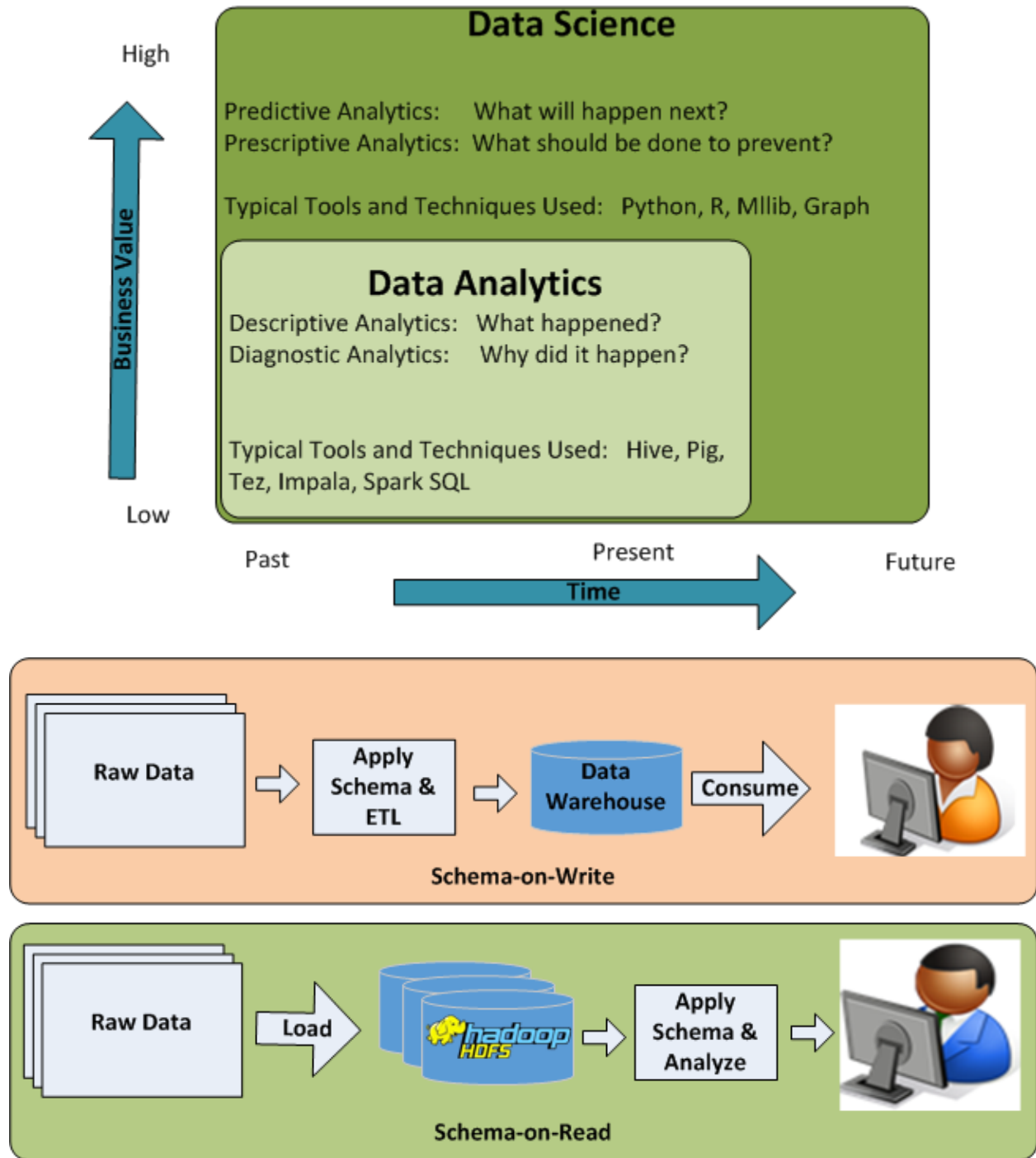
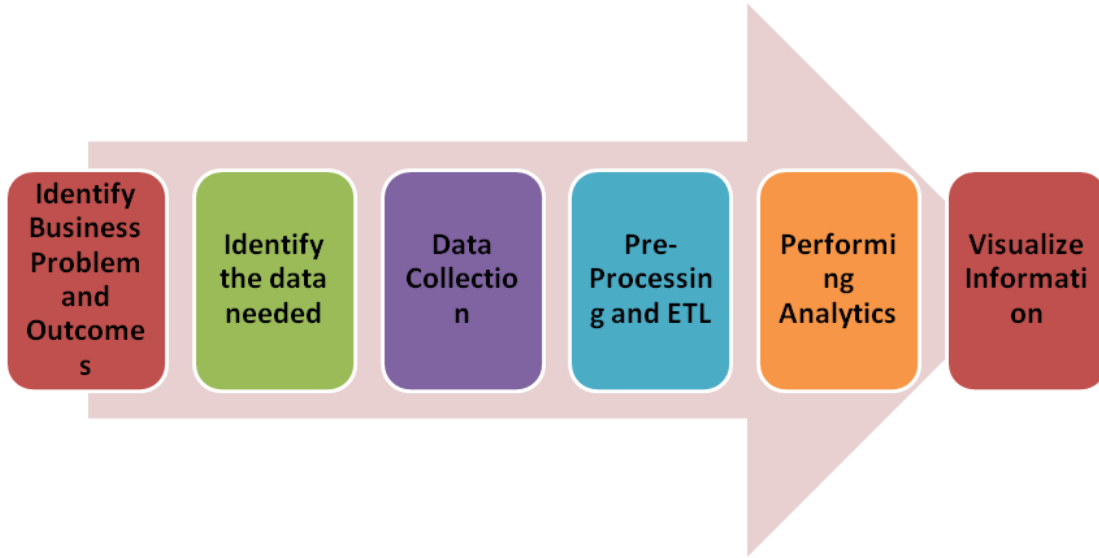
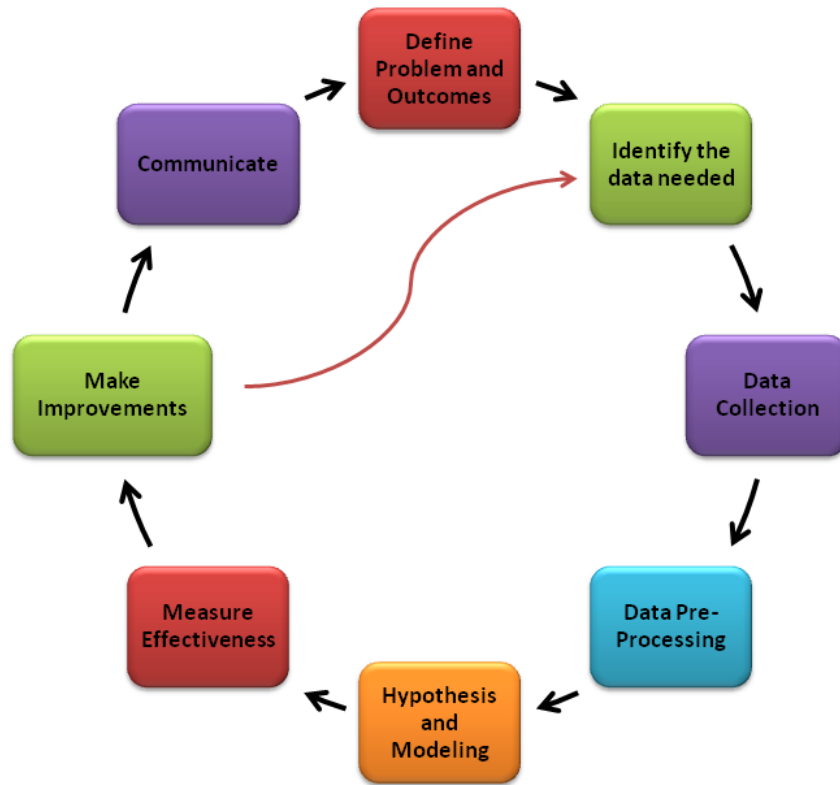


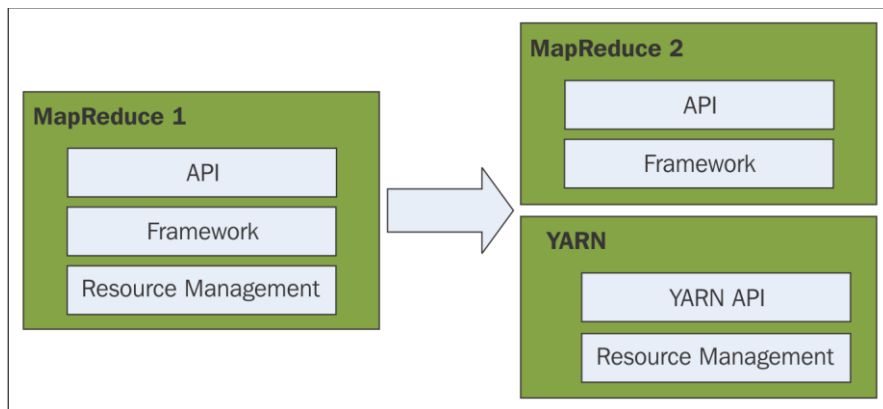
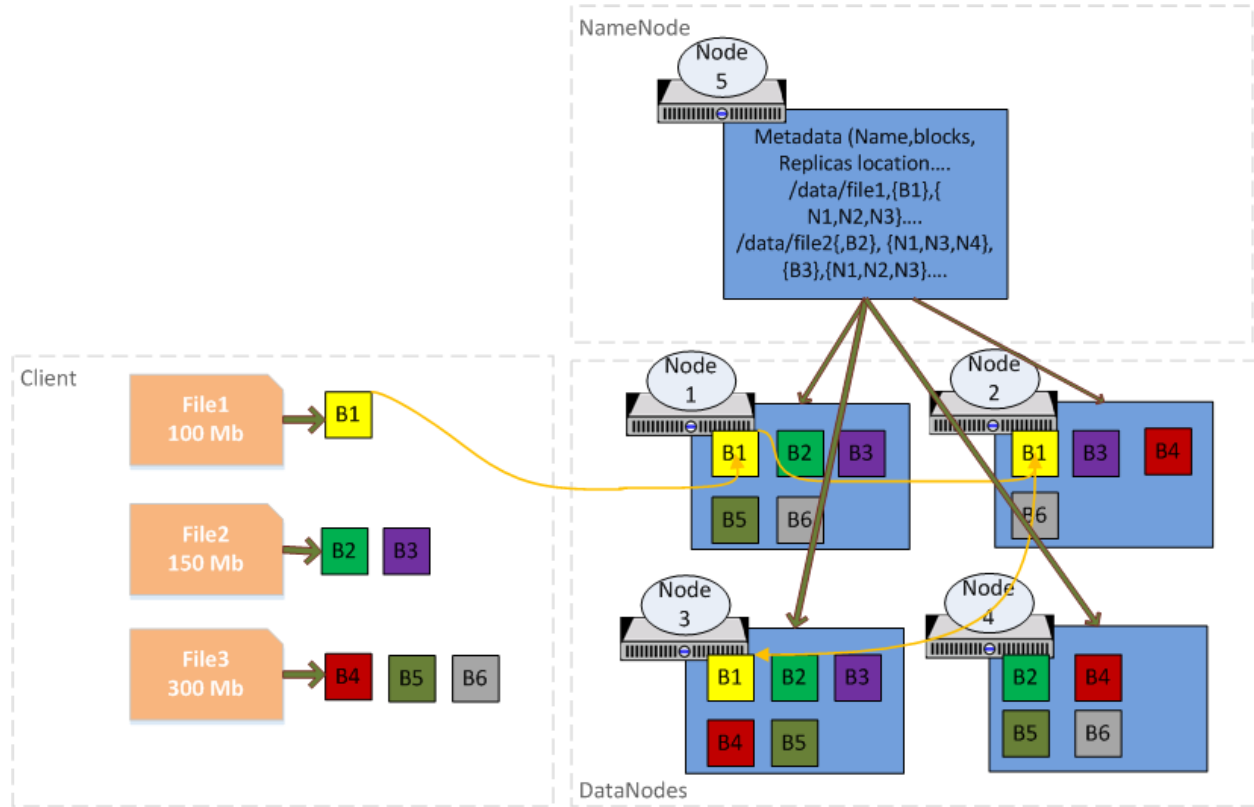
# Chapter 1: Big Data Analytics at a 10,000-Foot View

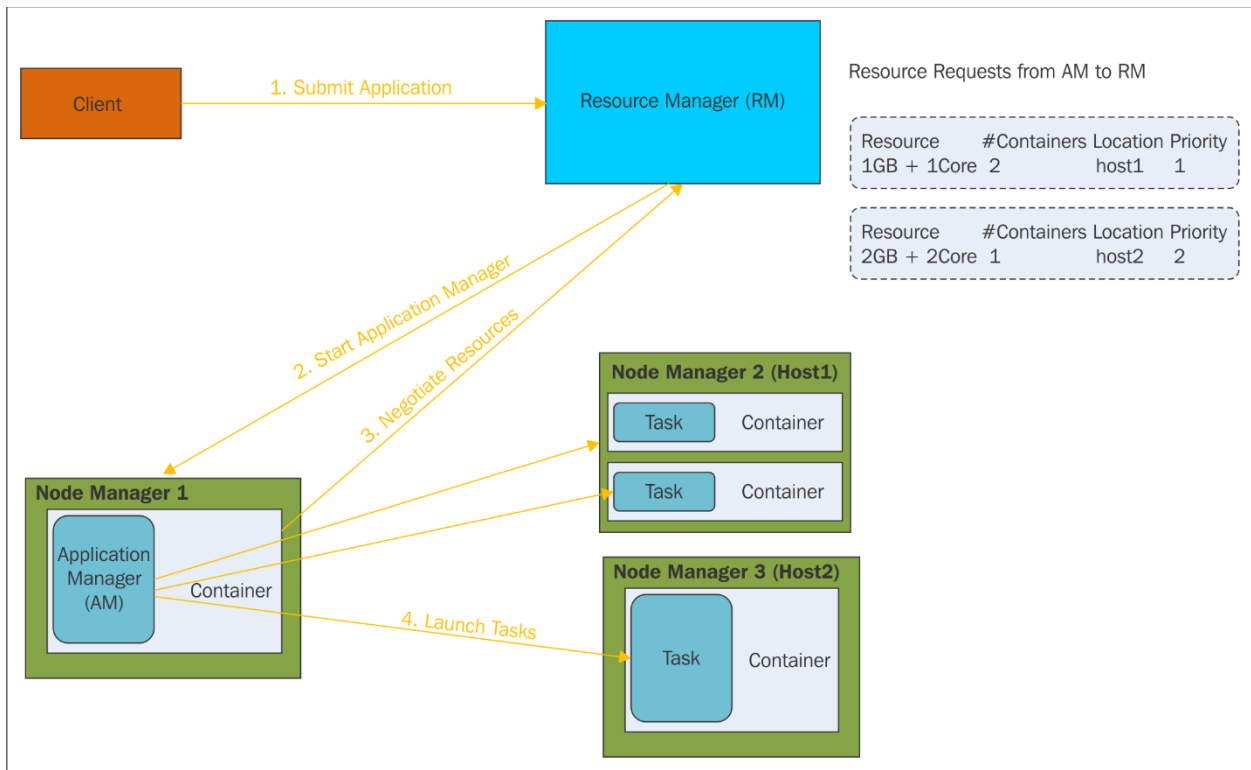
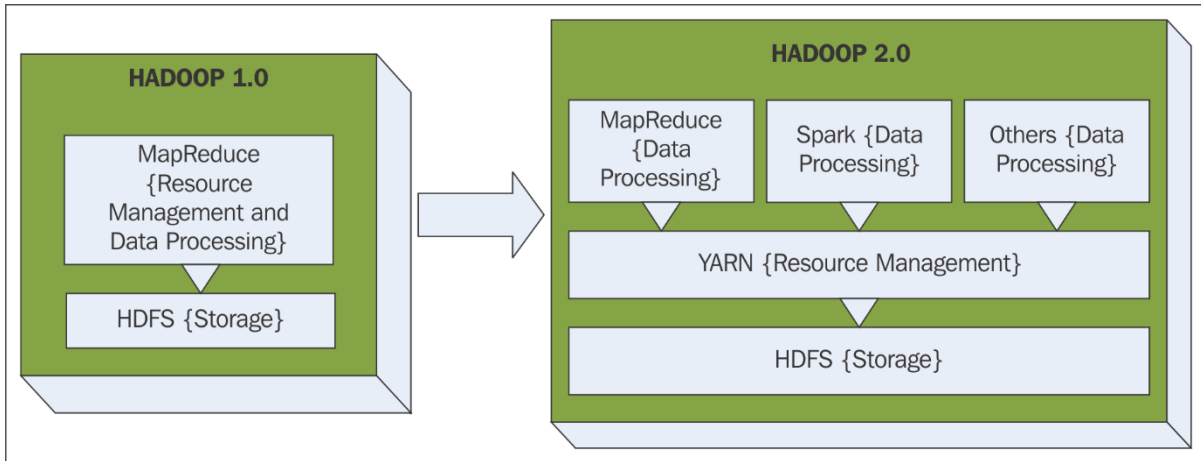


