

splunk> .conf2017

# Advanced Analytics With Splunk Using Apache Spark Machine Learning And Spark Graph

Raanan Dagan | Architect

September 25, 2017 | Washington, DC

# Forward-Looking Statements

During the course of this presentation, we may make forward-looking statements regarding future events or the expected performance of the company. We caution you that such statements reflect our current expectations and estimates based on factors currently known to us and that actual events or results could differ materially. For important factors that may cause actual results to differ from those contained in our forward-looking statements, please review our filings with the SEC.

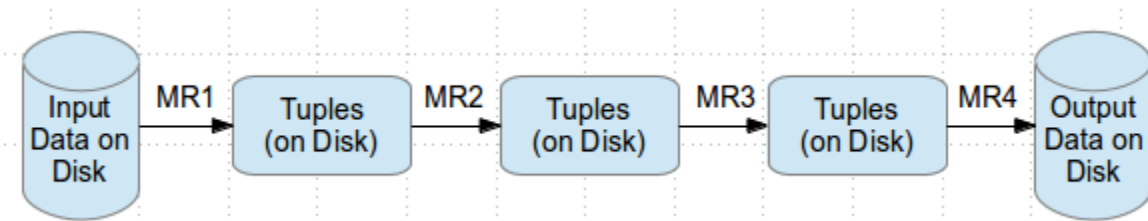
The forward-looking statements made in this presentation are being made as of the time and date of its live presentation. If reviewed after its live presentation, this presentation may not contain current or accurate information. We do not assume any obligation to update any forward looking statements we may make. In addition, any information about our roadmap outlines our general product direction and is subject to change at any time without notice. It is for informational purposes only and shall not be incorporated into any contract or other commitment. Splunk undertakes no obligation either to develop the features or functionality described or to include any such feature or functionality in a future release.

Splunk, Splunk>, Listen to Your Data, The Engine for Machine Data, Splunk Cloud, Splunk Light and SPL are trademarks and registered trademarks of Splunk Inc. in the United States and other countries. All other brand names, product names, or trademarks belong to their respective owners. © 2017 Splunk Inc. All rights reserved.

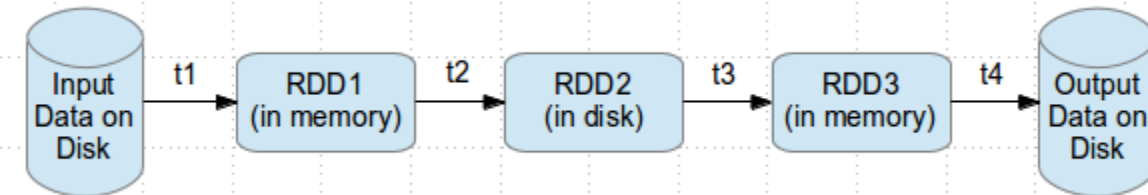
# Why Spark?

- ▶ Most of machine learning algorithms are iterative because each iteration can improve the results
- ▶ With disk based approach each iteration's output is written to disk, making it slow

## Hadoop execution flow



## Spark execution flow



<http://www.wiziq.com/blog/hype-around-apache-spark/>

# About Apache Spark



- ▶ Initially started at UC Berkeley in 2009
- ▶ Fast and general purpose cluster computing system
- ▶ 10x (on disk) - 100x (In-Memory) faster
- ▶ Most popular for running *Iterative Machine Learning Algorithms*.
- ▶ Provides high level APIs in
  - Java, Scala, Python
- ▶ Integration with Hadoop and its ecosystem and can **read existing data**

<http://spark.apache.org/>





# Spark Core

---



# Resilient Distributed Dataset (RDD)

- ▶ Resilient Distributed Dataset (RDD) is a basic abstraction in Spark
- ▶ Immutable, partitioned collection of elements that can be operated in parallel
- ▶ Basic Operations
  - map
  - filter
  - persist
- ▶ Multiple Implementation
  - [PairRDDFunctions](#) : RDD of Key-Value Pairs, groupByKey, Join
- ▶ RDD main characteristics:
  - A list of partitions
  - A function for computing each split

```
130.60.4 - - [07/Jan 18:10:57:153] "GET /category.screen?category_id=GIFTS&JSESSIONID=5D15L9FF10ADFF10 HTTP 1.1" 404 720 "http://buttercup-shopping.com/cart.do?action=view&itemId=EST-6&product_id=F1-5W-01" "Mozilla/5.0 (Macintosh; Intel Mac OS X 10_11_2; rv:53.0) Gecko/20100801 Firefox/53.0"
128.241.220.82 - - [07/Jan 18:10:57:123] "GET /product.screen?product_id=FL-DSH-01&JSESSIONID=5D35L7FF6ADFF0 HTTP 1.1" 404 3322 "http://buttercup-shopping.com/cart.do?action=purchase&itemId=EST-26&product_id=K0-CU-01" "Mozilla/5.0 (Macintosh; Intel Mac OS X 10_11_2; rv:53.0) Gecko/20100801 Firefox/53.0"
317.27.160.0.0 - - [07/Jan 18:10:56:156] "GET /oldlink?item_id=EST-26&JSESSIONID=5D55L9FF1ADFF3 HTTP 1.1" 200 1318 "http://buttercup-shopping.com/cart.do?action=changequantity&itemId=EST-18&product_id=AV-CB-01&JSESSIONID=5D55L9FF1ADFF3" "Mozilla/5.0 (Macintosh; Intel Mac OS X 10_11_2; rv:53.0) Gecko/20100801 Firefox/53.0"
10.0.0.0 - - [07/Jan 18:10:55:187] "GET /category.screen?category_id=FLOWERS&JSESSIONID=5D55L9FF1ADFF3 HTTP 1.1" 200 3865 "http://buttercup-shopping.com/cart.do?action=remove&itemId=EST-1" "Mozilla/5.0 (Macintosh; Intel Mac OS X 10_11_2; rv:53.0) Gecko/20100801 Firefox/53.0"
10.0.0.0 - - [07/Jan 18:10:55:189] "GET /category.screen?category_id=FLOWERS&JSESSIONID=5D55L9FF1ADFF3 HTTP 1.1" 200 3865 "http://buttercup-shopping.com/cart.do?action=remove&itemId=EST-1" "Mozilla/5.0 (Macintosh; Intel Mac OS X 10_11_2; rv:53.0) Gecko/20100801 Firefox/53.0"
```



