# Using the Latest Features from the Splunk Machine Learning Toolkit to Create Your Own Custom Models

Adam J. Oliner | Director of Engineering

Harsh Keswani | Product Manager - Machine Learning

October 2018

# 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> .conf18

# Speakers

**Adam J. Oliner**

Director of Engineering

**Harsh Keswani**

Product Manager: Machine Learning

splunk> .conf18

# Outline

▶ Splunk Machine Learning Toolkit

▶ Platform Extensions: ML-SPL, etc.

▶ Experiments: Guided Machine Learning

▶ Demo

▶ What's New

▶ Customer Success

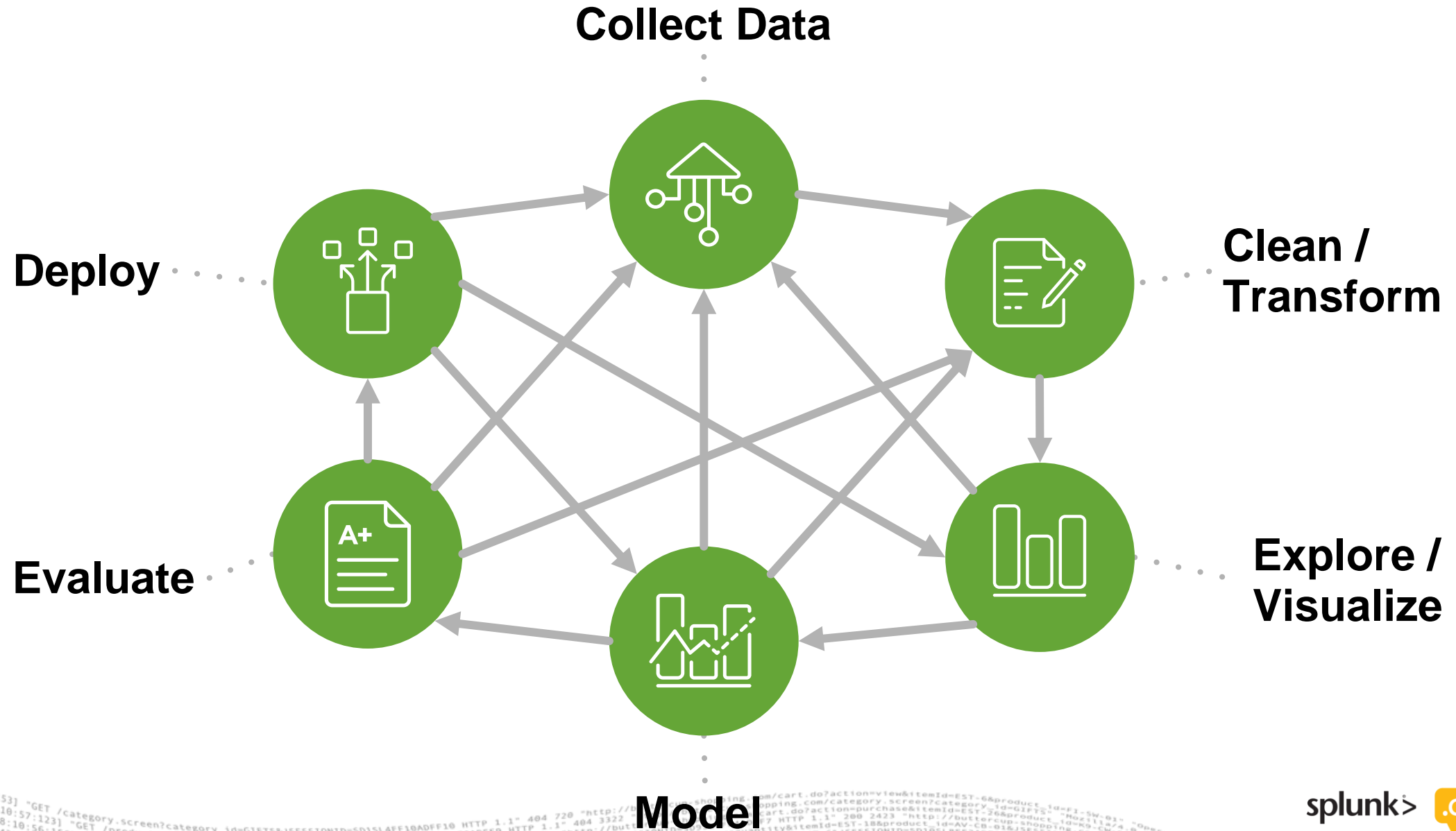# Splunk Machine Learning Toolkit
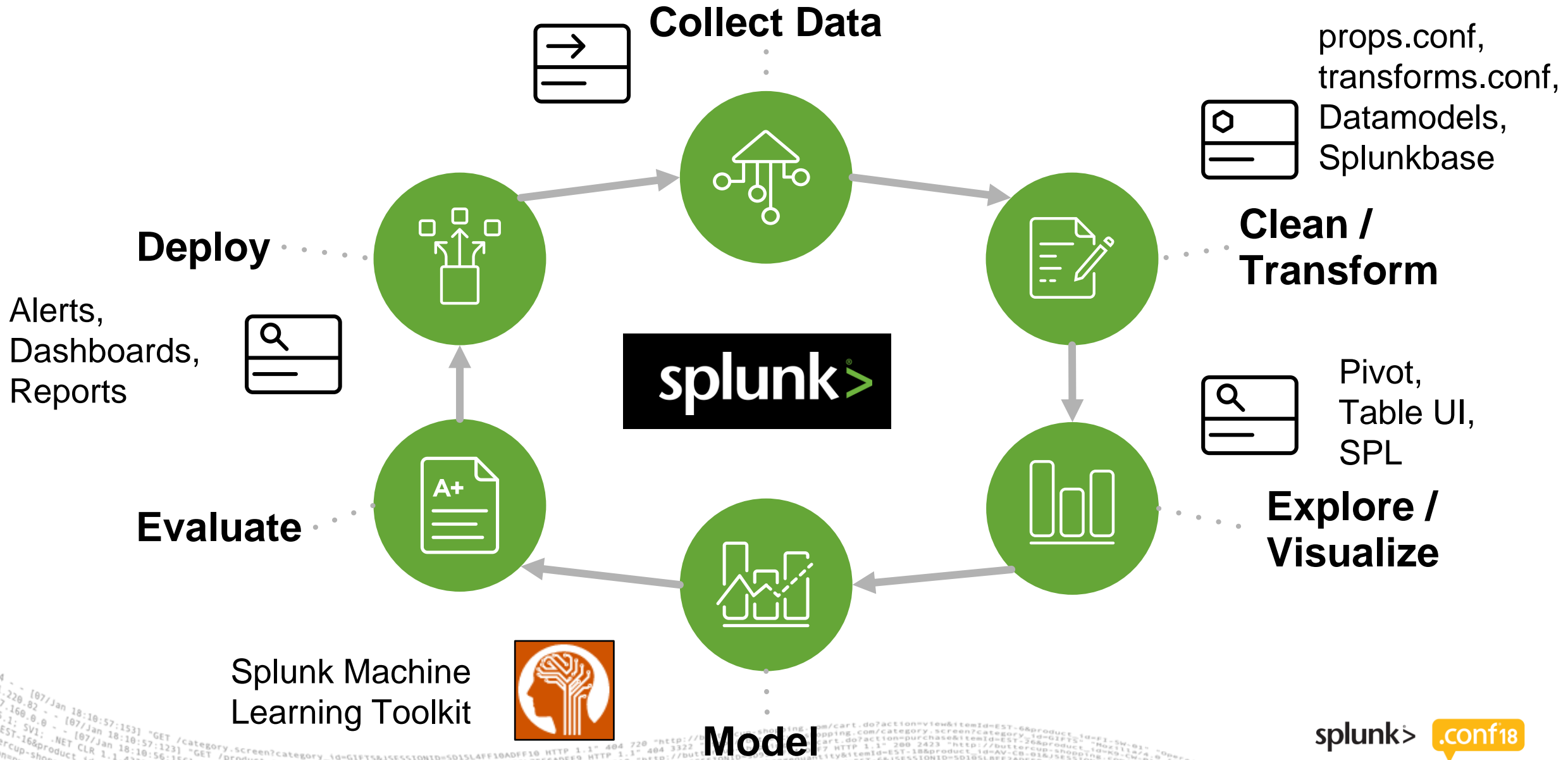
platform extensions and guided modeling dashboards

