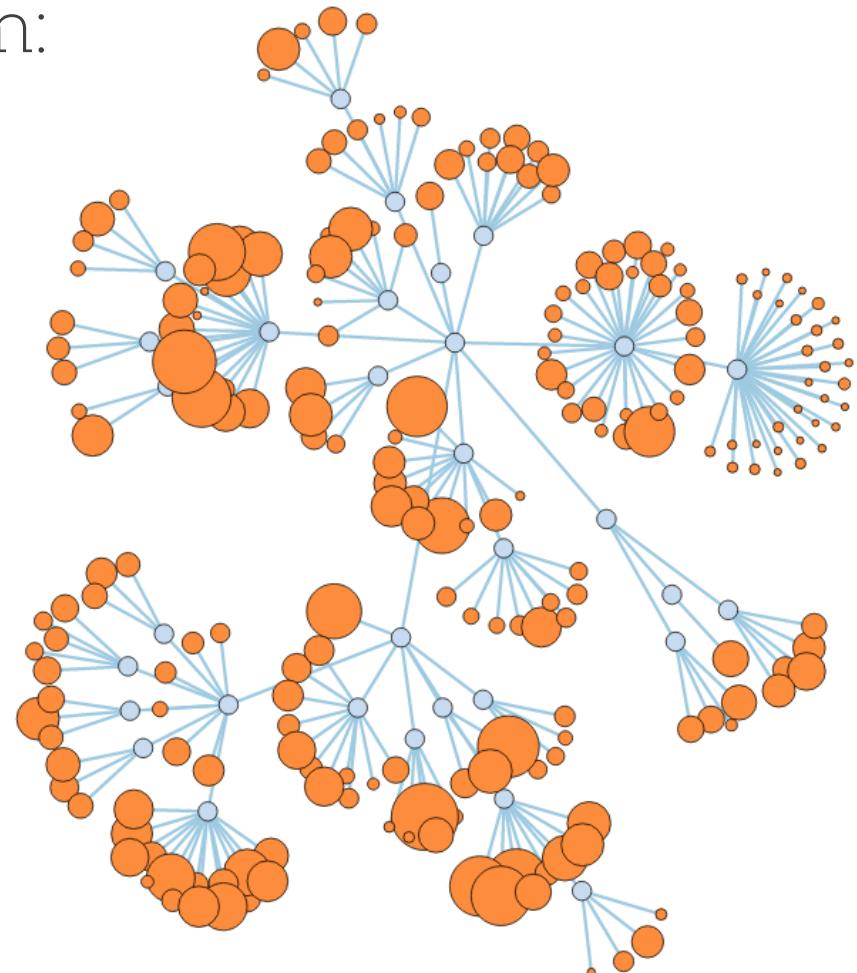
By Ellery Wulczyn
Data Scientist @ The Wikimedia Foundation



Wikipedia Clickstream

GraphX representation

Subgraph from clickstream:

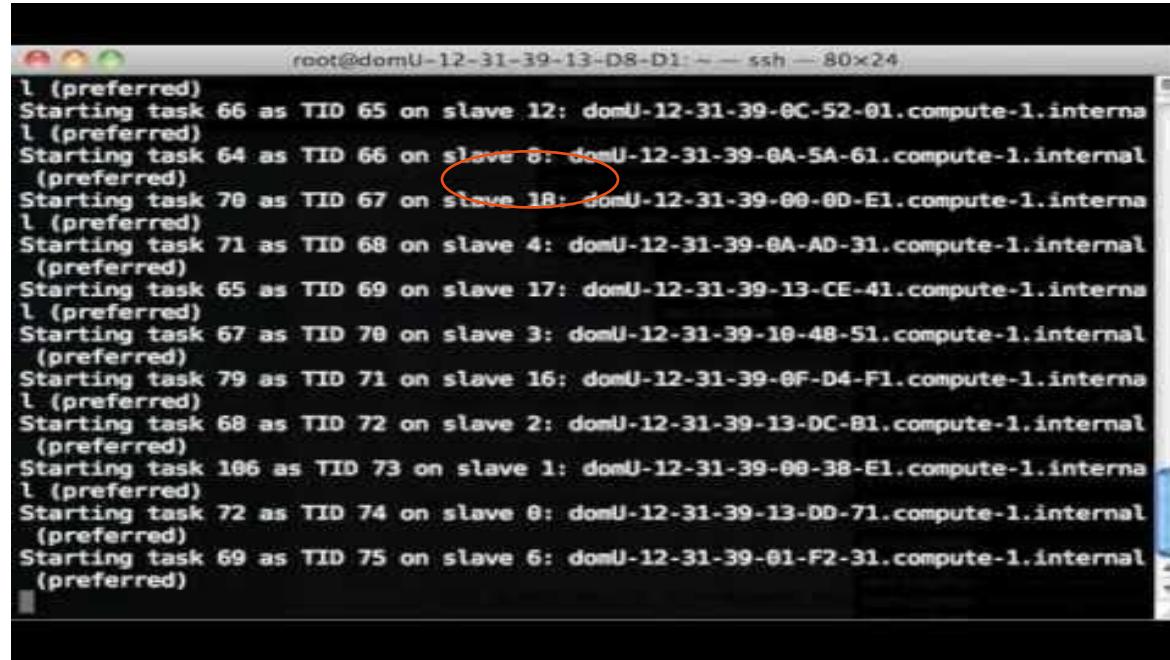




Machine Learning

Spark Demo

June 2011



```
root@domU-12-31-39-13-DB-D1: ~ — ssh — 80x24
l (preferred)
Starting task 66 as TID 65 on slave 12: domU-12-31-39-0C-52-01.compute-1.internal
l (preferred)
Starting task 64 as TID 66 on slave 8: domU-12-31-39-0A-5A-61.compute-1.internal
  (preferred)
Starting task 70 as TID 67 on slave 18: domU-12-31-39-00-00-E1.compute-1.internal
l (preferred)
Starting task 71 as TID 68 on slave 4: domU-12-31-39-0A-AD-31.compute-1.internal
  (preferred)
Starting task 65 as TID 69 on slave 17: domU-12-31-39-13-CE-41.compute-1.internal
l (preferred)
Starting task 67 as TID 70 on slave 3: domU-12-31-39-10-48-51.compute-1.internal
  (preferred)
Starting task 79 as TID 71 on slave 16: domU-12-31-39-0F-D4-F1.compute-1.internal
l (preferred)
Starting task 68 as TID 72 on slave 2: domU-12-31-39-13-DC-B1.compute-1.internal
  (preferred)
Starting task 106 as TID 73 on slave 1: domU-12-31-39-00-38-E1.compute-1.internal
l (preferred)
Starting task 72 as TID 74 on slave 0: domU-12-31-39-13-00-71.compute-1.internal
  (preferred)
Starting task 69 as TID 75 on slave 6: domU-12-31-39-01-F2-31.compute-1.internal
  (preferred)
```

<https://www.youtube.com/watch?v=Jw9yNTJI8iM>

“We start by loading some fields from a **tab-separated Wikipedia file** into a distributed collection of Article objects, then we perform queries on it. Queries on the **on-disk** data take **20 seconds**, but asking Spark to cache the Article objects in **memory** reduces the query latency to **1-2 seconds**. ”

- Matei



Machine Learning for Humans

- Make practical ML scalable and easy
- Inspired by scikit-learn

w/ Joseph Bradley @ Facebook

Algorithm coverage

Classification

- Logistic regression w/ elastic net
- Naive Bayes
- Streaming logistic regression
- Linear SVMs
- Decision trees
- Random forests
- Gradient-boosted trees
- Multilayer perceptron
- One-vs-rest

Regression

- Least squares w/ elastic net
- Isotonic regression
- Decision trees
- Random forests
- Gradient-boosted trees
- Streaming linear methods

Recommendation

- Alternating Least Squares

Frequent itemsets

- FP-growth
- Prefix span

Feature extraction & selection

- Binarizer
- Bucketizer
- Chi-Squared selection
- CountVectorizer
- Discrete cosine transform
- ElementwiseProduct
- Hashing term frequency
- Inverse document frequency
- MinMaxScaler
- Ngram
- Normalizer
- One-Hot Encoder
- PCA
- PolynomialExpansion
- RFormula
- SQLTransformer
- Standard scaler
- StopWordsRemover
- StringIndexer
- Tokenizer
- StringIndexer
- VectorAssembler
- VectorIndexer
- VectorSlicer
- Word2Vec

Clustering

- Gaussian mixture models
- K-Means
- Streaming K-Means
- Latent Dirichlet Allocation
- Power Iteration Clustering

Statistics

- Pearson correlation
- Spearman correlation
- Online summarization
- Chi-squared test
- Kernel density estimation

Linear algebra

- Local dense & sparse vectors & matrices
- Distributed matrices
 - Block-partitioned matrix
 - Row matrix
 - Indexed row matrix
 - Coordinate matrix
- Matrix decompositions

Model import/export

Pipelines

List based on Spark 1.5

High-level functionality

Learning tasks

- Classification
- Regression
- Recommendation
- Clustering
- Frequent itemsets

Data utilities

- Feature extraction & selection
- Statistics
- Linear algebra

Workflow utilities

- Model import/export
- Pipelines
- DataFrames
- Cross validation

Spark MLlib Components

`spark.mllib`
(package)

vs.

`spark.ml`
(package)

- Original “lower-level” API
- Built on top of RDDs
- Will go to maintenance mode starting Spark 2.0

- Newer “higher-level” API for constructing workflows
- Built on top of DataFrames

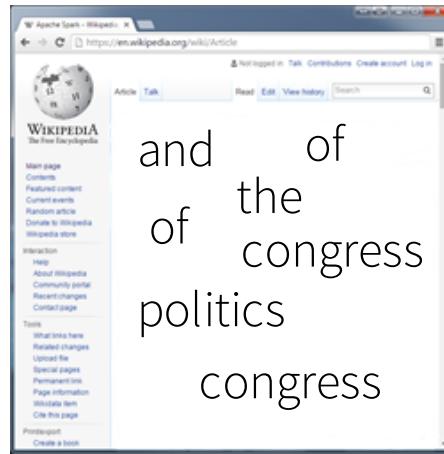
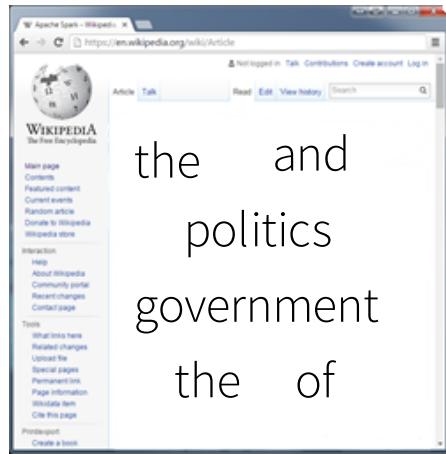
Term Frequency – Inverse Document Frequency

An information retrieval algorithm used to rank how important a word is to a collection of documents

TF: If a word appears frequently in a doc, it's important

IDF: But if a word appears in many docs (the, and, of), the word is not meaningful, so lower its score

TF-IDF

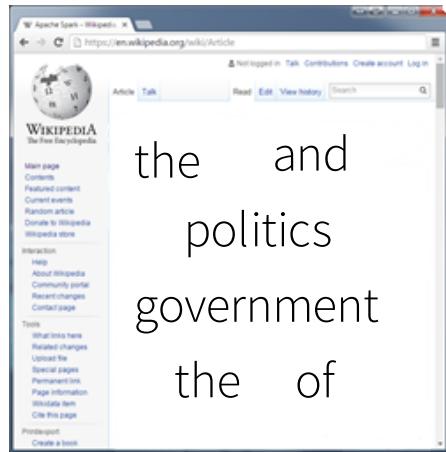


TF:

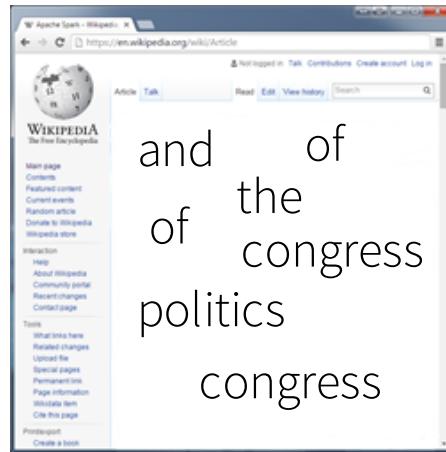
of times word X appears in a document

Total # of words in the document

TF-IDF



the: 2
and: 1
politics: 1
government: 1
of: 1

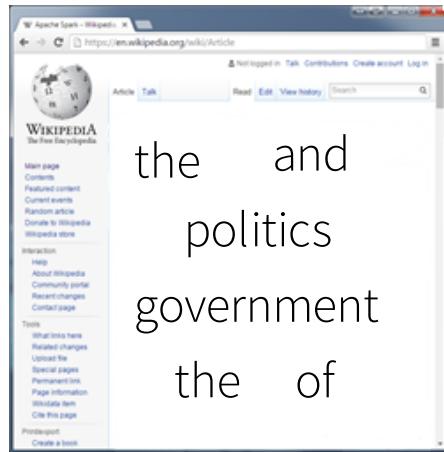


of: 2
congress: 2
and: 1
the: 1
politics: 1



sports: 2
the: 1
and: 1
games: 1
olympics: 1

TF-IDF



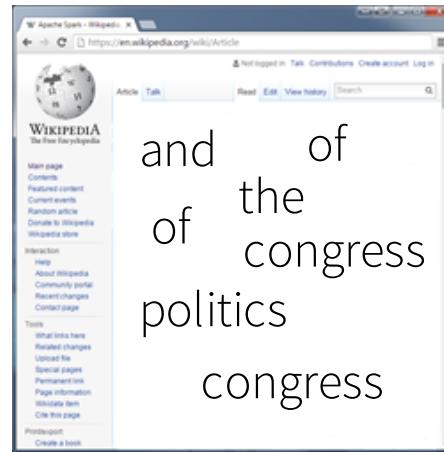
the: $2 / 6 = 0.3$

and: $1 / 6 = 0.16$

politics: $1 / 6 = 0.16$

government: $1 / 6 = 0.16$

of: $1 / 6 = 0.16$



of: $2 / 7 = 0.28$

congress: $2 / 7 = 0.28$

and: $1 / 7 = 0.14$

the: $1 / 7 = 0.14$

politics: $1 / 7 = 0.14$



sports: $2 / 6 = 0.3$

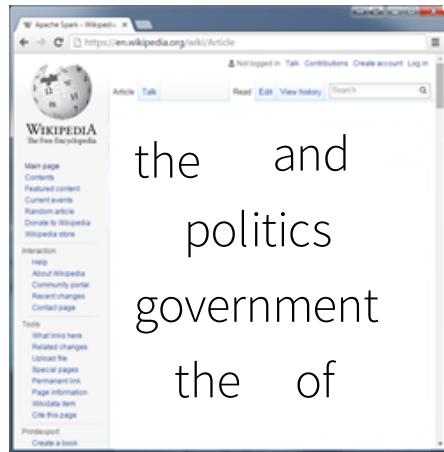
the: $1 / 6 = 0.16$

and: $1 / 6 = 0.16$

games: $1 / 6 = 0.16$

olympics: $1 / 6 = 0.16$

TF-IDF



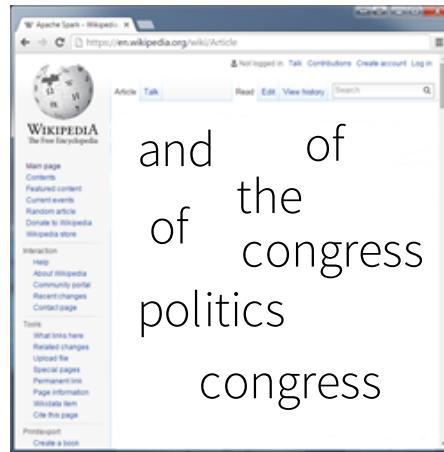
the: $2 / 6 = 0.3$

and: $1 / 6 = 0.16$

politics: $1 / 6 = 0.16$

government: $1 / 6 = 0.16$

of: $1 / 6 = 0.16$



of: $2 / 7 = 0.28$

congress: $2 / 7 = 0.28$

and: $1 / 7 = 0.14$

the: $1 / 7 = 0.14$

politics: $1 / 7 = 0.14$



sports: $2 / 6 = 0.3$

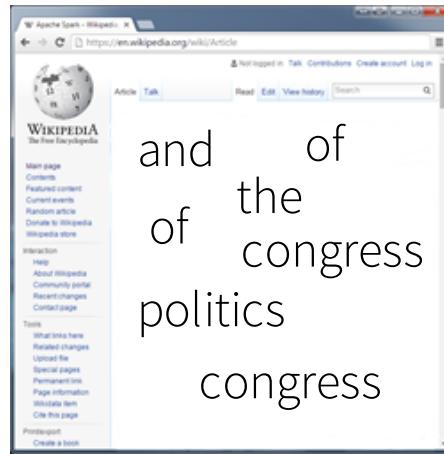
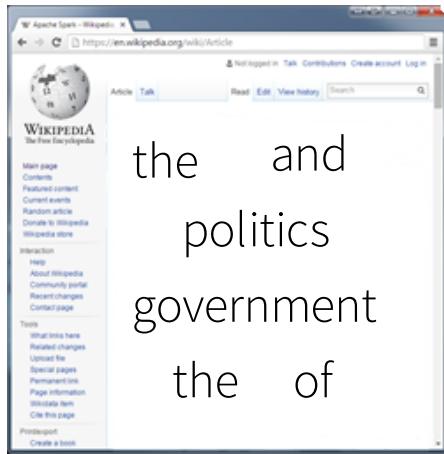
the: $1 / 6 = 0.16$

and: $1 / 6 = 0.16$

games: $1 / 6 = 0.16$

olympics: $1 / 6 = 0.16$

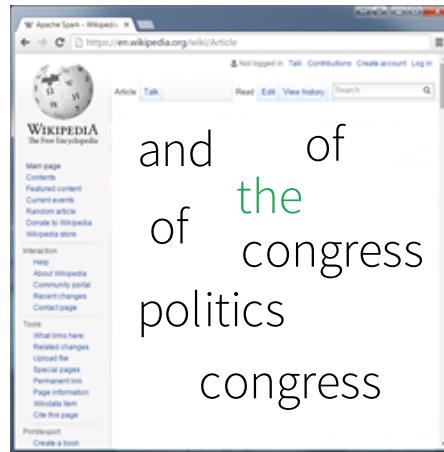
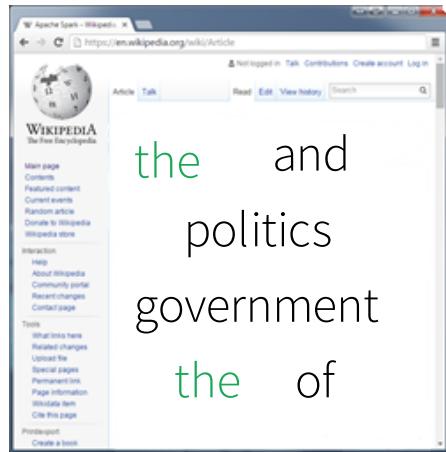
TF-IDF



IDE:

$$\log \left[\frac{\text{Total # of documents}}{\text{\# of documents with word X in it}} \right]$$

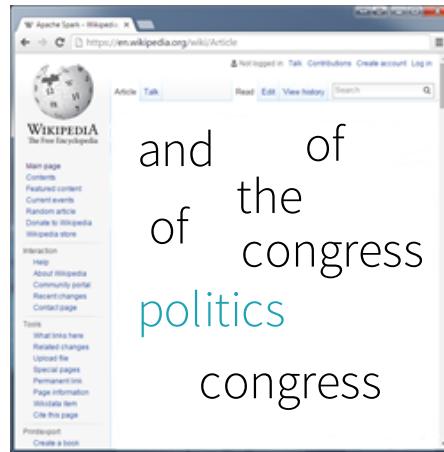
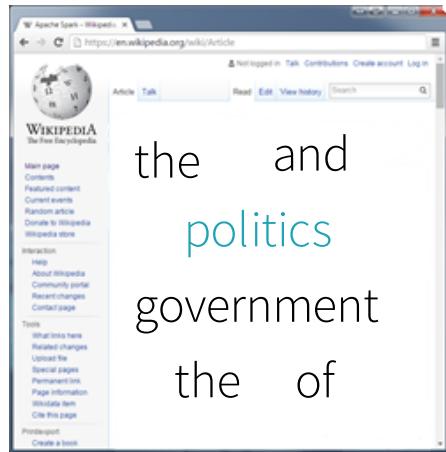
TF-IDF



IDE:

$$\log \left[\frac{3}{3} \right] = 0$$

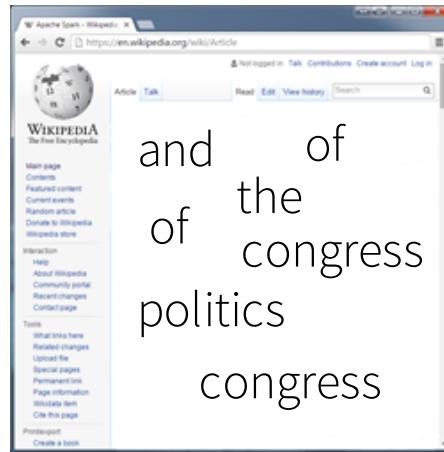
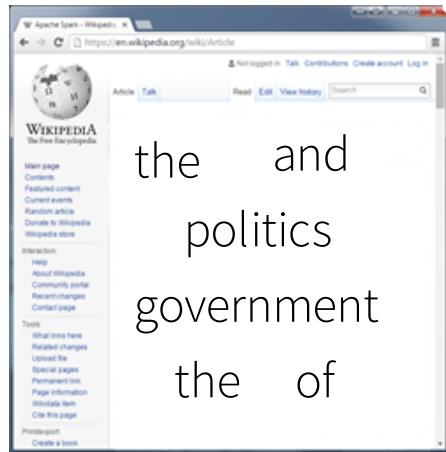
TF-IDF



IDF:

$$\log \left[\frac{3}{2} \right] = 0.18$$

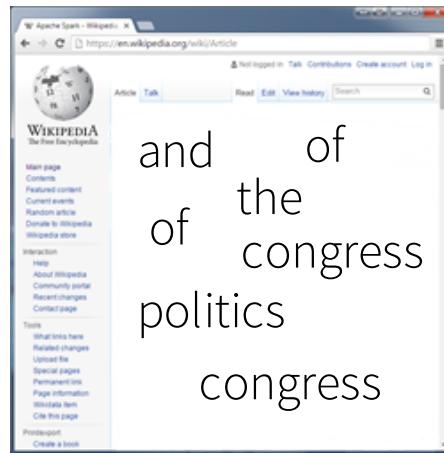
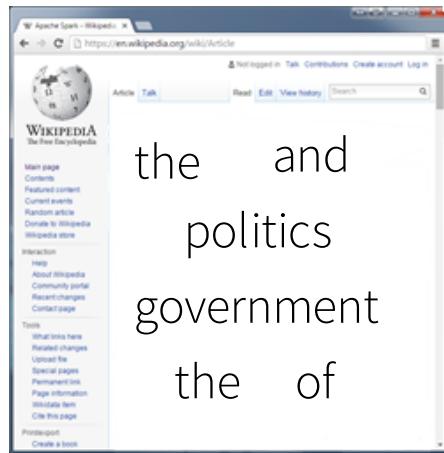
TF-IDF



IDF:

$$\log \left[\frac{3}{1} \right] = 0.48$$

TF-IDF



$$\text{government: } 0.16 * 0.48 = 0.08$$

$$\text{politics: } 0.16 * 0.18 = 0.03$$

$$\text{of: } 0.16 * 0.18 = 0.03$$

$$\text{the: } 0.3 * 0 = 0$$

$$\text{and: } 0.16 * 0 = 0$$

$$\text{congress: } 0.28 * 0.48 = 0.08$$

$$\text{of: } 0.28 * 0.18 = 0.05$$

$$\text{politics: } 0.14 * 0.18 = 0.03$$

$$\text{the: } 0.14 * 0 = 0$$

$$\text{and: } 0.14 * 0 = 0$$

$$\text{sports: } 0.3 * 0.48 = 0.144$$

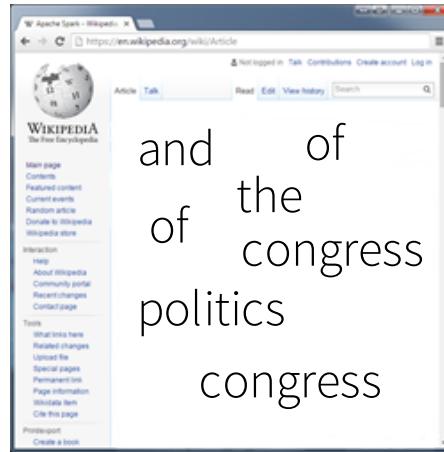
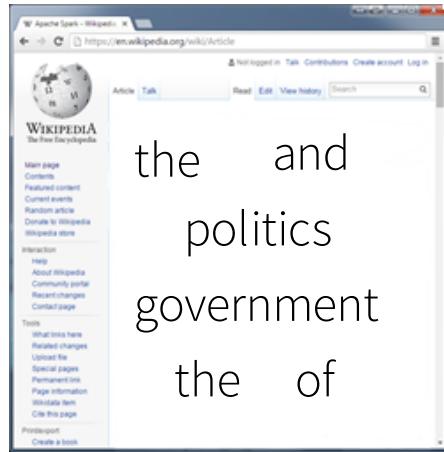
$$\text{olympics: } 0.16 * 0.48 = 0.08$$

$$\text{games: } 0.16 * 0.48 = 0.08$$

$$\text{the: } 0.16 * 0 = 0$$

$$\text{and: } 0.16 * 0 = 0$$

TF-IDF

$$\text{government: } 0.16 * 0.48 = 0.08$$

$$\text{politics: } 0.16 * 0.18 = 0.03$$

$$\text{of: } 0.16 * 0.18 = 0.03$$

$$\text{congress: } 0.28 * 0.48 = 0.08$$

$$\text{of: } 0.28 * 0.18 = 0.05$$

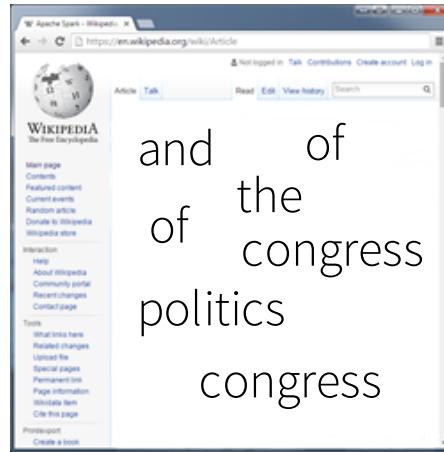
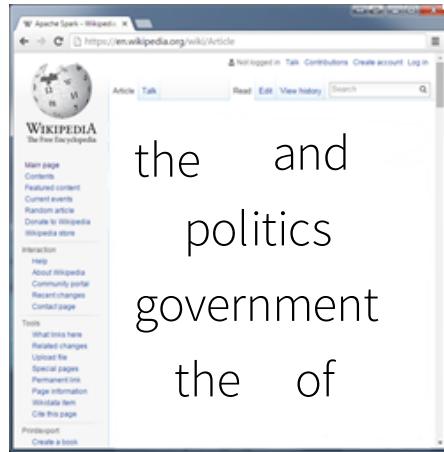
$$\text{politics: } 0.14 * 0.18 = 0.03$$

$$\text{sports: } 0.3 * 0.48 = 0.144$$

$$\text{olympics: } 0.16 * 0.48 = 0.08$$

$$\text{games: } 0.16 * 0.48 = 0.08$$

TF-IDF

 Search

$$\text{government: } 0.16 * 0.48 = 0.08$$

$$\text{politics: } 0.16 * 0.18 = 0.03$$

$$\text{of: } 0.16 * 0.18 = 0.03$$

$$\text{congress: } 0.28 * 0.48 = 0.08$$

$$\text{of: } 0.28 * 0.18 = 0.05$$

$$\text{politics: } 0.14 * 0.18 = 0.03$$

$$\text{sports: } 0.3 * 0.48 = 0.144$$

$$\text{olympics: } 0.16 * 0.48 = 0.08$$

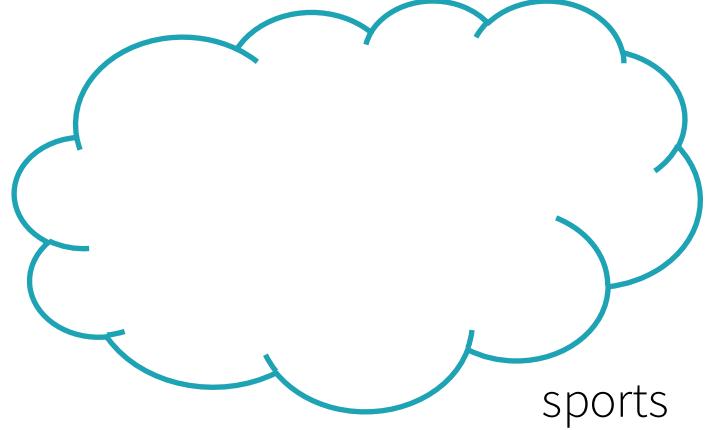
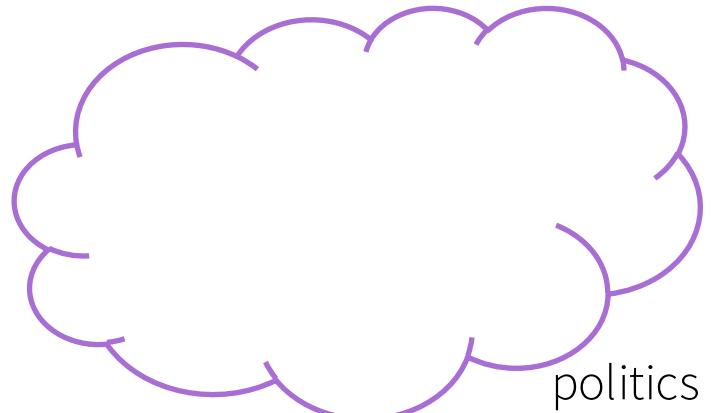
$$\text{games: } 0.16 * 0.48 = 0.08$$

K-means clustering

An unsupervised ML algorithm for classifying items into k groups.