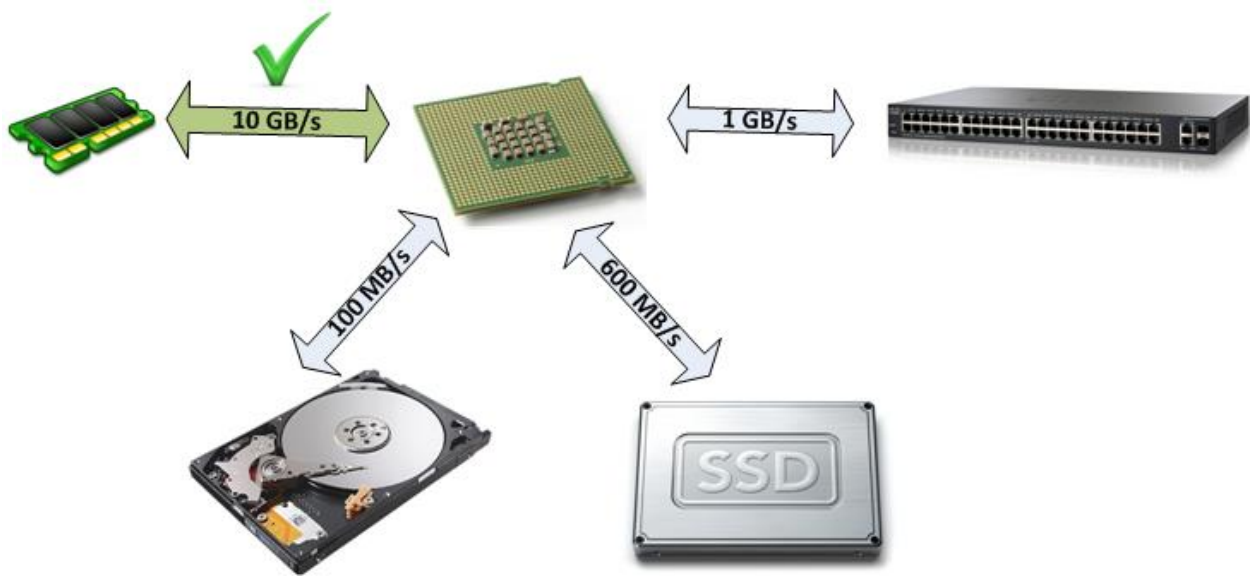
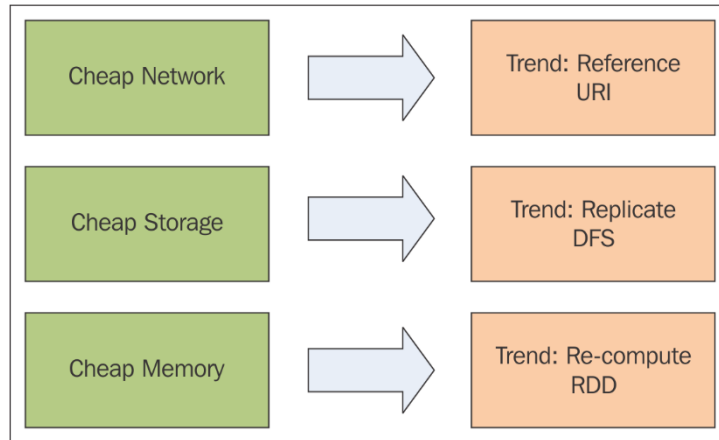




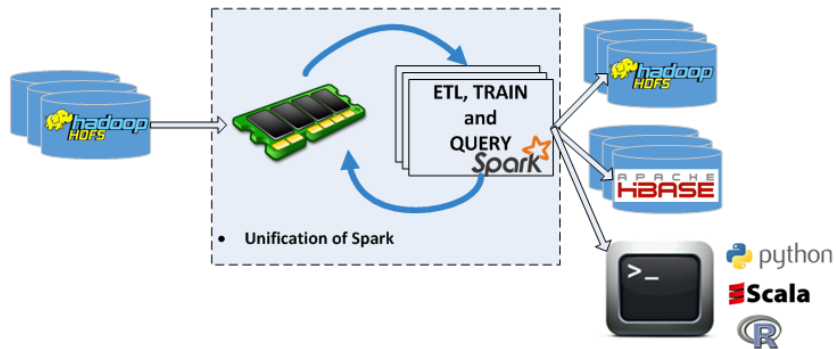
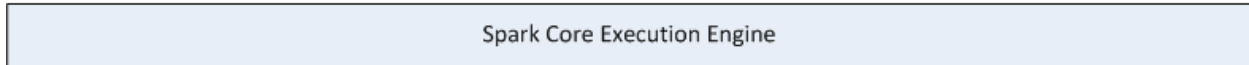
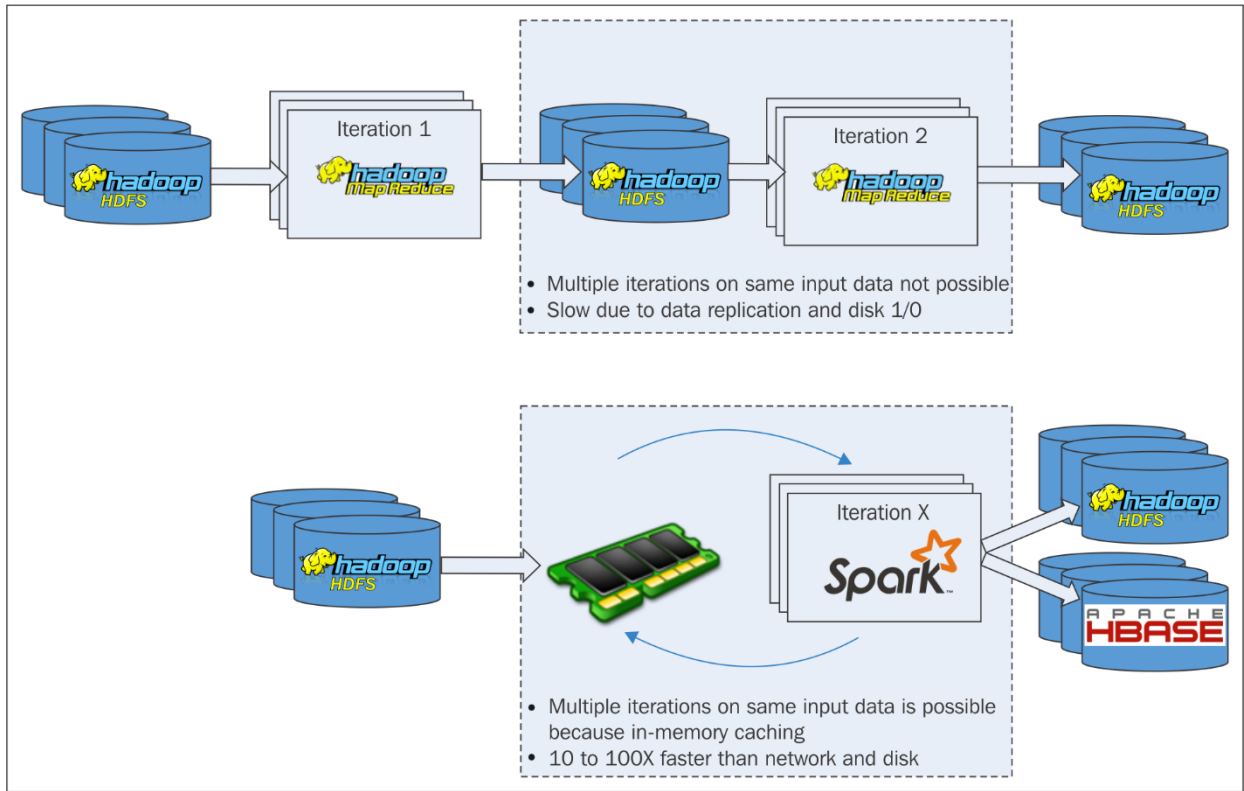
# Chapter 2: Getting Started with Apache Hadoop and Apache Spark

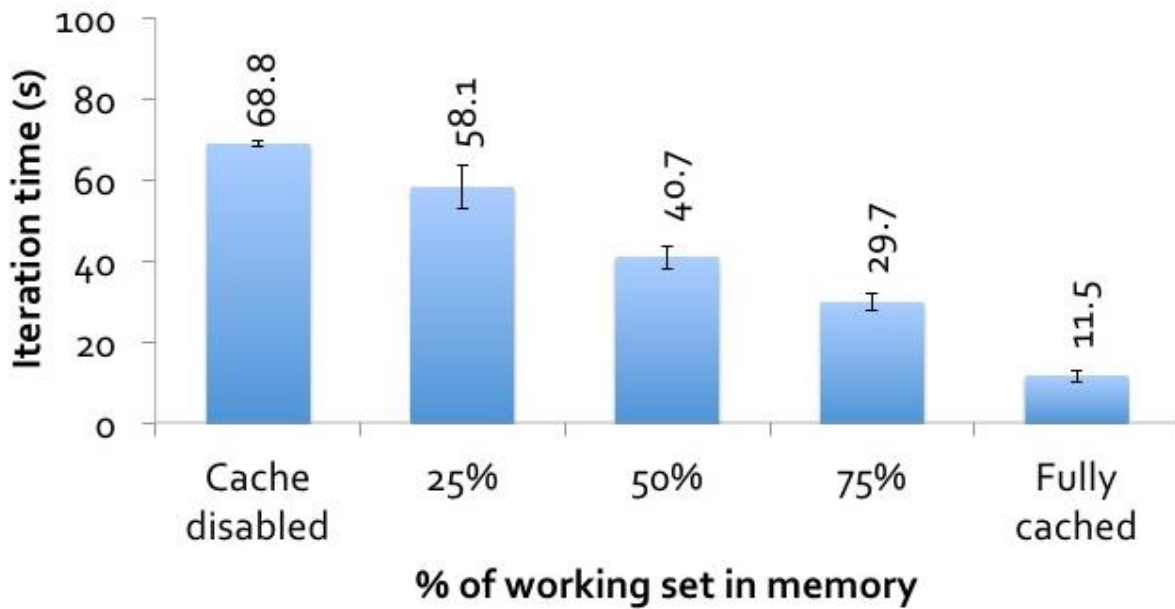
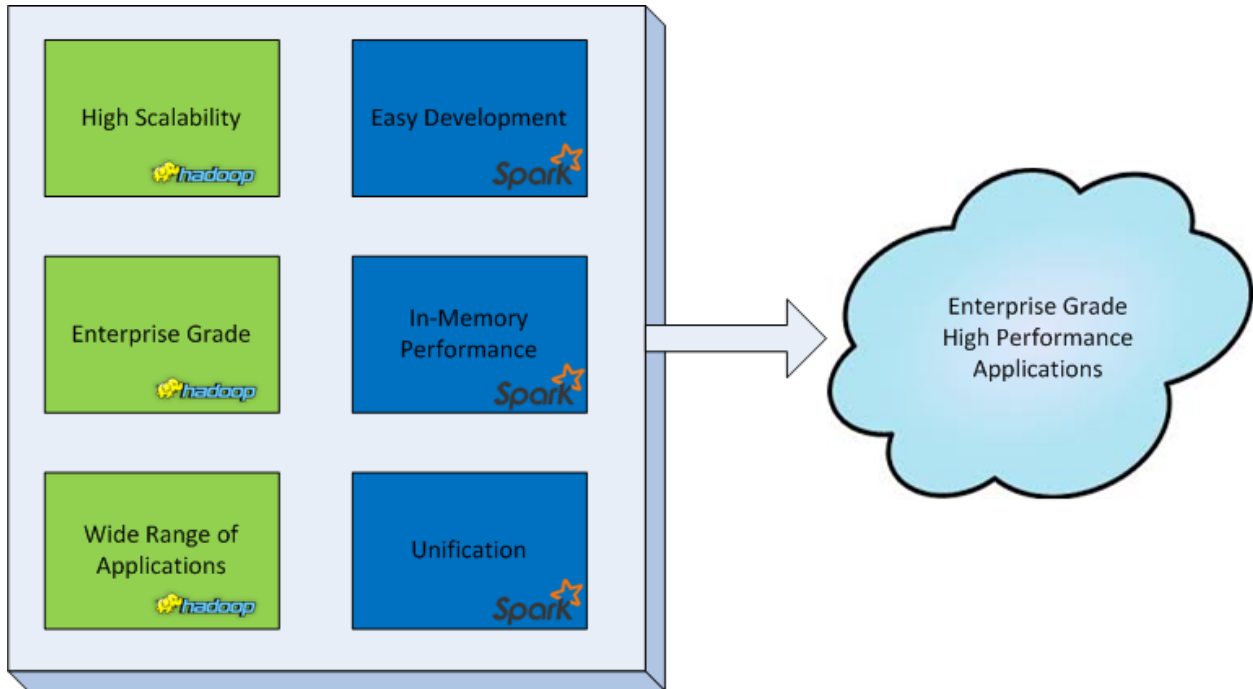






