



Using Spark to Tune Spark

Adrian Popescu, Shivnath Babu

#AI7SAIS

Meet the speakers



Adrian Popescu

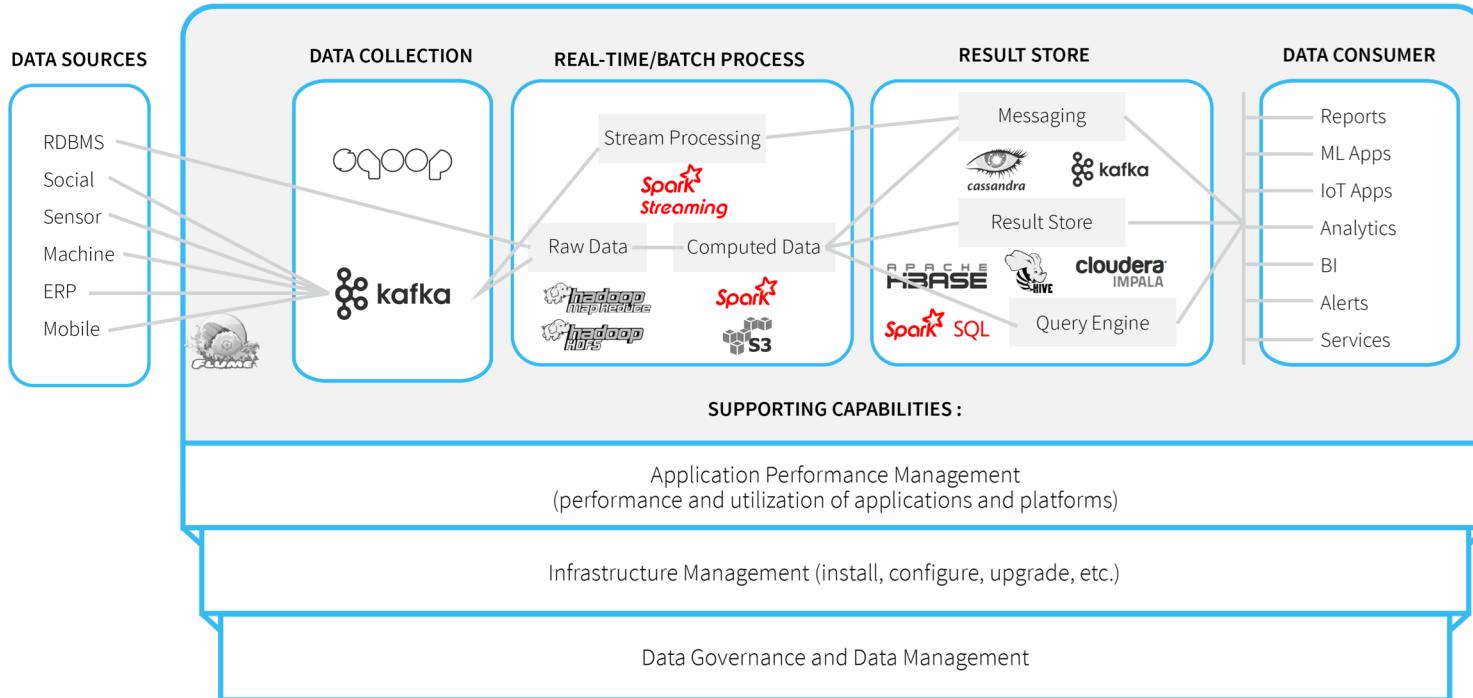
- Data engineer at Unravel
- PhD from EPFL, Switzerland
- 8+ years of experience in performance monitoring & modeling of data management systems
- Focusing on tuning and optimization of Big Data apps



Shivnath Babu

- Cofounder and CTO at Unravel, Adjunct Professor at Duke University
- Focusing on ease-of-use and manageability of data-intensive systems
- Recipient of US National Science Foundation CAREER Award, three IBM Faculty Awards, HP Labs Innovation Research Award

Many apps are being built in Spark

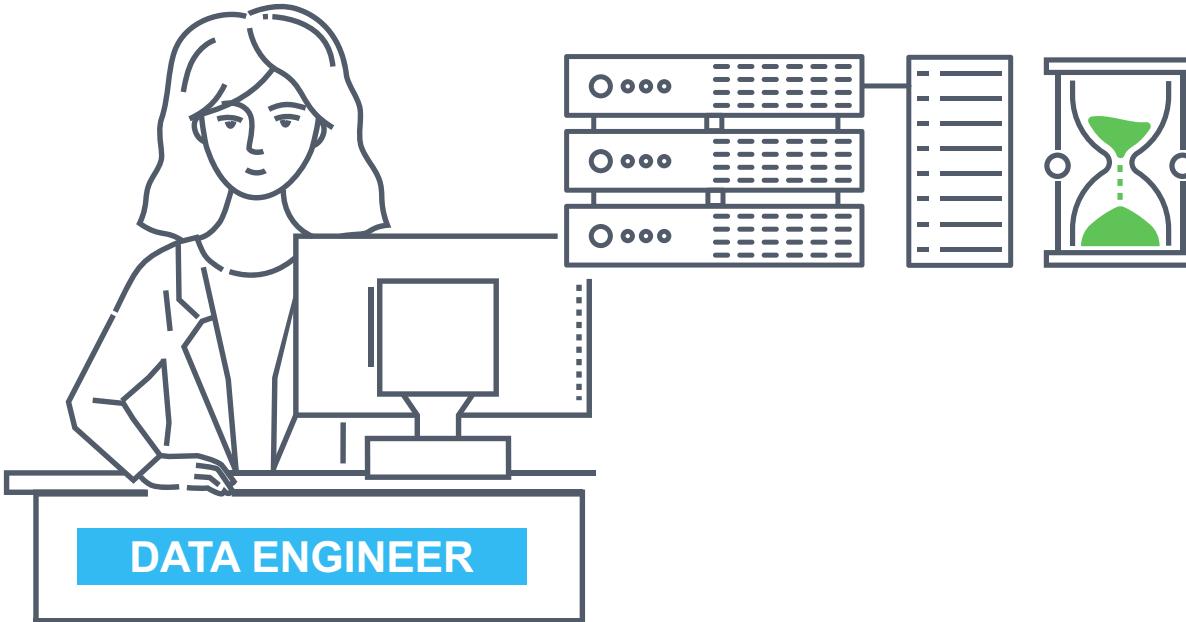


**But, let us face it:
Running Spark apps in
production is hard**

My app often fails with Out of Memory...



My app is too slow...



DATA ENGINEER

My app is missing SLA...



