



Productizing ML For Behavior Modeling and Security

Janet He | Chief Solution Architect at SAIC, Inc.

Marios Iliofotou | Principal Data Science Engineer UBA

September 27, 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.

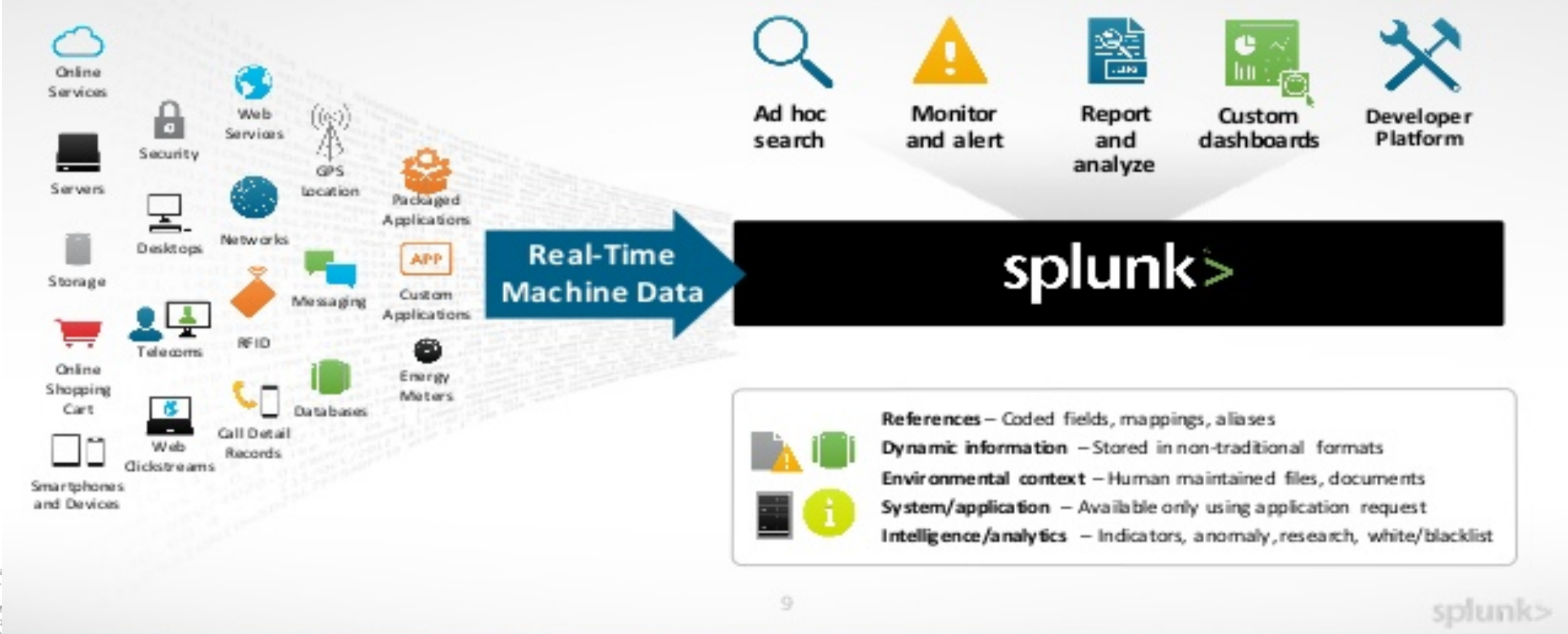
The forward-looking statements made in this presentation are being made as of the time and date of its live presentation. If reviewed after its live presentation, this presentation may not contain current or accurate information. We do not assume any obligation to update any forward looking statements 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.