splunk> .conf18

# Machine Learning

▸ A process for generalizing from examples
▸ Examples

- A, B, … → #                    (regression)

- A, B, ... → a                    (classification)

- $X_{past} \rightarrow X_{future}$                    (forecasting)

- like with like                    (clustering)

- $|X_{predicted} - X_{actual}| >> 0$                    (anomaly detection)

splunk> .conf18

# Machine Learning Process



**Collect Data**

**Deploy**

**Clean / Transform**

**Evaluate**

**Explore / Visualize**

**Model**

# Machine Learning Process with Splunk



**Collect Data**

props.conf,
transforms.conf,
Datamodels,
Splunkbase

**Clean /
Transform**

**Deploy**

Alerts,
Dashboards,
Reports

Pivot,
Table UI,
SPL

**Evaluate**

**Explore /
Visualize**

Splunk Machine
Learning Toolkit

**Model**

splunk> .conf18

# Data Gathering and Prep

## Source: CrowdFlower



### What data scientists spend the most time doing

- Building training sets: 3%
- Cleaning and organizing data: 60%
- Collecting data sets; 19%
- Mining data for patterns: 9%
- Refining algorithms: 4%
- Other: 5%

**Want to learn more about data prep? Download the slides and recording for the following session.**

Getting Your Data Ready for Machine Learning

**Speakers**

**Kristal Curtis**, Software Engineer, Machine Learning, Splunk
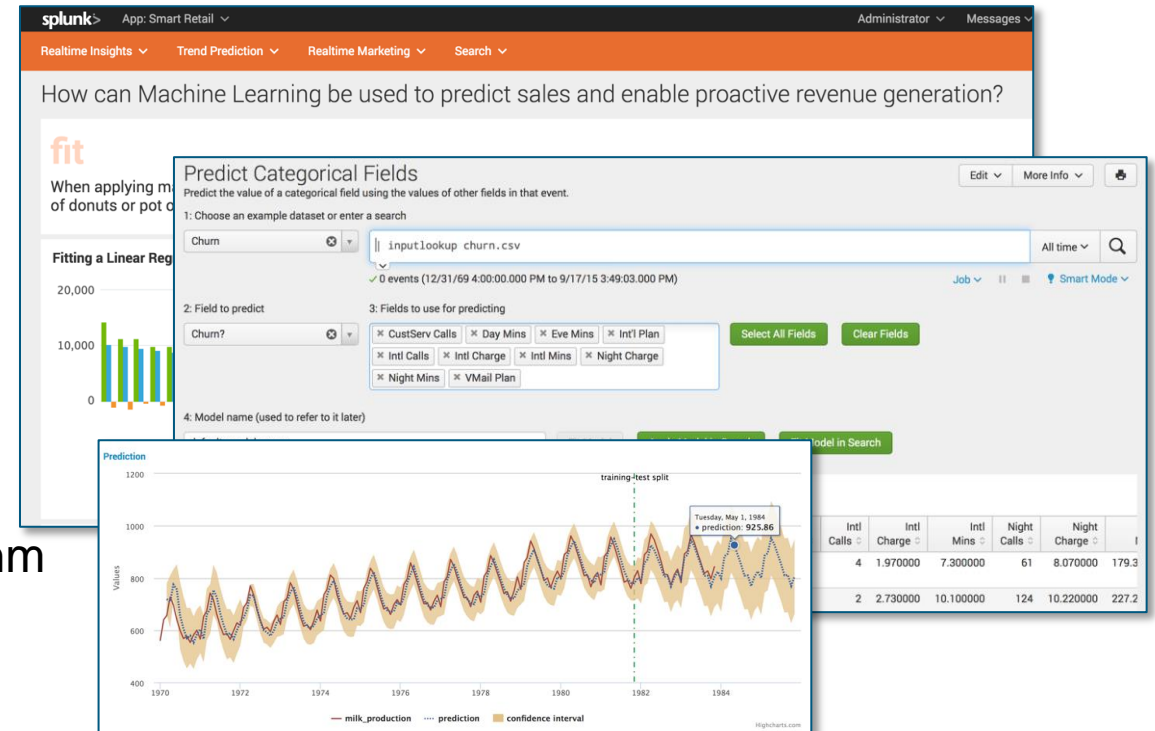**Adam J. Oliner**, Director of Engineering, Splunk

Wednesday, Oct 03, 12:45 p.m. - 1:30 p.m.

splunk> .conf18

# Splunk Machine Learning Toolkit

## extends Splunk with new tools and guided modeling

- **Experiments:** Guided model building, testing, and deployment for common objectives

- **Showcases:** Interactive examples for typical IT, security, business, and IoT use cases

- **Algorithms:** 30 standard algorithms

  (supervised & unsupervised)

- **ML Commands:** New SPL commands to fit, test and operationalize models

- **ML-SPL API**  Extensibility to easily import any algorithm

  (proprietary / open source)

- **Python for Scientific Computing Library:** Access to 300+ open source algorithms

# Platform Extensions

custom search commands for machine learning

splunk> .conf18

# SPL, Macros, & Viz

## Oh, my!

▶ **Commands (ML-SPL)**
- fit
- apply
- summary
- listmodels
- deletemodel
- sample
- score

▶ **Macros**
- regressionstatistics
- classificationstatistics
- classificationreport
- confusionmatrix
- forecastviz
- histogram
- modvizpredict
- splitby(1-5)