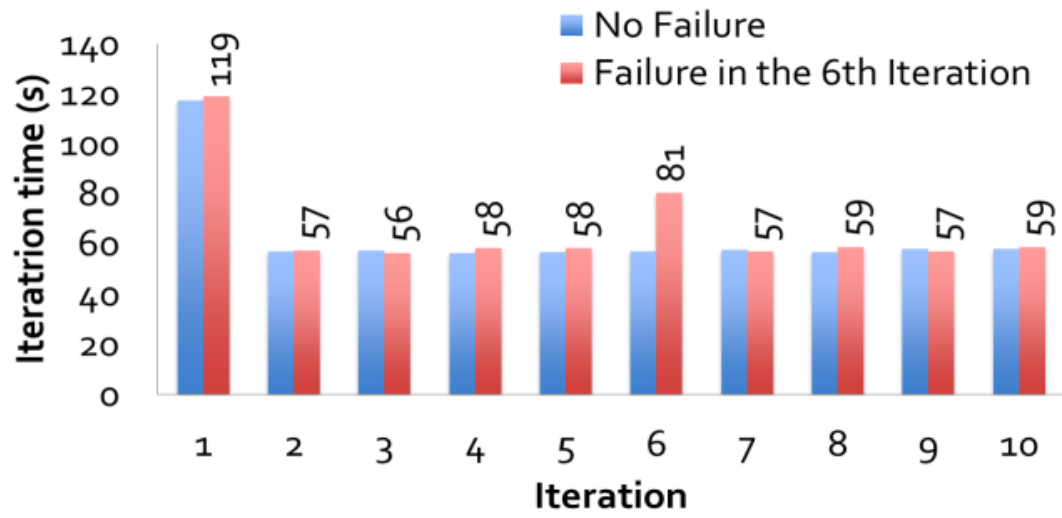
<p><b>2010</b></p> <ul style="list-style-type: none"> <li>• Spark Paper released</li> <li>• Open Sourced</li> </ul>	<p><b>2013</b></p> <ul style="list-style-type: none"> <li>• Spark goes to Apache Software Foundation</li> </ul>	<p><b>2014</b></p> <ul style="list-style-type: none"> <li>• Top level project at Apache Software Foundation.</li> <li>• Databricks setup 100 TB Sort new record in 23 Mins with Spark.</li> <li>• Spark v 1.0 released.</li> </ul>	<p><b>2015 &amp; 2016</b></p> <ul style="list-style-type: none"> <li>• Spark v 1.5 released</li> <li>• Spark v 1.6 released</li> <li>• Spark v 2.0 released</li> <li>• New APIs, Features and Performance Improvements</li> </ul>
---	---	--	---





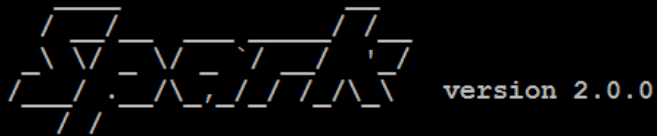


### Fault Recovery Address



## Chapter 3: Deep Dive into Apache Spark

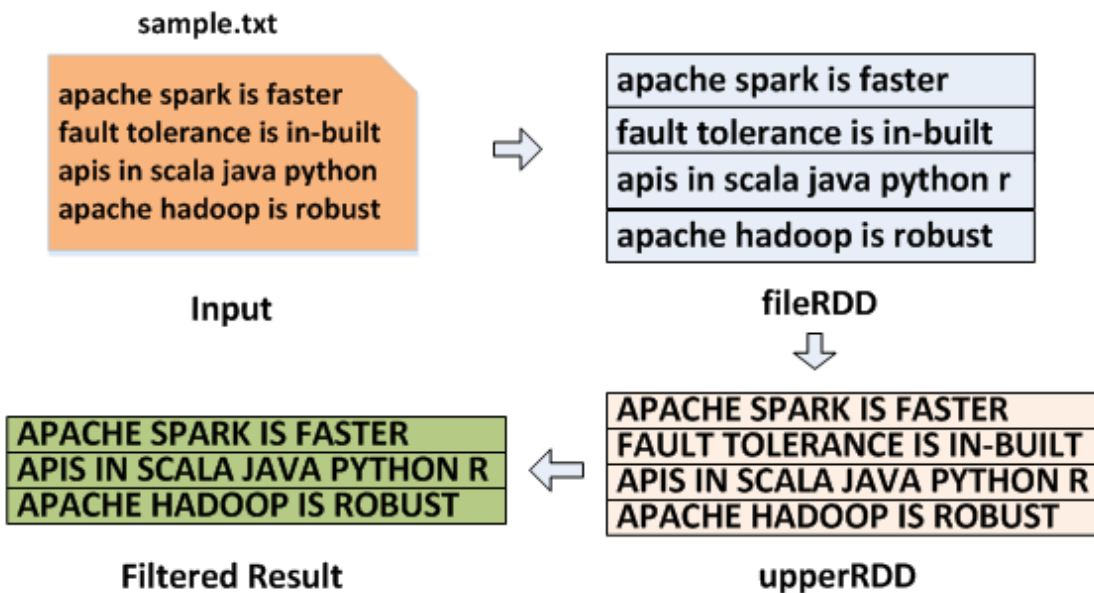
```
[cloudera@quickstart spark-2.0.0-bin-hadoop2.7]$ bin/spark-shell
Spark context Web UI available at http://192.168.139.175:4040
Spark context available as 'sc' (master = local[*], app id = local-1470527900398).
Spark session available as 'spark'
Welcome to
```

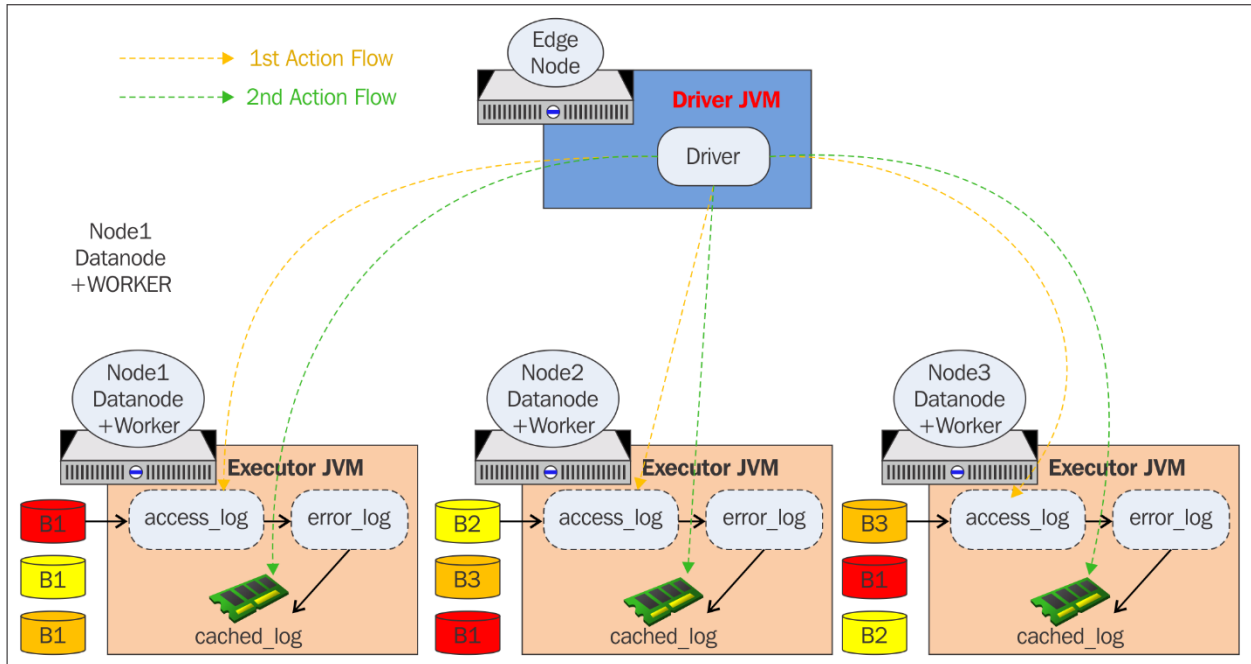


```
Using Scala version 2.11.8 (Java HotSpot(TM) 64-Bit Server VM, Java 1.7.0_67)
Type in expressions to have them evaluated.
Type :help for more information.
```

```
scala> █
```

```
fileRDD = sc.textFile("/data/sample.txt")
upperRDD = fileRDD.map(lambda line: line.upper())
upperRDD.filter(lambda line: line.startswith('A')).collect()
```





## Spark Jobs (?)

Total Uptime: 1.2 min  
 Scheduling Mode: FIFO  
 Completed Jobs: 1

▶ [Event Timeline](#)

### Completed Jobs (1)

Job Id	Description	Submitted	Duration	Stages: Succeeded/Total	Tasks (for all stages): Succeeded/Total
0	count at <stdin>:1	2016/04/07 04:01:06	1 s	1/1	4/4

### Summary Metrics for 4 Completed Tasks

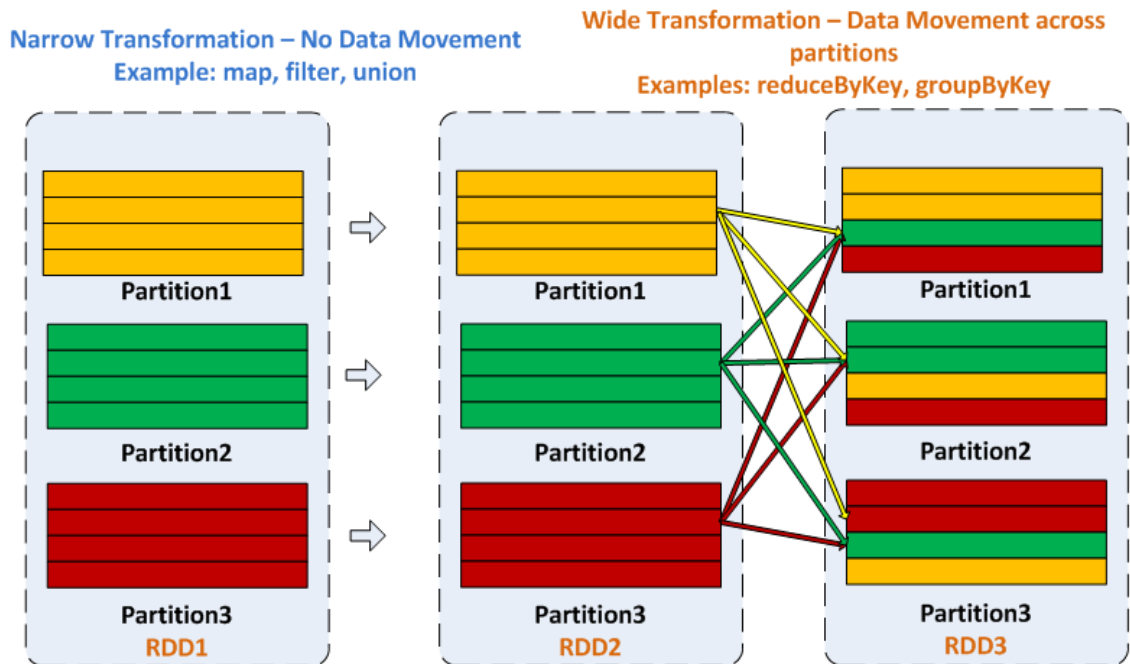
Metric	Min	25th percentile	Median	75th percentile	Max
Duration	0.4 s	0.4 s	0.4 s	0.4 s	0.4 s
GC Time	0.3 s	0.3 s	0.3 s	0.3 s	0.3 s

### Aggregated Metrics by Executor

Executor ID ^	Address	Task Time	Total Tasks	Failed Tasks	Succeeded Tasks
0	192.168.139.175:43785	5 s	4	0	4

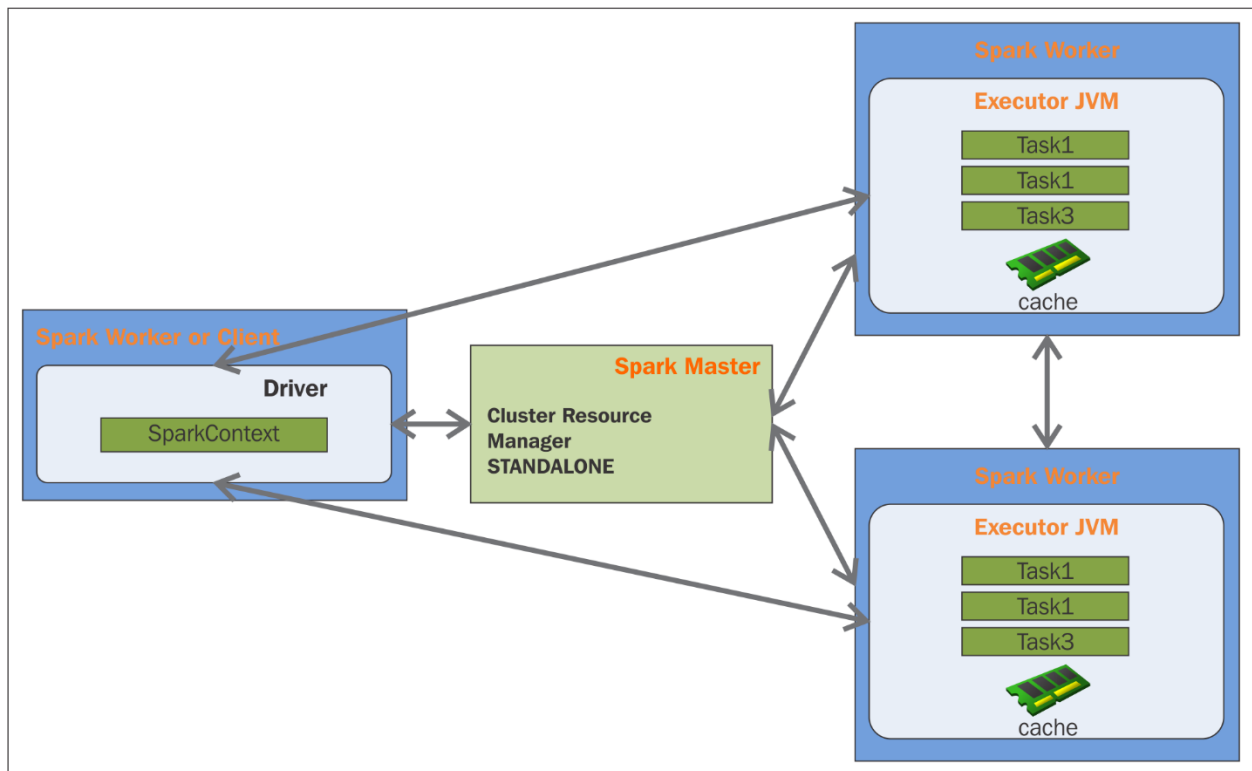
### Tasks

Index ^	ID	Attempt	Status	Locality Level	Executor ID / Host	Launch Time	Duration	GC Time	Errors
0	0	0	SUCCESS	PROCESS_LOCAL	0 / 192.168.139.175	2016/08/06 18:06:19	0.4 s	0.3 s	
1	1	0	SUCCESS	PROCESS_LOCAL	0 / 192.168.139.175	2016/08/06 18:06:19	0.4 s	0.3 s	
2	2	0	SUCCESS	PROCESS_LOCAL	0 / 192.168.139.175	2016/08/06 18:06:19	0.4 s	0.3 s	
3	3	0	SUCCESS	PROCESS_LOCAL	0 / 192.168.139.175	2016/08/06 18:06:19	0.4 s	0.3 s	



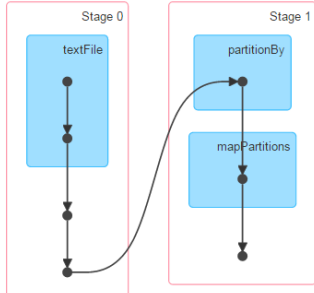
## Tasks

Index ▲	ID	Attempt	Status	Locality Level	Executor ID / Host
0	0	0	SUCCESS	PROCESS_LOCAL	0 / 192.168.139.164
1	1	0	SUCCESS	PROCESS_LOCAL	0 / 192.168.139.164
2	2	0	SUCCESS	PROCESS_LOCAL	0 / 192.168.139.164
3	3	0	SUCCESS	PROCESS_LOCAL	0 / 192.168.139.164