# Spark SQL DataFrames

---





# SQLContext

- ▶ Most powerful way to use Spark SQL is inside a Spark application
- ▶ Load data and query it with SQL while simultaneously combining it with “regular” program code utilizing SQLContext or HiveContext

```
// SQL Imports
// Import Spark SQL. If you can't have
the
// hive dependencies

import org.apache.spark.sql.SQLContext

// Construct SQL Context
val sqlContext = new SQLContext(...)
```

```
// SQL Imports
// Import Spark SQL

import
org.apache.spark.sql.hive.HiveContext

// Construct Hive Context
val hiveContext = new HiveContext(...)
```



# DataFrames

- ▶ Offers rich relational/procedural integration within Spark programs
- ▶ DataFrames:
  - Collections of structured records that can be manipulated using Spark's procedural API or new relational API
  - Perform relational operations on DataFrames using a domain-specific language (DSL) similar to R data frames and Python Pandas
  - Pass Scala, Java or Python functions through DataFrames to build a logical plan
  - Create directly from Spark's distributed objects
  - Enable relational processing in existing Spark programs
- ▶ Automatically store data in a columnar format
- ▶ Go through a relational optimizer, *Catalyst*
- ▶ Standard data representation in a new "ML pipeline" API for machine learning

# Query Federation To External Databases

- ▶ Data pipelines often combine data from heterogeneous sources
- ▶ Spark SQL data sources leverage Catalyst to push predicates down into the data sources whenever possible

## Example: Use JDBC data source and JSON data source to join two tables together

```

▶ CREATE TEMPORARY TABLE users USING jdbc
OPTIONS(driver "mysql" url "jdbc:mysql://userDB/users ")

▶ CREATE TEMPORARY TABLE logs
USING json OPTIONS (path "logs.json")

▶ SELECT users.id,users.name,logs.message
FROM users JOIN logs WHERE users.id=logs.userId
AND users.registrationDate > "2015-01-01"