k : The # of pre-chosen groups

K-means clustering



K-means clustering

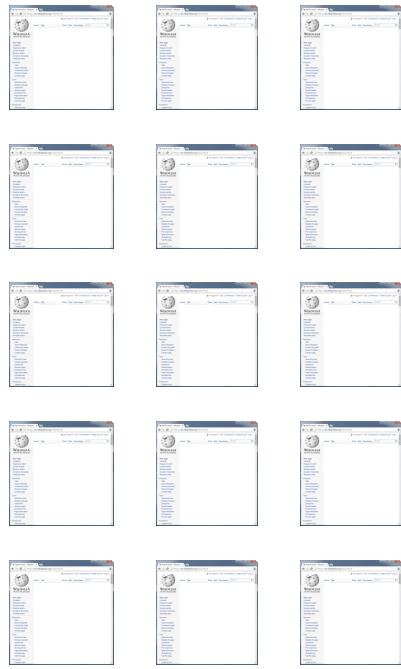


politics

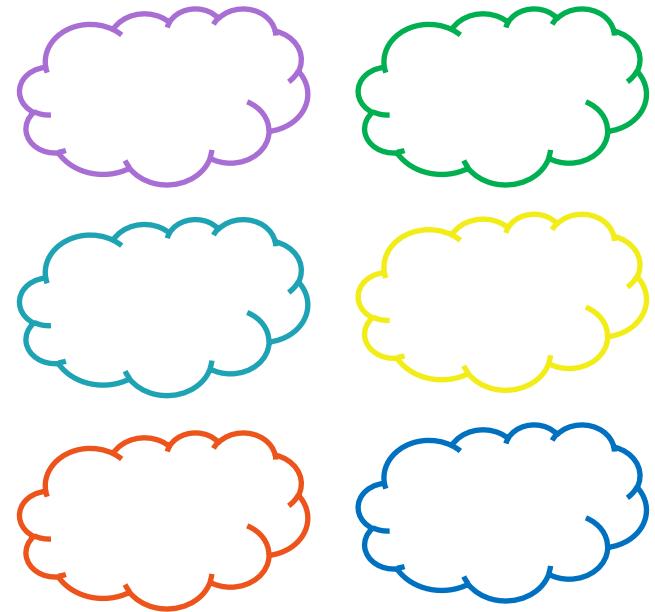


sports

K-means clustering

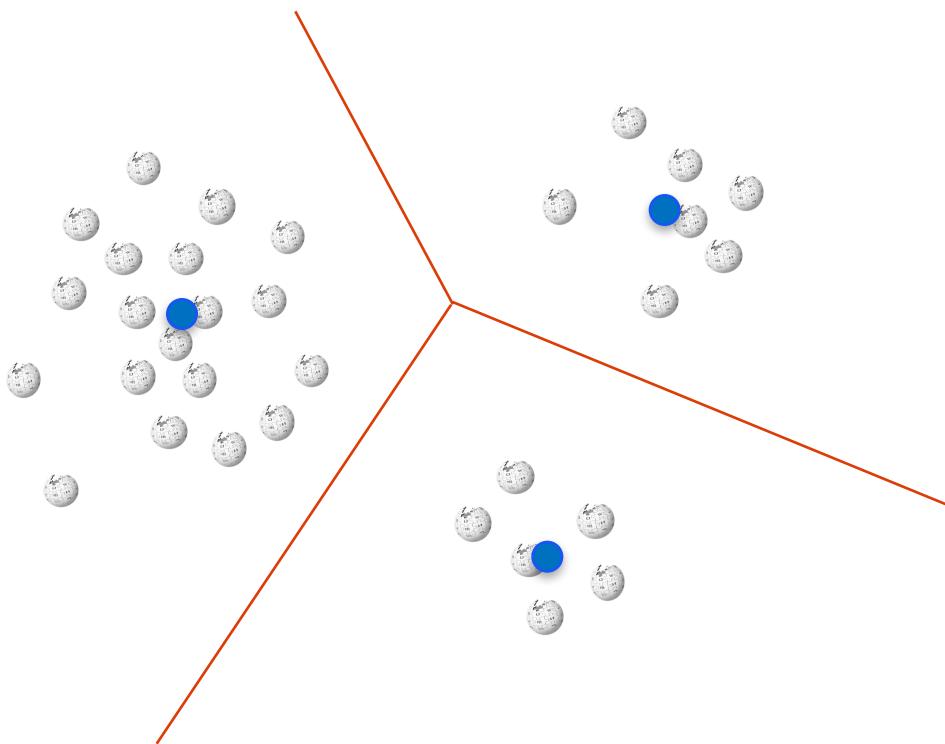


5 million articles



100 clusters

K-means clustering

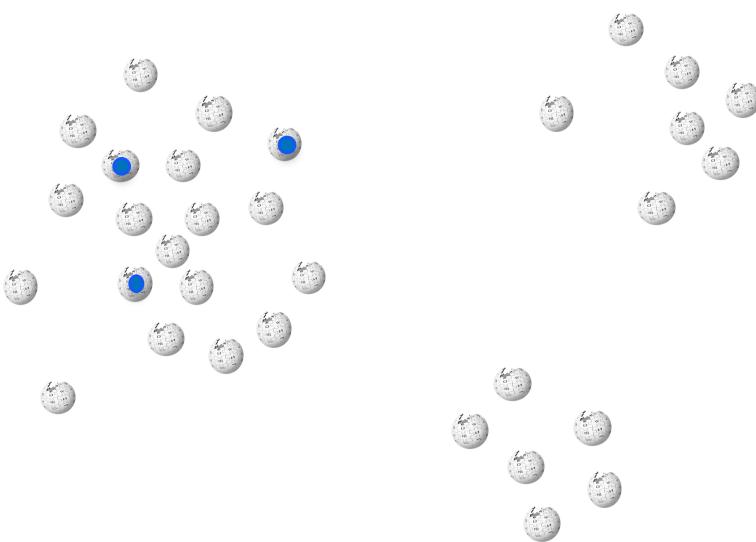


black dots: Wikipedia Articles (items)

red lines: Cluster Partitions

blue dots: Cluster Centroids

K-means clustering: Initialization



Approach #1:

Randomly select centroids to start

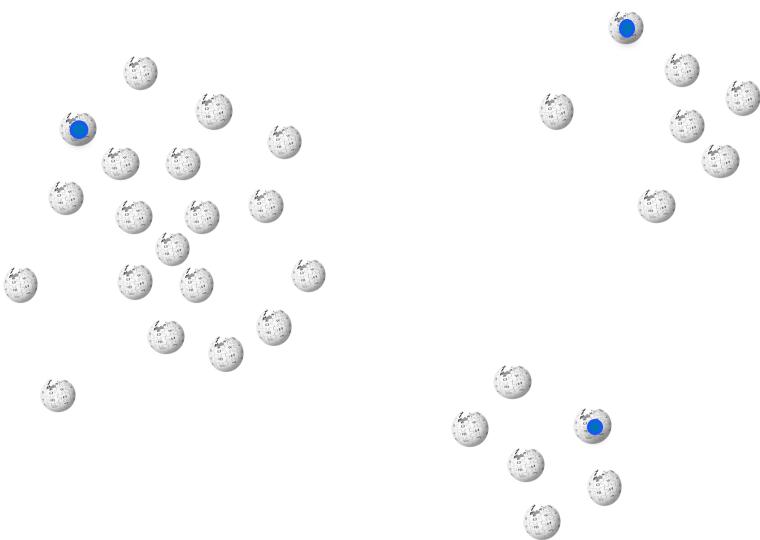
Pro:

Does not require full iteration

Con:

If a cluster is large compared to neighbors, initial centroids will likely end up in large cluster.

K-means clustering: Initialization



Approach #2: k-means++

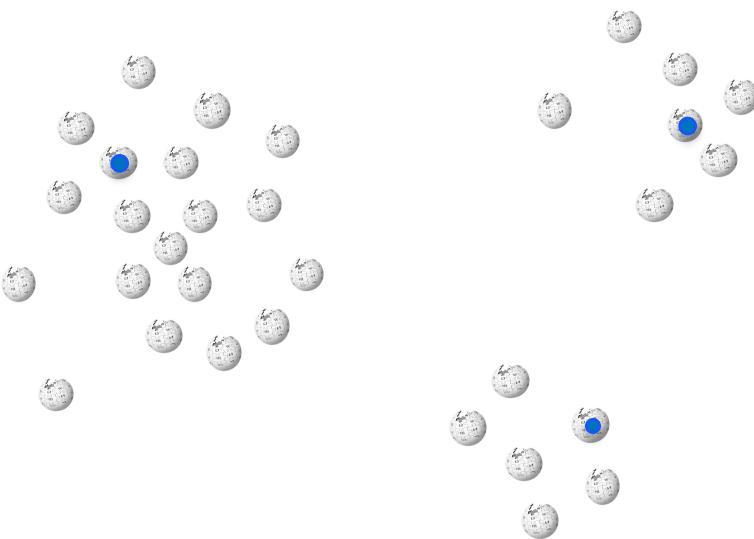
Pro:

Avoids the poor clusterings found by standard algorithm

Con:

Needs k passes over the data, so does not scale well to large data sets

K-means clustering: Initialization



Approach #3: k-means||

Pro:

Also avoids the poor clusterings found by standard algorithm

Can select multiple candidate centroids in first iteration

K-means clustering: Initialization

(expert-only) Parameters

A list of advanced, expert-only (hyper-)parameter keys this algorithm can take. Users can set and get the parameter values through setters and getters, respectively.

▼ **final val initMode: Param[String]**

Param for the initialization algorithm. This can be either "random" to choose random points as initial cluster centers, or "k-means||" to use a parallel variant of k-means++ (Bahmani et al., Scalable K-Means++, VLDB 2012). Default: k-means||.

Definition Classes KMeansParams

Annotations @Since("1.5.0")

► **final val initSteps: IntParam**

Param for the number of steps for the k-means|| initialization mode.

Source: <https://spark.apache.org/docs/1.6.0/api/scala/index.html#org.apache.spark.ml.clustering.KMeans>

K-means clustering

$k: 2$

Article	Word X	Word Y
A	1.0	1.0
B	1.5	2.0
C	3.0	4.0
D	5.0	7.0
E	3.5	5.0
F	4.5	5.0
G	3.5	4.5

Source: <http://mnemstudio.org>

K-means clustering

k : 2

	Article	Mean Vector (<i>centroid</i>)
Group 1	A	(1.0, 1.0)
Group 2	D	(5.0, 7.0)

K-means clustering

k : 2

	Article	Mean Vector (<i>centroid</i>)
Cluster 1	A, B, C	(1.8, 2.3)
Cluster 2	D, E, F, G	(4.1, 5.4)

K-means clustering

$k: 2$

Article	Distance to mean (<i>centroid</i>) of Cluster 1	Distance to mean (<i>centroid</i>) of Cluster 2
Cluster 1	A	1.5
	B	0.4
	C	2.1
Cluster 2	D	5.7
	E	3.2
	F	3.8
	G	2.8

K-means clustering

k : 2

	Article	Mean Vector (<i>centroid</i>)
Cluster 1	A, B	(1.3, 1.5)
Cluster 2	C, D, E, F, G	(3.9, 5.1)

K-means clustering

```
204  /**
205  * :: Experimental ::
206  * K-means clustering with support for k-means|| initialization proposed by Bahmani et al.
207  *
208  * @see [[http://dx.doi.org/10.14778/2180912.2180915 Bahmani et al., Scalable k-means++.]]
209  */
210 @Since("1.5.0")
211 @Experimental
212 class KMeans @Since("1.5.0") (
213     @Since("1.5.0") override val uid: String
214     extends Estimator[KMeansModel] with KMeansParams with DefaultParamsWritable {
215
216     setDefault(
217         k -> 2,
218         maxIter -> 20,
219         initMode -> MLlibKMeans.K_MEANS_PARALLEL,
220         initSteps -> 5,
221         tol -> 1e-4)
222
223     @Since("1.5.0")
224     override def copy(extra: ParamMap): KMeans = defaultCopy(extra)
225
226     @Since("1.5.0")
227     def this() = this(Identifiable.randomUUID("kmeans"))
228
229     /** @group setParam */
230     @Since("1.5.0")
231     def setFeaturesCol(value: String): this.type = set(featuresCol, value)
```

Source: <https://github.com/apache/spark/blob/master/mllib/src/main/scala/org/apache/spark/ml/clustering/KMeans.scala#L216>

K-means clustering

$k: 2$

	Cluster 1	Cluster 2
	Dist to (1, 1)	Dist to (5, 7)
A	(1, 1)	
B		D
C		D
D		D, E (5, 7)
		D, E, F
mean	(1.1, 1.2)	(1, 2)



Spark ML Pipelines

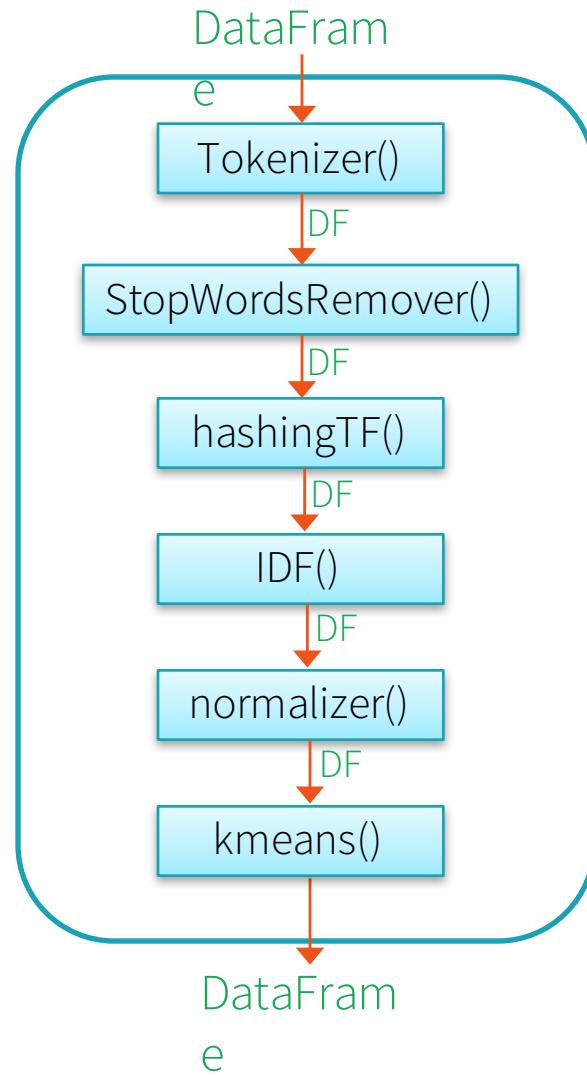
DataFrame: uses DF from Spark SQL as a ML dataset. Different columns can store text, feature vectors, true labels and predictions

Transformer: an algorithm which can transform one DataFrame into another DataFrame
(example: ML model is a transformer that transforms a DF with features into a DF with predictions)