SAIC Approach

Technology Platform – Splunk

Solution: Splunk, The Engine For Machine Data

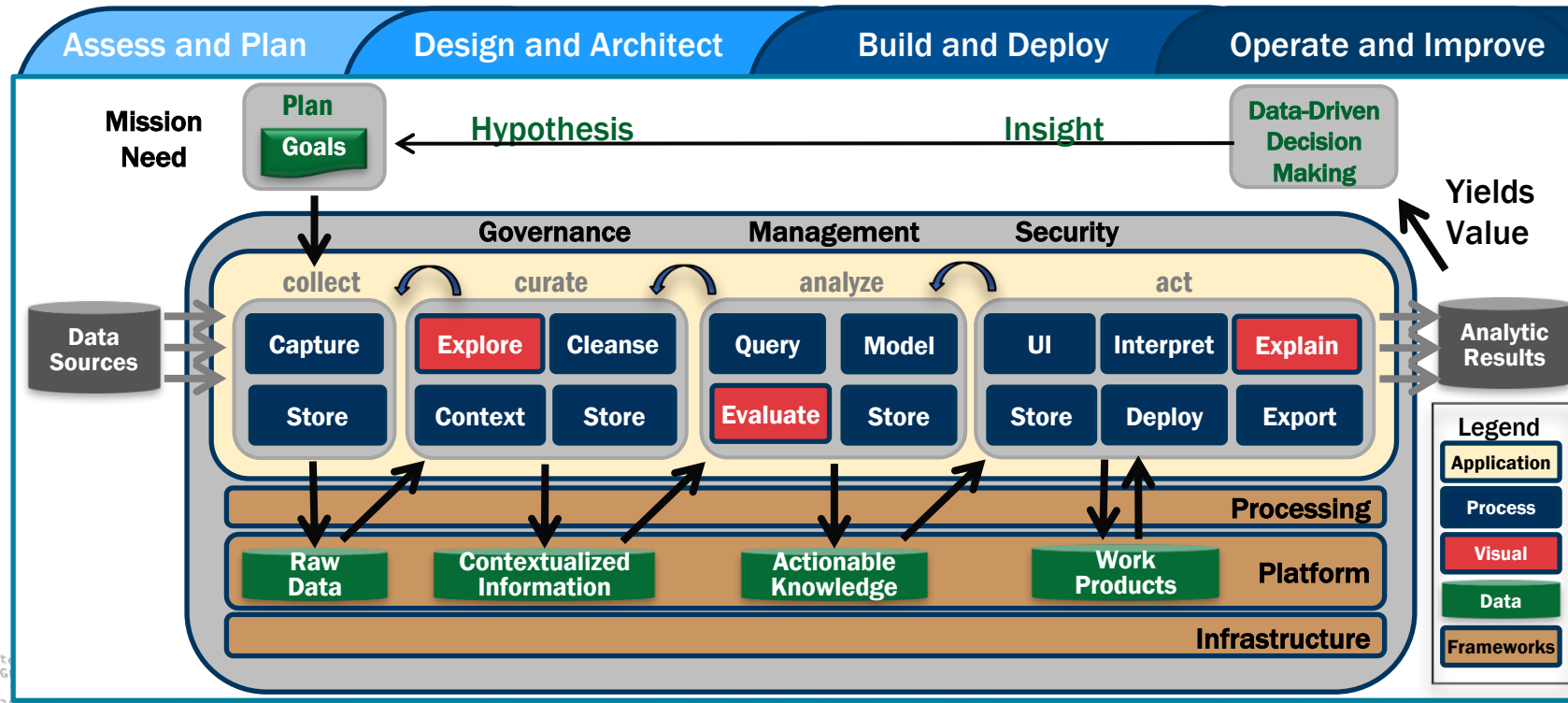


SAIC Approach

Data Management Methodology – SAIC DSE™

SAIC's Data Management Model enables innovation of analysis

- ▶ Data Science Edge™ (DSE) is SAIC's proprietary data lifecycle model geared toward the efficient planning and execution of enterprise data planning and analytics
- ▶ Model includes four phases of execution; Assess, Design, Build, and Improve. DSE Improve focuses on the performance and optimization of existing data and analytic systems
- ▶ SAIC has successfully used this process model to design a big data lake for our clients, and perform real-world testing of airport check-in biometrics devices



SAIC Approach

Data Protection – SAIC CSE™

CyberSecurity Edge Three Phase Methodology

1

Discover offers highly trained objective experts to identify real-world security risk and validate the implementation and effectiveness of an organization's existing security controls against industry recognized best practices and adversarial threats.

2

Mitigation is a highly tailored offering designed to help a customer design, plan, and implement solutions to meet specific goals and improve overall cybersecurity.

3

Manage provides cost efficient, low risk options for ongoing and continuous monitoring support by certified cybersecurity experts.

- Three options include managed, staff, and hybrid.

Advantages of SAIC's Approach

Verified | Recognized

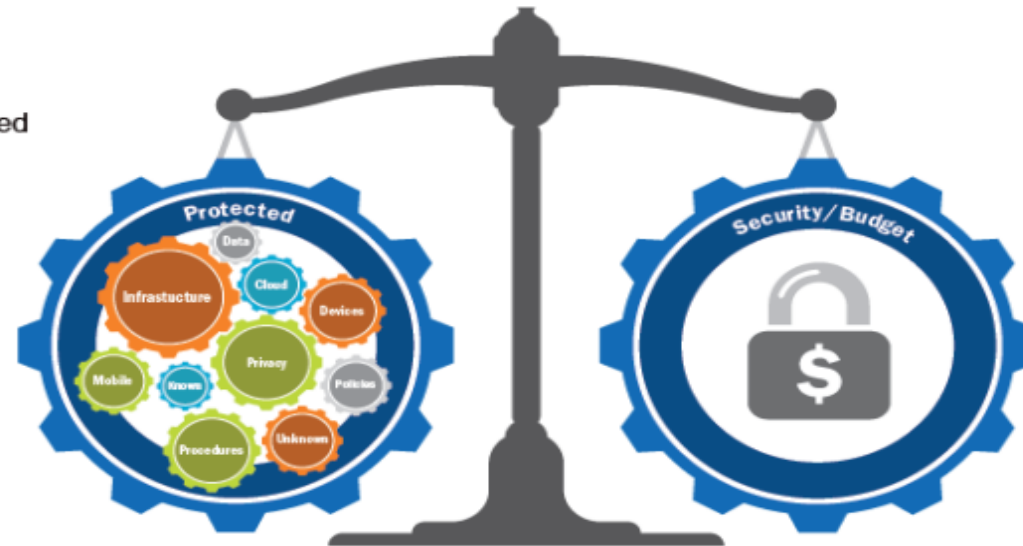
- Proven methodologies that have been developed and refined over countless engagements.

Automated | Optimized | Balanced | Tailored

- Offers customer-tailored solutions without the customization price tag.
- Optimizes current customer toolset.
- Fills gaps to strengthen ecosystem.
- Automates information assurance tasks.
- Balances tools, risk tolerance, and budget.

Packaged | Defined

- Clearly defines scope across all three phases with a fixed-priced model.



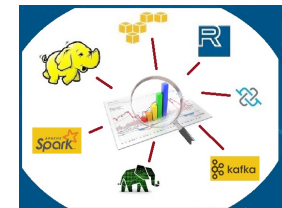
SAIC Approach

SAIC Big Data Analytics Solutions



Repeatable Solutions

- **Big Data Assessment and Roadmap** – templates and processes to assess an organization’s big data maturity and devise a roadmap.
- **Big Data Platform Accelerator** – reference architecture, blueprints, conops and security guidance to accelerate development and deployment of a big data platform.
- **Big Data Analytics Sandbox** – an SAIC cloud-based platform enabling client organizations to “play” with big data tools and technologies and develop advanced analytic products. Augmented for Deep Learning tools.
- **Big Data as a Service** – a scalable “as a service” offering allowing streaming analysis, batch analytics and data exploration in a secure fashion. Augmented for logical data analytics solution to handle the *Variety* problem of big data.



Solution Examples

SAIC Internal Splunk UBA Implementation

- ▶ SAIC Splunk UBA implementation based on machine data and our existing SIEM infrastructure.
- ▶ Objectives
 - Detect hidden security threats
 - Monitor networking, system, application, user and device anomalous behavior
 - Provide threat visualization
 - Increase SOC response to threats efficiently and effectively

Integrating with currently Splunk infrastructure include ES, ITSI, and SAIC capabilities in DSE, CSE and big data analytics services.