DATA PIPELINE OWNER

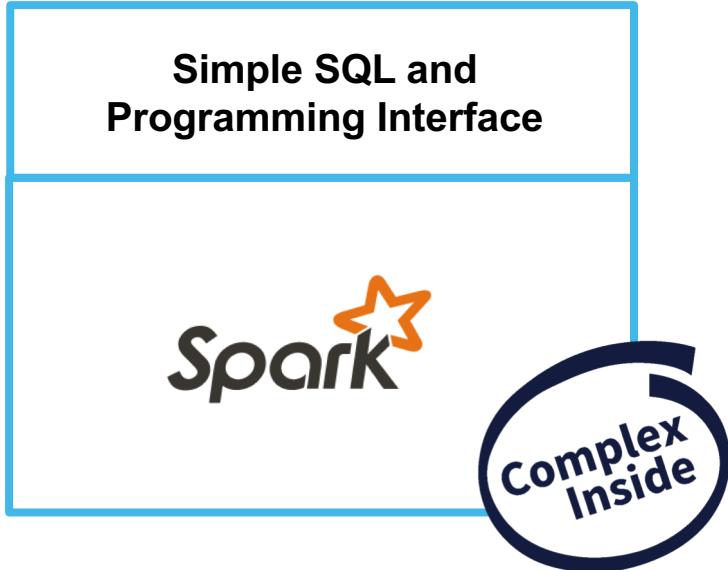
This rogue app is wasting resources and reducing cluster throughput



Many factors affect app performance



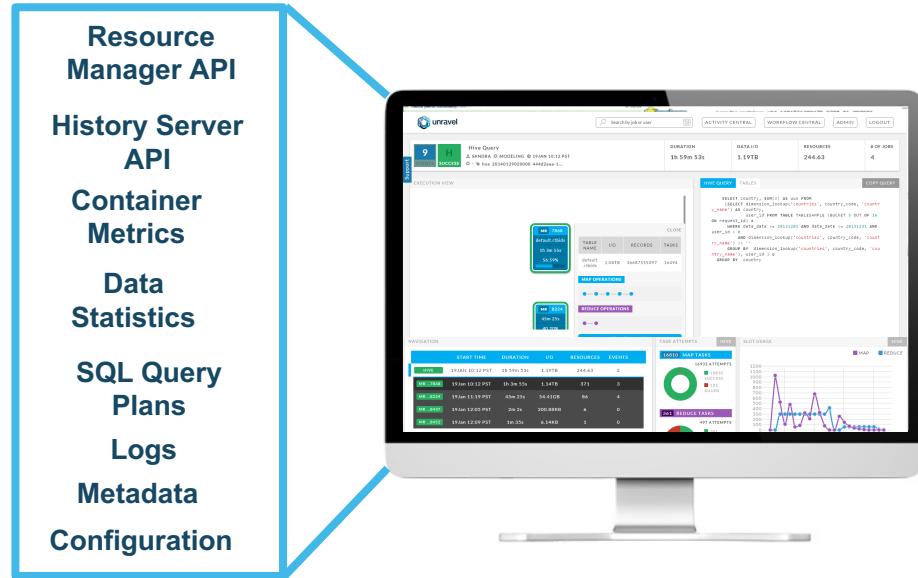
To add to Spark's complexity



- **Many types of Spark apps**
 - SQL
 - Streaming
 - AI/ML
 - Graph
 - Scala/Python/R
- **Many app submission methods in Spark**
 - CLI
 - Thrift Server
 - Notebooks like Zeppelin, Jupyter, Hue
 - ETL tools like Informatica, Pentaho, and Talend
 - Schedulers like Airflow, Autosys, Control M, Oozie, Tidal, TWS
- **Many infrastructure choices for Spark**
 - On-premises multi-tenant clusters
 - Transient cloud clusters
 - Auto-scaling clusters
 - Containerized deployments

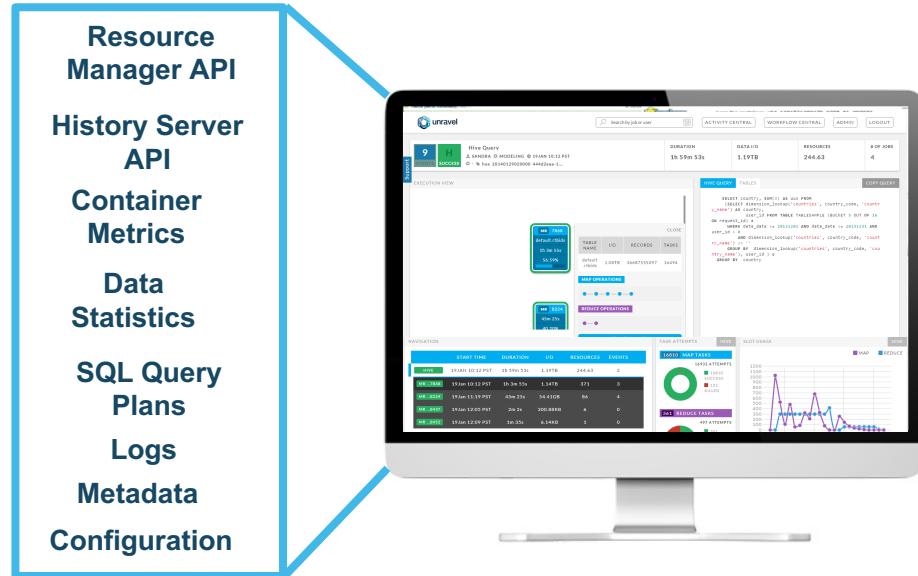
**Can we convert this problem
into a data problem?**

First: Bring all monitoring data to a single platform

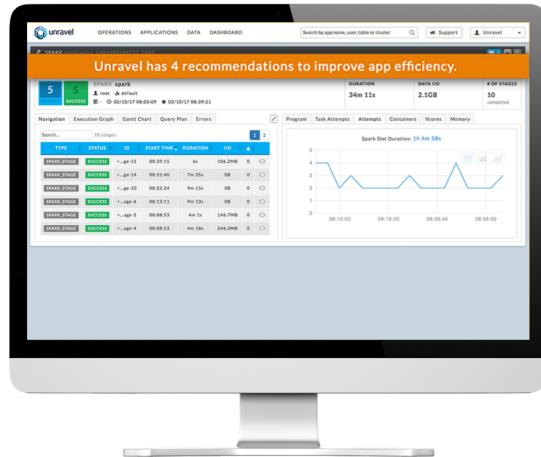


One complete correlated view.

Then: Apply intelligent algorithms to analyze the data automatically



One complete correlated view.