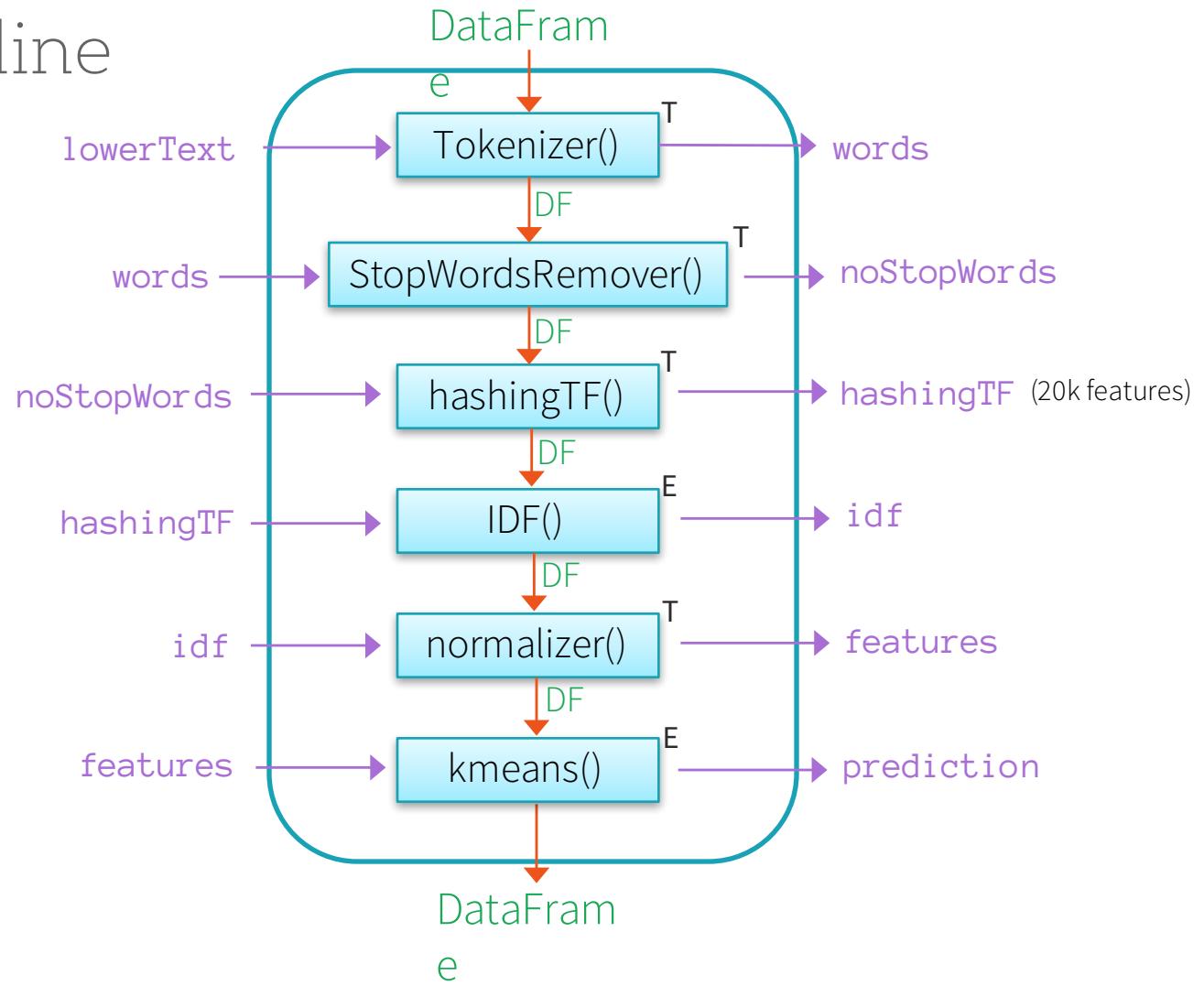
Estimator: an algorithm which can be fit on a DF to produce a Model
(example: a learning algorithm is an Estimator which trains on a DF and produces a model)

Pipeline: chains multiple Transformers and Estimators together to specify a ML workflow

Spark ML Pipeline



Spark ML Pipeline





Chat w/ Vida

Compression types

What are the benefits/disadvantages of various compression types?

- Non-splittable: gzip, zip
- Splittable: Lzo, lz4, bzip2, snappy
- Columnar: Parquet



Advantages & Pitfalls of File Types



Data Storage Options

- SQL Databases/Data Warehouses



- NoSQL Stores



- Specialized Engines



Want to Learn More?

Not your father's database:

How to use Apache Spark properly in your big data
architecture

Wednesday, 4:20pm–5:00pm

Enterprise Adoption Track





Spark Streaming

Motivation

Stability condition for any streaming app

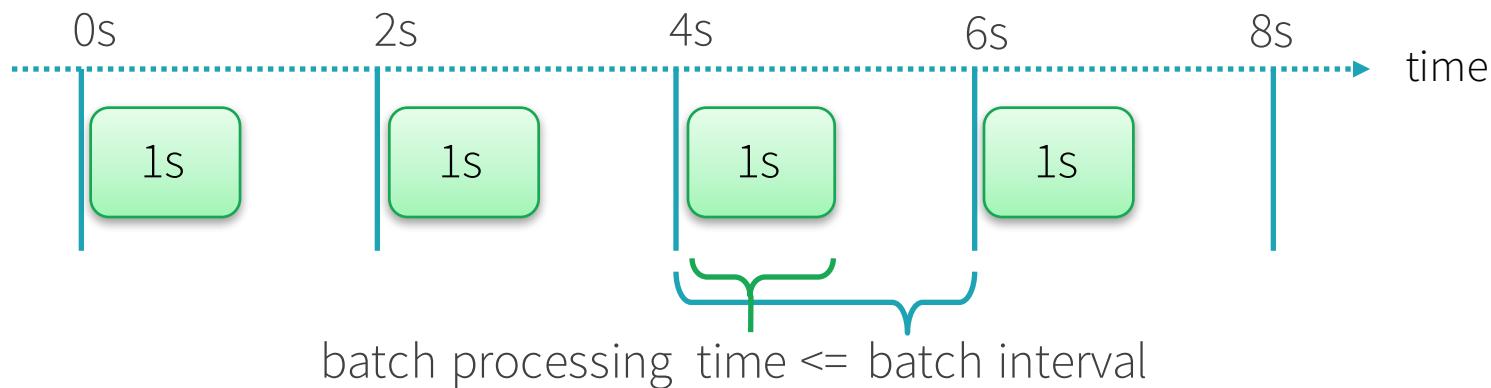
Receive data only as fast as the system can process it

Stability condition for Spark Streaming's “micro-batch” model

Finish processing previous batch before next one arrives

Stable micro-batch operation

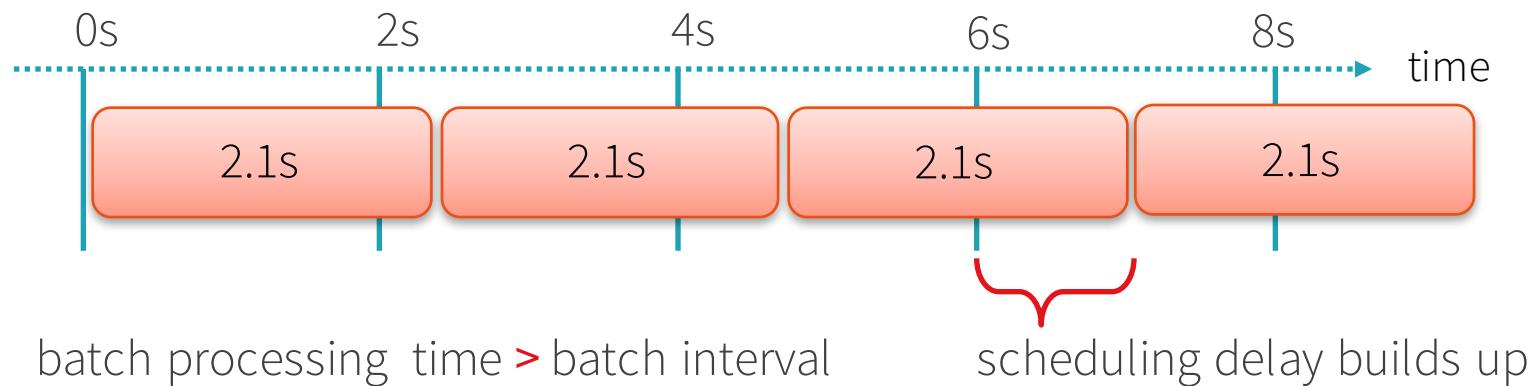
Spark Streaming runs micro-batches at fixed *batch intervals*



Previous batch is processed before next one arrives => **stable**

Unstable micro-batch operation

Spark Streaming runs micro-batches at fixed *batch intervals*



Batches continuously gets delayed and backlogged => **unstable**

6 Edit Streams:



English : en



52.89.53.194:9002



German : de



54.68.10.240:9003



French : fr



54.68.10.240:9004



Russian : ru



54.68.10.240:9005



Italian : it



54.68.10.240:9006



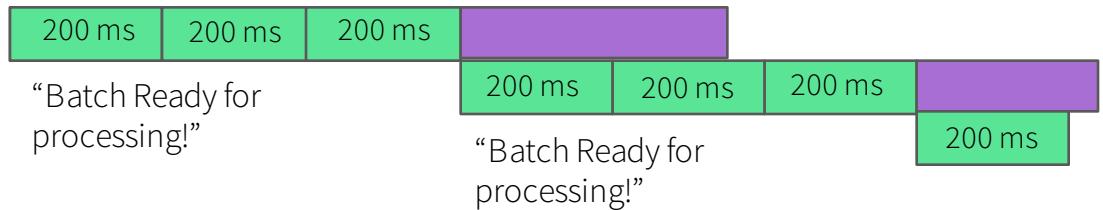
Spanish : es



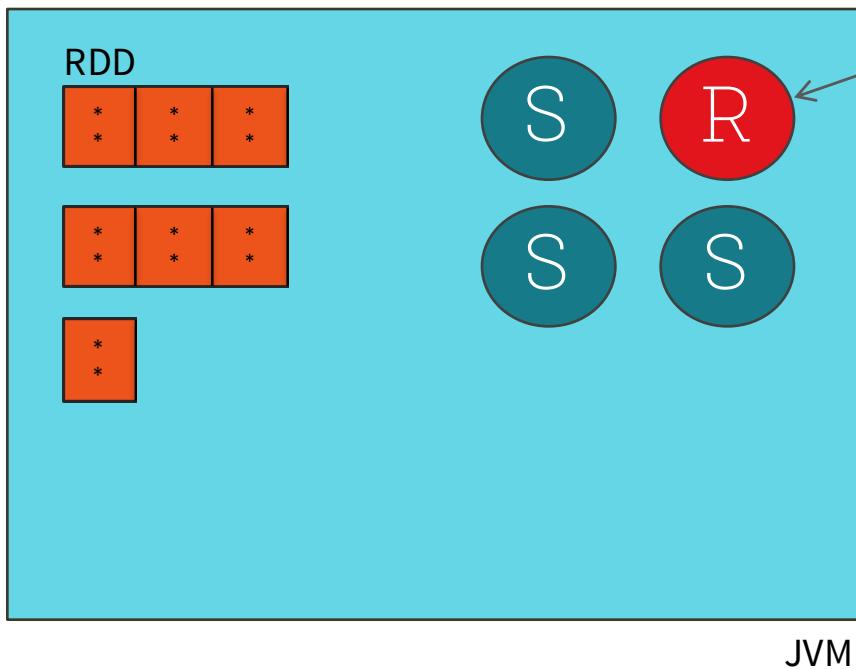
54.68.10.240:9007



Batch Interval: 600 ms



English Edits



Batch Interval: 600 ms

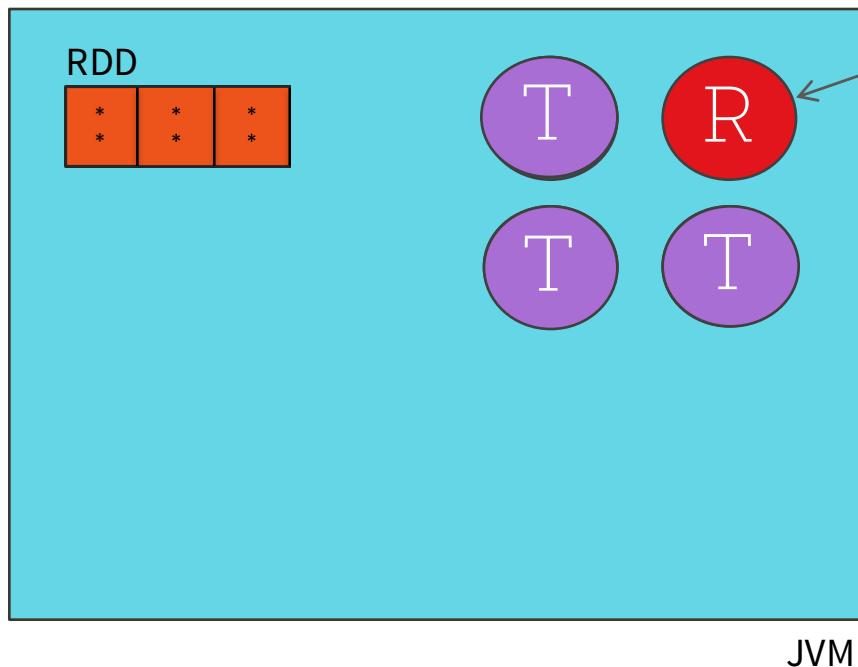


"Batch Ready for
processing!"