```

# Spark MLlib

---



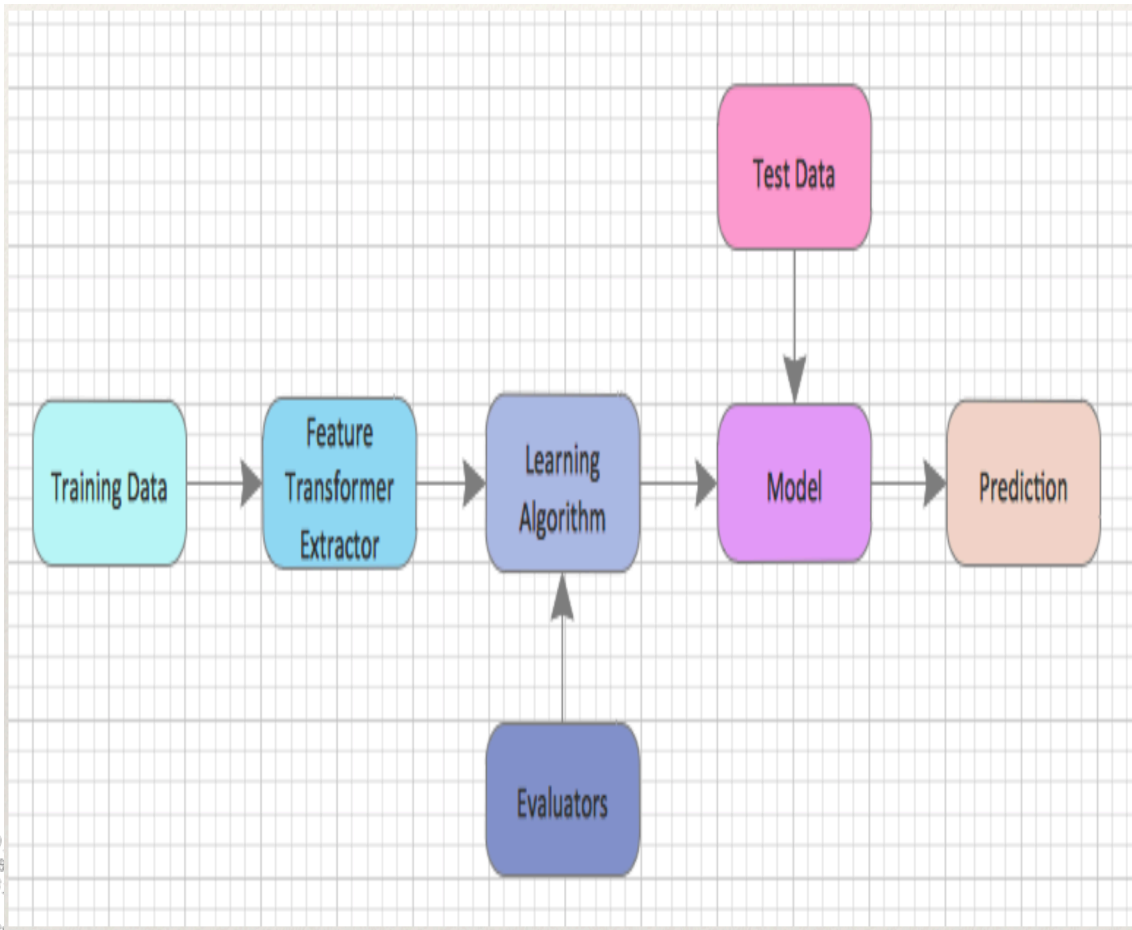
# Spark Machine Learning Basics

## ML algorithms include:

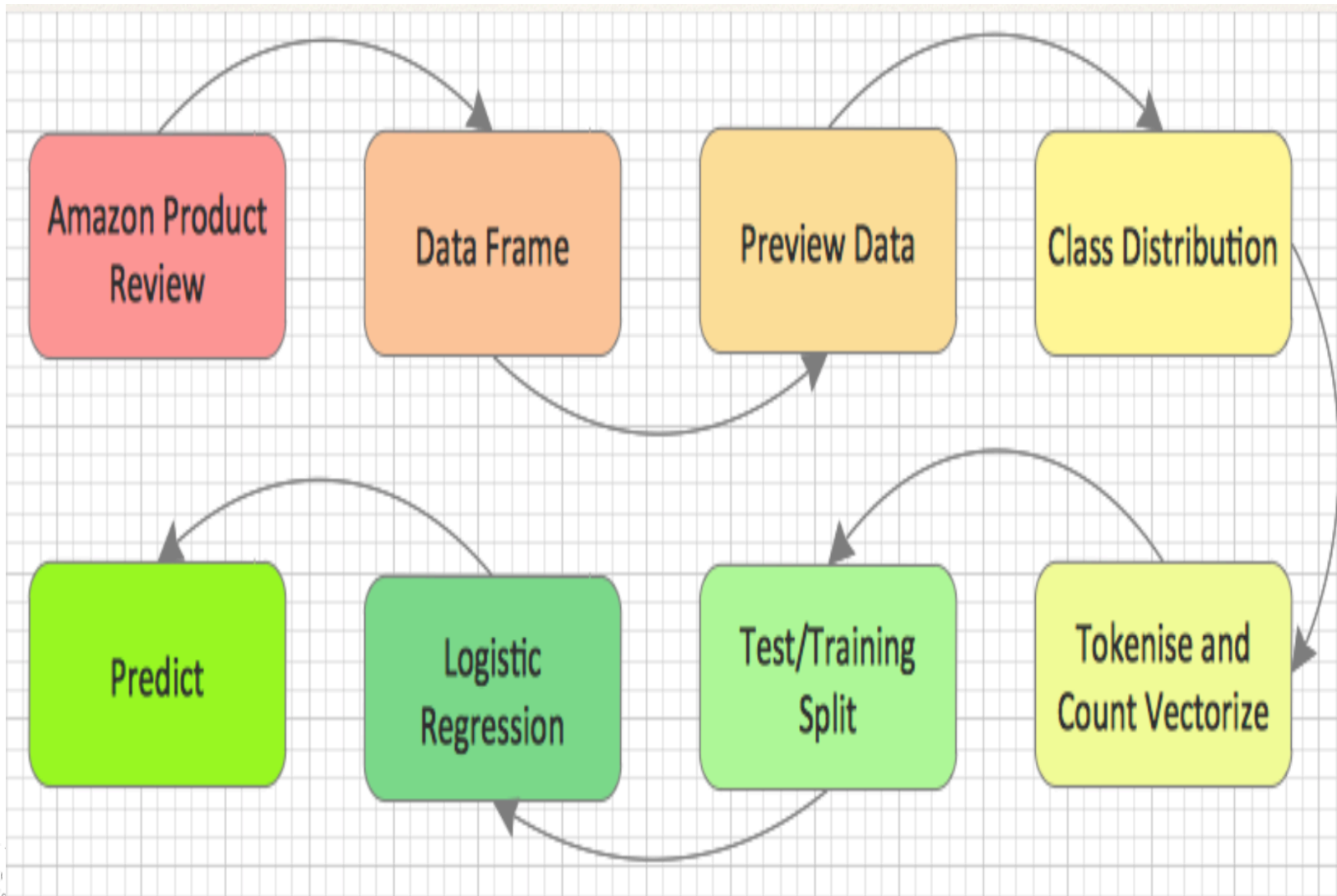
- ▶ Classification: logistic regression, naive Bayes,...
- ▶ Regression: generalized linear regression, survival regression...
- ▶ Decision trees, random forests, and gradient-boosted trees
- ▶ Recommendation: alternating least squares (ALS)
- ▶ Clustering: K-means, Gaussian mixtures (GMMs),...
- ▶ Topic modeling: latent Dirichlet allocation (LDA)
- ▶ Frequent itemsets, association rules, and sequential pattern mining

## ML workflow utilities include:

- ▶ Feature transformations: standardization, normalization, hashing,...
- ▶ ML Pipeline construction
- ▶ Model evaluation and hyper-parameter tuning
- ▶ ML persistence: saving and loading models and Pipelines
- ▶ Distributed linear algebra: SVD, PCA,...
- ▶ Statistics: summary statistics, hypothesis testing,...