Aaron Bishop, @SAICinc

A CISO's Perspective on User Behavior Analytics: Setting the Right Expectations for All Stakeholders

Solution Example

Machine Learning – Threat Analysis

Need:

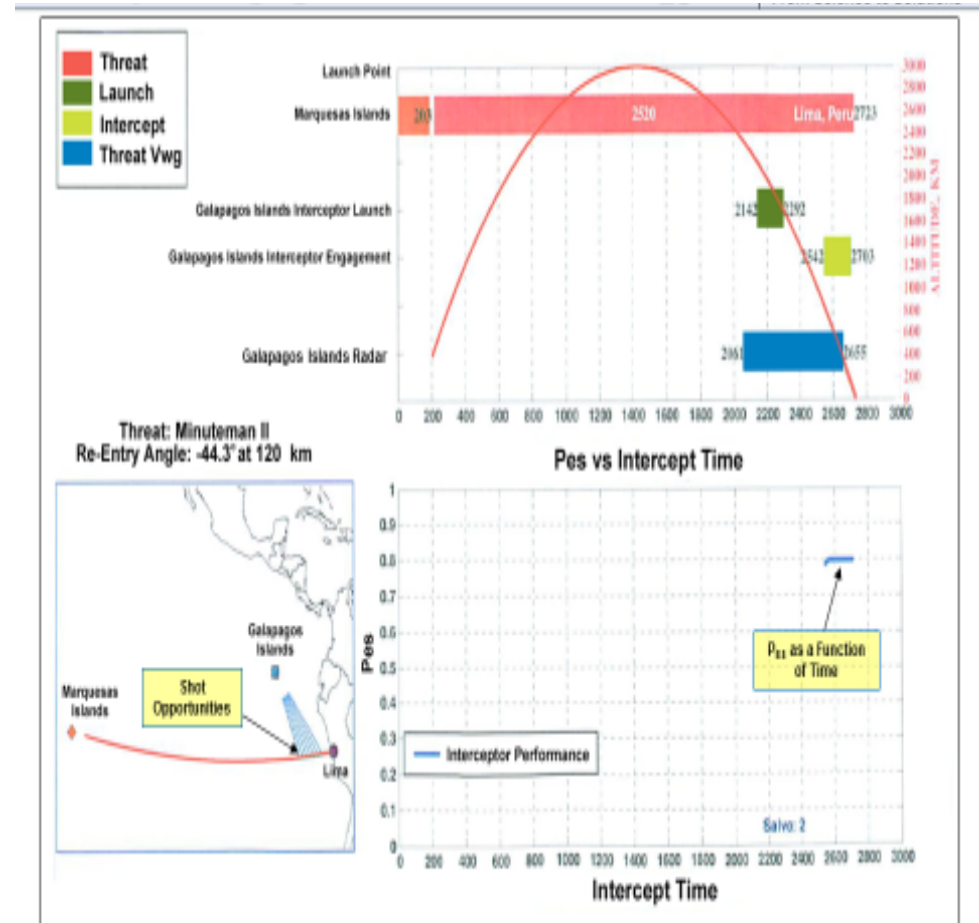
- Support studies on emerging threats and impacts. Manage and coordinate over 30 TB of raw data and processed products between multiple sites

Solution:

- Use of SQL databases as well as NoSQL (MongoDB)
- Developed and modeled advanced threat discrimination algorithms using Neural Networks, and Bayesian classifiers
- Automated tools to run simulations, generate KPIs, and create briefings
- Variety of tools used for visual displays including GIS and 3-D plots

Benefits:

- Eliminates laborious manual effort on part of analysts
- Provides frequent insights to leadership



Contact

- ▶ For additional information please contact us
- ▶ Splunk .conf 2017 SAIC booth M36

Sanjay Sardar
 VP | Advanced Analytics
 Advanced Analytics, Simulation and Training
 Cell: 703.861.5620 | Desk: 703.676.5028
 email: sanjay.sardar@saic.com | @SAICinc

Janet He
 Chief Solution Architect | Advanced Analytics
 Advanced Analytics, Simulation and Training
 Cell: 301.366.2078 | Desk: 703.676.2378
 email: janet.he@saic.com | @SAICinc



Overview

- ▶ Introduction
- ▶ Challenges
- ▶ Platform
- ▶ Programmability (SDK)
- ▶ Conclusions Q/A

```

130.60.4 - - [07/Jan 18:10:57:153] "GET /category.screen?category_id=GIFTS&JSESSIONID=5D15LAF10ADFF10 HTTP 1.1" 404 720 "http://buttercup-shopping.com/cart.do?action=view&itemId=EST-6&product_id=F1-5W-01"
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-268&product_id=K0-CW-01"
ows NT 5.1; SV1; .NET CLR 1.1.4322" 468 125.17 14.1.189] "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"
: //buttercup-shopping.com/oldlink?item_id=EST-26&JSESSIONID=5D55L9FF1ADFF3" 468 125.17 14.1.189] "GET /category.screen?category_id=FLOWERS&JSESSIONID=5D55L9FF1ADFF3 HTTP 1.1" 200 3885 "http://buttercup-shopping.com/cart.do?action=purchase&itemId=EST-26&product_id=K0-CW-01"
pping.com/cart.do?action=purchase&itemId=EST-26&product_id=K0-CW-01" 468 125.17 14.1.189] "GET /category.screen?category_id=FLOWERS&JSESSIONID=5D55L9FF1ADFF3 HTTP 1.1" 200 3885 "http://buttercup-shopping.com/cart.do?action=remove&itemId=EST-18&product_id=AV-CB-01"

```

Why use Machine Learning (ML)?