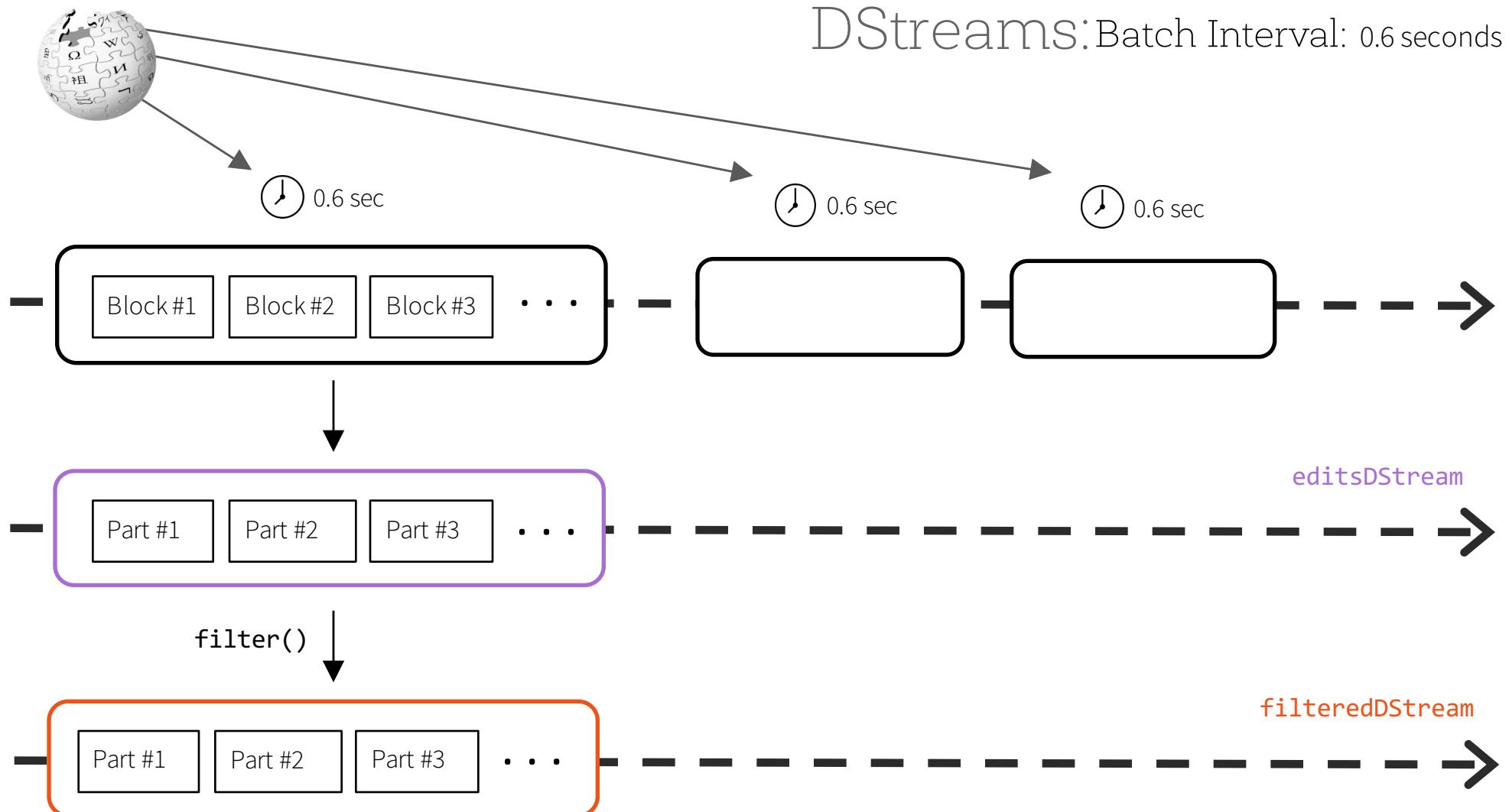


:9002

English Edits



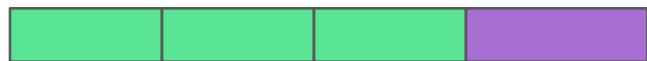
DStreams: Batch Interval: 0.6 seconds



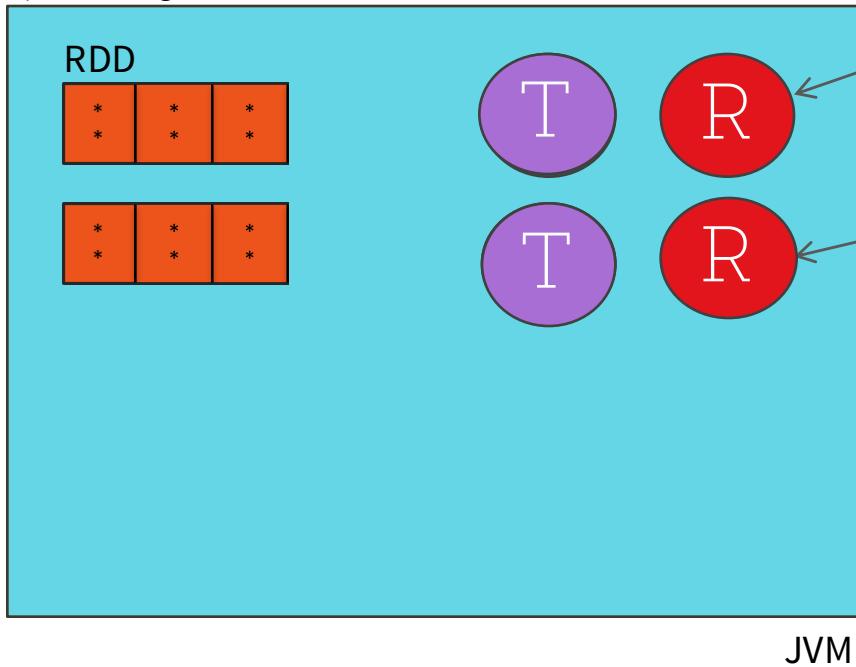
Batch Interval: **3 sec**



“English Batch Ready for processing!”



“Spanish Batch Ready for processing!”



English Edits

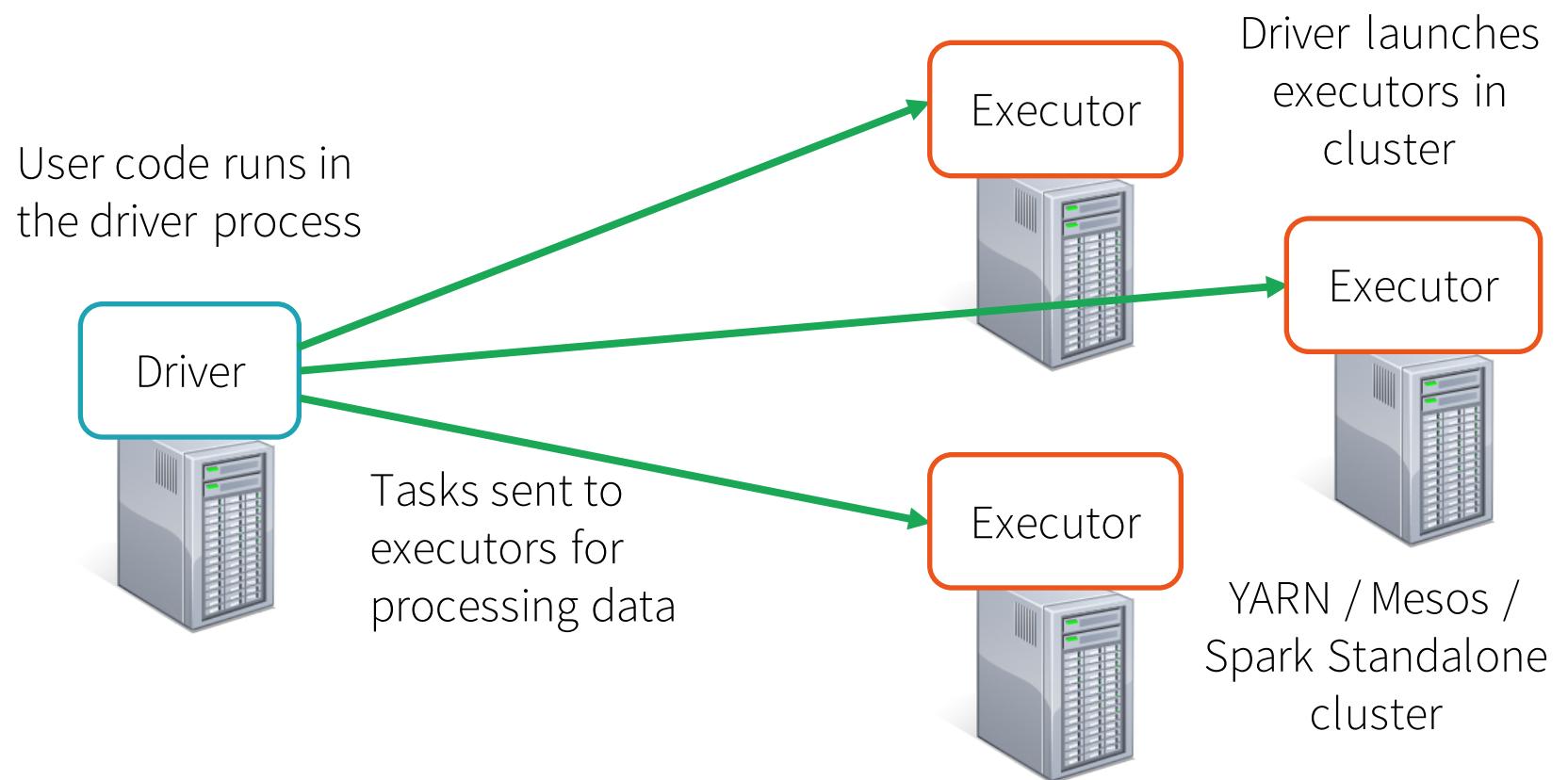


Spanish Edits

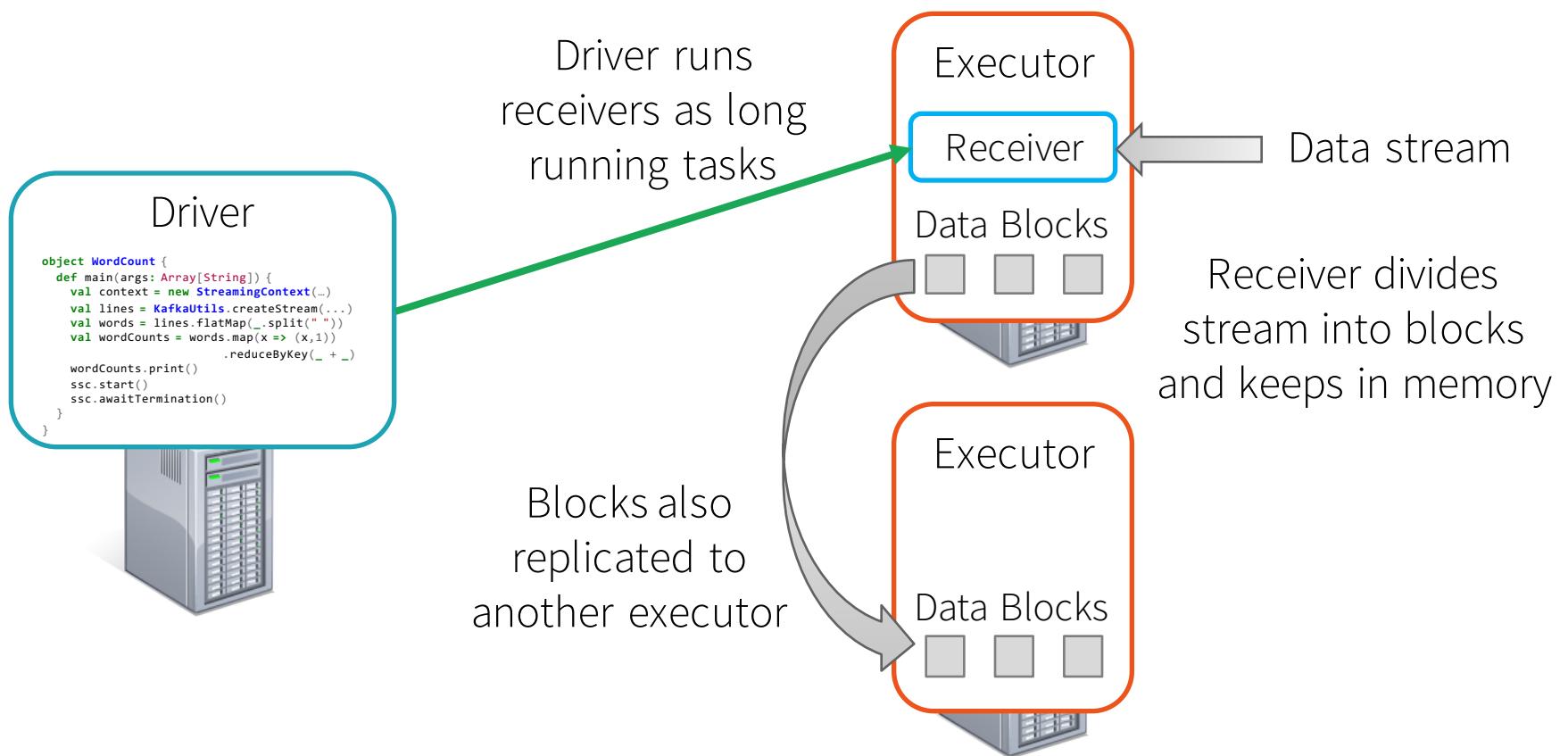
Transformations on DStreams:

map()	filter()	reduce()
flatMap()	repartition()	updateStateByKey()
join()	union()	countByValue()
cogroup()	ReduceByKey()	transform()
count()		