# Spark Classification ML Example



Supervised learning for predicting discrete labels

Multiple algorithms

- ▶ **logistic regression**
- ▶ Decision tree classifier
- ▶ Random forest classifier
- ▶ Gradient boosted tree classifier
- ▶ Multi-layer neural network classifier

# Spark Classification ML Code Example

1

Extract Fields

```
val trainingDataTable = sql("""
SELECT e.action
      u.age,
      u.latitude,
      u.logitude
FROM Users u
JOIN Events e
ON u.userId = e.userId""")
```

2

Build Model

```
val trainingData = trainingDataTable.map { row =>

val model =

  new LogisticRegressionWithSGD().run(trainingData)
```

3

Predict

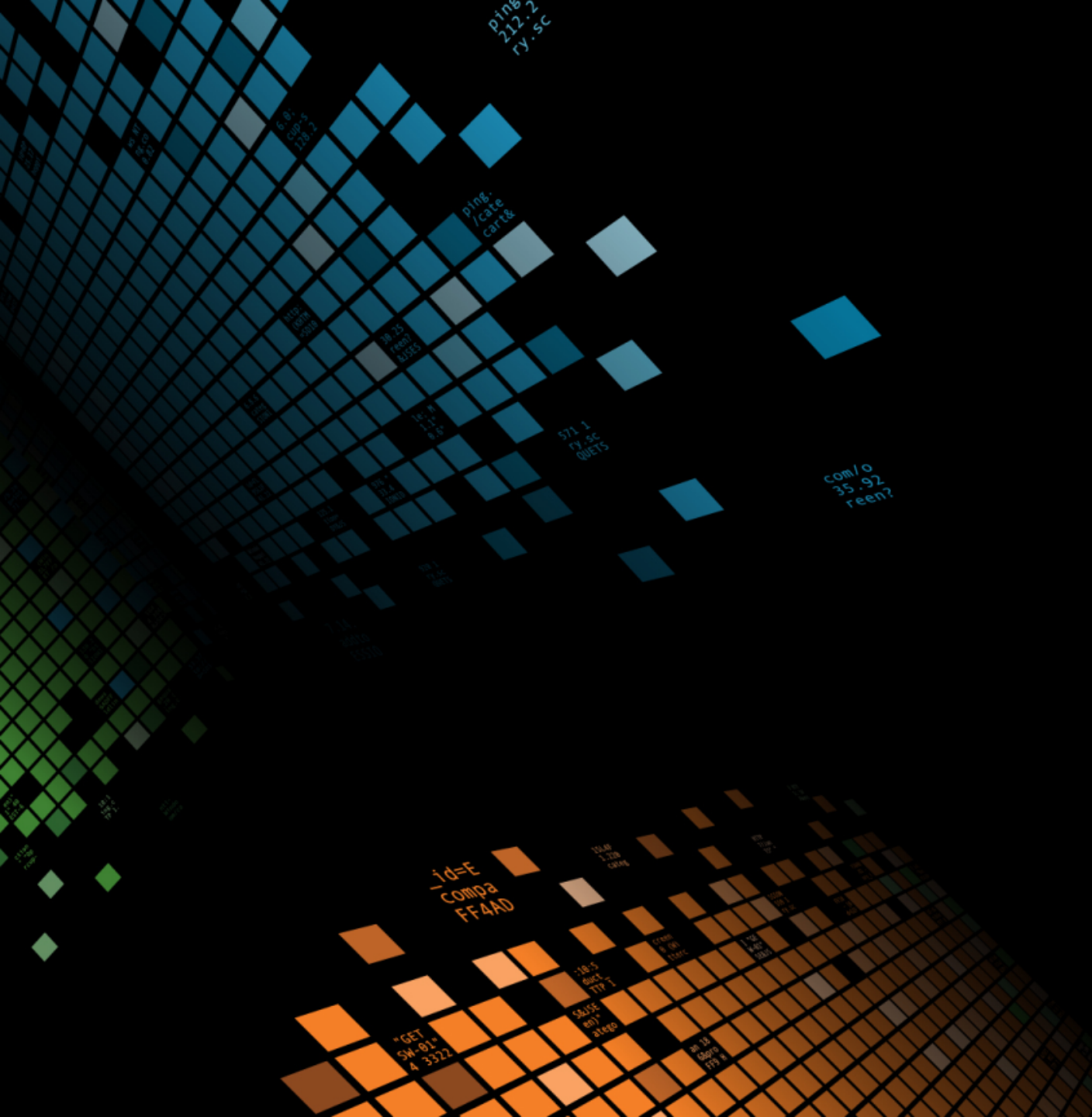
```
case class Score(userId: Int, score: Double)

val scores = allCandidates.map { row =>

  val features = Array[Double](row(1), row(2), row(3))

  Score(row(0), model.predict(features))

}
```