► You are probably trying to solve one of these problems:

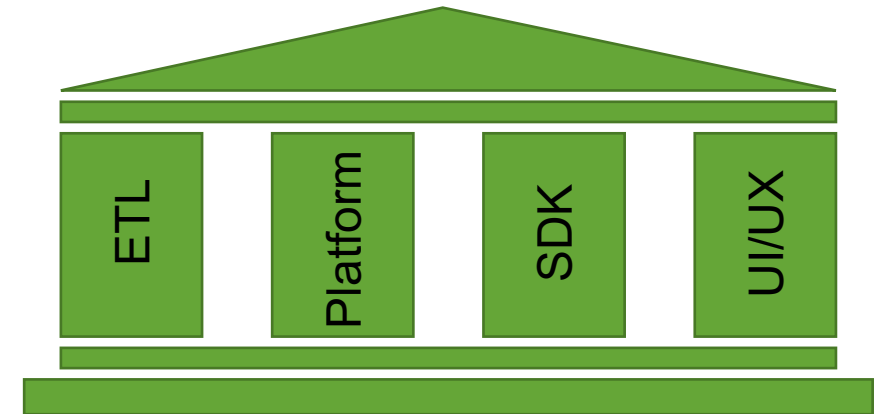
- Insider threats
- Malware/hackers
- Fraud

What they all have in common?
Simplistic solutions don't work!

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=FI-5W-03"
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-268product_id=KQ-CU-01"
137.27.160.0 - - [07/Jan 18:10:56:156] "GET /oldlink?item_id=EST-26&JSESSIONID=5D5L9FF1ADFF3 HTTP 1.1" 200 1318 "http://buttercup-shopping.com/cart.do?action=changequantity&itemId=EST-18&product_id=AV-CB-01&JSESSIONID=5D55LFF2ADFF3 HTTP 1.1" 200 3865 "http://buttercup-shopping.com/cart.do?action=remove&itemId=EST-189"
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=FI-5W-03"
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-268product_id=KQ-CU-01"
137.27.160.0 - - [07/Jan 18:10:56:156] "GET /oldlink?item_id=EST-26&JSESSIONID=5D5L9FF1ADFF3 HTTP 1.1" 200 1318 "http://buttercup-shopping.com/cart.do?action=changequantity&itemId=EST-18&product_id=AV-CB-01&JSESSIONID=5D55LFF2ADFF3 HTTP 1.1" 200 3865 "http://buttercup-shopping.com/cart.do?action=remove&itemId=EST-189"

Why Productizing a Solution is Hard?

1. ETL – Parsing, normalizing, cleaning
2. Platform – Scalability, performance, monitoring, orchestration
3. Programmability – Change/add new logic, test, develop (SDK)
4. Presentation – UI/UX, exploration/investigation



Any of the four pillars being weak the solution will fail!

Overview

► Introduction

► Challenges

► Platform

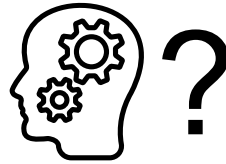
► Programmability (SDK)

► Conclusions Q/A

130.60.4 - - [07/Jan 18:10:57:153] "GET /category.screen?category_id=GIFTS&SESSIONID=5D5SLAFF10ADFF10 HTTP 1.1" 404 720 "http://buttercup-shopping.com/cart.do?action=view&itemId=EST-6&product_id=F1-5W-01"
128.241.220.82 - - [07/Jan 18:10:57:123] "GET /product.screen?product_id=FL-DSH-01&SESSIONID=5D5SL7FF6ADFF0 HTTP 1.1" 404 3322 "http://buttercup-shopping.com/cart.do?action=purchase&itemId=EST-268product_id=K9-CW-01"
10.0.0.1:51; SV1; .NET CLR 1.1.4322" 468 125.17 14.1.1 "GET /oldlink?item_id=EST-26&SESSIONID=5D5SL9FF1ADFF3 HTTP 1.1" 200 1318 "http://buttercup-shopping.com/cart.do?action=changequantity&itemId=EST-18&product_id=AV-CB-01&SESSIONID=5D5SL7FF6ADFF0 HTTP 1.1" 200 2423 "http://buttercup-shopping.com/cart.do?action=purchase&itemId=EST-268product_id=K9-CW-01"
10.0.0.1:51; SV1; .NET CLR 1.1.4322" 468 125.17 14.1.1 "GET /category.screen?category_id=FLOWERS&SESSIONID=5D5SL7FF6ADFF0 HTTP 1.1" 200 3865 "http://buttercup-shopping.com/cart.do?action=purchase&itemId=EST-268product_id=K9-CW-01"
10.0.0.1:51; SV1; .NET CLR 1.1.4322" 468 125.17 14.1.1 "GET /category.screen?category_id=FLOWERS&SESSIONID=5D5SL7FF6ADFF0 HTTP 1.1" 200 3865 "http://buttercup-shopping.com/cart.do?action=purchase&itemId=EST-268product_id=K9-CW-01"
10.0.0.1:51; SV1; .NET CLR 1.1.4322" 468 125.17 14.1.1 "GET /category.screen?category_id=FLOWERS&SESSIONID=5D5SL7FF6ADFF0 HTTP 1.1" 200 3865 "http://buttercup-shopping.com/cart.do?action=purchase&itemId=EST-268product_id=K9-CW-01"
10.0.0.1:51; SV1; .NET CLR 1.1.4322" 468 125.17 14.1.1 "GET /category.screen?category_id=FLOWERS&SESSIONID=5D5SL7FF6ADFF0 HTTP 1.1" 200 3865 "http://buttercup-shopping.com/cart.do?action=purchase&itemId=EST-268product_id=K9-CW-01"
10.0.0.1:51; SV1; .NET CLR 1.1.4322" 468 125.17 14.1.1 "GET /category.screen?category_id=FLOWERS&SESSIONID=5D5SL7FF6ADFF0 HTTP 1.1" 200 3865 "http://buttercup-shopping.com/cart.do?action=purchase&itemId=EST-268product_id=K9-CW-01"
10.0.0.1:51; SV1; .NET CLR 1.1.4322" 468 125.17 14.1.1 "GET /category.screen?category_id=FLOWERS&SESSIONID=5D5SL7FF6ADFF0 HTTP 1.1" 200 3865 "http://buttercup-shopping.com/cart.do?action=purchase&itemId=EST-268product_id=K9-CW-01"