▶ **Viz**
- Outliers Chart
- Forecast Chart
- Scatter Line Chart
- Histogram Chart
- Downsampled Line Chart
- Scatterplot Matrix
- Box Plot Chart

130.60.4 - - [07/Jan 18:10:57:153] "GET /category.Screen?category_id=GIFTS&JSESSIONID=SD1SL4FF10ADFF10 HTTP 1.1" 404 720 "http://buttercup-shopping.com/cart.do?action=view&itemId=EST-6&product_id=FL-SW-01" "Mozilla/5.0
128.241.220.82 - - [07/Jan 18:10:57:123] "GET /product.Screen?product_id=FL-DSH-01&JSESSIONID=SD5SL7FF6ADFF9 HTTP 1.1" 404 3322 "http://buttercup-shopping.com/cart.do?action=purchase&itemId=EST-26&product_id=GIFTS" "Mozilla/4.0
317 27.160.0.0 - - [07/Jan 18:10:56:156] "GET /oldlink?item_id=EST-18&product_id=AV-CB-01&JSESSIONID=SD5SL9FF1ADFF3 HTTP 1.1" 200 1318 "http://JSESSIONID=SD9SL4FF4ADFF7 HTTP 1.1" 200 2423

splunk> .conf18

# ML-SPL Commands

- Fit (i.e., train) a model from search results

  ```
  … | fit <ALGORITHM> <TARGET> from <VARIABLES …>
            <PARAMETERS> into <MODEL>
  ```

- Apply a model to obtain predictions from (new) search results

  ```
  … | apply <MODEL>
  ```

- Inspect a model (e.g., display coefficients)

  ```
  | summary <MODEL>
  ```

- Score the prediction results

  ```
  … | score <SCORE_METHOD> <ACTUAL> ~ <PREDICTED>
  ```

# ML-SPL Commands: fit

… | fit <ALGORITHM> <TARGET> from <VARIABLES> <PARAMETERS> into <MODEL>

*optional*

Examples:

```
… | fit LinearRegression
        system_temp from cpu_load fan_rpm
        into temp_model
… | fit KMeans k=10
        downloads purchases posts days_active visits_per_day
        into user_behavior_clusters
```

splunk> .conf18

# ML-SPL Algorithms

- 30 algorithms OotB
  - prediction, clustering, forecasting, feature engineering
- Extensibility API for 300+ more
- Pipeline for advanced use cases

```
…      | fit TFIDF message
       | fit StandardScaler files bytes
       | fit KMeans message_tfidf_* SS_* k=5
       | fit PCA message_tfidf_* k=2
       | …
```

# ML-SPL Commands: apply

```
… | apply <MODEL>
```

Examples:

```
    … | apply temp_model
    … | apply user_behavior_clusters
```

# ML-SPL Commands: score

… | score <SCORE_METHOD> <ACTUAL> ~ <PREDICTED>

Examples:

    … | score accuracy_score vehicleType ~ LR_prediction
DT_prediction
    … | score confusion_matrix actual=vehicleType
predicted=pred_type