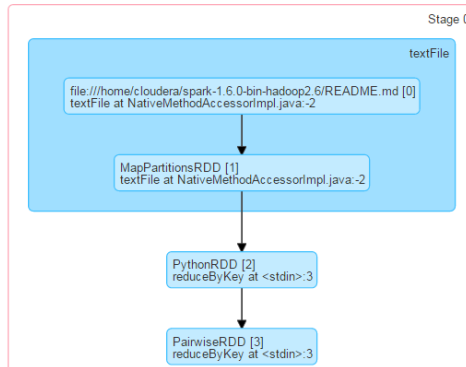
### Details for Job 0

Status: SUCCEEDED  
 Completed Stages: 2  
 Event Timeline  
 DAG Visualization



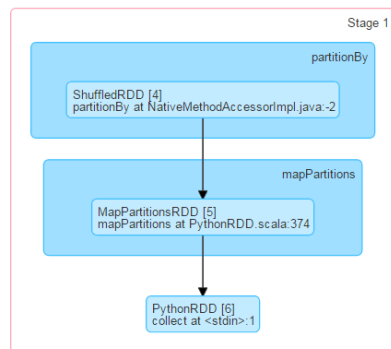
### Details for Stage 0 (Attempt 0)

Total Time Across All Tasks: 2 s  
 Locality Level Summary: Process local: 2  
 Shuffle Write: 4.9 KB / 16  
 DAG Visualization



### Details for Stage 1 (Attempt 0)

Total Time Across All Tasks: 0.2 s  
 Locality Level Summary: Node local: 2  
 Shuffle Read: 4.9 KB / 16  
 DAG Visualization

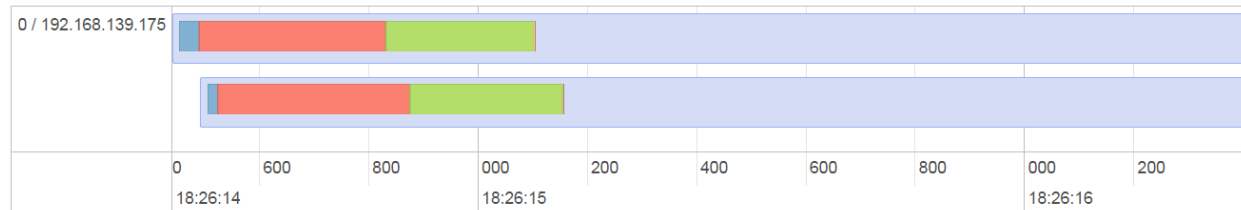


### Details for Stage 0 (Attempt 0)

Total Time Across All Tasks: 2 s  
 Locality Level Summary: Process local: 2  
 Input Size / Records: 1914.0 B / 99  
 Shuffle Write: 5.0 KB / 16

- DAG Visualization
- Show Additional Metrics
- Event Timeline
- Enable zooming

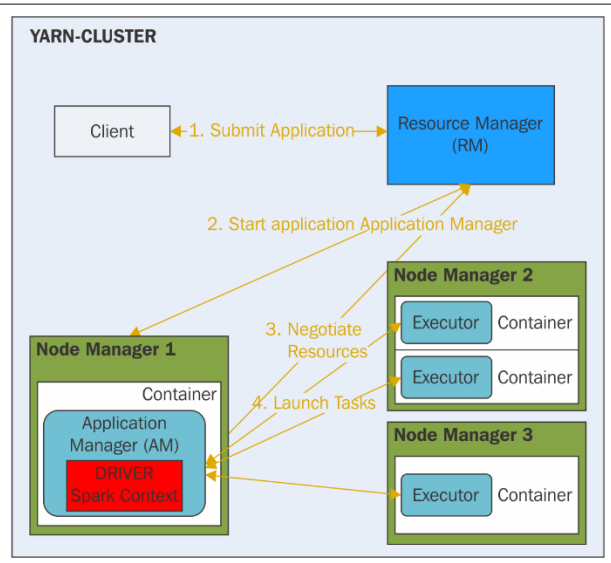
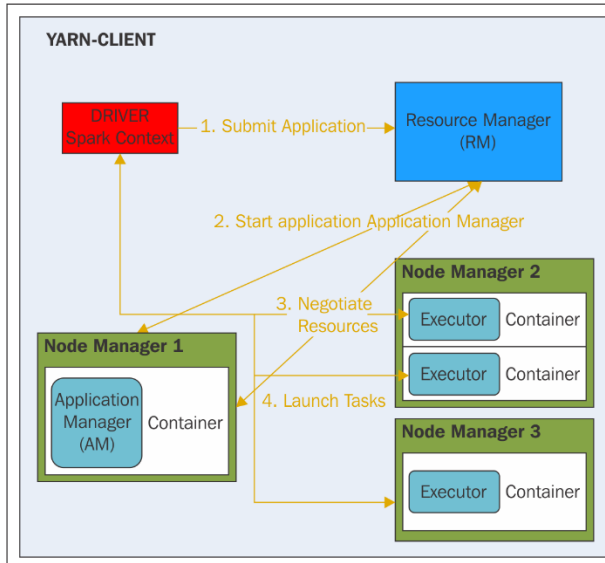
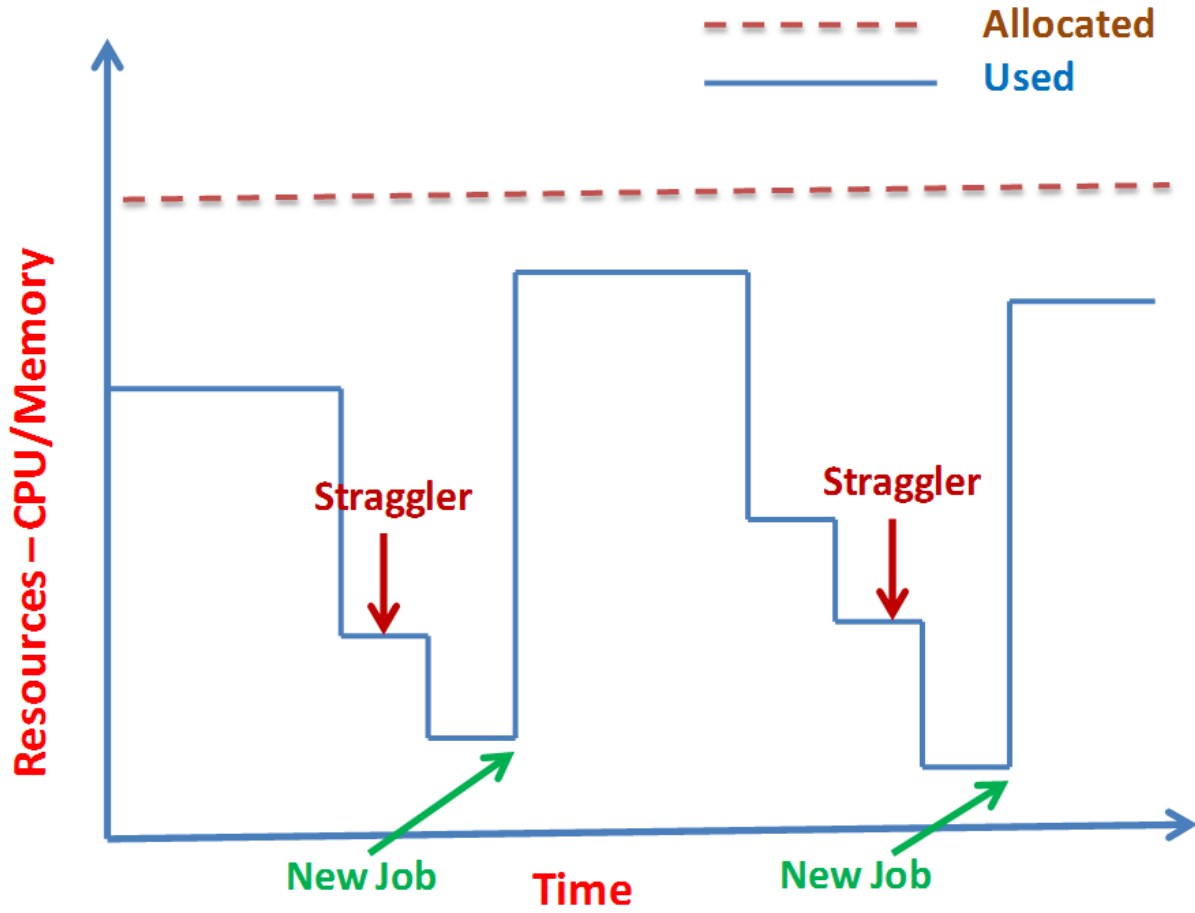
- Scheduler Delay
- Task Deserialization Time
- Shuffle Read Time
- Executor Computing Time
- Shuffle Write Time
- Result Serialization Time
- Getting Result Time



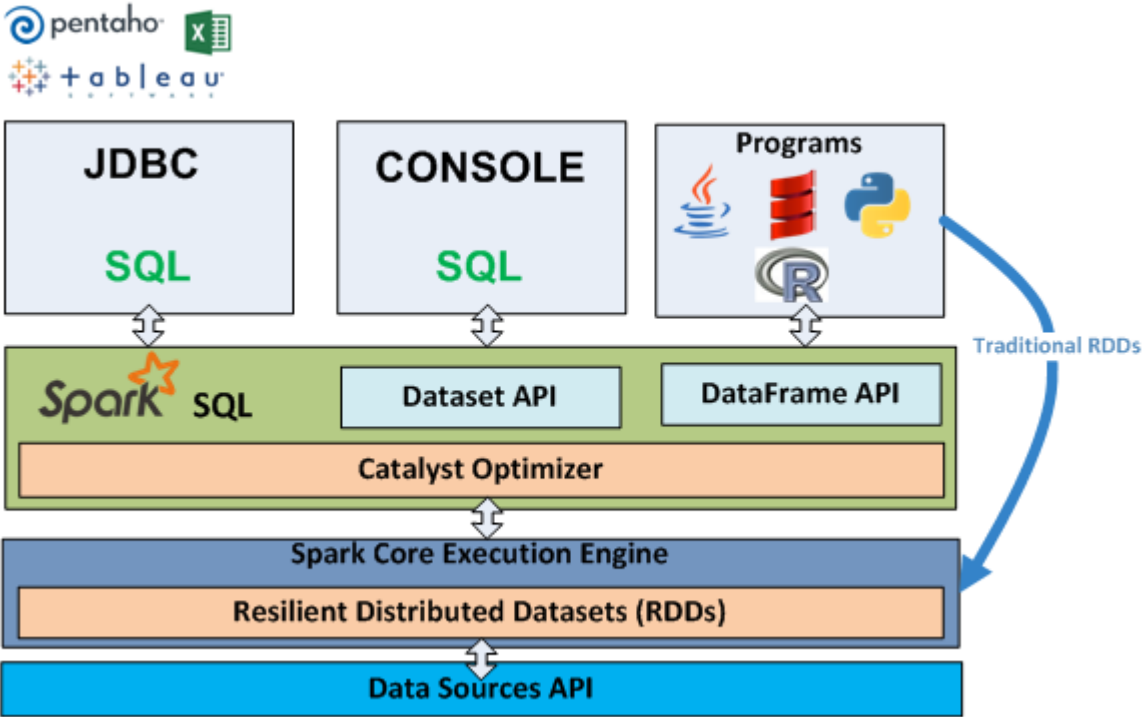
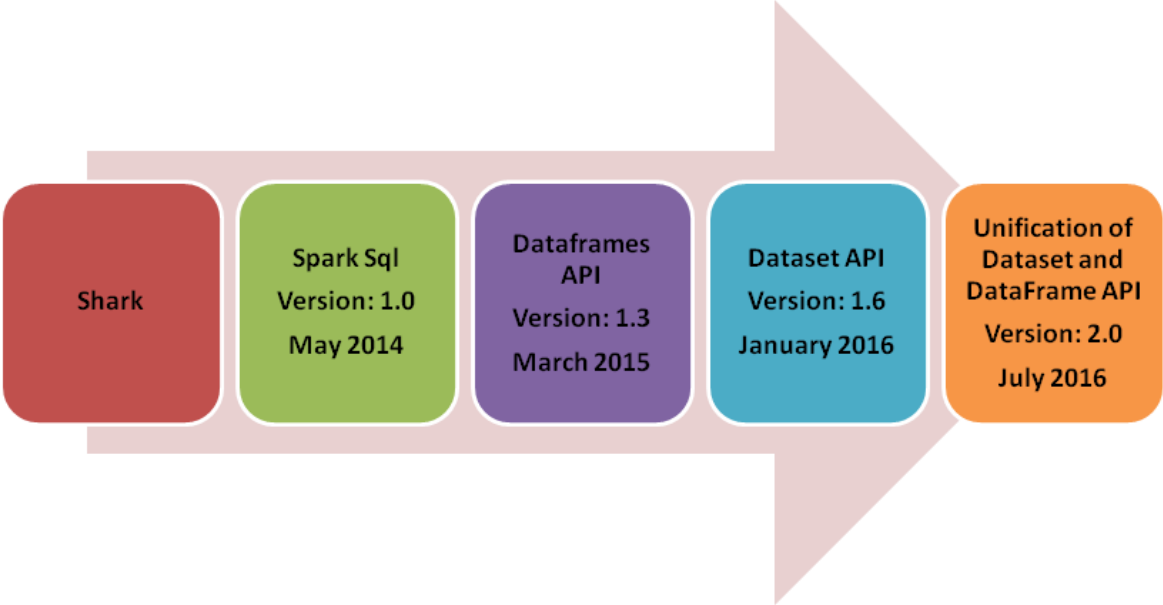
### Storage