Goal: Manage Multiple ML Models

Model = “Training and scoring of ML models plus utility tasks”



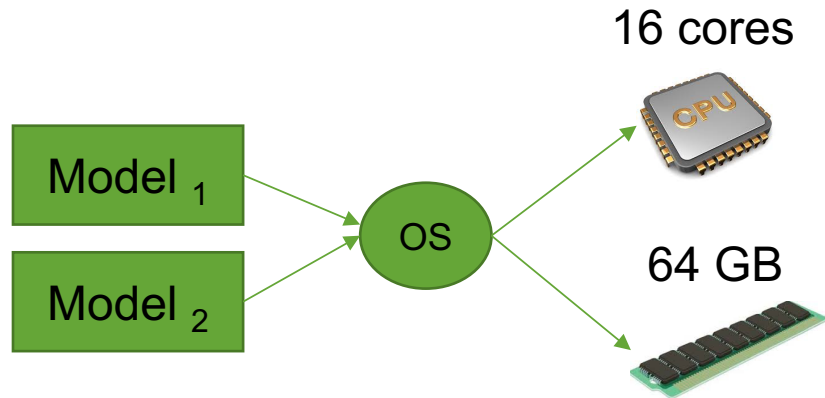
► Isolate models (processes)

- Out of memory
- Out of disk space
- High CPU usage

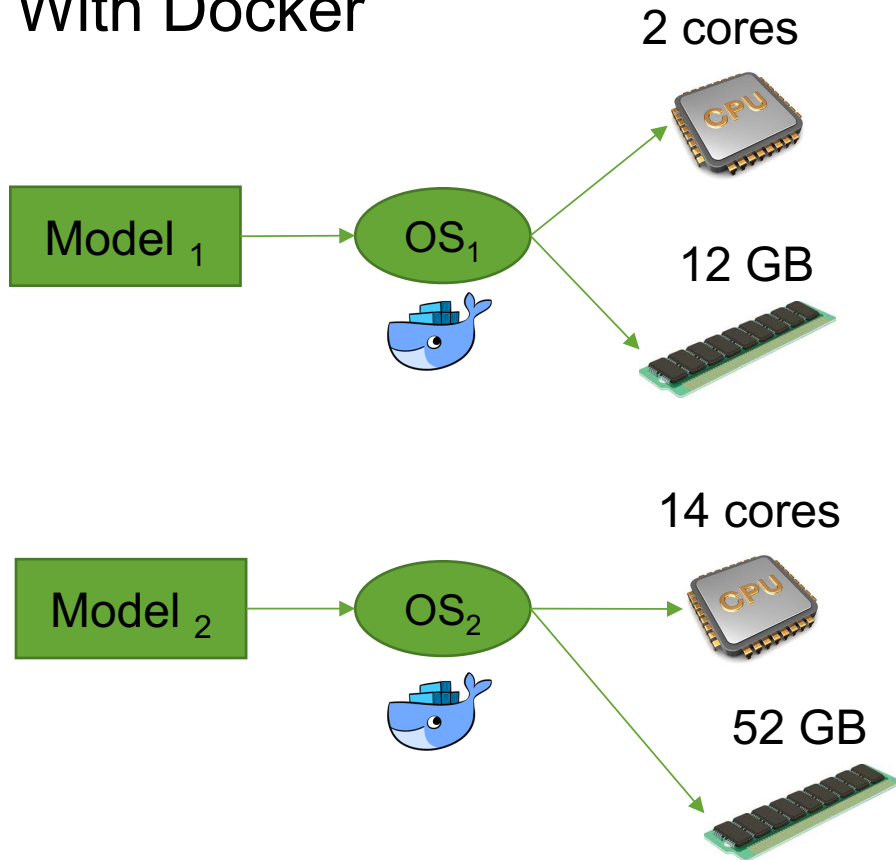


Model Isolation With Docker

Without Docker



With Docker



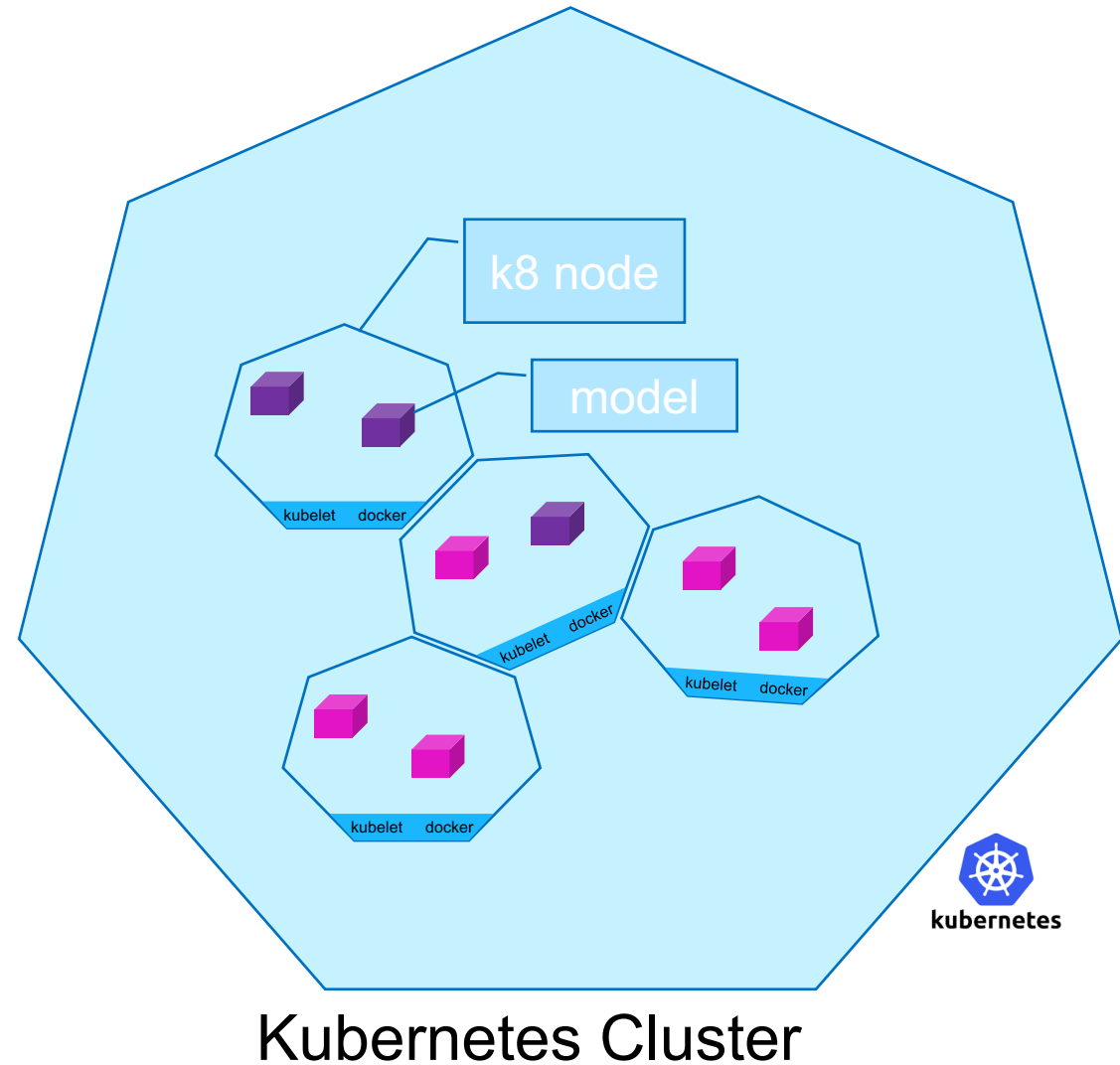
If "Model 2" takes all the resource, "Model 1" is not affected

130.60.4 - - [07/Jan 18:10:57:153] "GET /category.screen?category_id=GIFTS&JSESSIONID=5D15LAF10ADFF10 HTTP 1.1" 404 720 "http://buttercup-shopping.com/cart.do?action=view&itemId=EST-6&product_id=FI-SW-03" Moz/1.12.0
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-268product_id=KQ-CW-01" Moz/1.12.0
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" Moz/1.12.0
:/buttercup-shopping_id=RP-LI-02" 468 125.17.14.105:1871 "GET /category.screen?category_id=FLOWERS&JSESSIONID=5D55L9FF1ADFF3 HTTP 1.1" 200 3865 "http://buttercup-shopping.com/cart.do?action=purchase&itemId=EST-189" Moz/1.12.0
:/buttercup-shopping_id=RP-LI-02" 468 125.17.14.105:1871 "GET /category.screen?category_id=FLOWERS&JSESSIONID=5D55L9FF1ADFF3 HTTP 1.1" 200 3865 "http://buttercup-shopping.com/cart.do?action=purchase&itemId=EST-189" Moz/1.12.0

Multi Node Containers with Kubernetes

► Scenarios

- A model is struggling
 - Spin more instances of the model and balance the load
- New model is added
 - Create new containers and assign them to a node
- When cluster gets overloaded
 - Add extra nodes