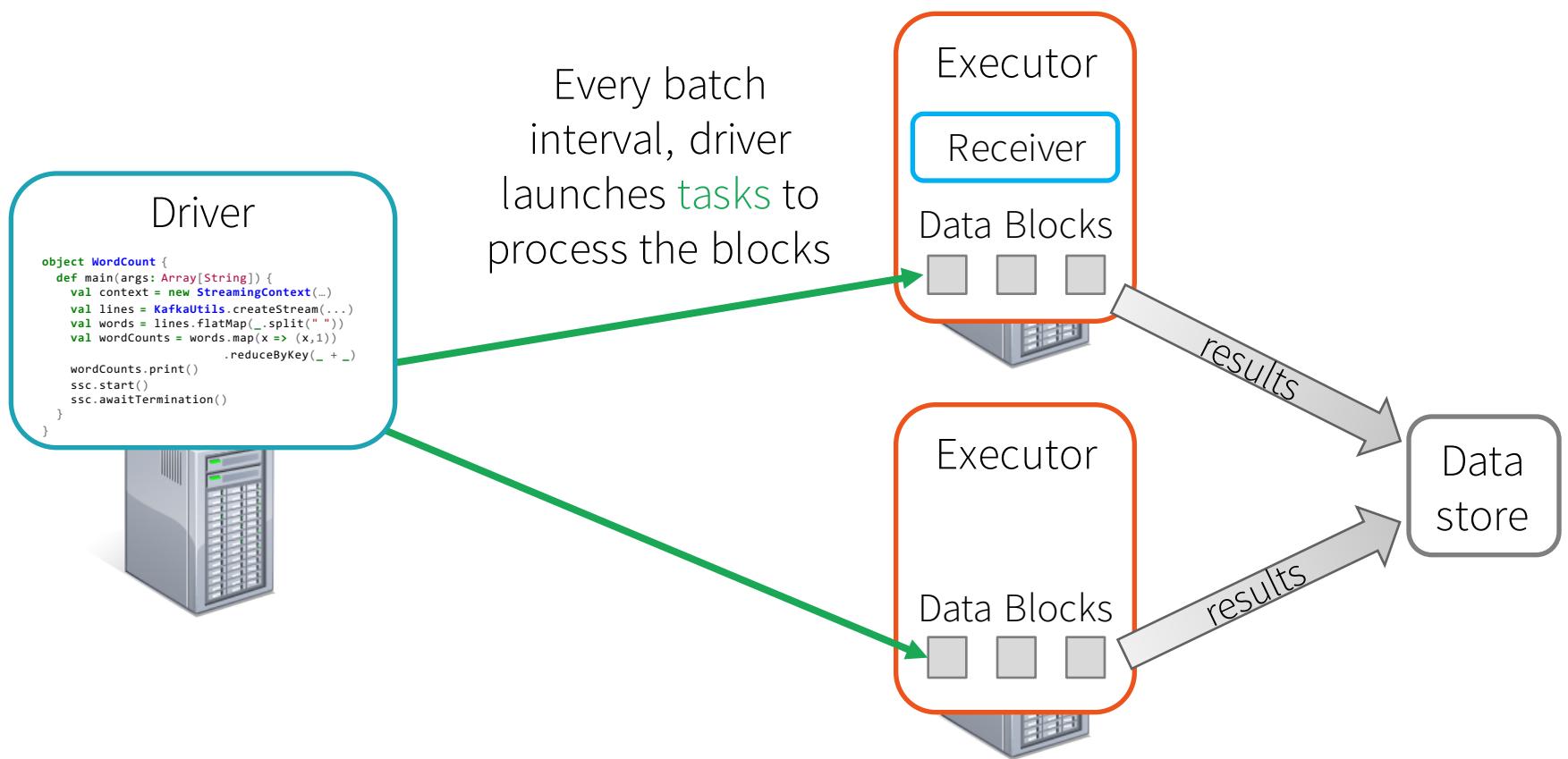
Any Spark Application



Spark Streaming Application: Receive data

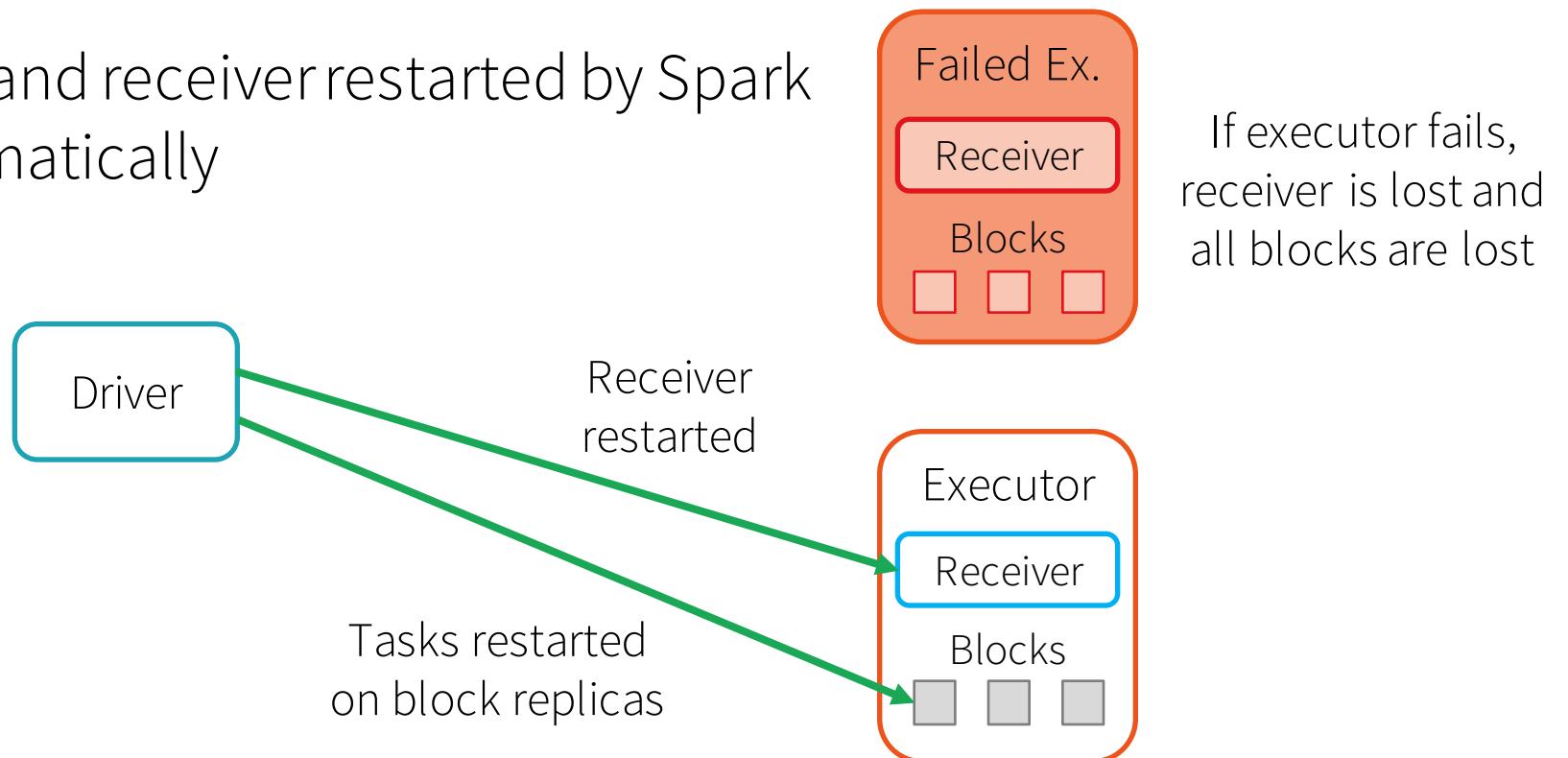


Spark Streaming Application: Process data



What if an executor fails?

Task and receiver restarted by Spark automatically



What if the driver fails?

Failed
Driver

When the driver fails, all the executors fail

All computation, all received blocks are lost

How to recover?

Failed Ex.

Receiver

Blocks



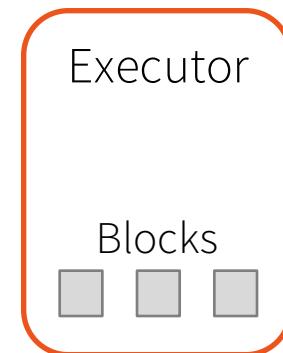
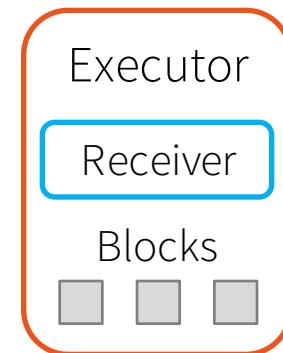
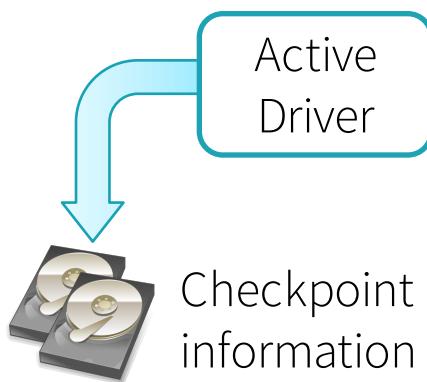
Failed
Executor

Blocks



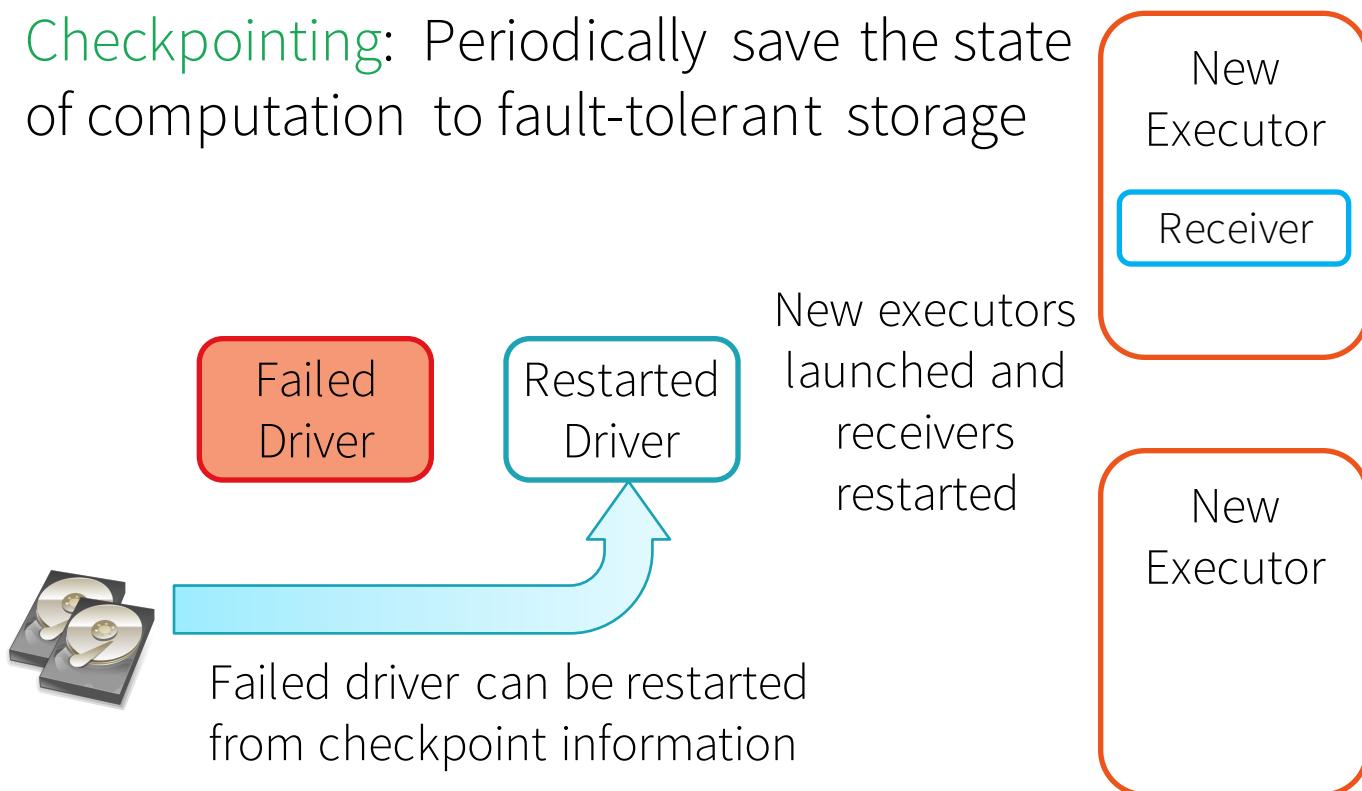
Recovering Driver with Checkpointing

Checkpointing: Periodically save the state of computation to fault-tolerant storage

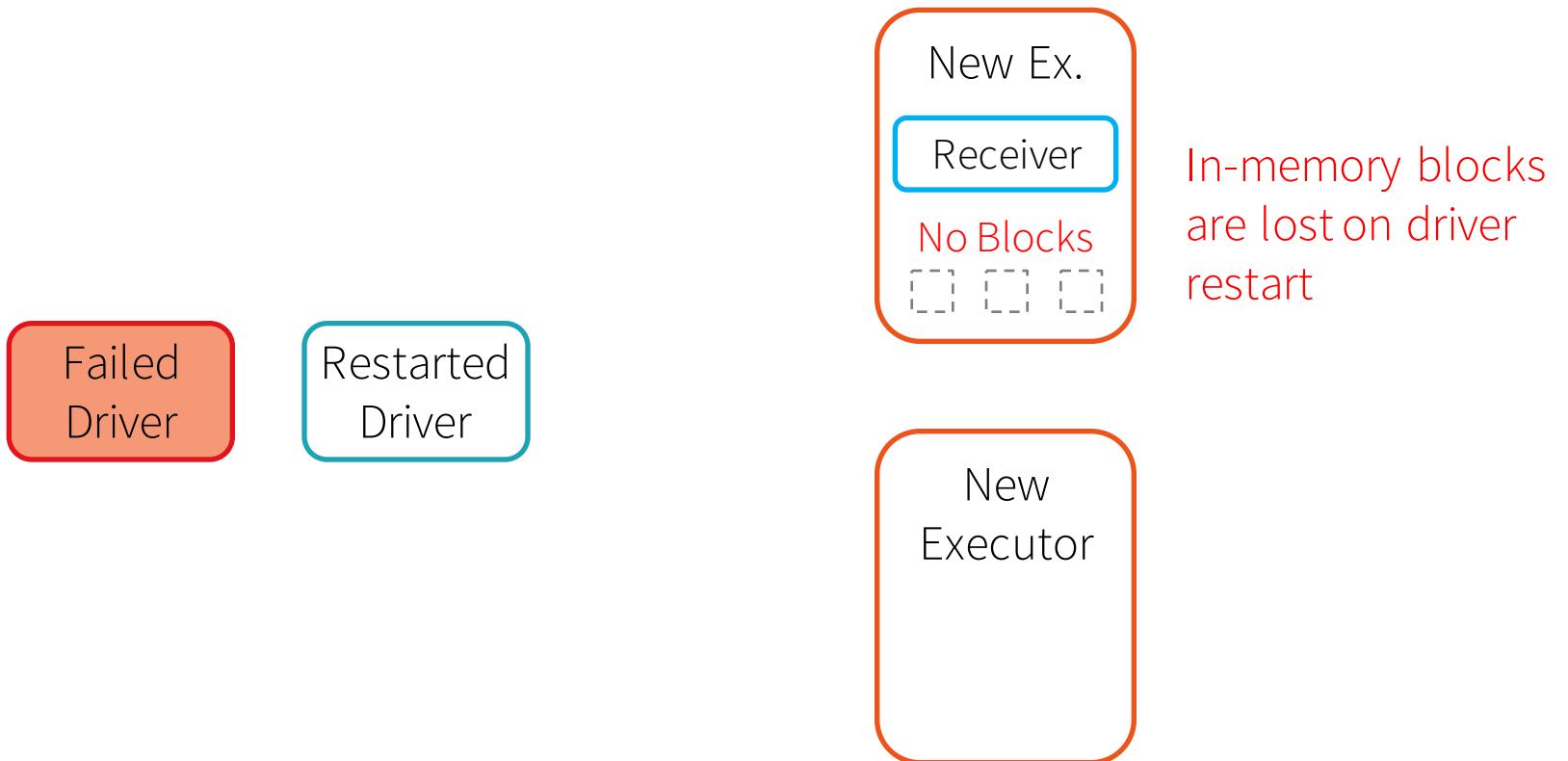


Recovering Driver with Checkpointing

Checkpointing: Periodically save the state of computation to fault-tolerant storage

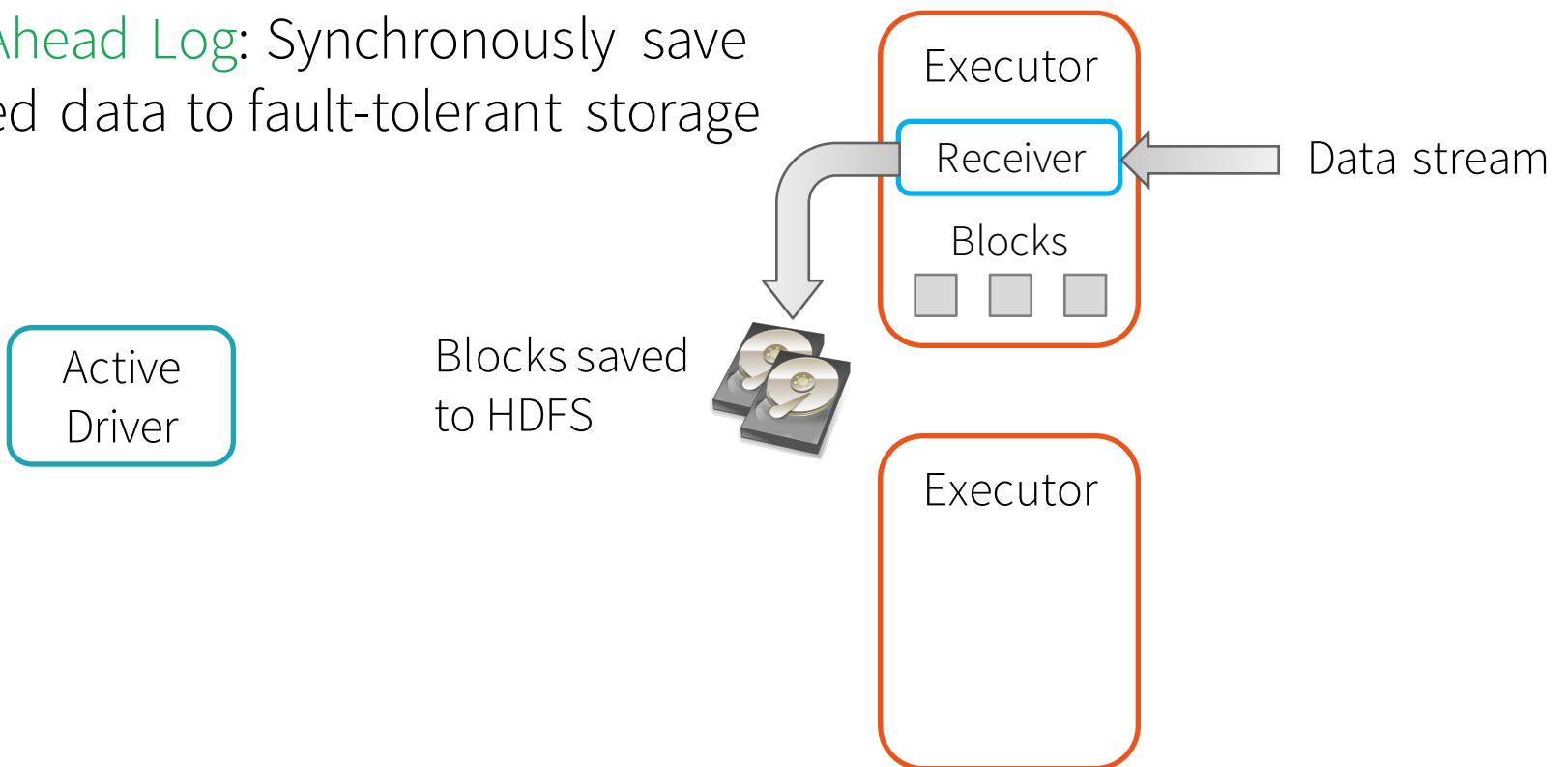


Received blocks lost on Restart!