#### RDDs

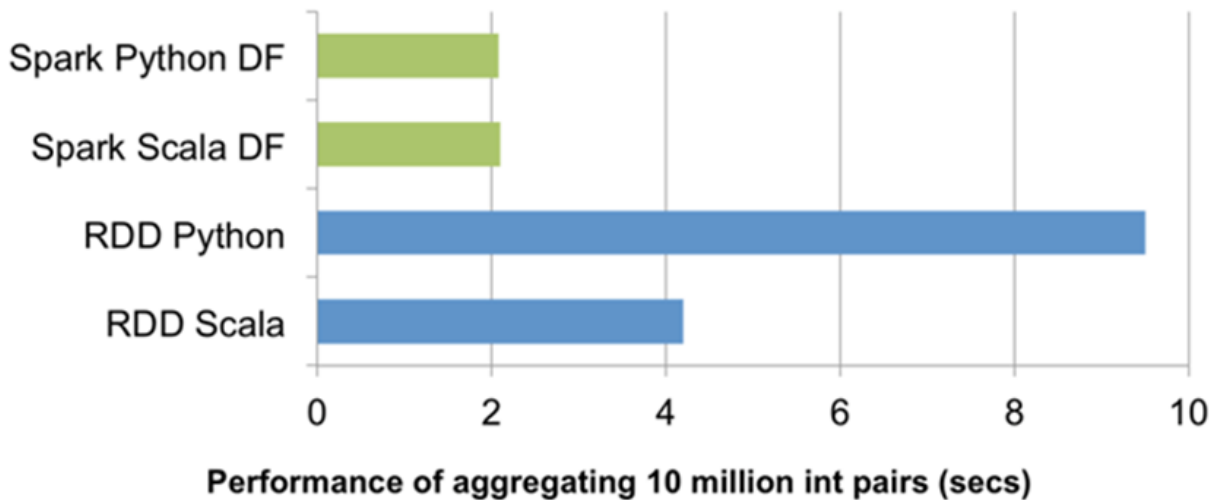
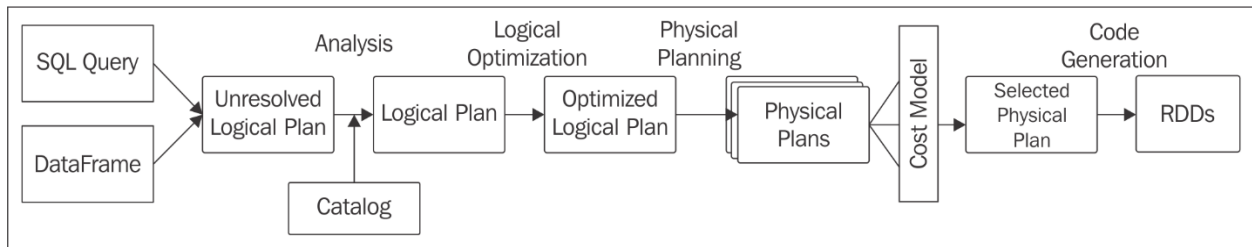
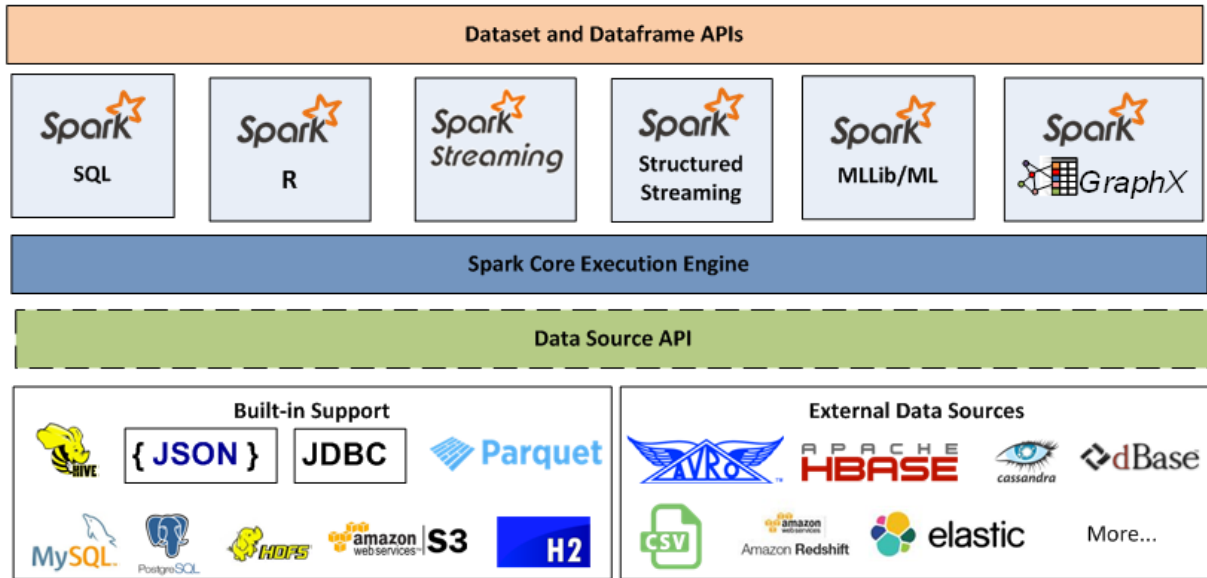
RDD Name	Storage Level	Cached Partitions	Fraction Cached	Size in Memory	Size on Disk
PythonRDD	Memory Serialized 1x Replicated	2	100%	3.9 KB	0.0 B



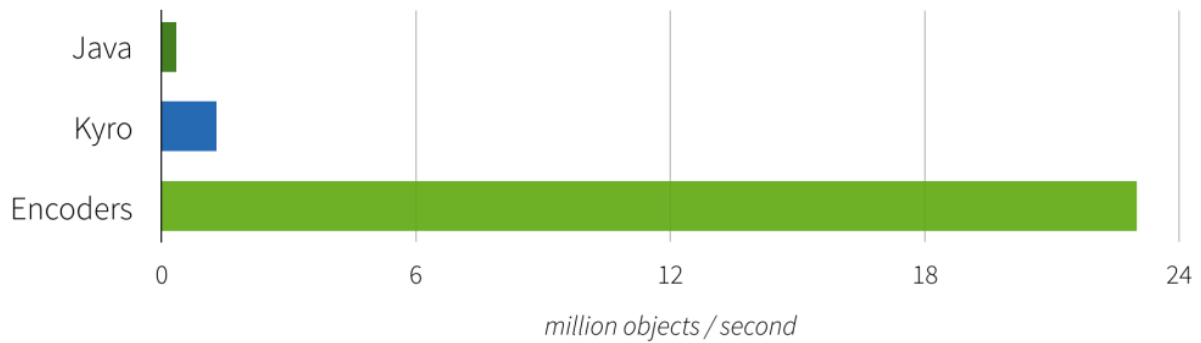
# Chapter 4: Big Data Analytics with Spark SQL, DataFrames, and Datasets







## Serialization / Deserialization Performance



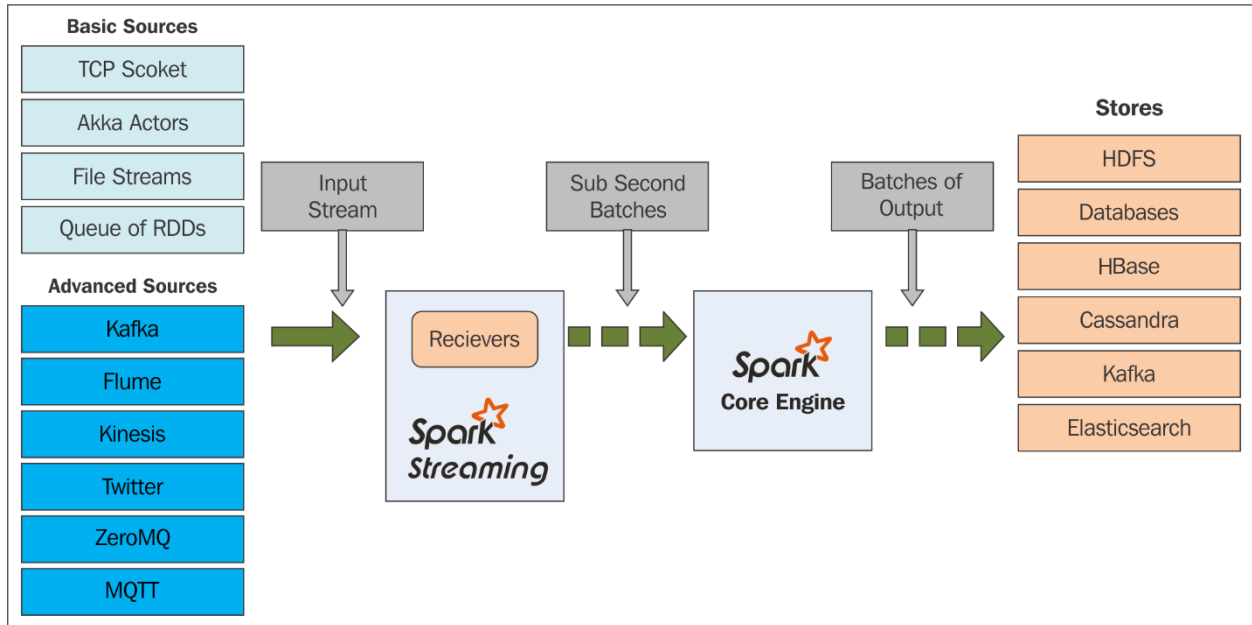
```
>>> mylist = [(50, "DataFrame"), (60, "pandas")]
>>> myschema = ['col1', 'col2']
>>> df = sc.parallelize(mylist).toDF(myschema)
>>> df.printSchema()
root
 |-- col1: long (nullable = true)
 |-- col2: string (nullable = true)

>>> df.show()
+----+-----+
|col1|    col2|
+----+-----+
|  50|DataFrame|
|  60|   pandas|
+----+-----+
```

```
>>> df = spark.read.format('jdbc').options(url='jdbc:mysql://localhost:3306/retail_db?user=root&password=cloudera', dbtable='departments').load()
>>> df.show()
+-----+-----+
|department_id|department_name|
+-----+-----+
|          2|      Fitness|
|          3|    Footwear|
|          4|    Apparel|
|          5|        Golf|
|          6|    Outdoors|
|          7|    Fan Shop|
+-----+-----+

>>> df2rdd = df.rdd
>>> df2rdd.take(2)
[Row(department_id=2, department_name=u'Fitness'), Row(department_id=3, department_name=u'Footwear')]
```