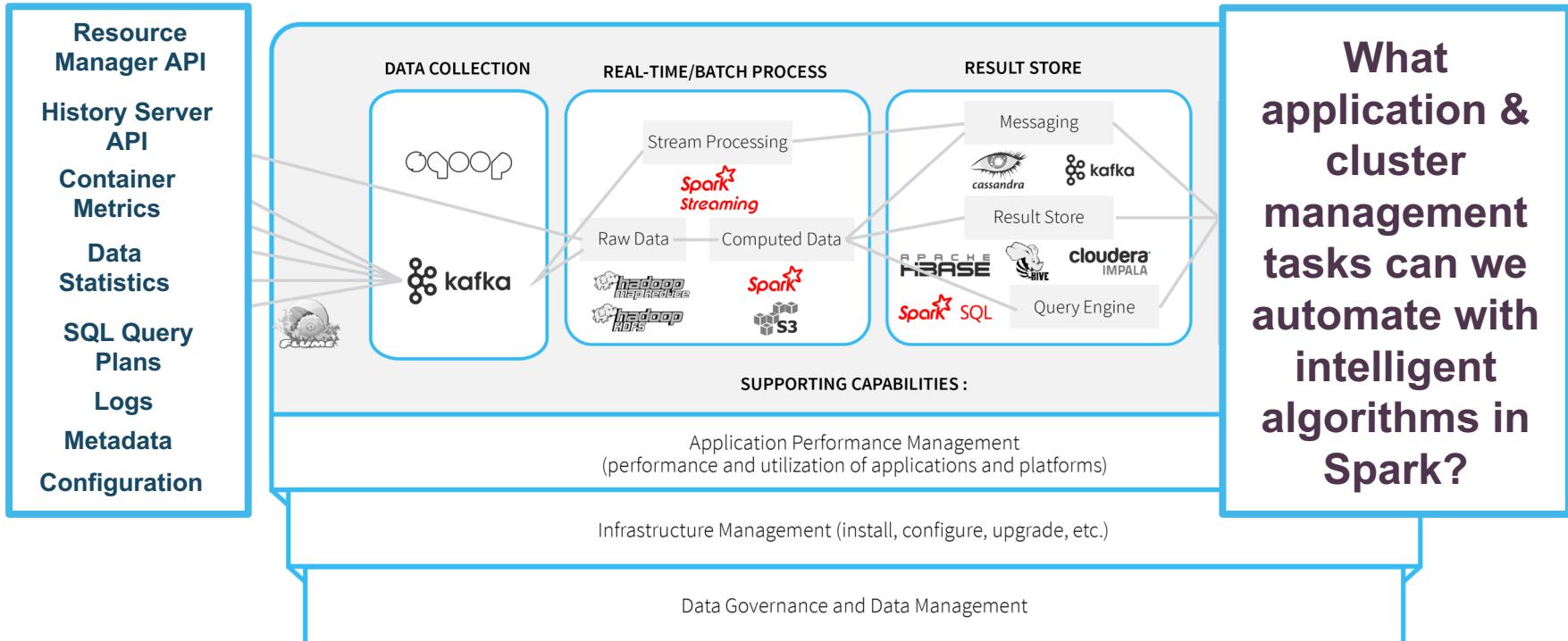


Built-in intelligence.



USE THE
RIGHT TOOL
FOR THE JOB

Why not use Spark itself?



Let us take three (hard) tasks

- **Failures in Spark**
- SLA management for real-time data pipelines
- Application autotuning

Manual Root Cause Analysis of Spark Failures

Typical Failure in Spark

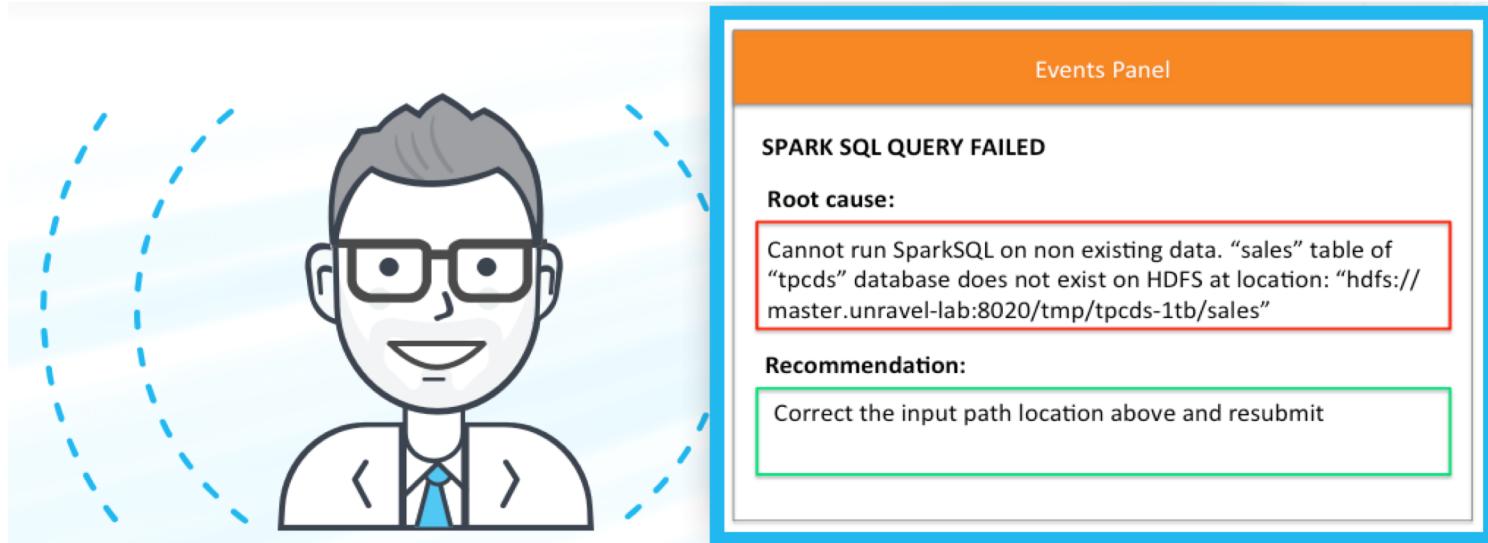


```
...  
.spark.sql.execution.Converter.toSafe, doExecute(rowFormatConverters.scala:56)  
.spark.sql.execution.SparkPlan$anonfun$execute$5.apply(SparkPlan.scala:132)  
.spark.sql.execution.SparkPlan$anonfun$execute$5.apply(SparkPlan.scala:130)  
.spark.rdd.RDDOperationScopes.withScope(RDDOperationScope.scala:150)  
.spark.execution.SparkPlan.execute(SparkPlan.scala:130)  
.spark.execution.Exchange.prepareShuffleDependency(Exchange.scala:164)  
.spark.sql.execution.Exchange$anonfun$doExecute$1.apply(Exchange.scala:254)  
.spark.sql.execution.Exchange$anonfun$doExecute$1.apply(Exchange.scala:248)  
.spark.sql.catalyst.errors.package$.attachTree(package.scala:48)  
  
.spark.sql.catalyst.errors.package$$TreeNodeException: execute, tree:  
parJoining(nation#71,o_year#72,...,None  
[key=[nation#71,o_year#72], functions=[(sum(amount#73),mode=Partial,isDistinct=false)], output=[nation#71,o_year#72,sum#78])  
:#50 AS nation#71,year(cast(_orderdate#20 as date)) AS o_year#72,((l_extendedprice#5 * (1.0 - l_discount#6)) - (ps_supplycost#28 * l_quantity#4)) AS  
in [_orderkey#16L, [_orderkey#0L]  
orderkey#16L ASC], false, @  
tenExchange hashpartitioning(_orderkey#16L,200), None  
ject [_orderkey#16L,_orderdate#20]  
Scan ExistingRDD[_orderkey#16L,_custkey#17L,_orderstatus#18,_totalprice#19,_orderdate#20,_orderpriority#21,_clerk#22,_shippriority#23L,_comm  
orderkey#0L ASC], false, @  
tenExchange hashpartitioning(_orderkey#0L,200), None  
ject [_extendedprice#5,_discount#6,_quantity#4,_orderkey#0L,n_name#50,ps_supplycost#28]  
SortMergeJoin[_partkey#30L,[_partkey#1L]  
:- Sort [_partkey#30L ASC], false, @  
:+ TungstenExchange hashpartitioning(_partkey#30L,200), None  
:+ Project [_partkey#30L]  
:+ Filter Contains(p_name#31,ghost)  
:+ Scan ExistingRDD[_partkey#30L,p_name#31,p_mfgr#32,p_brand#33,p_type#34,p_size#35L,p_container#36,p_retailprice#37,p_comment#38]  
:+ Sort [_partkey#1L ASC], false, @  
:+ TungstenExchange hashpartitioning(_partkey#1L,200), None  
:+ Project [_extendedprice#5,_discount#6,_quantity#4,_partkey#1L,_orderkey#0L,n_name#50,ps_supplycost#28]  
:+ SortMergeJoin[_suppkey#26L,_partkey#25L,[_suppkey#2L,_partkey#1L]  
:- Sort [_suppkey#26L ASC,ps_partkey#25L ASC], false, @  
:+ TungstenExchange hashpartitioning(_suppkey#26L,ps_partkey#25L,200), None
```