Recovering data with Write Ahead Logs

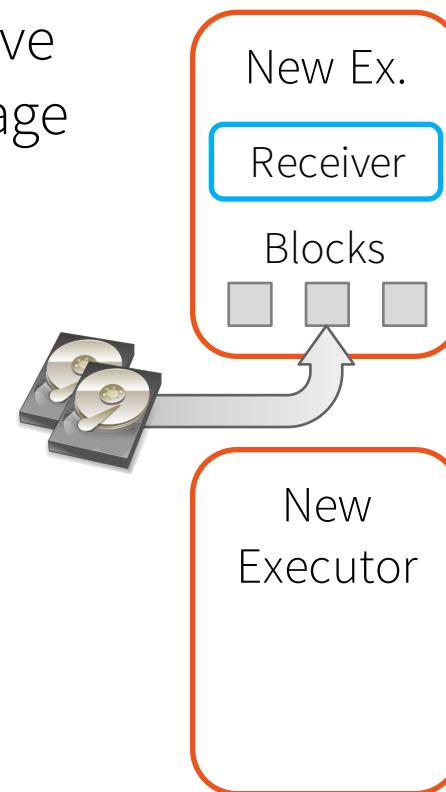
Write Ahead Log: Synchronously save received data to fault-tolerant storage



Received blocks lost on Restart!

Write Ahead Log: Synchronously save received data to fault-tolerant storage

Failed Driver Restarted Driver



Blocks recovered
from Write Ahead Log

Kafka:

Approach 1: Receiver-based Approach

This approach uses a Receiver to receive the data.

Approach 2: Direct Approach (No Receivers)

New receiver-less “direct” approach since Spark 1.3

Periodically queries Kafka for the latest offsets in each topic + partition
and accordingly defines the offset ranges to process in each batch

```
val directKafkaStream = KafkaUtils.createDirectStream[  
    [key class], [value class], [key decoder class], [value decoder class]](  
    streamingContext, [map of Kafka parameters], [set of topics to consume])
```

MLlib Operations on DStreams:

Streaming Linear Regression

Streaming Logistic Regression

Streaming KMeans

Online hypothesis testing (A/B testing)



Streaming DataFrames



Umbrella ticket to track what's needed to make streaming DataFrame a reality:

<https://issues.apache.org/jira/browse/SPARK-8360>

Streaming DataFrames

- Easier-to-use APIs (batch, streaming, and interactive)

- And optimizations:
 - - Tungsten backends
 - - native support for out-of-order data
 - - data sources and sinks

```
val stream = read.kafka("...")  
stream.window(5 mins, 10 secs)  
  .agg(sum("sales"))  
  .write.jdbc("mysql://...")
```

Spark 1.3
Static DataFrames

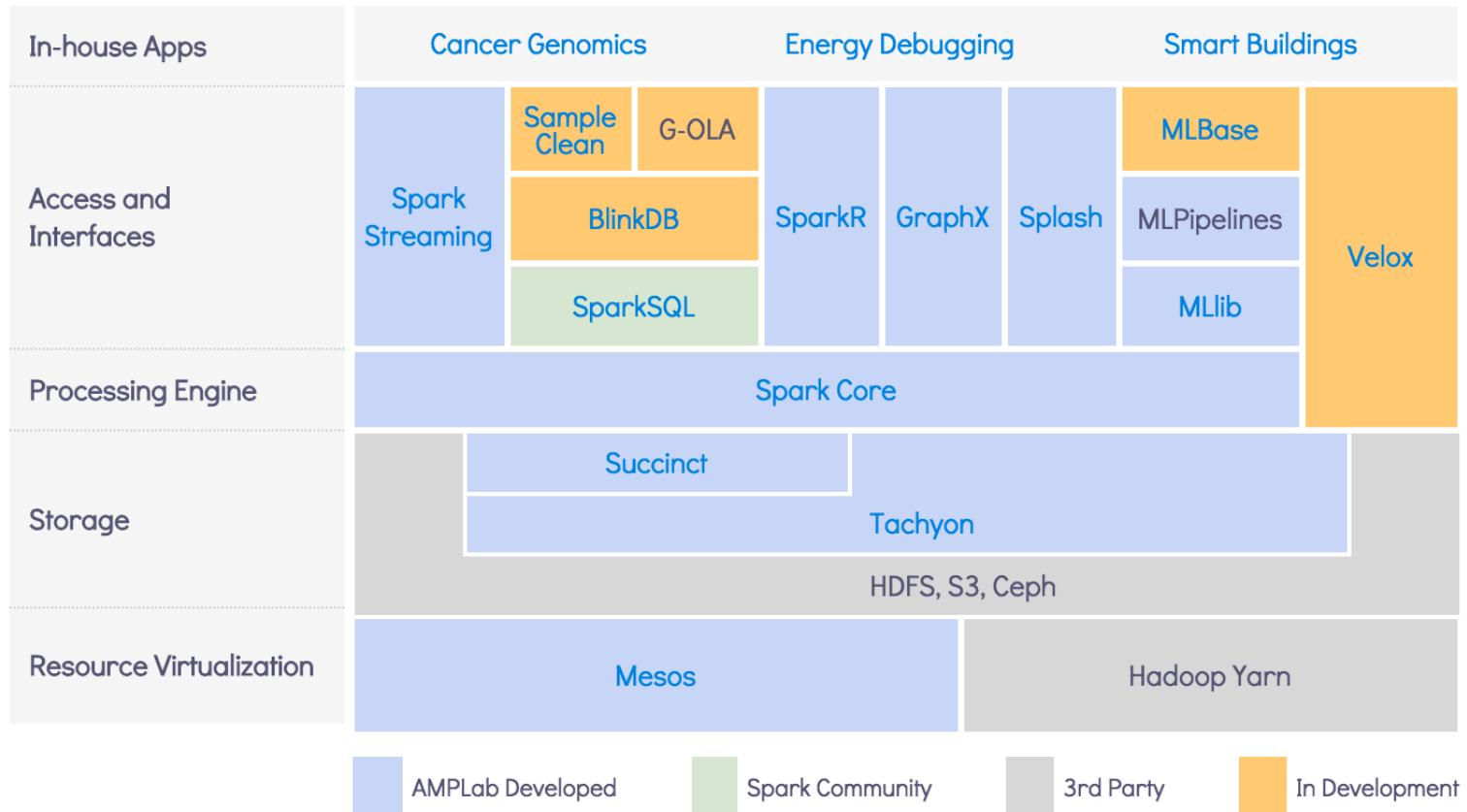
Spark 2.0
Continuous DataFrames

Single API !



Closing Remarks

BDAS: Berkeley Data Analytics Stack



Hardware Trends

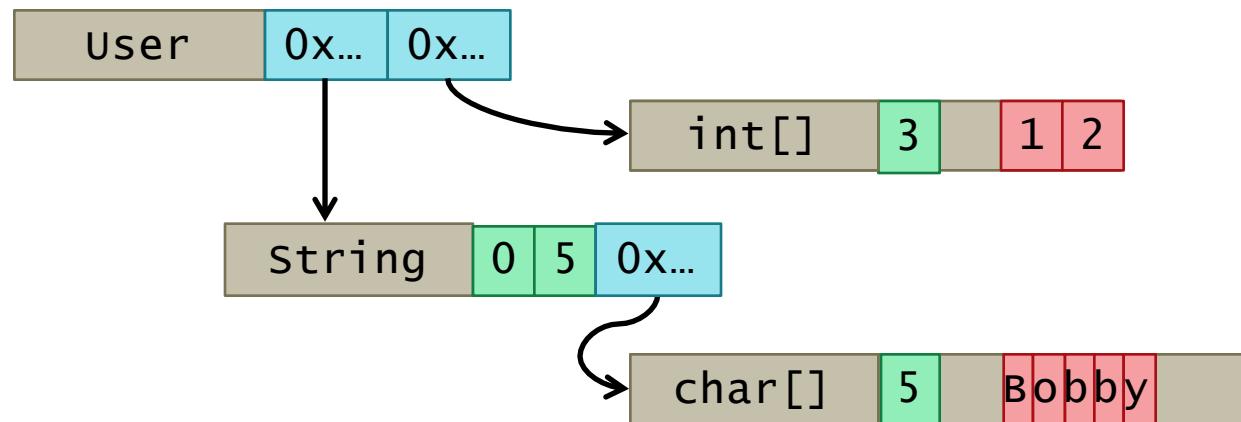
	2010	2015	
Storage	50+MB/s (HDD)	500+MB/s (SSD)	10X
Network	1Gbps	10Gbps	10X
CPU	~3GHz	~3GHz	:(

Project Tungsten

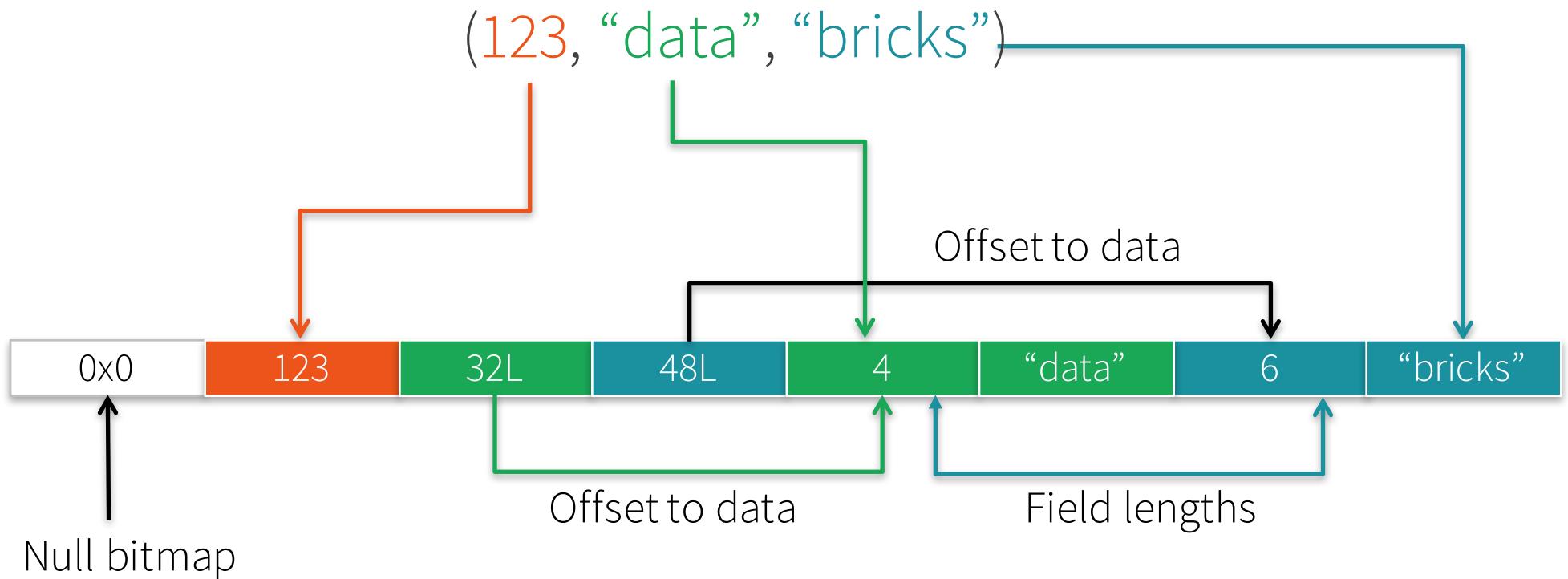
- Substantially speed up execution by optimizing CPU efficiency, via:
 - (1) Runtime code generation
 - (2) Exploiting cache locality
 - (3) Off-heap memory management

Challenge: Data Representation

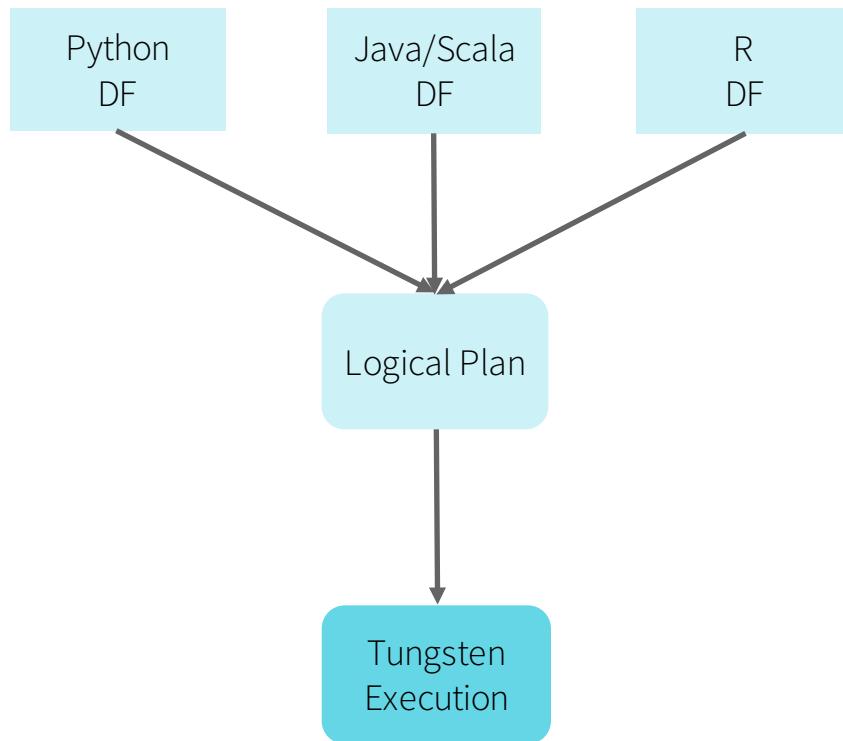
- Java objects often many times larger than underlying fields
- `class User(name: String, friends: Array[Int])`
- `User("Bobby", Array(1, 2))`



Tungsten's Compact Encoding



From DataFrame to Tungsten



Initial phase in Spark 1.5

More work coming in 2016

BLOG > ANNOUNCEMENTS , VIRTUAL MACHINES

Largest VM in the Cloud

THURSDAY, JANUARY 8, 2015



DREW McDANIEL
Principal Program Manager, Azure

G-Series Size Details

VM Size	Cores	RAM
Standard_G1	2	2
Standard_G2	4	5
Standard_G3	8	11
Standard_G4	16	22
Standard_G5	32	448 GiB
		6596 GB
		64

AWS Announces X1 Instances For EC2 With 2TB Of Memory, Launching Next Year

Posted Oct 8, 2015 by [Frederic Lardinois \(@fredericl\)](#)

926 SHARES



Above: Amazon today announced a massive new instance type for its AWS EC2 compute service. The

CrunchBase

Amazon

FOUNDED
1994

OVERVIEW
Amazon is an e-commerce retailer formerly known as Amazon.com. It offers consumers with products in thousands of categories. It also purchases for resale from vendors and third-party sellers. Operating in North America and International markets, Amazon provides services through websites such as amazon.com, amazon.ca. It also enables authors, music filmmakers, ...