# Spark GraphX

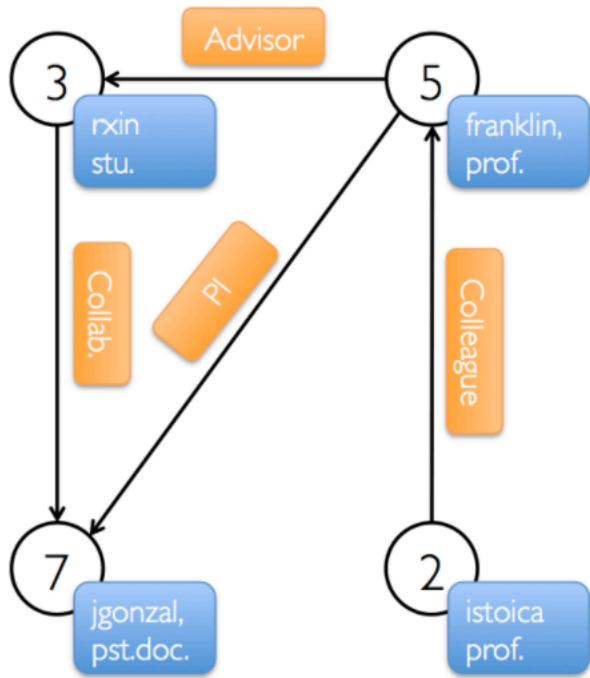
---





# Spark GraphX Example

## Property Graph



## Vertex Table

Id	Property (V)
3	(rxin, student)
7	(jgonzal, postdoc)
5	(franklin, professor)
2	(istoica, professor)

## Edge Table

Srcld	Dstld	Property (E)
3	7	Collaborator
5	3	Advisor
2	5	Colleague
5	7	PI

```
// Assume the SparkContext has already been constructed
```

```
val sc: SparkContext
```

```
// Create an RDD for the vertices
```

```
val users: RDD[(VertexId, (String, String))] =
  sc.parallelize(Array((3L, ("rxin", "student")), (7L, ("jgonzal", "postdoc")),
    (5L, ("franklin", "prof")), (2L, ("istoica", "prof"))))
```

```
// Create an RDD for edges
```

```
val relationships: RDD[Edge[String]] =
  sc.parallelize(Array(Edge(3L, 7L, "collab"), Edge(5L, 3L, "advisor"),
    Edge(2L, 5L, "colleague"), Edge(5L, 7L, "pi")))
```

```
// Define a default user in case there are relationship with missing user
```

```
val defaultUser = ("John Doe", "Missing")
```

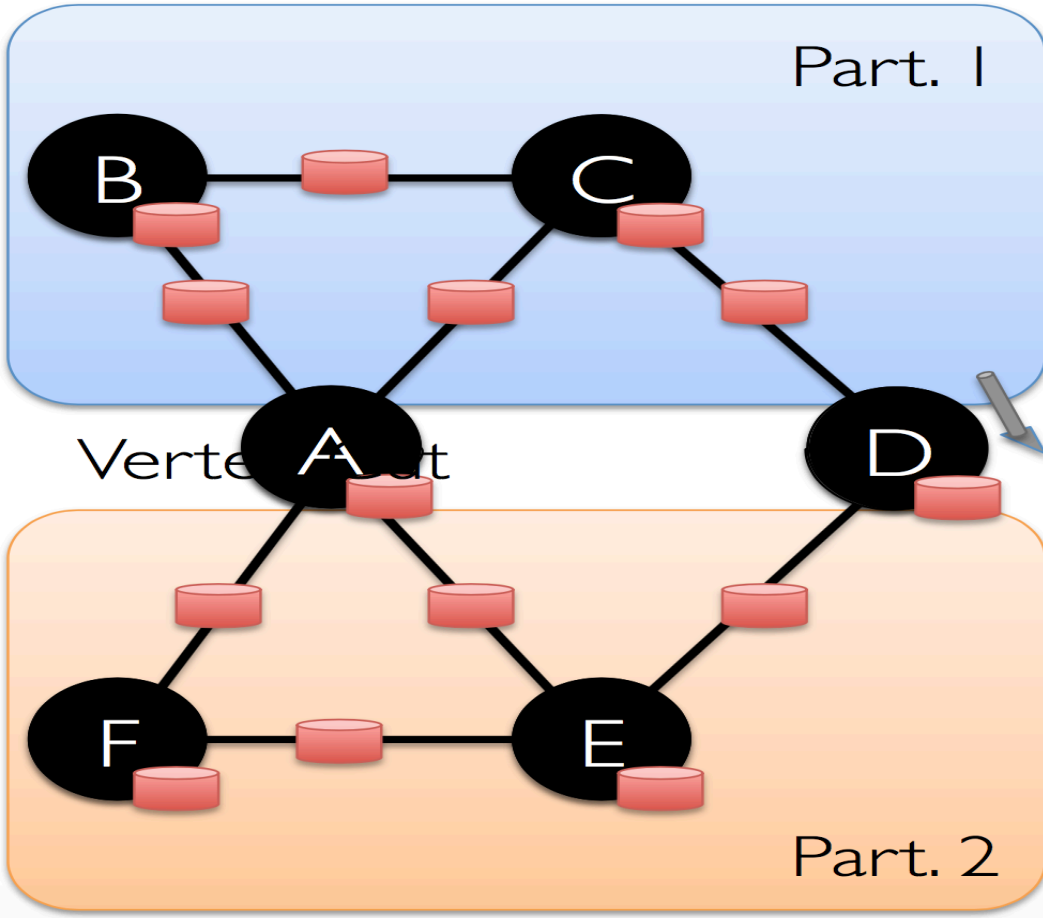
```
// Build the initial Graph
```

```
val graph = Graph(users, relationships, defaultUser)
```