# Chapter 5: Real-Time Analytics with Spark Streaming and Structured Streaming

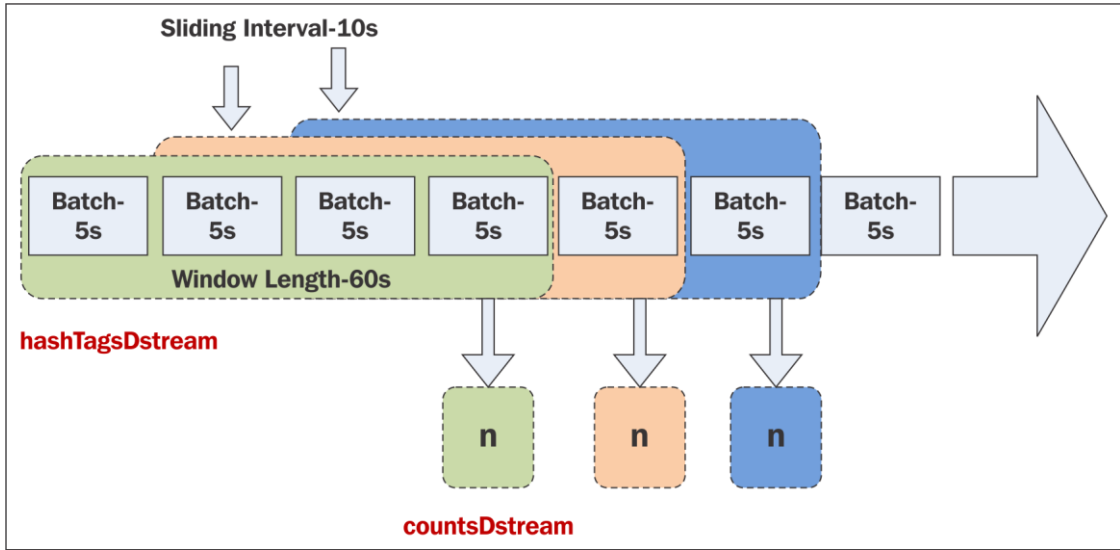


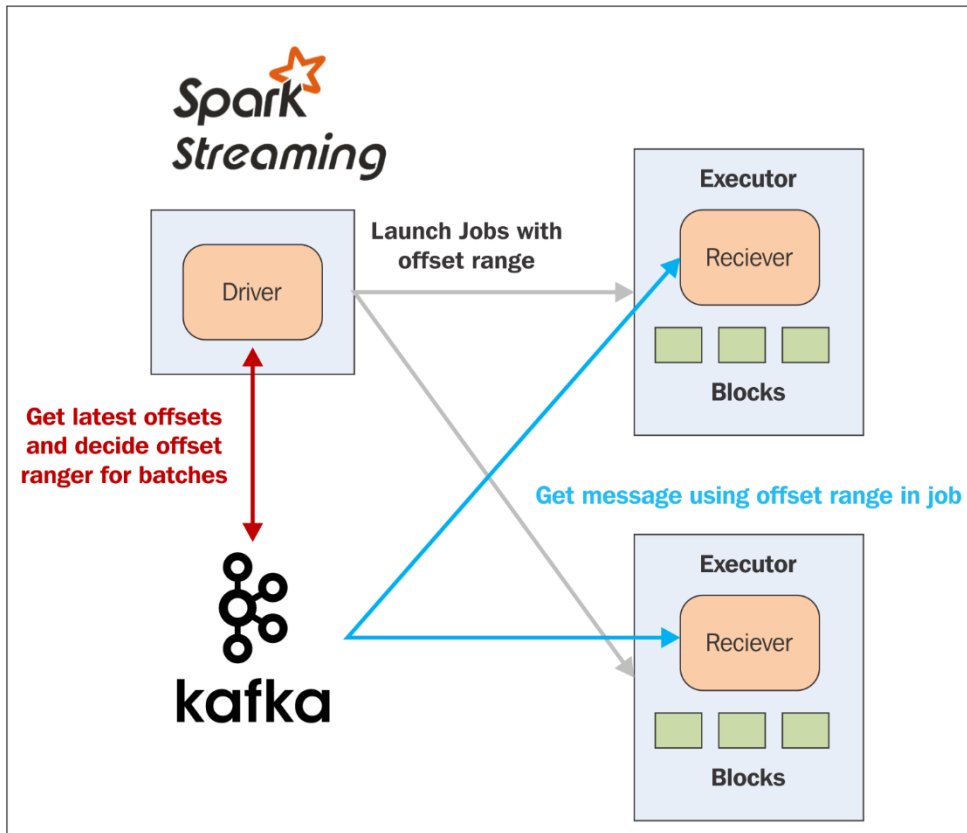
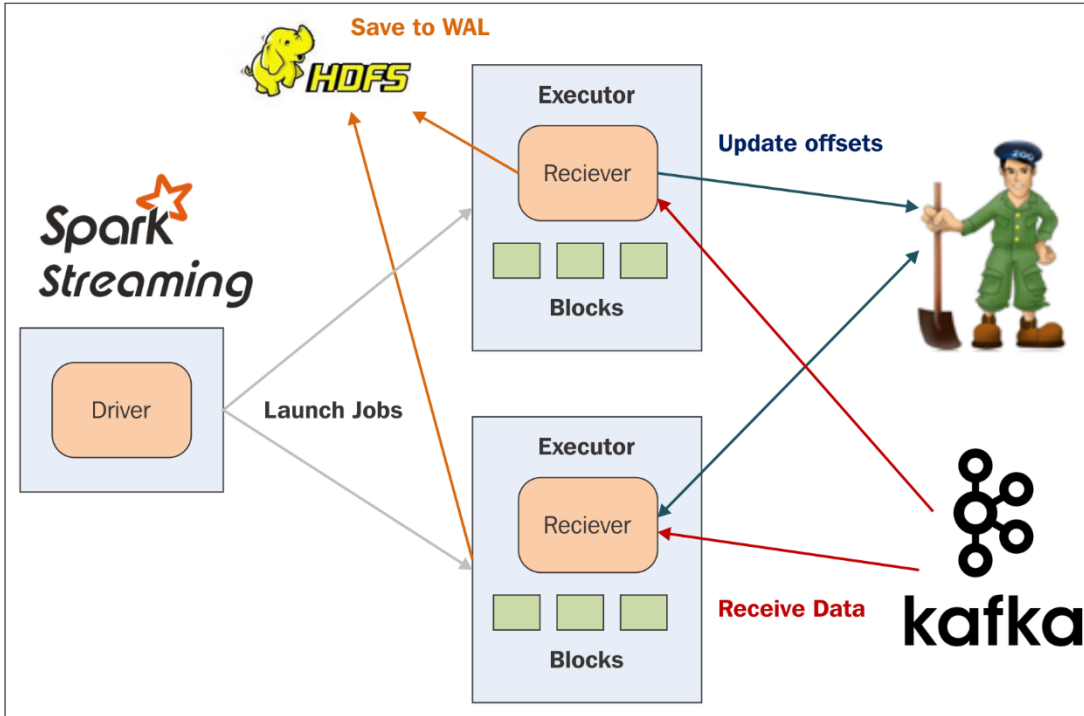
```
[root@quickstart ~]# nc -l 9999
spark
spark
spark
hadoop
hadoop
hadoop

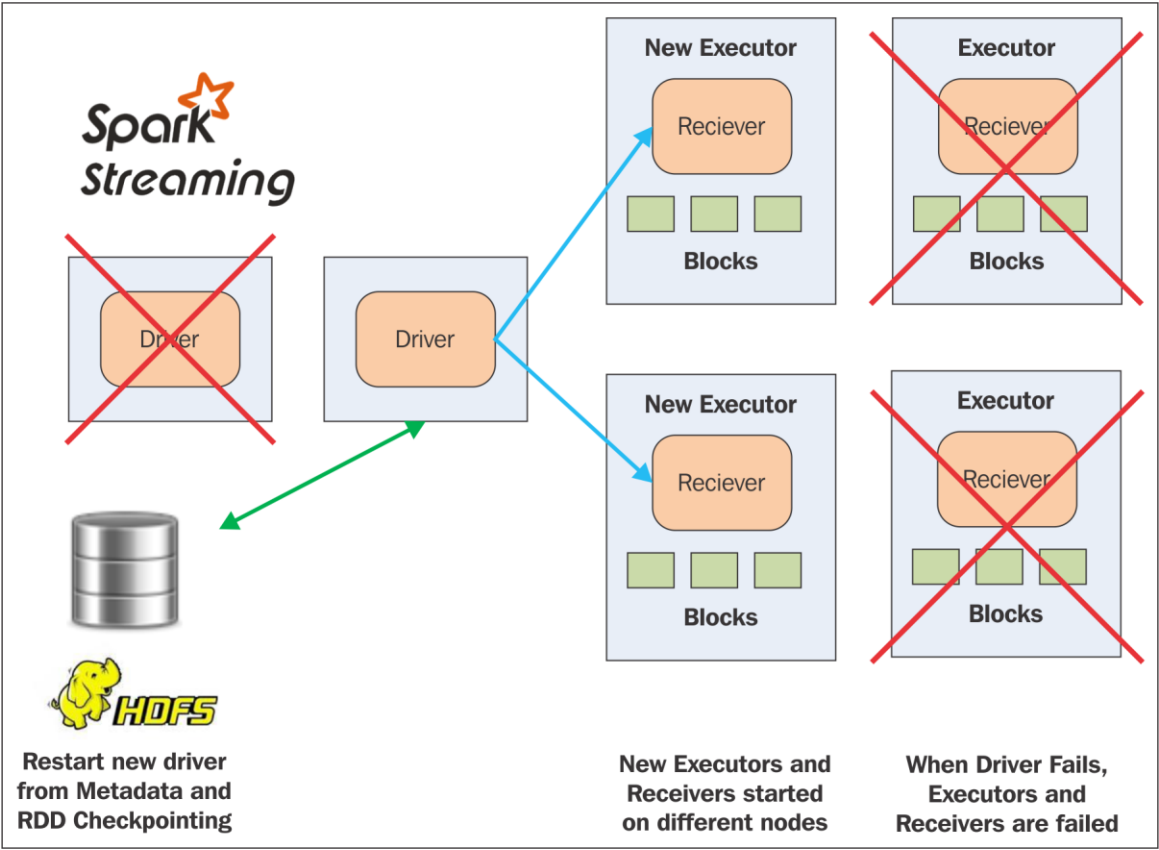
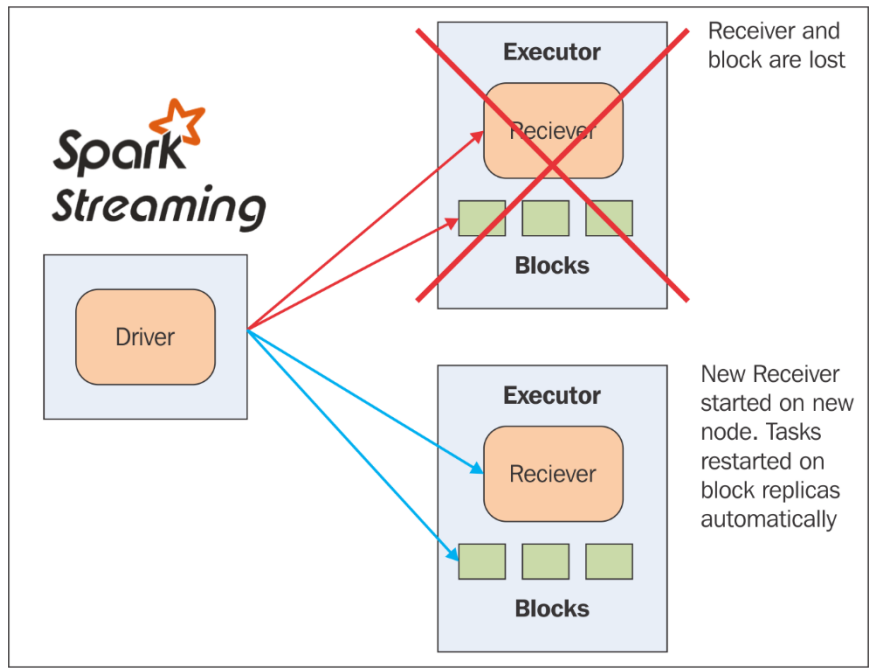
spark
spark
spark
hadoop
hadoop
hadoop
█

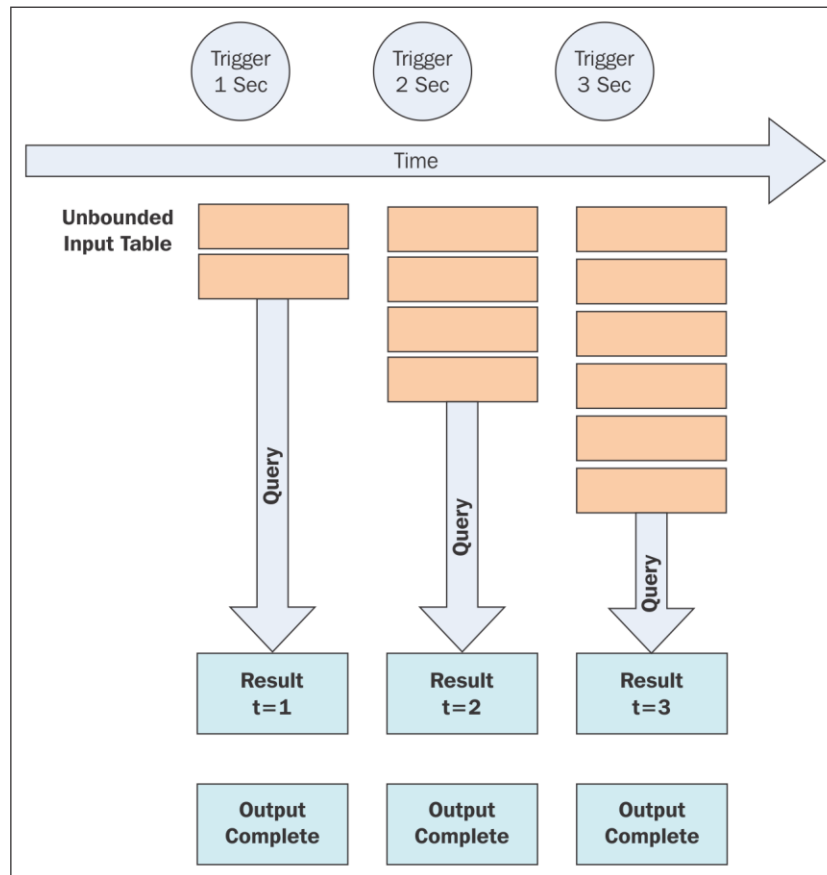
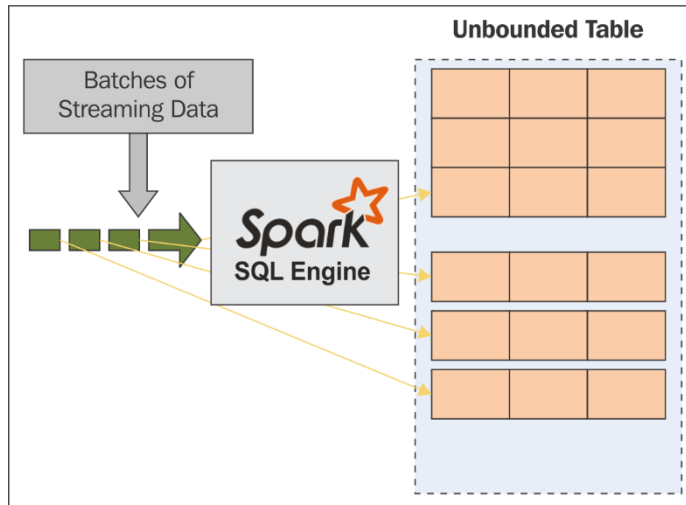
-----
Time: 2015-12-19 06:40:06
-----
(u'', 1)
(u'hadoop', 3)
(u'spark', 3)
-----
Time: 2015-12-19 06:40:07
-----
(u'hadoop', 3)
(u'spark', 3)
-----
Time: 2015-12-19 06:40:08
-----
Time: 2015-12-19 06:40:09
-----
```

```
[root@quickstart ~]# nc -l k 9999
spark
spark
spark
hadoop
hadoop
hadoop
spark
spark
spark
hadoop
hadoop
hadoop
Time: 2015-12-19 06:44:54
-----
(u'hadoop', 3)
(u'spark', 3)
-----
Time: 2015-12-19 06:44:55
-----
(u'hadoop', 5)
(u'spark', 6)
-----
Time: 2015-12-19 06:44:56
-----
(u'hadoop', 6)
(u'spark', 6)
-----
Time: 2015-12-19 06:44:57
-----
(u'hadoop', 6)
(u'spark', 6)
```











# Chapter 6: Notebooks and Dataflows with Spark and Hadoop

192.168.139.170:8889/tree#notebooks

**jupyter**

Files Running Clusters

Select items to perform actions on them. Upload New ↕

bin

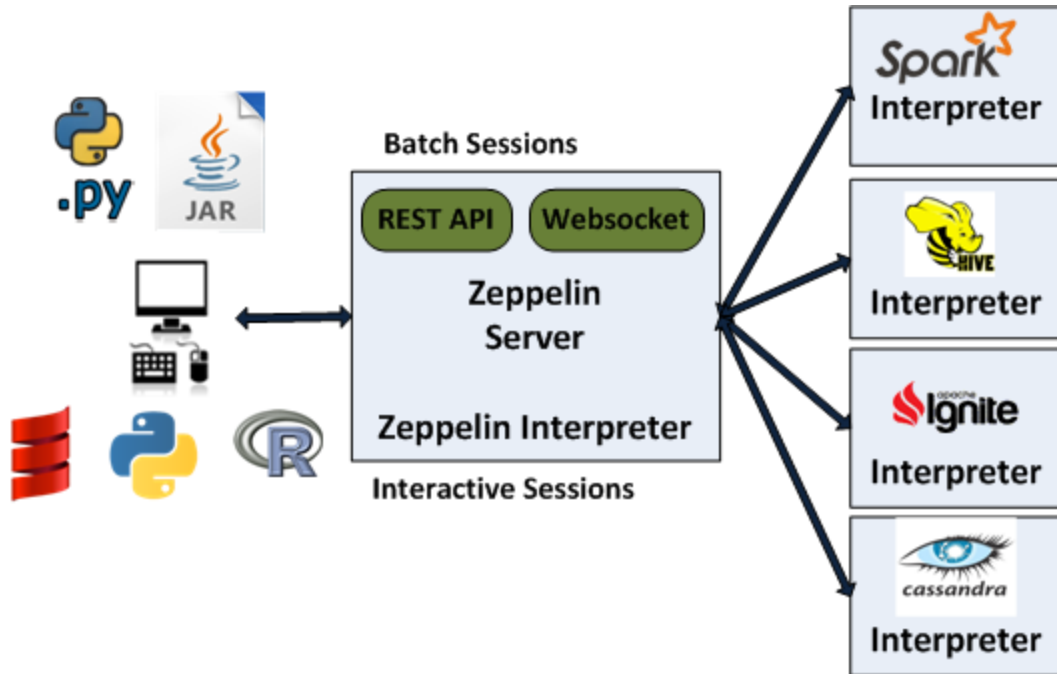
**jupyter** Chapter 6 Notebook (autosaved)

File Edit View Insert Cell Kernel Help

CellToolbar

```
In [7]: plot(counts.collect())
```

Category	Count
spark	3.0
hadoop	2.0
interactive	1.0
notebook	1.0
mapreduce	1.0
analytics	1.0
jupyter	1.0
ipython	1.0



Chapter6ZepNote1

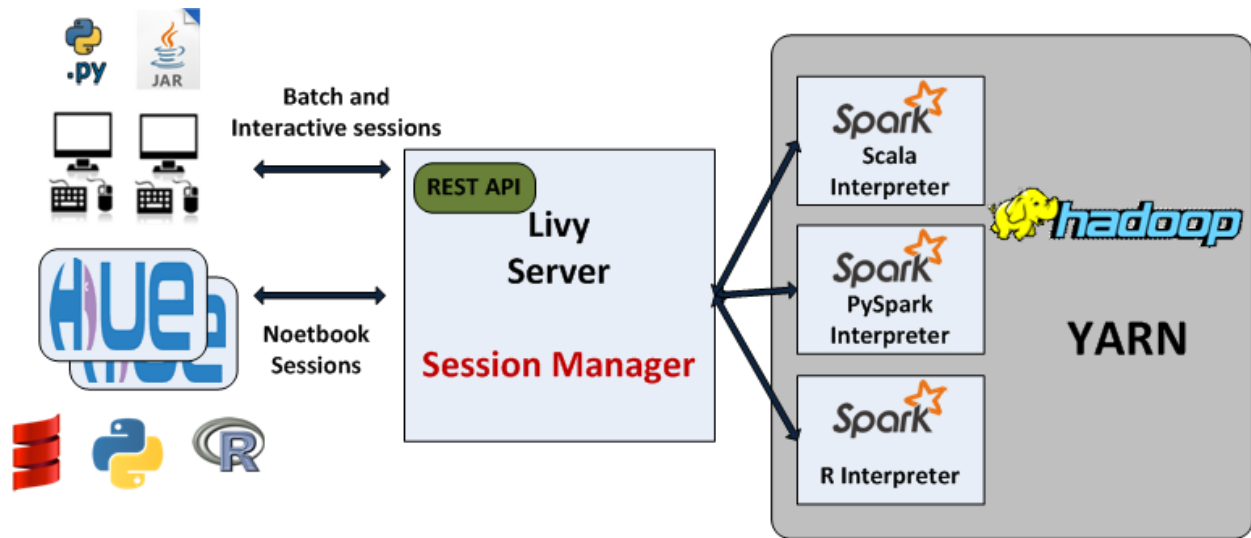
```
%sql select word, max(count) from words group by word
```

FINISHED



word	_c1
json	62
nodemanager	15
during	25
port	3
so	21
flume	178
configurations	149
nimbus	13
removed	60

Took 16 seconds (outdated)



This screenshot shows the Hue web interface. The top navigation bar includes "HUE", "Query Editors", "Notebooks", "Data Browsers", "Workflows", "Search", and "Security". The "Notebooks" menu is open, showing options to "Add Notebook" and "View Notebooks". A "Sample Notebook" is selected. The main content area displays a snippet editor with the text "Add a snippet to start your r". Below this, a circular menu offers various languages: Hive, Impala, PySpark, Scala, and R.