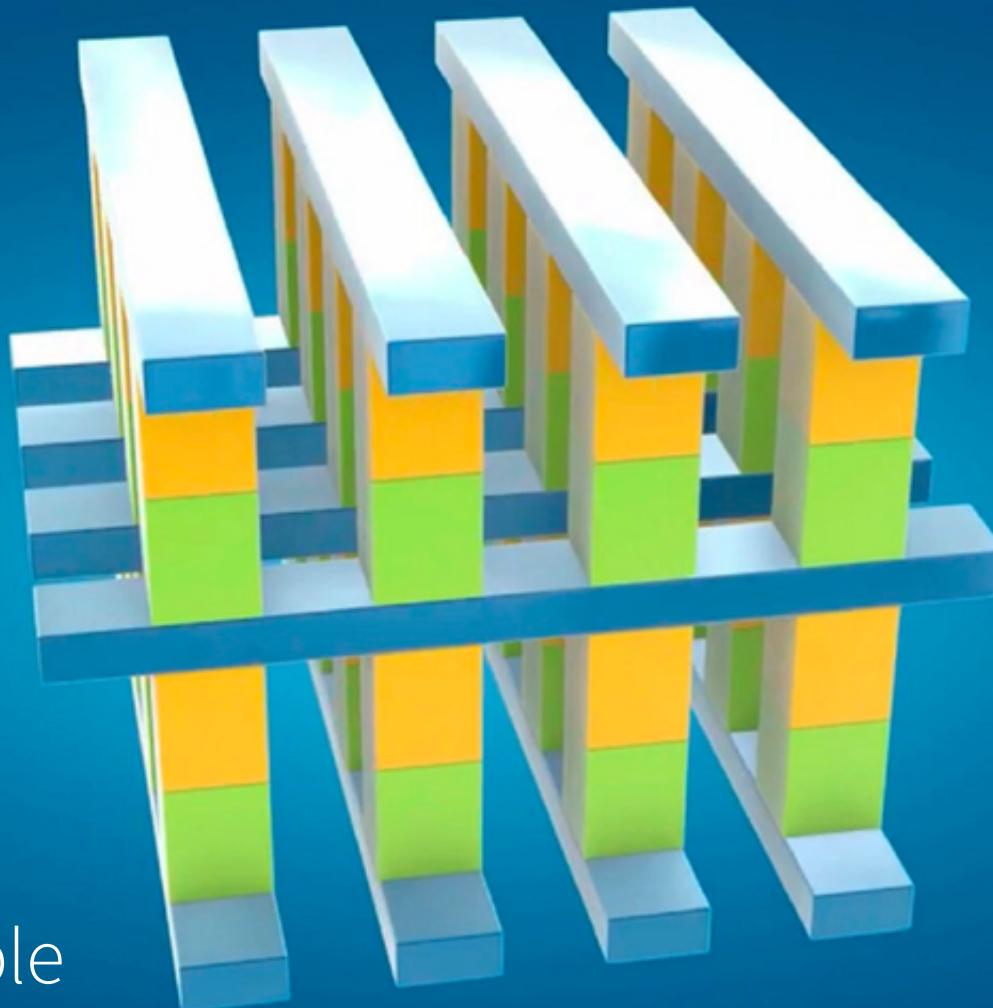
LOCATION
Seattle, WA

CATEGORIES
E-Commerce, Crowdsourcing, Groceries, Consumer Goods, Delivery, Software, Retail, Internet

FOUNDERS
Jeff Bezos

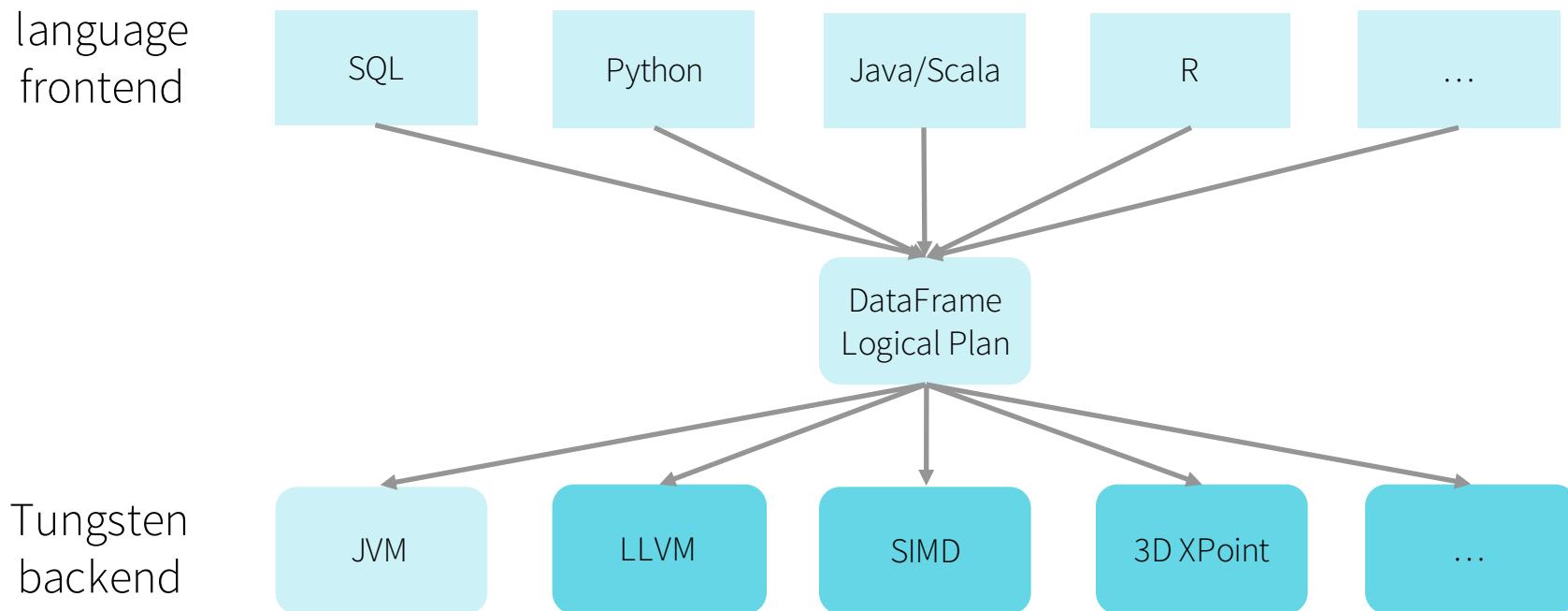
3D XPoint

- DRAM latency
- SSD capacity
- Byte addressible



Learn more: https://en.wikipedia.org/wiki/3D_XPoint

Unified API, One Engine, Automatically Optimized



SQL

Python

R

Streaming

Advanced
Analytics

DataFrame (& Dataset)

Tungsten Execution



Q) How many articles are needed to cover the sum of all human knowledge?

104 million!

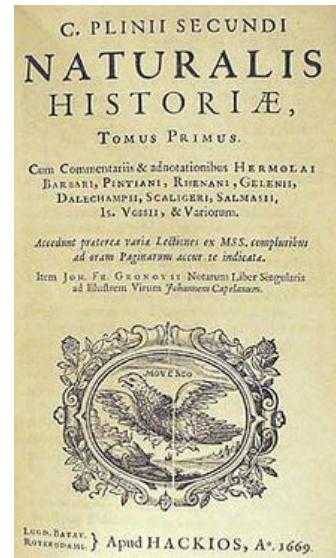
Source: https://en.wikipedia.org/wiki/User:Emijrp/All_human_knowledge

Library of Alexandria



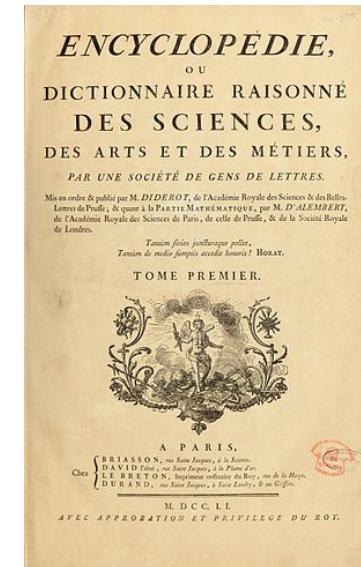
- From Egypt
- 3rd century B.C.
- 400,000 scrolls
- Burned down by Julius Caesar

Naturalis Historia



- From Roman Empire
- 77 A.D.
- 37 books

L'Encyclopédie



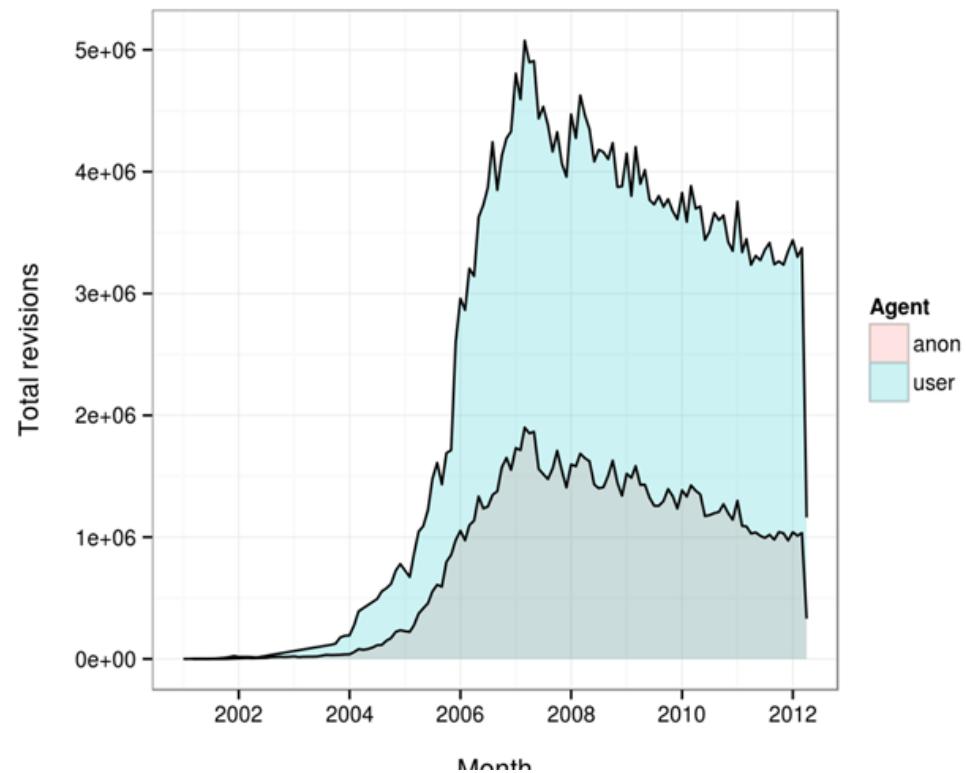
- French
- 1751
- 17 volumes
- 18,000 pages of text
- 75,000 entries

Upcoming Challenges



- Rise of smartphones
- Sock puppets, bias, paid editors...

Drop in # of Edits:



Source:
https://meta.wikimedia.org/wiki/Research:Anonymous_editor_acquisition/Volume_and_impact



CADE MEIZ BUSINESS 12.01.15 7:00 AM

WIKIPEDIA DEPLOYS AI TO EXPAND ITS RANKS OF HUMAN EDITORS



EVAN MILLS/WIRED

AARON HALFAKER JUST built an artificial intelligence engine designed to automatically analyze changes to Wikipedia.

Wikipedia is the online encyclopedia anyone can edit. In crowdsourcing the creation of an encyclopedia, the not-for-profit website forever changed the way we get information. It's among the ten most-visited sites on the Internet, and it has swept tomes like World Book and Encyclopedia Britannica into the dustbin of history. But it's not without flaws. If anyone can edit Wikipedia, anyone can mistakenly add bogus information. And anyone can vandalize the site, purposefully adding bogus information. Halfaker, a senior research scientist at the Wikimedia Foundation, the organization that oversees Wikipedia, built his AI engine as a way of identifying such vandalism.

O'REILLY®



Advanced Analytics with **Spark**

PATTERNS FOR LEARNING FROM DATA AT SCALE

Sandy Ryza, Uri Laserson,
Sean Owen & Josh Wills

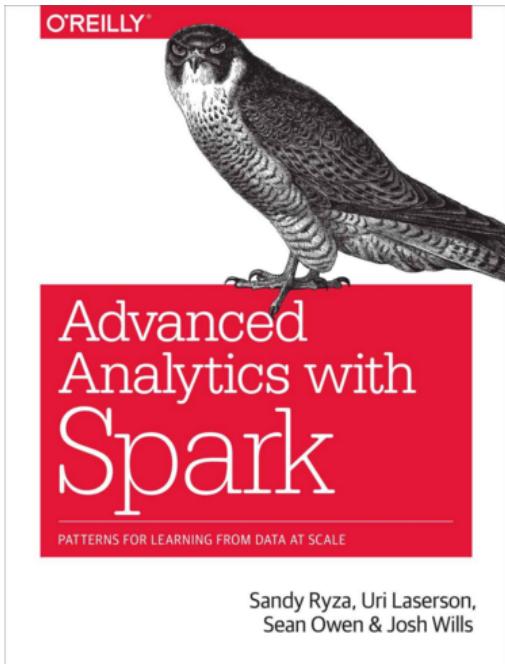
O'REILLY®



Learning **Spark**

LIGHTNING FAST DATA ANALYSIS

Holden Karau, Andy Konwinski,
Patrick Wendell & Matei Zaharia



Chapter 6: Understanding Wikipedia with Latent Semantic Analysis

- The Term-Document Matrix
- Getting the Data
- Parsing and Preparing the Data
- Lemmatization
- Computing the TF-IDFs
- Singular Value Decomposition
- Finding Important Concepts
- Querying and Scoring with the Low-Dimensional Representation
- Term-Term Relevance
- Document-Document Relevance
- Term-Document Relevance
- Multiple-Term Queries

You will have access to Databricks and these datasets via Community Edition after class...

Good luck!