Overview

▶ Introduction

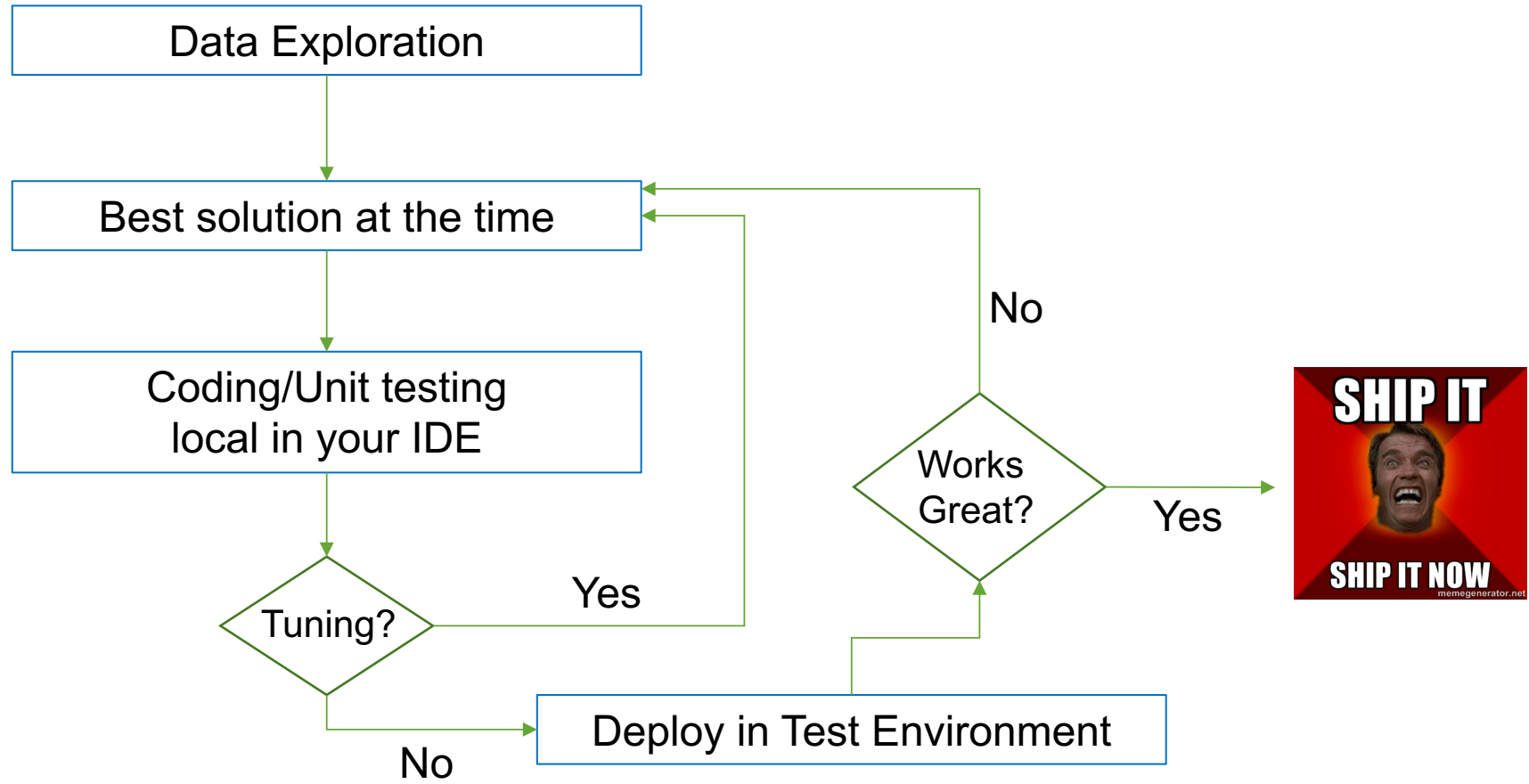
▶ Challenges

▶ Platform

▶ **Programmability (SDK)**

▶ Conclusions Q/A

Realistic Model Development Life-Cycle



SDK should support all of these steps!



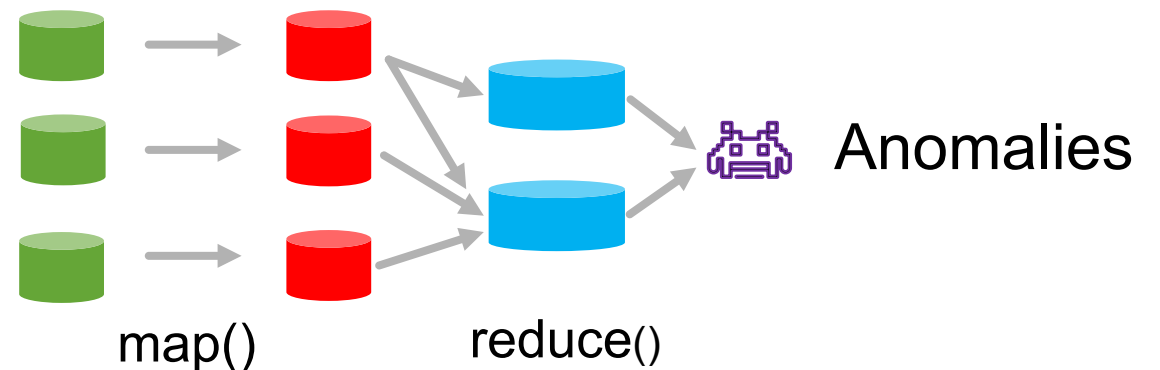
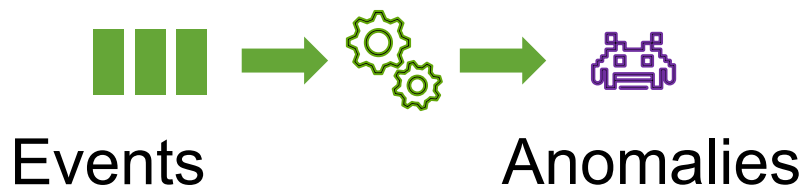
Different Use Cases Require Different Model Types

▶ Streaming

- Single pass over the data
- Quick response to events
- Run continuously

▶ Batch

- Multiple passes over the data
- Can run expensive correlations (joins)
- Run at scheduled intervals
(think Linux Cron jobs)



splunk> UBA: Streaming Model APIs

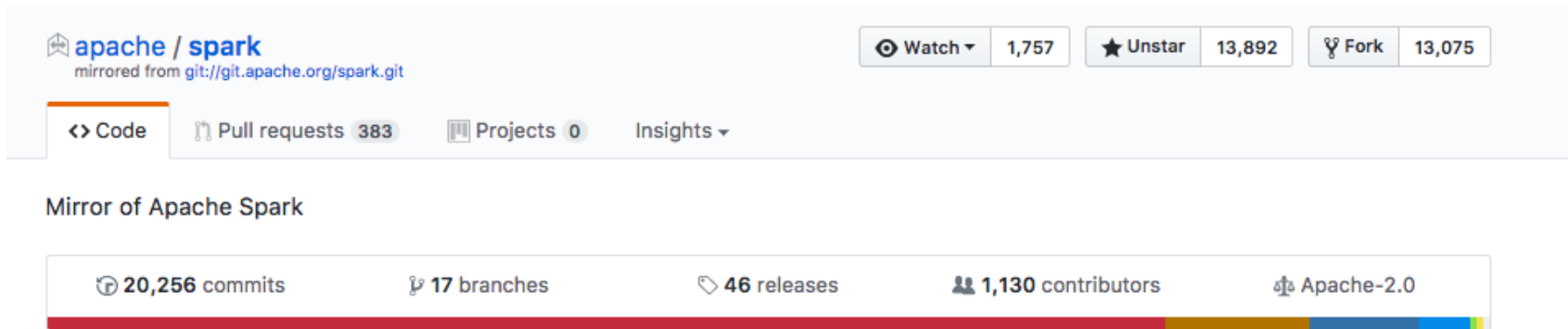
```
Option [Anomaly] analyzeData (DataEvent currentEvent)
```