5 levels of stack
traces of this form

- Many levels of correlated stack traces
- Identifying the root cause is hard and time consuming

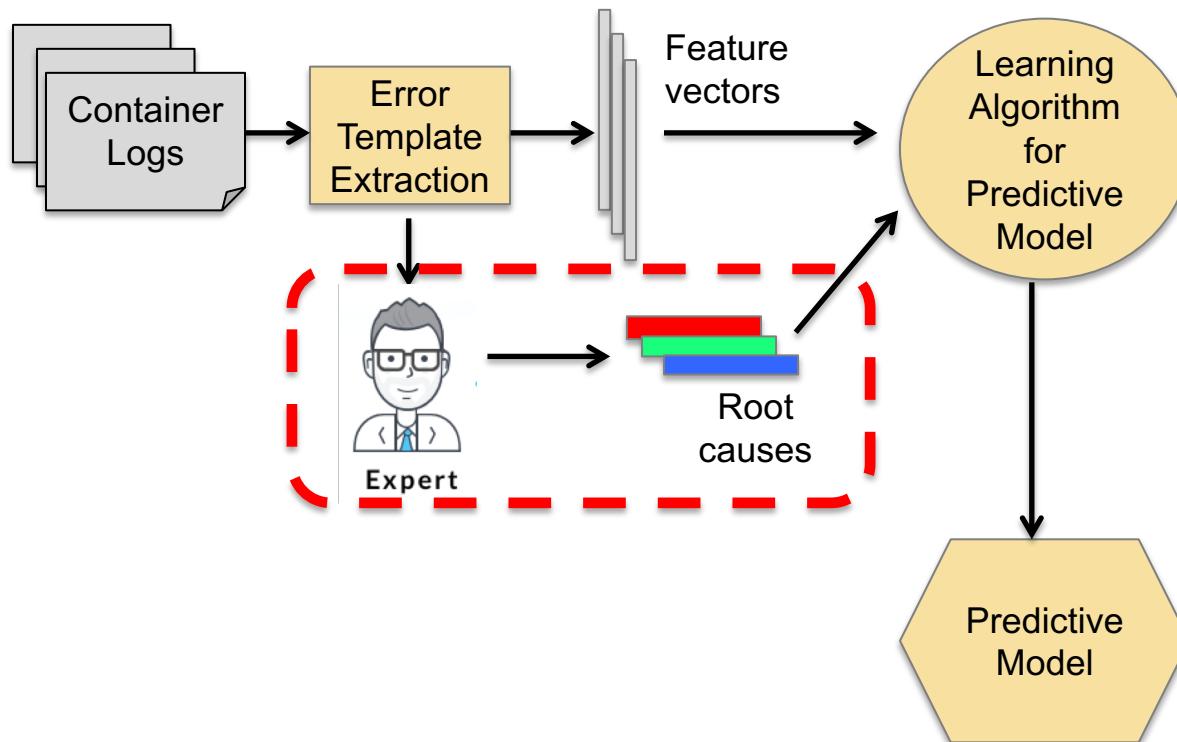
Automated Root Cause Analysis of Spark Failures