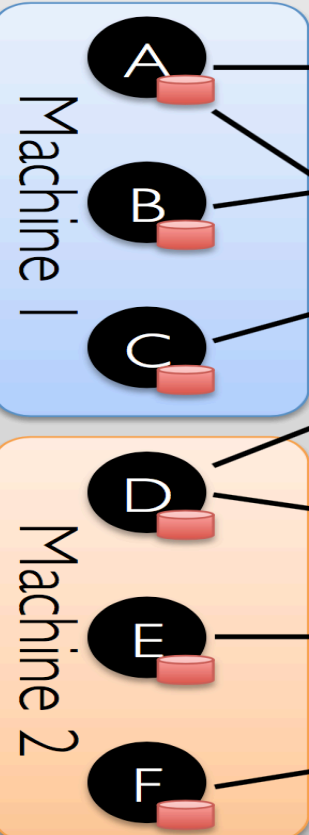


# Spark GraphX Architecture

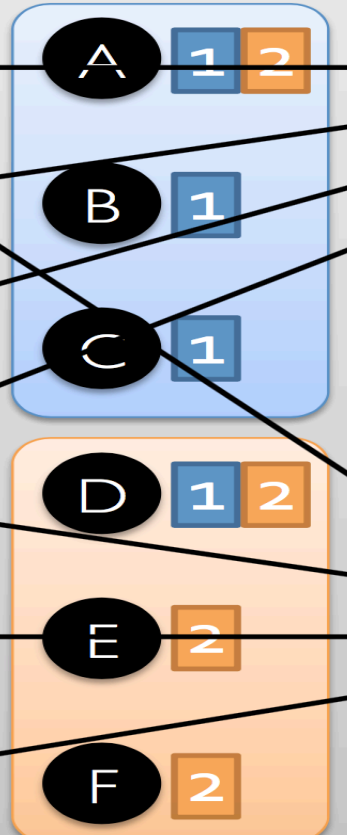
## Property Graph



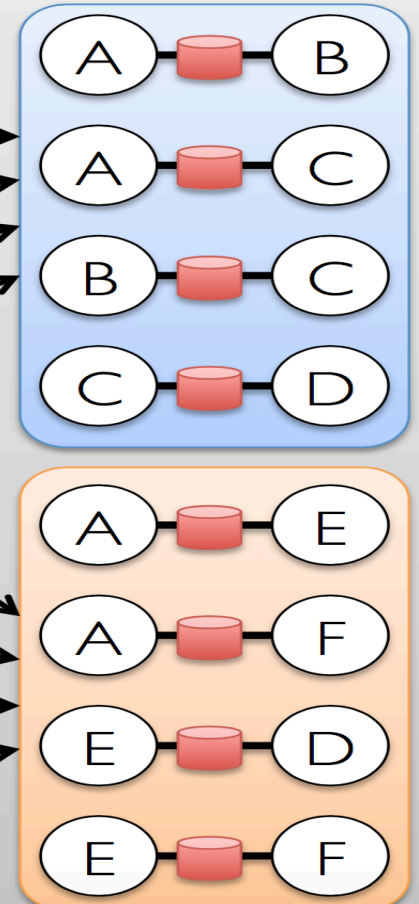
## Vertex Table (RDD)