splunk> .conf18

# ML-SPL Commands: summary

```
… | summary <MODEL>
```

Examples:

```
… | summary temp_model
… | summary user_behavior_clusters
```

# ML-SPL Commands

```
| listmodels
| deletemodel <MODEL>
```
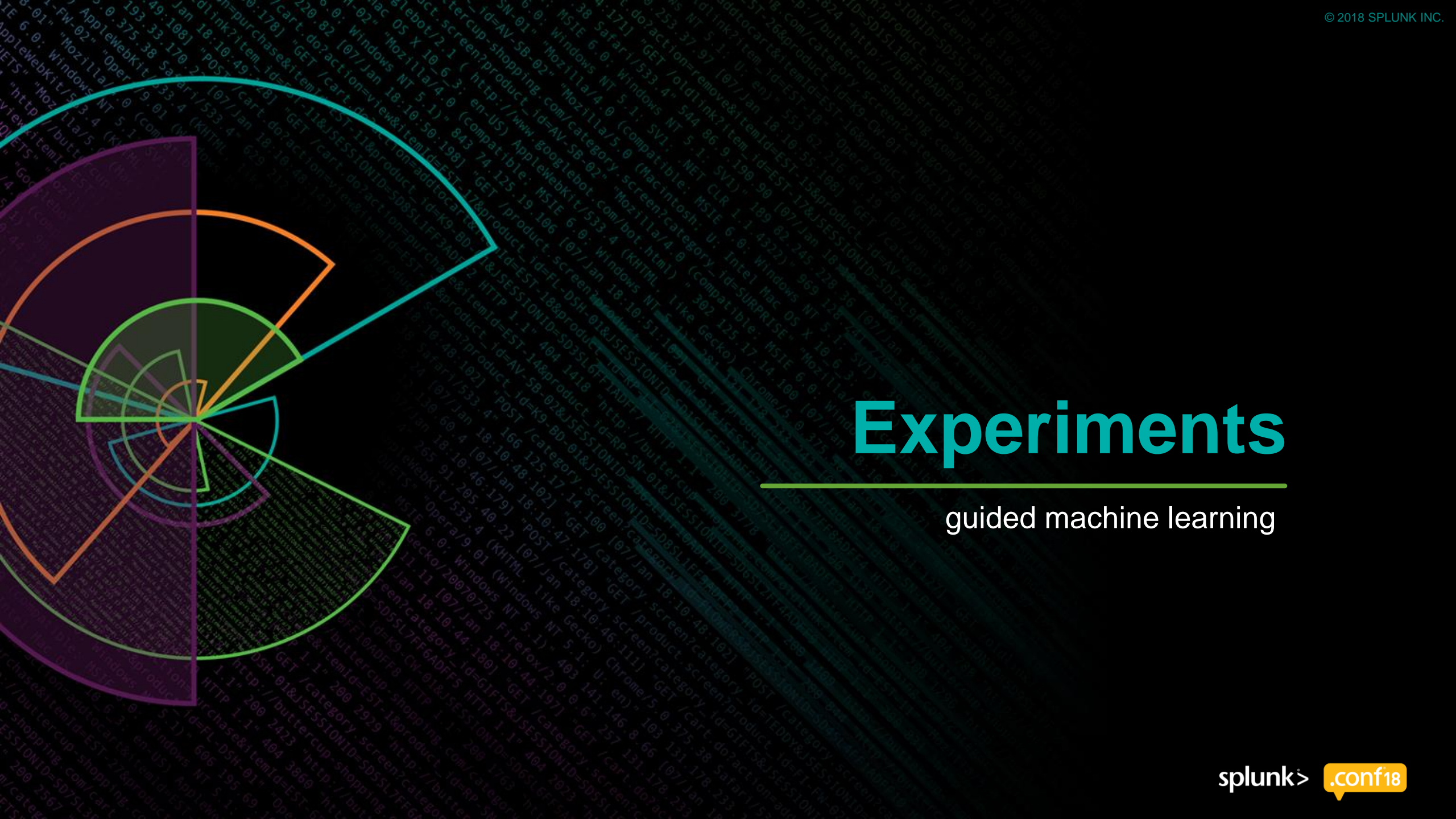
# ML-SPL Commands: `sample`

- ► Randomly sample or partition events

… | sample <PARAMETERS>

- ► Four modes
  - • Ratio                           … | sample 0.01
  - • Count                           … | sample 20
  - • Proportional        … | sample proportional="some_field"
  - • Partition             … | sample partitions=10

splunk> .conf18

# Experiments

guided machine learning

# Guided ML with Experiments

▸ Guides you through an analysis
▸ Automatically generates all the relevant SPL

Fit a model on all your data in search ↗                                    ✕

```
| inputlookup server_power.csv

| fit StandardScaler "total*" with_mean=true        // apply preprocessing steps
with_std=true into
example_server_power_StandardScaler_0

| fit PCA "SS*" k=2 into
example_server_power_PCA_1

| fit LinearRegression fit_intercept=true           // fit and save a model using the entire dataset
"ac_power" from "SS*" into "example_server_power"    and provided parameters
```

splunk> .conf18

# Experiments: Fit

**Algorithm**

RandomForestRegress... ▼

**Field to predict**

ac_power ▼

**Fields to use for predicting**

SS* (1) ▼

**Split for training / test: 50 / 50**

──────○──────

**N Estimators**

(optional)

**Max Depth**

(optional)

**Max Features**

(optional)

**Min Samples Split**

(optional)

**Max Leaf Nodes**

(optional)