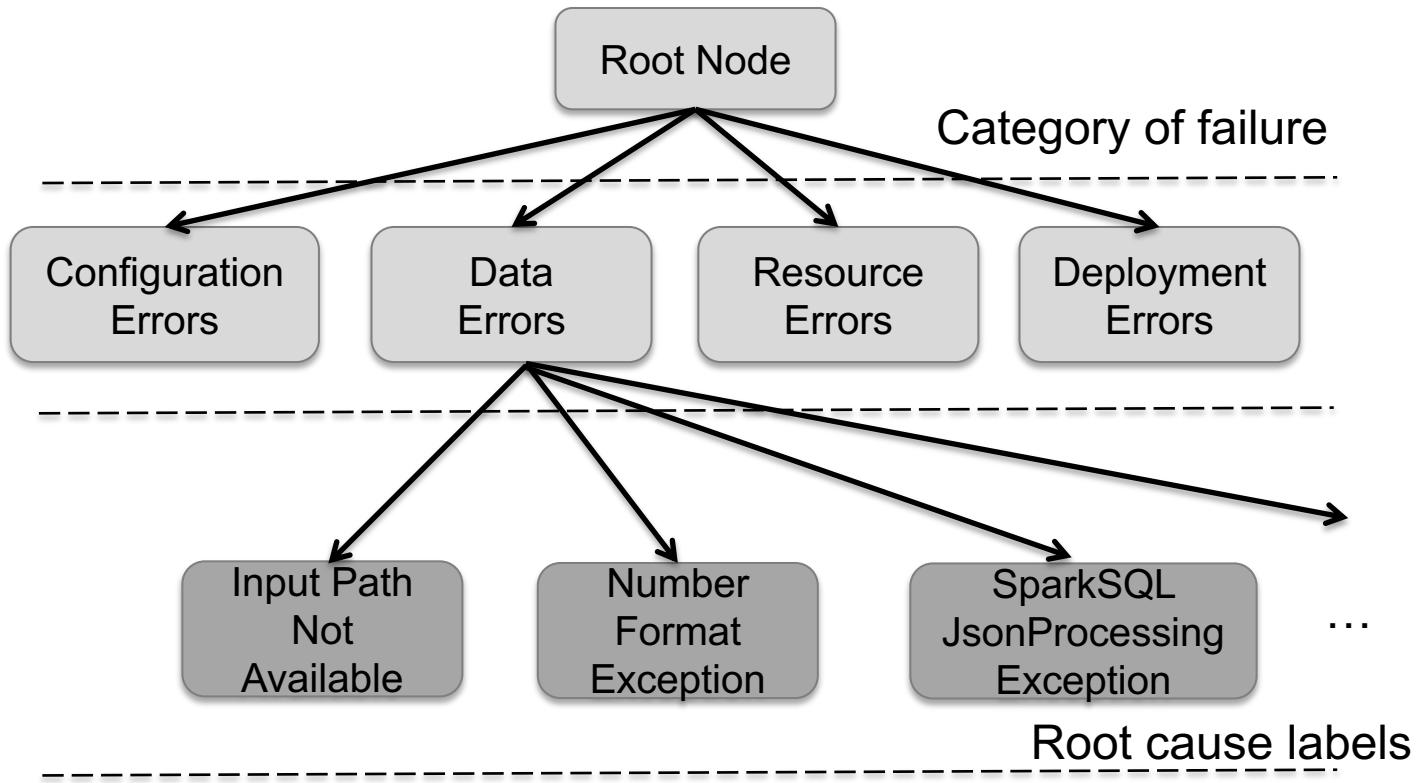
The illustration shows a data scientist with short grey hair and glasses, wearing a white shirt and a blue tie. He is looking towards the right side of the screen. On the right side, there is a rectangular window titled "Events Panel". Inside the panel, the text "SPARK SQL QUERY FAILED" is displayed in bold. Below it, the section "Root cause:" is shown in bold, followed by a red-bordered box containing the error message: "Cannot run SparkSQL on non existing data. "sales" table of "tpcds" database does not exist on HDFS at location: "hdfs://master.unravel-lab:8020/tmp/tpcds-1tb/sales"". Underneath this, the section "Recommendation:" is shown in bold, followed by a green-bordered box containing the text "Correct the input path location above and resubmit".

- **Reduce troubleshooting time from days to seconds**
- Improve productivity of data scientists and analysts

Automatic Root Cause Analysis



We have created a Failure Taxonomy



Two Ways to get Root-Cause Labels

- Manual diagnosis by a domain expert
- Automatic injection of the root cause

Unravel's Large-scale Lab Framework for Automatic Root Cause Analysis

Environment:

- Lab created on demand on cloud or on-premises
- Workloads are run and failures are injected

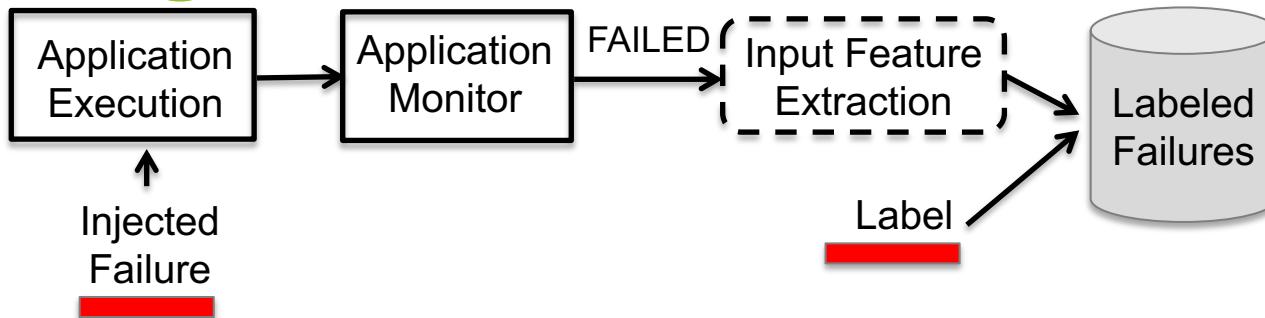
Spark and multi-tenant Workloads:

- Variety of workloads: Batch, ML, SQL, Streaming, etc.

Failures:

- Large set of root causes learned from customers & partners. Constantly updated
- Continuously inject these root causes to train & test models for root-cause prediction

Injecting Failures



Injected failure examples:

- Invalid input
- Invalid memory configuration
- OOME: Java heap space
- OOME: GC overhead limit
- Container killed by YARN
- Runtime incompatibility
- No space left on device
- Transformations inside other transformations
- Runtime error
- Arithmetic error
- Invalid configuration settings

Extracting Input Features from Logs

```
java.lang.OutOfMemoryError: Java heap space
    at
scala.reflect.ManifestFactory$$anon$9.newArray(Manifest.scala:114)
    at
scala.reflect.ManifestFactory$$anon$9.newArray(Manifest.scala:112)
    at ...
```

- Extracting stack traces and error messages
- Tokenize by class names and words
- Create a vocabulary of words from all words collected

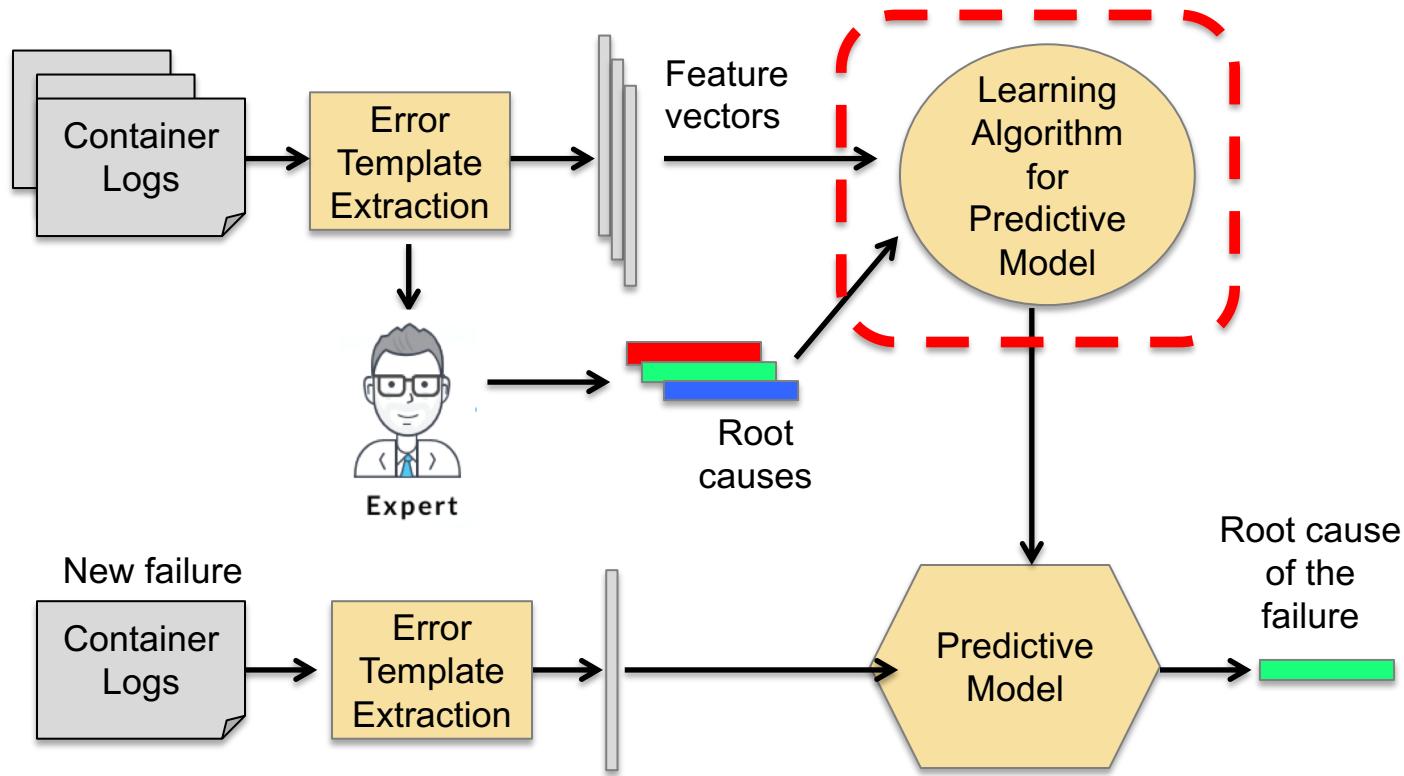
Tokens example:

```
java.lang.OutOfMemoryError Java heap space at
scala.reflect.ManifestFactory$$anon$9.newArray(Manife
st.scala:114)
```