## Routing Table (RDD)



## Edge Table (RDD)

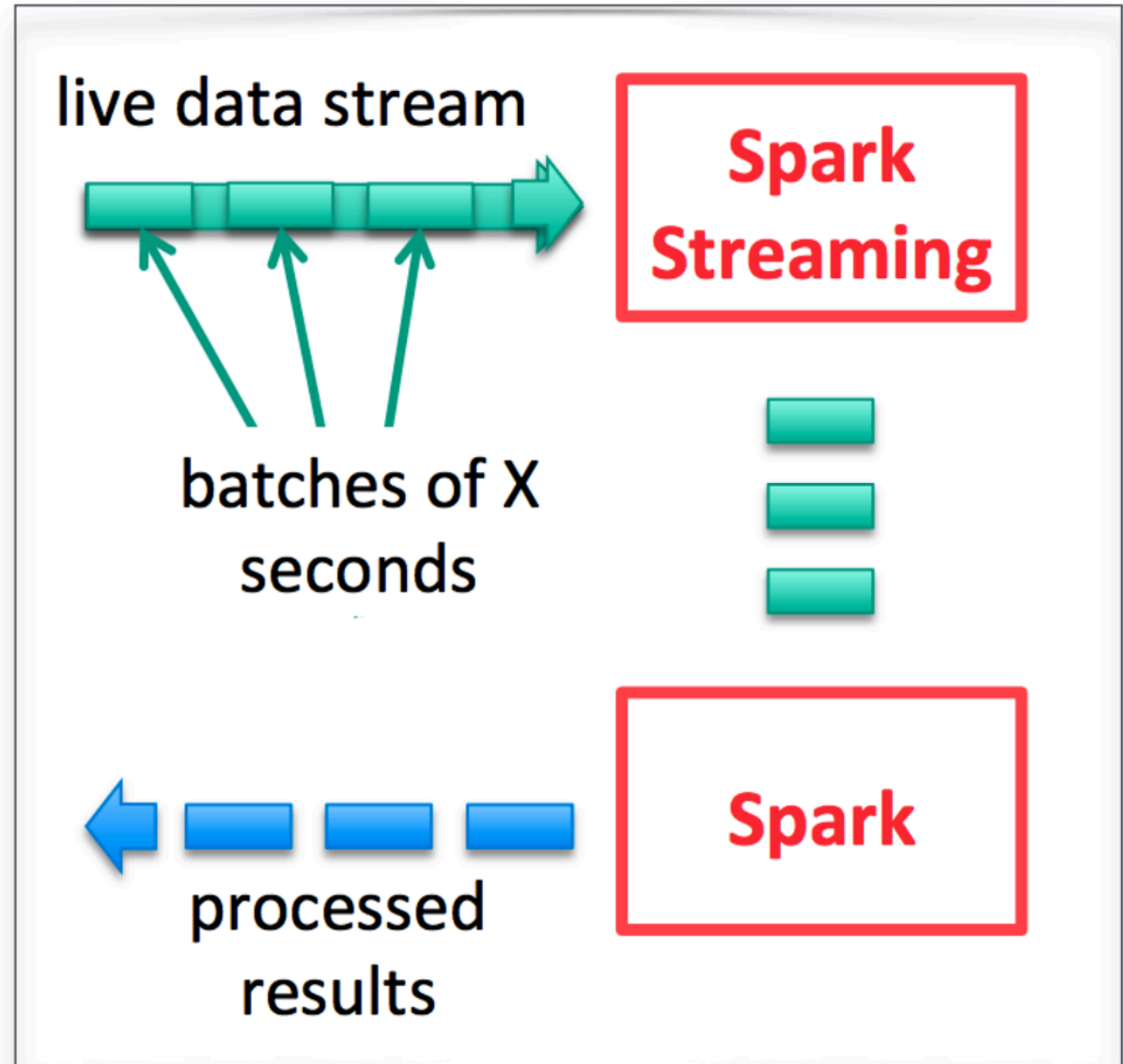


# Spark Stream

---

# Spark Stream

- *Chop up the live stream into batches of X seconds*
- *Spark treats each batch of data as RDDs and processes them using RDD operations*
- *Finally, the processed results of the RDD operations are returned in batches*



	<b>Kafka Streams</b>	<b>Storm</b>	<b>Spark Streaming</b>	<b>Flink</b>
<b>Integration</b>	Easy	Difficult	Difficult	Difficult
<b>Development</b>	Easy, flexible	Difficult	Difficult	Difficult
<b>Operations</b>	Easy	Difficult (Clustering)	Difficult (Clustering)	Difficult (Clustering)
<b>Infrastructure</b>	Small	Large (Clustering)	Large (Clustering)	Large (Clustering)
<b>Delivery</b>	At least once	At least once	Exactly Once	Exactly Once
<b>Latency</b>	Milliseconds	Seconds	Milliseconds	Milliseconds
<b>Fault Tolerance</b>	Yes	Yes	Yes	Yes
<b>Scalability</b>	Yes	No	Yes	No

# Document Classification With Splunk And Spark

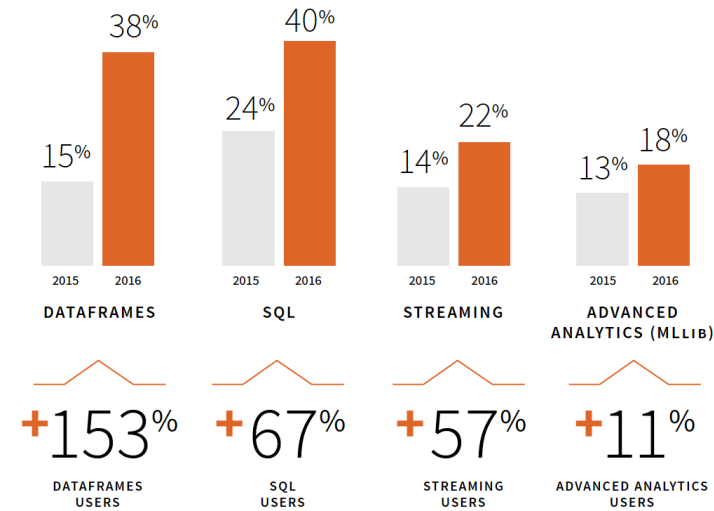
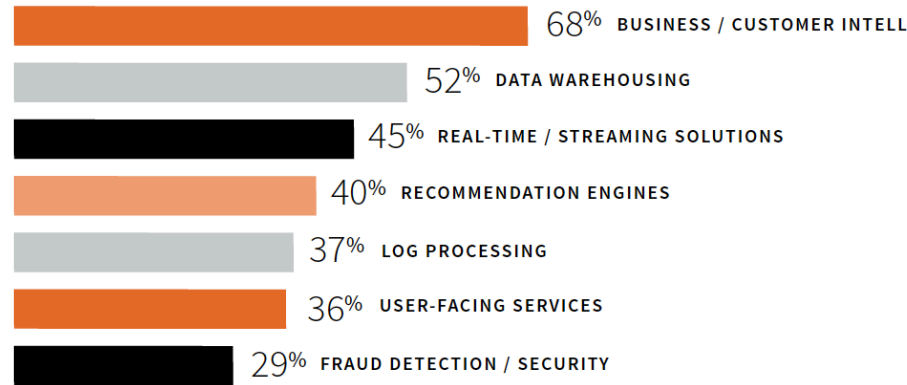
---