**Save the model as**

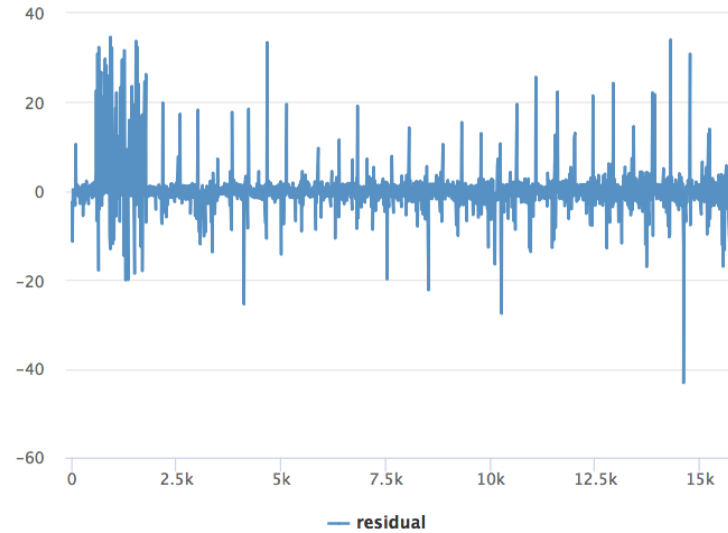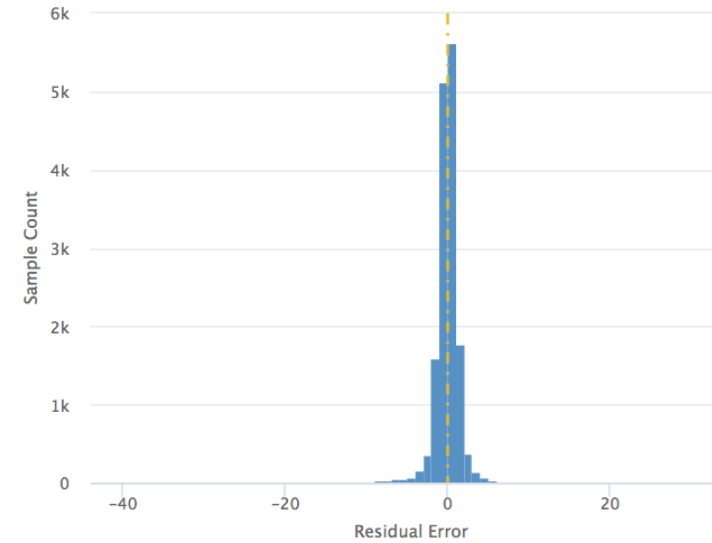server_power

[ Fit Model ]  [ Schedule Training ]  [ Open in Search ]  [ Show SPL ]

# Experiments: Validate

# Experiments: Deploy

# Experiments

- ▸ Predict Numeric Fields
- ▸ Predict Categorical Fields
- ▸ Detect Numeric Outliers
- ▸ Detect Categorical Outliers
- ▸ Forecast Time Series
- ▸ Cluster Numeric Events

130.60.4 - - [07/Jan 18:10:57:153] "GET /category.screen?category_id=GIFTS&JSESSIONID=SD15L4FF10ADFF10 HTTP 1.1" 404 720 "http://buttercup-shopping.com/cart.do?action=view&itemId=EST-6&product_id=FL-SW-01" ...
128.241.220.82 - - [07/Jan 18:10:57:123] "GET /product.screen?product_id=FL-DSH-01&JSESSIONID=SD5SL7FF6ADFF9 HTTP 1.1" 404 3322 "http://buttercup-shopping.com/cart.do?action=purchase&itemId=EST-26&product_id=GIFTS" ...
317 27.160.0.0 - - [07/Jan 18:10:56:156] "GET /oldlink?item_id=EST-26&JSESSIONID=SD5SL9FF1ADFF3 HTTP 1.1" 200 1318 "http://JSESSIONID=SD9SL4FF4ADFF7 HTTP 1.1" 200 2423 "http://buttercup-shopping..." ...

# Let's Build a Custom Model!

# What's New?

since last .conf

# Major Highlights

(since .conf 2017)

**Splunk Machine Learning Toolkit Updates**

Includes new features for the Experiment Framework, algorithms, pre-processing steps, validation options etc.

**Python for Scientific Computing 1.3 Update**

Updated libraries giving you access to new and modified algorithms and its parameters.

**Splunk MLTK Connector for Apache Spark™**

Massive model building with MLlib directly from Splunk and SPL, No Scala skills required. (Limited Availability Release)

**GitHub MLTK Community**

Leverage and share algorithms collaboratively with the broader MLTK community

**Splunk MLTK Container for Tensor Flow**

Container based neural networks, leveraging GPUs/CPUs.

splunk> .conf18

# Splunk Machine Learning Toolkit Updates

- **Experiment Management Framework:** A unified UI that provides the ability to:

  → Set roles based access control on experiments
  → Browsing and filtering pre-built models
  → Monitoring and scheduling alerts and searches
  → Getting history statistics about experiment's previous runs and alerts

- **Score Command:** A new command for validating models and statistical tests for any use case, shipping with N algorithms today.

- **K-fold Cross-validation:** A popular and powerful way to quickly reduce model overfitting.

- **UI for MLSPL.CONF:** A interface to give user the power to change the safe settings if required for app level mlspl configuration.