Input Feature extraction

- **Bag of Words with TF-IDF**
 - Computes a vocabulary of words
 - Uses TF-IDF to reflect importance of words in a document
- **Doc2Vec**
 - Maps words, paragraphs, or documents to multi-dimensional vectors
 - Evaluates the placement of words wrt neighboring words
 - Uses a 3-layer neural network

System Architecture



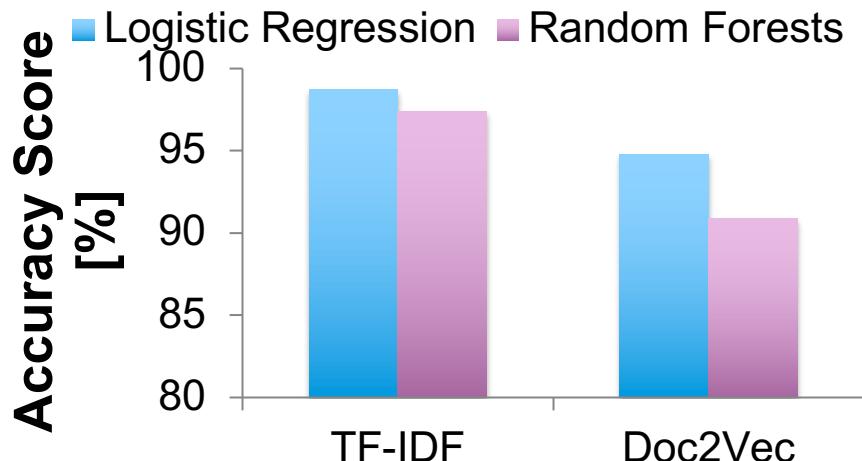
Predictive Models

- **Shallow Learning**
 - Logistic Regression
 - Random forests
- **Deep Learning**
 - Neural networks

Very easy to
implement these
in Spark

Predicting the Root Cause of Failures

- **Training** and **testing** with injected failures
- Test to train data set ratio **75% to 25%**
- **Models**: logistic regression, random forests