This screenshot shows the Hue web interface with a notebook open. The browser address bar shows "quickstart.cloudera:8888/notebook/notebook". The notebook title is "Chapter 6 Hue Notebook". The main content area shows a snippet titled "My Snippet" with the following Scala code:
 

```
sc.parallelize(range(1000)).map(lambda x: x * x).take(10)
```

 Below the code, the output is displayed as a list:
 

```
[0, 1, 4, 9, 16, 25, 36, 49, 64, 81]
```

 The left sidebar shows a list of tables: customers, sample\_07, sample\_08, and web\_logs.

My Notebook

REGION

country\_code3

VALUE

count(\*)



Hortonworks DataFlow

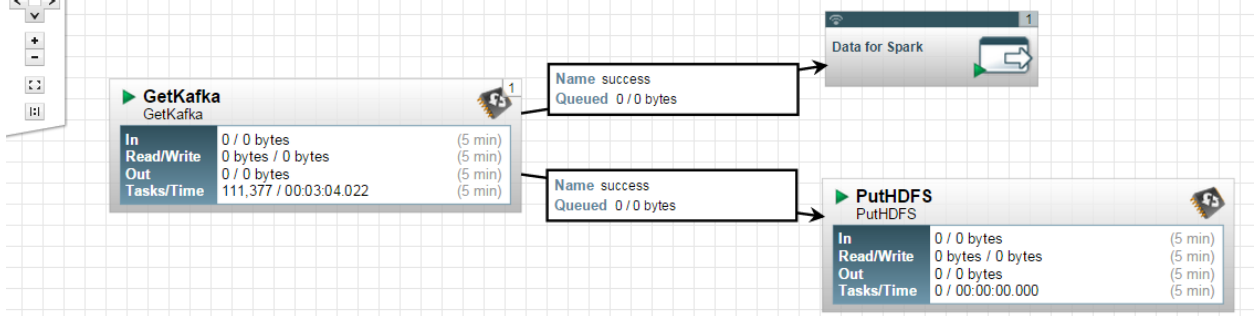
Nifi in Ambari

POWERED

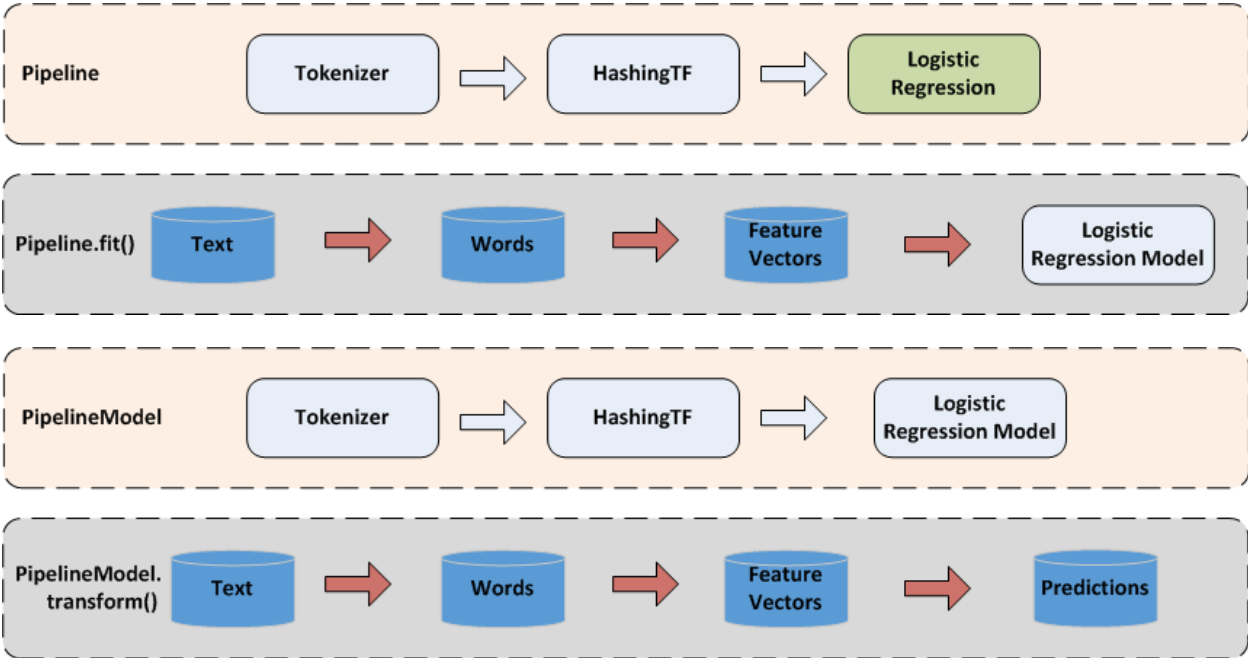
Navigation icons and search bar for the NiFi interface.

NiFi Flow

Active threads: 2 Queued: 0 / 0 bytes Stats last refreshed: 01:21:37 UTC



# Chapter 7: Machine Learning with Spark and Hadoop

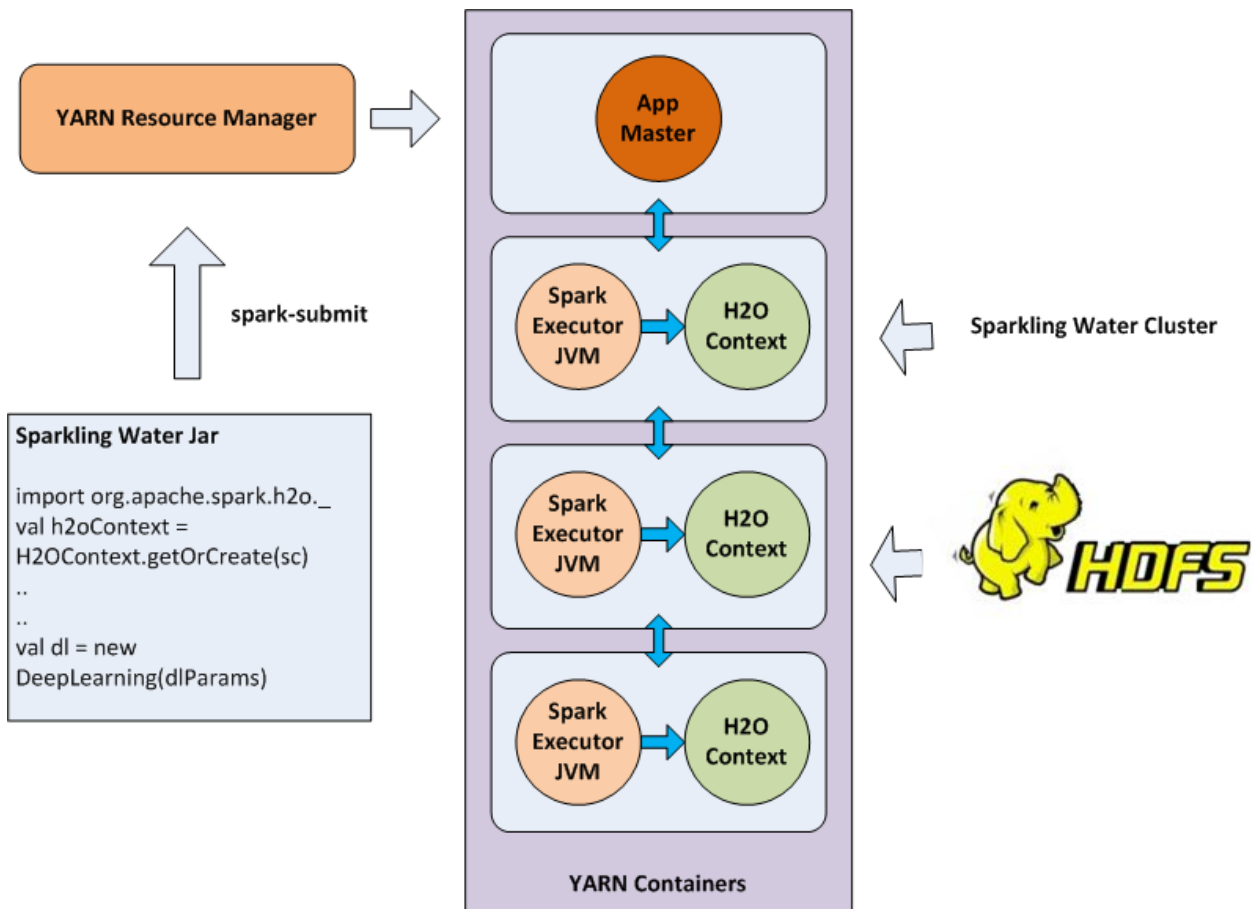


```
+---+-----+-----+-----+
| id|                text|label|
+---+-----+-----+-----+
|  0|apache spark rdd ...|  1.0|
|  1|      mllib pipeline|  0.0|
|  2|      hadoop mahout |  1.0|
|  3|mapreduce iterative|  0.0|
+---+-----+-----+-----+
```

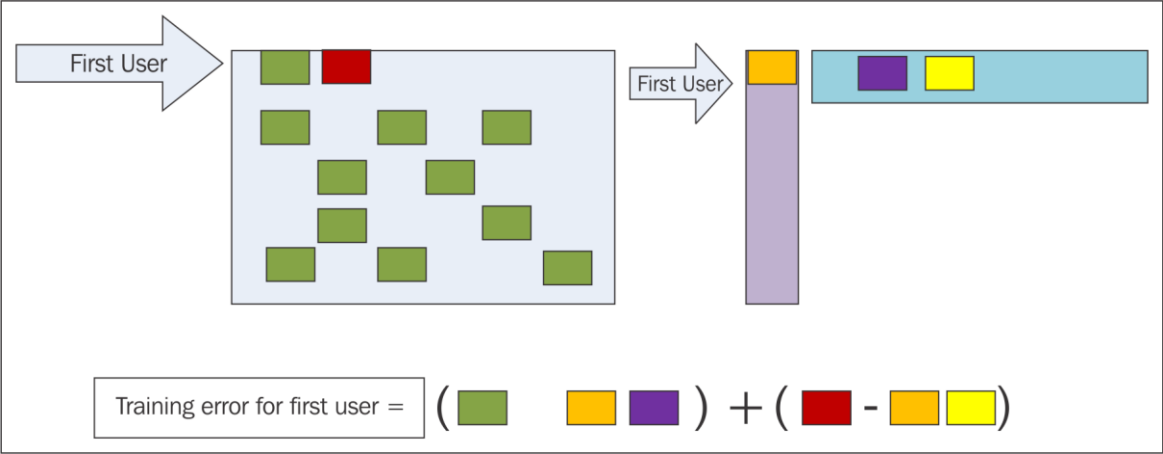
```

+-----+
| id|      text|
+-----+
|  4|    p q r|
|  5|mllib pipeline|
|  6|    x y z|
|  7|  hadoop mahout|
+-----+

```



# Chapter 8: Building Recommendation Systems with Spark and Mahout



```

+-----+-----+-----+-----+
|          name | maxrtng | minrtng | cntusr |
+-----+-----+-----+-----+
| American Beauty (... |    5.0 |    1.0 | 3428 |
| Star Wars: Episod... |    5.0 |    1.0 | 2991 |
| Star Wars: Episod... |    5.0 |    1.0 | 2990 |
| Star Wars: Episod... |    5.0 |    1.0 | 2883 |
| Jurassic Park (1993) |    5.0 |    1.0 | 2672 |
+-----+-----+-----+-----+
only showing top 5 rows

```

```

+-----+-----+
|userid| cnt|
+-----+-----+
| 4169|2314|
| 1680|1850|
| 4277|1743|
| 1941|1595|
| 1181|1521|
+-----+-----+
only showing top 5 rows

```

```

+-----+-----+-----+-----+
|userid|movieid|rating| name|
+-----+-----+-----+-----+
| 4169| 1231| 5.0|Right Stuff, The ...|
| 4169| 232| 5.0|Eat Drink Man Wom...|
| 4169| 3632| 5.0|Monsieur Verdoux ...|
| 4169| 1233| 5.0|Boat, The (Das Bo...|
| 4169| 1834| 5.0|Spanish Prisoner,...|
+-----+-----+-----+-----+
only showing top 5 rows

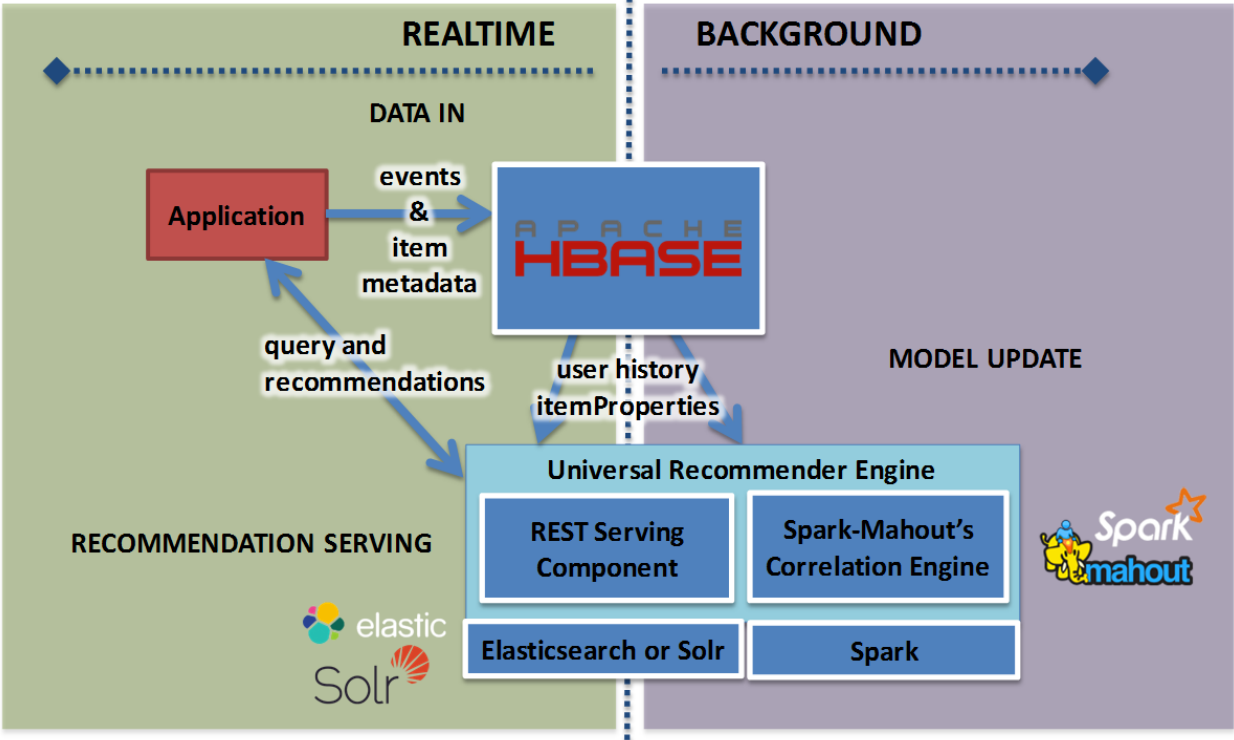
```

```

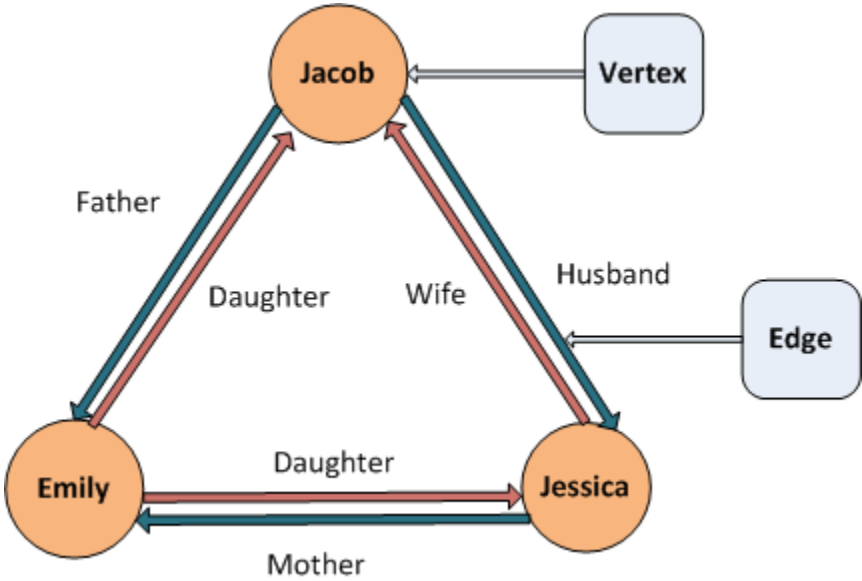
>>> dir(model)
['_class_', '_del_', '_delattr_', '_dict_', '_doc_', '_format_
_', '_getattr_', '_hash_', '_init_', '_module_', '_new_'
, '_reduce_', '_reduce_ex_', '_repr_', '_setattr_', '_sizeof_'
, '_str_', '_subclasshook_', '_weakref_', '_java_loader_class_',
java_model', '_load_java', '_sc', 'call', 'load', 'predict', 'predictAll
', 'productFeatures', 'rank', 'recommendProducts', 'recommendUsers', 'sa
ve', 'userFeatures']