# Splunk Machine Learning Toolkit Updates

- ## New out-of-the-box algorithms

  - **Local Outlier Factor** : Unsupervised anomaly detection.
  - **Multi-layer Perceptron Classifier :** Neural network-based supervised classifier.
  - **Robust Scaler** : Re-scaling algorithm that is robust to outliers.
  - **X-Means:** Unsupervised clustering.

- ## New pre-processing steps

  - **Term Frequency-Inverse Document Frequency** : Feature extraction on unstructured text.
  - **Field Selector:** Feature selection.

splunk> .conf18

# MLTK 4.0 - Python for Scientific Computing 1.3

# Splunk MLTK Connector for Apache Spark™ (Limited Availability Release)



| sfit [y from x* into "model"]

BYO APACHE Spark™

Industrial Assets

Consumer and Mobile Devices

| sapply "model"

persisted model

IoT

Real Time **splunk>** Search

Alert

Visualize

Not on Cloud

Cloud

splunk> .conf18

# MLTK 4.0 - *Splunk Community for MLTK Algorithms on GitHub*

**A community github for sharing algorithm files**

*"The creation of the Splunk Community for MLTK Algorithms on GitHub will help us find new functionality within the catalog at a much faster rate, which will allow us to get even more use out of the Splunk Machine Learning Toolkit,"*
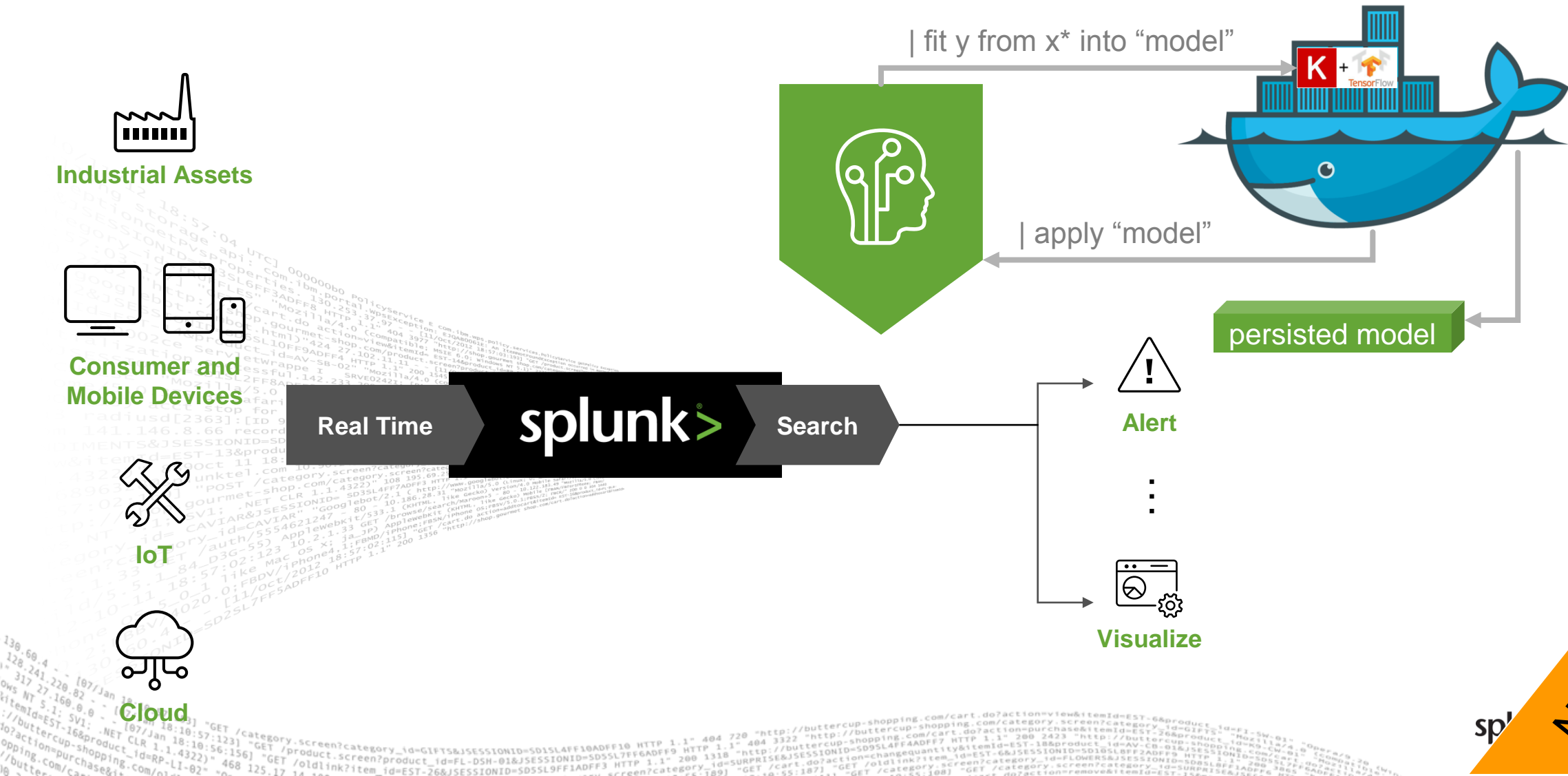*said **Nathan Worsham, IS Security Administrator, Pinnacol***