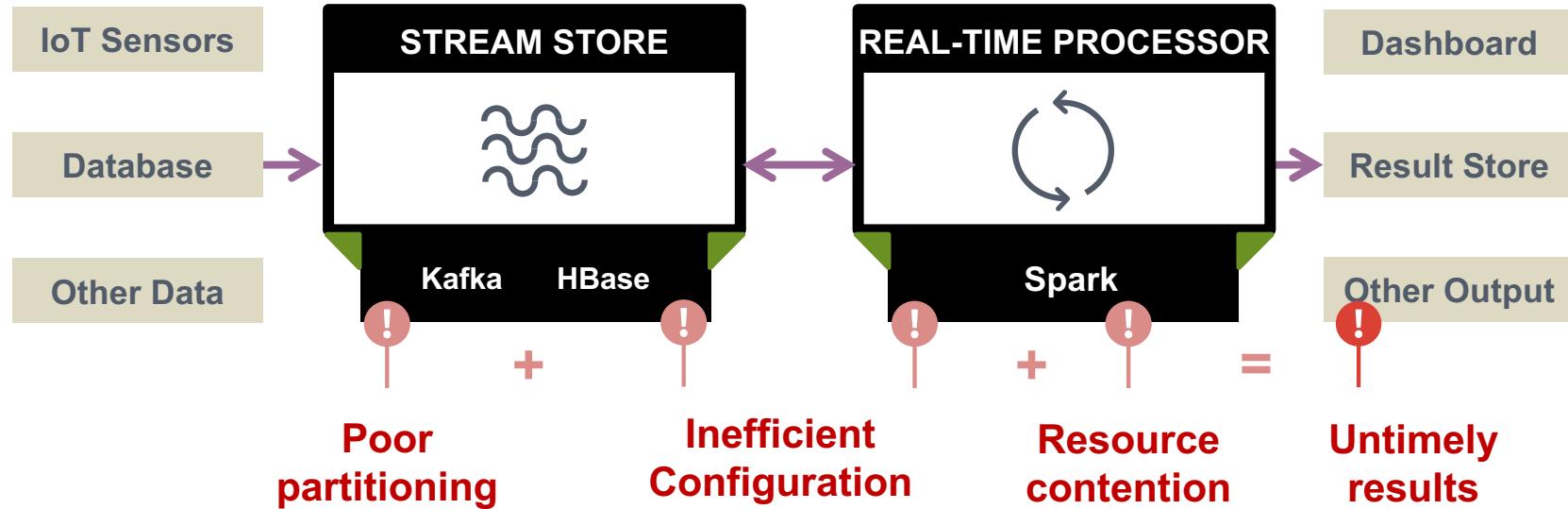
Work with
deep learning
is in progress

See our talk at
Strata, NY 2017
for more details

Let us take three (hard) tasks

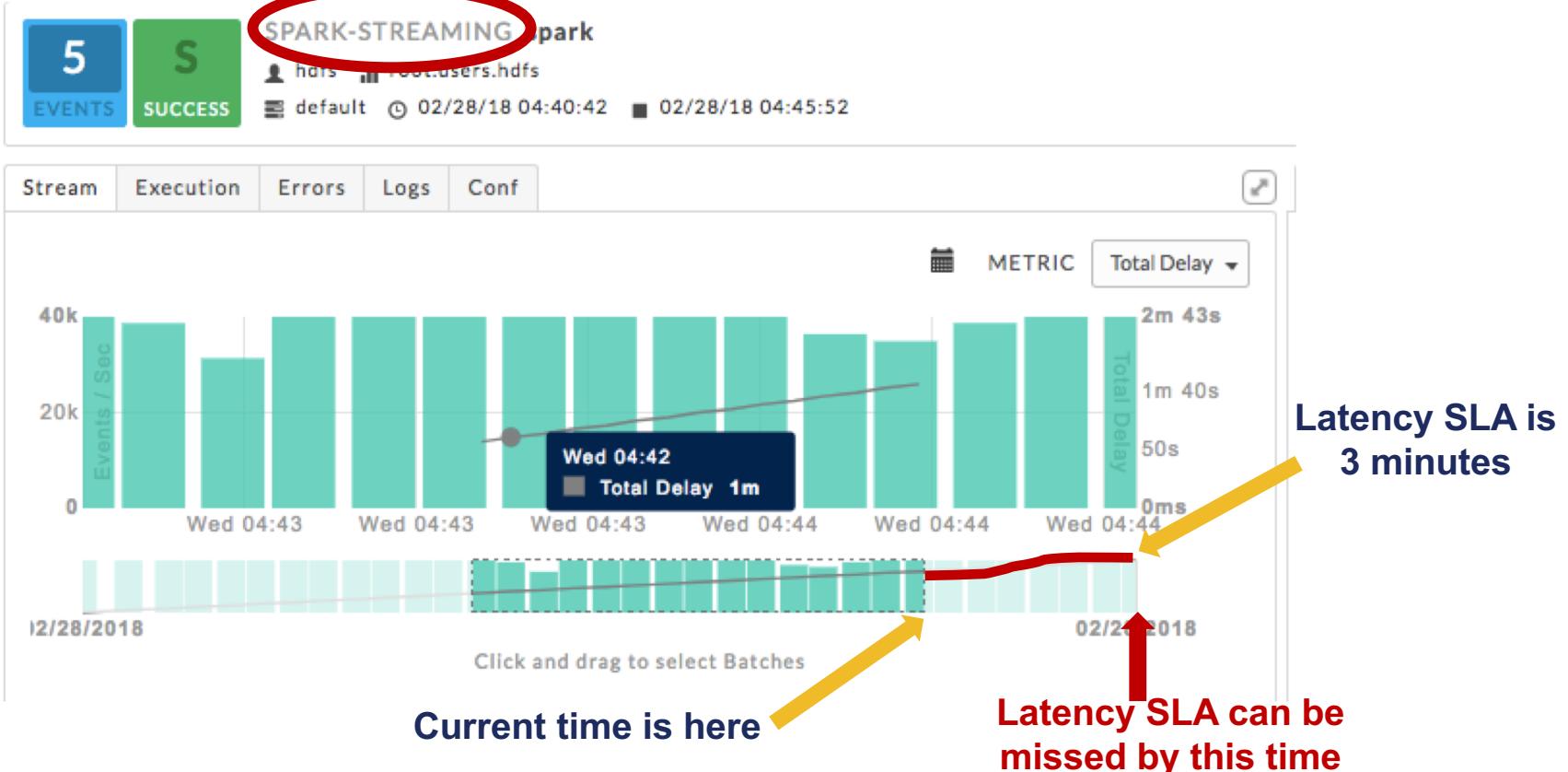
- Failures in Spark
- **SLA management for real-time data pipelines**
- Application autotuning

Many problems can cause unreliable performance in streaming pipelines

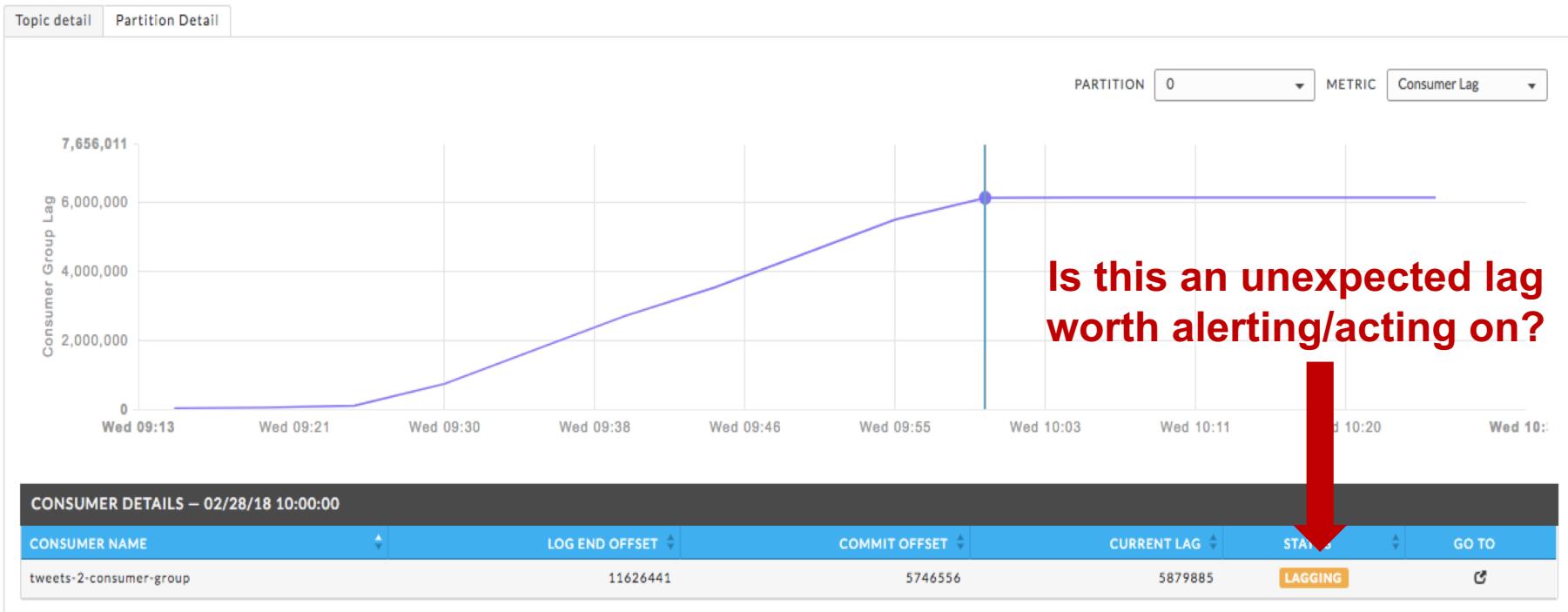


**Forecasting & Anomaly
Detection are important for
streaming pipelines**

Predicting when SLAs are in danger of being missed



Detecting anomalies is tricky



What is an anomaly?