```

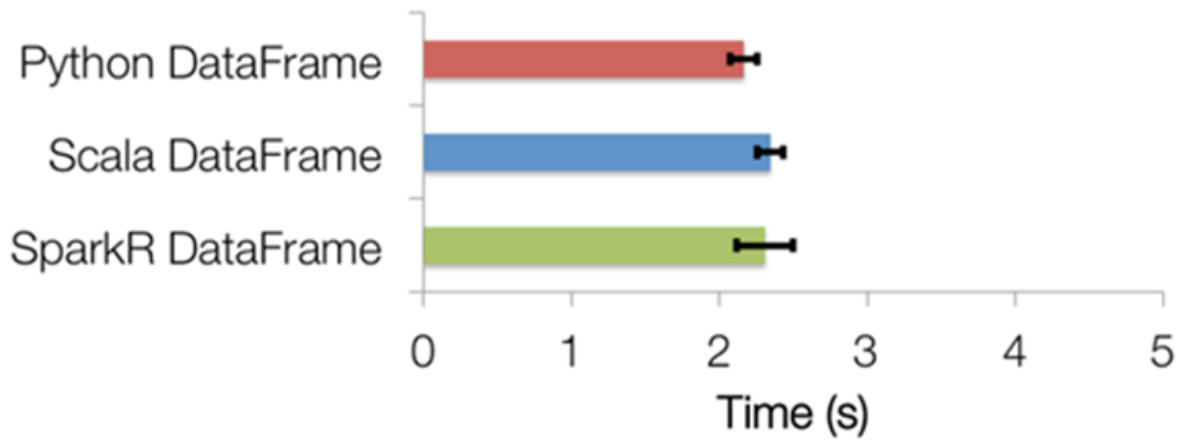
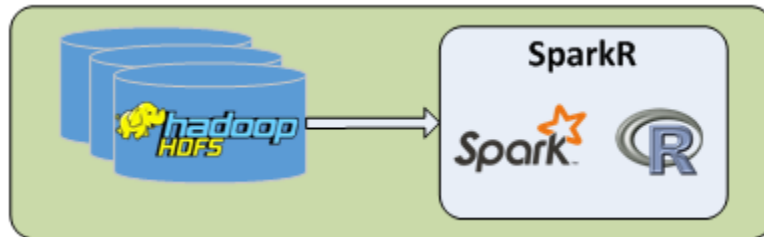
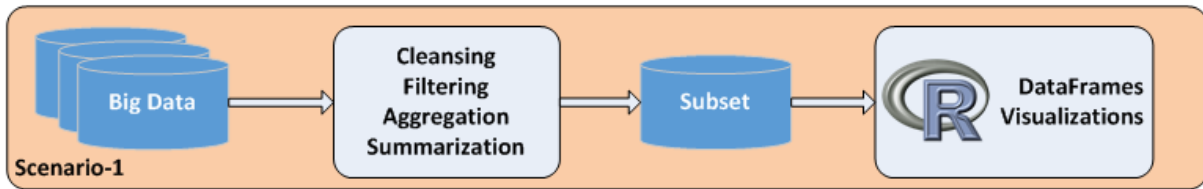


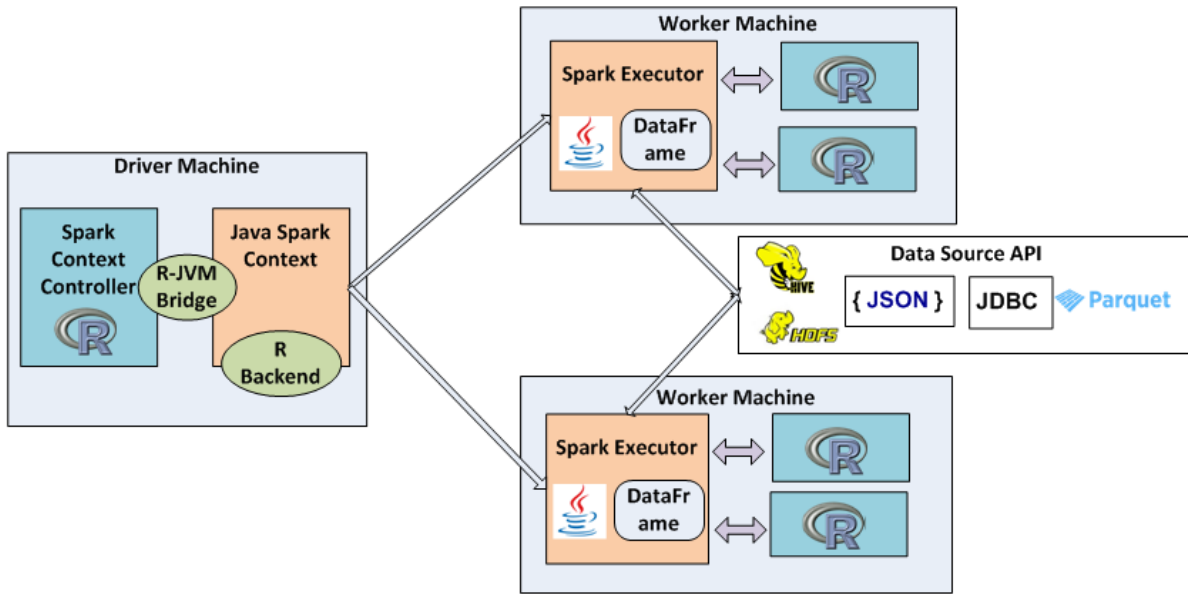


# Chapter 9: Graph Analytics with GraphX



# Chapter 10: Interactive Analytics with SparkR





```

Launching java with spark-submit command /home/cloudera/spark-2.0.0-bin-hadoop2.7/bin/spark-submit --master "local[2]" "sparkr-shell" /tmp/RtmpQ6DC4T/backend_port50fe71dbe06d
Setting default log level to "WARN".
To adjust logging level use sc.setLogLevel(newLevel).

Welcome to

Spark version 2.0.0

SparkSession available as 'spark'.
>

```

```

Launching java with spark-submit command /home/cloudera/spark-2.0.0-bin-hadoop2.7/bin/spark-submit --master "spark://quickstart.cloudera:7077" "sparkr-shell" /tmp/RtmpRknKgY/backend_port566f765d5f12
Setting default log level to "WARN".
To adjust logging level use sc.setLogLevel(newLevel).

Welcome to

Spark version 2.0.0

SparkSession available as 'spark'.
>

```

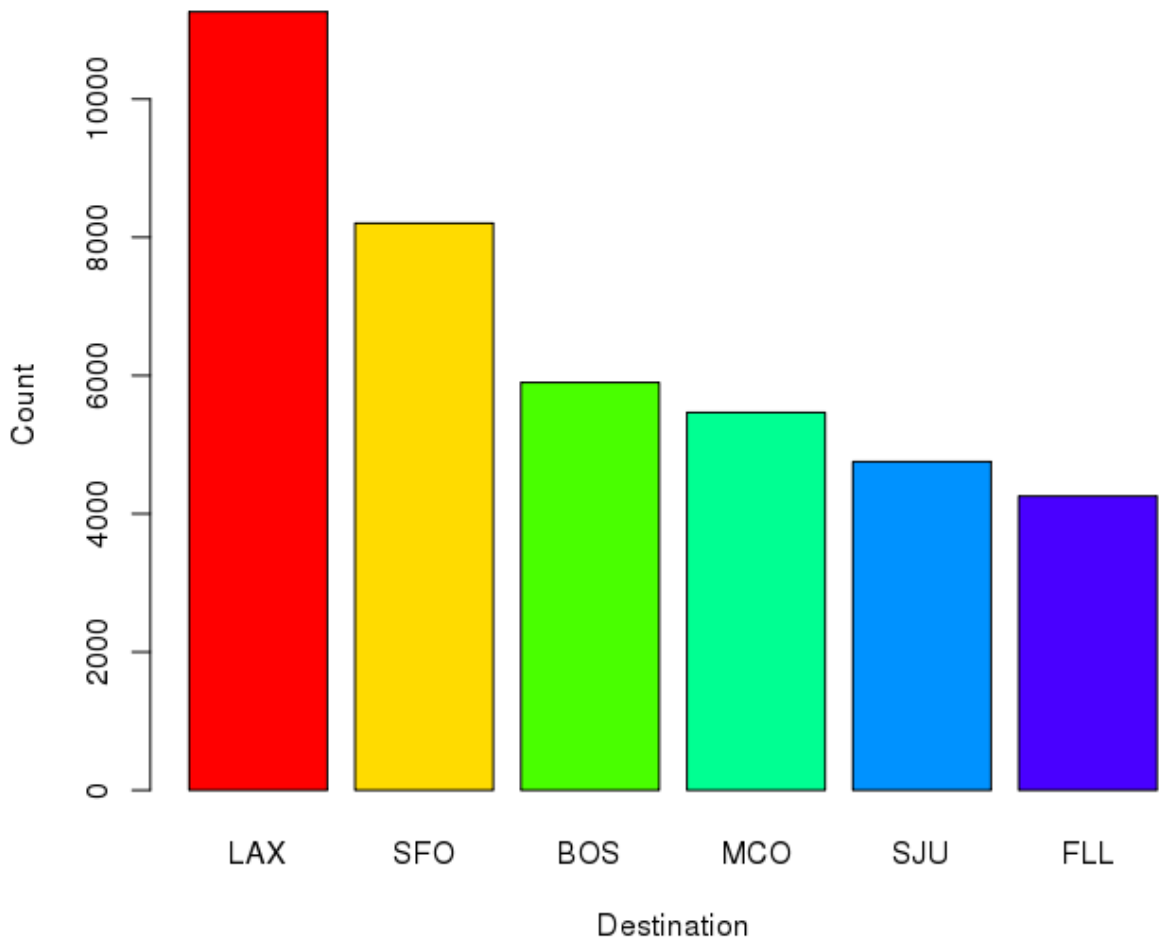
```
Launching java with spark-submit command /home/cloudera/spark-2.0.0-bin-hadoop2.7/bin/spark-submit --master "yarn" "sparkr-shell" /tmp/RtmprecjdS/backend_port1bc3213922d9
Setting default log level to "WARN".
To adjust logging level use sc.setLogLevel(newLevel).

Welcome to

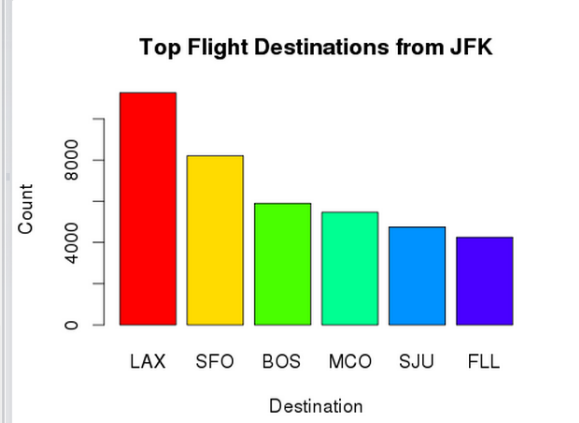
  Spark version 2.0.0

SparkSession available as 'spark'.
>
```

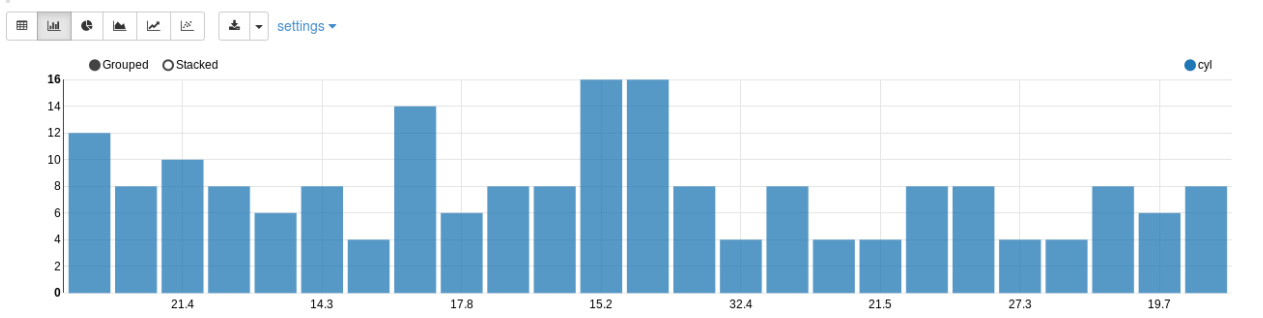
### Top Flight Destinations from JFK



```
Source
Console ~/
+ group_by(flights$dest) %>%
+   summarize(count = n(flights$dest))
+
> top_dests <- head(arrange(jfk_dest, desc(jfk_dest$count)))
[Stage 1:>
1) / 2][Stage 1:=====>
(1 + 1) / 2][Stage 2:>
(0 + 0) / 200][Stage 2:====>
(14 + 1) / 200][Stage 2:====>
(22 + 1) / 200][Stage 2:====>
(32 + 1) / 200][Stage 2:====>
(42 + 1) / 200][Stage 2:====>
(53 + 1) / 200][Stage 2:====>
(67 + 1) / 200][Stage 2:====>
(82 + 1) / 200][Stage 2:====>
(95 + 1) / 200][Stage 2:====>
(109 + 1) / 200][Stage 2:====>
(127 + 1) / 200][Stage 2:====>
(144 + 2) / 200][Stage 2:====>
(159 + 1) / 200][Stage 2:====>
(175 + 1) / 200][Stage 2:====>
(196 + 1) / 200][Stage 2:====>
> barplot(top_dests$count, names.arg = top_dests$dest,col=rainbow(7),main="Top Flight Destinations from JFK", xlab = "Destination", ylab= "Count", beside=TRUE )
```



%sql select \* from carstable FINISHED



%sql select \* from carstable FINISHED