- ▶ State is checkpointed internally
 - Serialization (Protocol Buffers, Kryo)
- ▶ Streaming models choose
 - a) Pivot (e.g., per user or device)
 - b) Input types (e.g., HTTP traffic data)

```
130.60.4 - - [07/Jan 18:10:57:153] "GET /category.screen?category_id=GIFTS&JSESSIONID=5D15LAF10ADFF10 HTTP 1.1" 404 720 "http://buttercup-shopping.com/cart.do?action=view&itemId=EST-6&product_id=FI-SW-03" "Mozilla/5.0 (Macintosh; Intel Mac OS X 10_11_7; 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=KQ-CU-01" "Mozilla/5.0 (Macintosh; Intel Mac OS X 10_11_7; rv:53.0) Gecko/20100801 Firefox/53.0"
ows NT 5.1; SV1; .NET CLR 1.1.4322)" 468 125.17 14 - - [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_7; rv:53.0) Gecko/20100801 Firefox/53.0"
itemId=EST-16&product_id=RP-LI-02" 404 125.17 14 - - [07/Jan 18:10:56:156] "GET /oldlink?item_id=EST-26&JSESSIONID=5D55L9FF1ADFF3 HTTP 1.1" 200 3865 "http://buttercup-shopping.com/cart.do?action=remove&itemId=EST-14&product_id=KQ-CU-01" "Mozilla/5.0 (Macintosh; Intel Mac OS X 10_11_7; rv:53.0) Gecko/20100801 Firefox/53.0"
buttercup-shopping.com/cart.do?action=purchase&itemId=EST-26&product_id=KQ-CU-01" 404 125.17 14 - - [07/Jan 18:10:56:156] "GET /oldlink?item_id=EST-26&JSESSIONID=5D55L9FF1ADFF3 HTTP 1.1" 200 3865 "http://buttercup-shopping.com/cart.do?action=remove&itemId=EST-14&product_id=KQ-CU-01" "Mozilla/5.0 (Macintosh; Intel Mac OS X 10_11_7; rv:53.0) Gecko/20100801 Firefox/53.0"
```

splunk> UBA: Batch Model APIs

▶ Apache Spark



apache / spark
mirrored from git://git.apache.org/spark.git

Watch 1,757 Unstar 13,892 Fork 13,075

Code Pull requests 383 Projects 0 Insights

Mirror of Apache Spark

20,256 commits 17 branches 46 releases 1,130 contributors Apache-2.0

▶ Apache 2.0.x full set of APIs are supported

- RDD, Dataset, DataFrames, Spark SQL

```
httpData.groupBy('userId').agg(sum('bytesOut'), unique('dstIP'))
```


Thank You

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

splunk® **.conf2017**

Splunk/UbaExample

   default ▾

df: Unit = ()

Took 25 sec. Last updated by anonymous at August 03 2017, 11:49:14 AM.

Show the top destinations in terms of number of events

FINISHED   

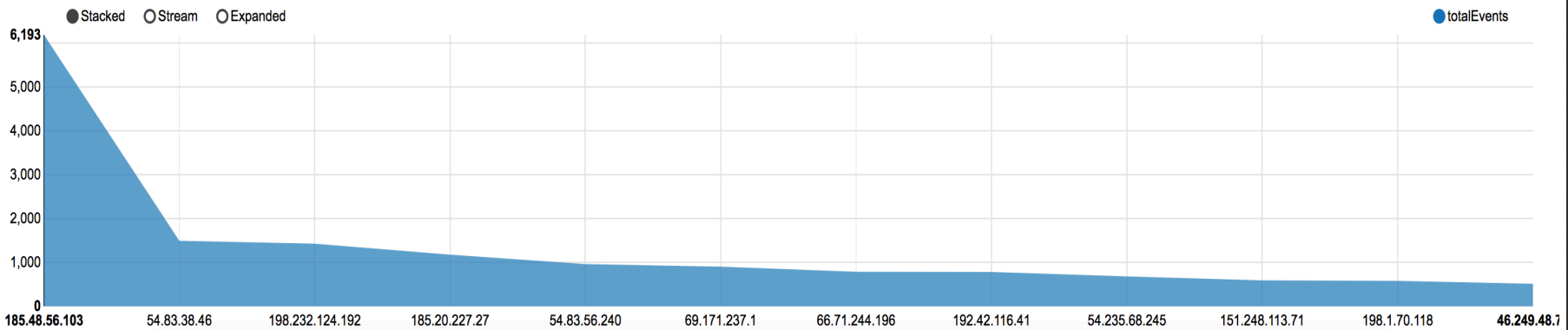
Took 0 sec. Last updated by anonymous at August 03 2017, 11:51:58 AM.

```
%sql
select destination, sum(numEvents) as totalEvents
from EventsPerDest
group by 1
order by totalEvents desc limit ${top=20}
```

FINISHED   

top

       settings ▾



Notebook Driven Development

► Zeppelin notebook example

The screenshot shows a Zeppelin notebook interface. At the top, there is a navigation bar with the Zeppelin logo, 'Notebook', and 'Job' tabs. A search bar contains 'Search your Notes' and a user profile 'anonymous'. Below this, the notebook title is 'Splunk/UbaExample'. The main content area shows a notebook cell with the title 'Show the top destinations in terms of number of events' and a status of 'FINISHED'. The cell contains a SQL query:

```
%sql
select destination, sum(numEvents) as totalEvents
from EventsPerDest
group by 1
order by totalEvents desc limit ${top}=20;
```

The query result is visualized as a bar chart. The y-axis represents the number of events, ranging from 0 to 6,193. The x-axis lists 12 IP addresses. The bars are blue and represent the 'totalEvents' for each destination. The top destination is 185.48.56.103 with approximately 6,193 events. The number of events decreases significantly for the second destination, 54.83.38.46, and continues to decrease for the remaining destinations.

Destination	totalEvents
185.48.56.103	6193
54.83.38.46	~1500
198.232.124.192	~1200
185.20.227.27	~1000
54.83.56.240	~800
69.171.237.1	~700
66.71.244.196	~600
192.42.116.41	~500
54.235.68.245	~400
151.248.113.71	~300
198.1.70.118	~200
46.249.48.7	~100

Different Model Types and Challenges

- ▶ Streaming: Quick response to events (Kafka)
 - State explosion
 - Slow EPS

- ▶ Batch: Stronger correlations (Apache Spark)
 - Job execution time (timeout)

- ▶ Common challenges
 - Scaling up/down
 - CPU/Memory/IO fairness (one task interfering with another)