Splunk MLTK Container for TensorFlow™ (via PS Whiteglove)

# Want to know more?

Download the slides and recordings for these sessions.

**FN1364 - Using Spark and MLLib for Large Scale Machine Learning With Splunk Machine Learning Toolkit**

(Thursday, Oct 04, 11:00 a.m. - 11:45 a.m.)

**Lin Ma**, Principal Software Engineer, Splunk
**Fred Zhang**, Principal Data Scientist, Splunk

**FN1409 - Thank You for Sharing: Expanding Machine Learning using Splunk MLTK GitHub Collaboration**

(Thursday, Oct 04, 11:00 a.m. - 11:45 a.m.)

**Gyanendra Rana**, Senior Product Manager, Splunk
**Nathan Worsham**, IS Security Administrator, Pinnacol Assurance

**FN1478 - Exciting, To-Be-Announced Platform Session**

Wednesday, Oct 03, 4:30 p.m. - 5:15 p.m.

**Phillipp Drieger**, Staff Machine Learning Architect, Splunk

splunk> .conf18

# Customer Success

splunk> .conf18

# ML Success Story

**Consumer Credit Reporting Agency**

**Acting on a Critical Customer Outages before the Customer Calls You**

RECURSION pharmaceuticals

**Many different machines are part of the drug discovery process, and machines acting abnormally mean a loss in efficiency and increased costs.**

TELUS
the future is friendly®

UE

NodeB

**Detect interference in cell towers Re-configure underperforming cells for optimal services levels**

Telco

**Improving cell tower uptime and reducing repair truck rolls with anomaly detection and root cause analysis**

splunk> .conf18

# ML Success Story



**Entertainment Company**

**Predicting and averting potential gaming outage conditions with finer-grained detection**

**Preventing fraud by Identifying malicious accounts and suspicious activities**



**Find errors in server pools, then prioritize actions and associate root cause**

**Online Retailer**



**Failed orders detected in real time to avoid lost revenue and unhappy customers**



**Predicting Student Achievement and taking action to improve grades**

130.60.4 - - [07/Jan 18:10:57:153] "GET /category.screen?category_id=GIFTS&JSESSIONID=SD1SL4FF10ADFF10 HTTP 1.1" 404 720 "http://buttercup-shopping.com/cart.do?action=view&itemId=EST-6&product_id=F1-SW-01 128.241.220.82 - - [07/Jan 18:10:57:123] "GET /product.screen?product_id=FL-DSH-01&JSESSIONID=SD5SL7FF6ADFF9 HTTP 1.1" 404 3322 "http://buttercup-shopping.com/category.screen?category_id=GIFTS ows NT 5.1; SV1; .NET CLR 1.1.4322) "GET /oldlink?item_id=EST-26&JSESSIONID=SD9SL4FF4ADFF7 HTTP 1.1" 200 2423 "http://buttercup-shopping itemId=EST-16&product_id=RP-LI-02- .NET CLR 1.1.4322) "GET /product.screen?product_id=FL-DSH-01&JSESSIONID=SD5SL9FF1ADFF3 HTTP 1.1" 200 1318 "http://buttercup-shopping.com/category.screen?category_id=SURPRISE&JSESSIONID=SD5SL9FF1ADFF /action=purchase&itemId=EST-26&JSESSIONID=SD5SL9FF1ADFF3 HTTP 1.1" 200 1318 "http://buttercup-shopping.com/oldlink?item_id=EST-26&JSESSIONID=SD5SL9FF1ADFF3

splunk> .conf18

**The latest release of Splunk Machine Learning Toolkit makes it significantly easier to process large amounts of data and find patterns to see what's right or wrong. Splunk's continued evolution of the Experiment Management Framework, including new tools to help validate our machine learning models, streamlines the complicated process of operationalizing machine learning.**

**Sundaresh Ramanathan, Director, IT Operations Analytics, Kinney Group, Inc.**

# Thank You

**Don't forget to rate this session
in the .conf18 mobile app**

.conf18

splunk>