# 2016 Spark Survey

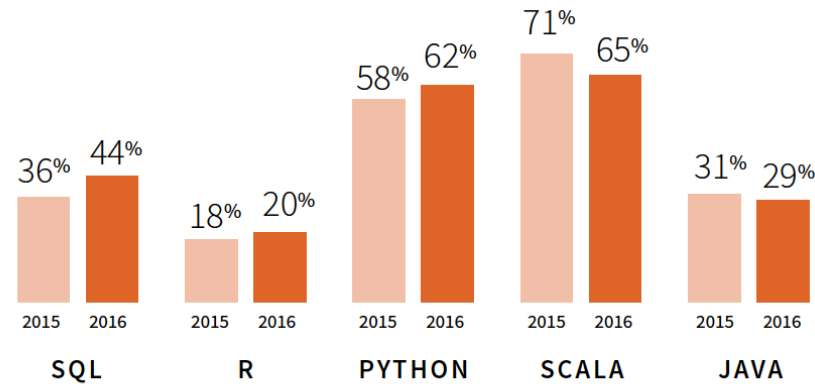
## TYPES OF PRODUCTS BUILT

% of respondents who use Spark to create each product (more than one product could be selected)



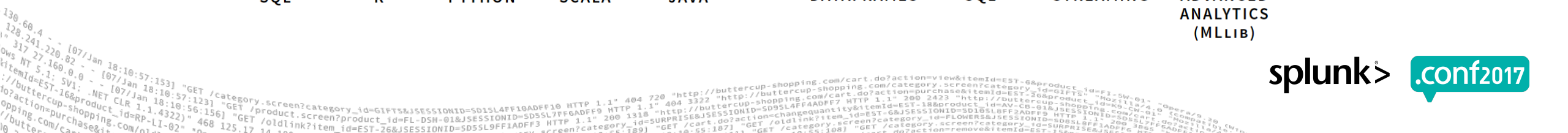
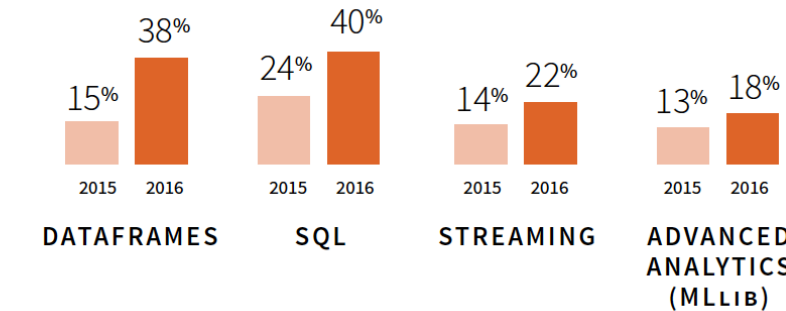
## LANGUAGES USED IN SPARK YEAR-OVER-YEAR

% of respondents who use each language (more than one language could be selected)



## SPARK COMPONENTS USED IN PRODUCTION YEAR-OVER-YEAR

% of respondents who use each component in production (more than one component could be selected)





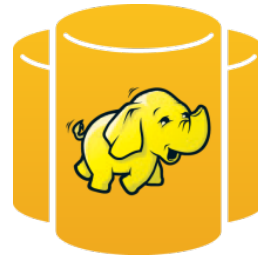


# Architecture



Splunk / Splunk Analytics  
for Hadoop Search Head

5



Hadoop  
Metadata

4

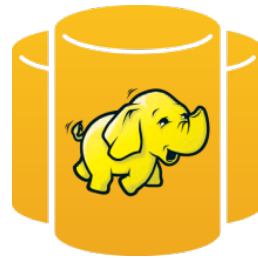


3



Machine  
Learning  
Extracts  
Metadata

2



Hadoop  
raw data



Transport

1



Documents





# Spark SQL And Splunk

splunk> App: Search & Reporting ▾ Administr

Search Pivot Reports Alerts Dashboards

## Spark DBXQuery

```
dbxquery query="SELECT * FROM `Spark`.`xademo`.`customer_details`" connection="spark_local_2" wrap=t
```

✓ 30 results (before 8/22/16 10:34:23.000 PM) No Event Sampling ▾

Events Patterns Statistics (30) Visualization

20 Per Page ▾ Format ▾ Preview ▾

(001)	(002)	(003)	(004)
customer_details.phone_number.STRING ▾	customer_details.plan.STRING ▾	customer_details.rec_date.STRING ▾	customer_details.status.STRING ▾
PHONE_NUM	PLAN	REC_DATE	STAUS





# Thank You

Don't forget to **rate this session** in the  
.conf2017 mobile app

splunk> .conf2017