1. An unexpected change that needs your attention
2. Smart alerts/actions:
 1. False positives should be minimal
 2. False negatives should be minimal

Time Series Analysis in Spark

Anomaly Detection with Apache Spark

A Gentle Introduction



Sean Owen // Director of Data Science



DATA SCIENCE AND ENGINEERING AT SCALE

AMSTERDAM • OCTOBER 27-29

Real-Time Anomaly Detection with Spark MLLib, Akka and Cassandra

Natalino Busa
Data Platform Architect at Ing

Time Series Analytics with Spark

Simon Ouellette
Faimdata



PySpark for Time Series Analysis

David Palaitis
Two Sigma Investments



Time Series Analysis in Spark

Anomaly
Detection with
Apache Spark

A Gentle Introduction

Sean Owen // Director of Data Science



DATA SCIENCE AND ENGINEERING AT SCALE

AMSTERDAM • OCTOBER 27-29

**See our talk at Strata, San Jose 2018
for more insights**

Time Series Analytics
with Spark

Simon Quellette



Real-Time Anomaly Detection
with Spark MLlib, Akka and
Cassandra

Natalino Busa
Data Platform Architect at Ing

PySpark for Time Series Analysis

David Palaitis
Two Sigma Investments

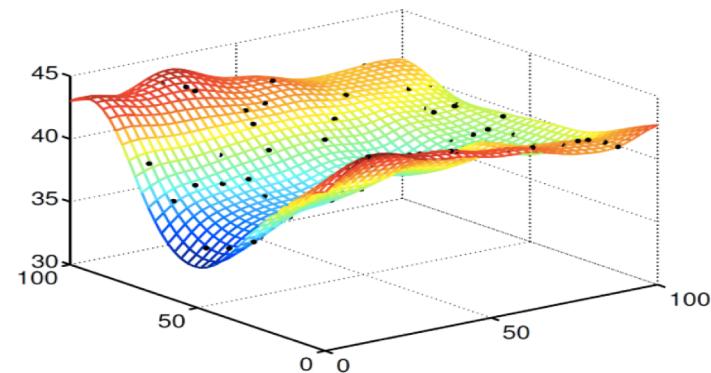


Let us take three (hard) tasks

- Failures in Spark
- SLA management for real-time data pipelines
- **Application autotuning**

| | |
|---|------|
| spark.driver.cores | 2 |
| spark.executor.cores | 10 |
| ... | |
| spark.sql.shuffle.partitions | 300 |
| spark.sql.autoBroadcastJoinThreshold | 20MB |
| ... | |
| SKEW('orders', 'o_custId') | true |
| spark.catalog.cacheTable("orders") | true |
| ... | |

PERFORMANCE



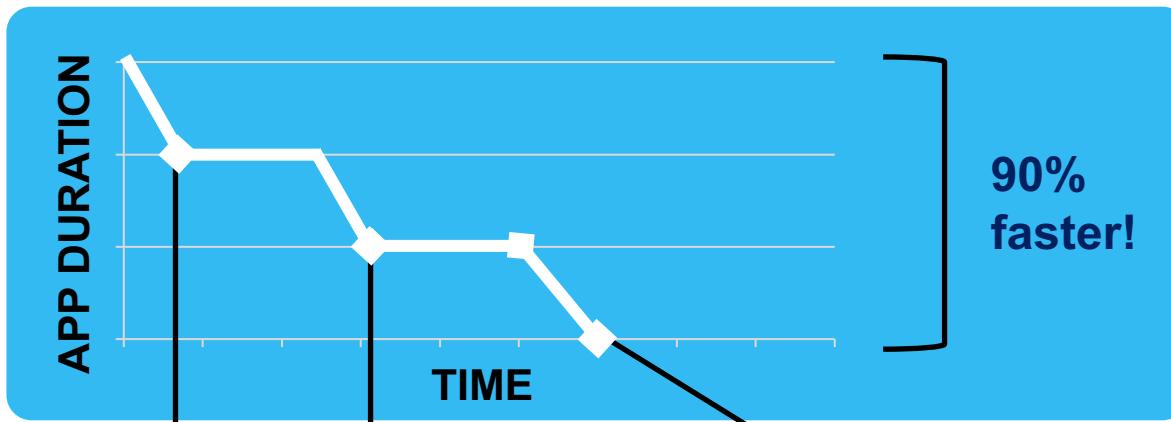
Today, tuning is often by trial-and-error

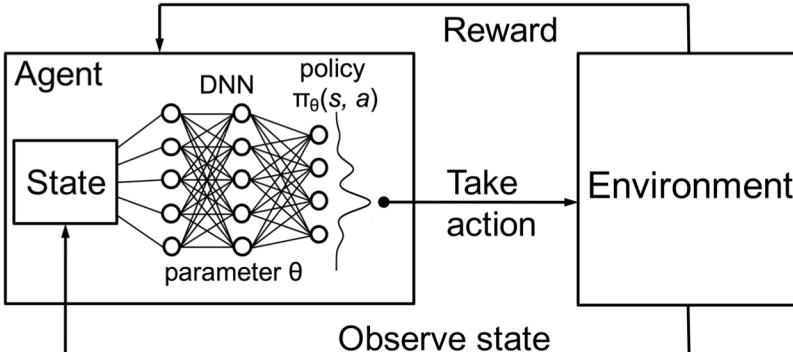
Imagine a new world



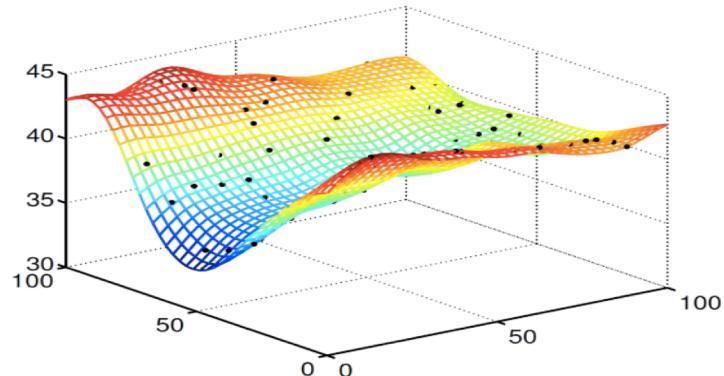
"I need to make this app faster"

Imagine a new world





Reinforcement Learning



Response Surface Methodology

Tuning Database Configuration Parameters with iTuned

Songyun Duan, Vamsidhar Thummala, Shivnath Babu*
 Department of Computer Science
 Duke University
 Durham, North Carolina, USA
 {syduan,vamsi,shivnath}@cs.duke.edu

ABSTRACT

Database systems have a large number of configuration parameters that control memory distribution, I/O optimization, costing of query plans, parallelism, many aspects of logging, recovery, and

Amy recalls that the database has *configuration parameters*. For lack of better understanding, she had set them to default values during installation. The parameters may need tuning, so Amy pulls out the 1000+ page database tuning manual. She finds many dozens

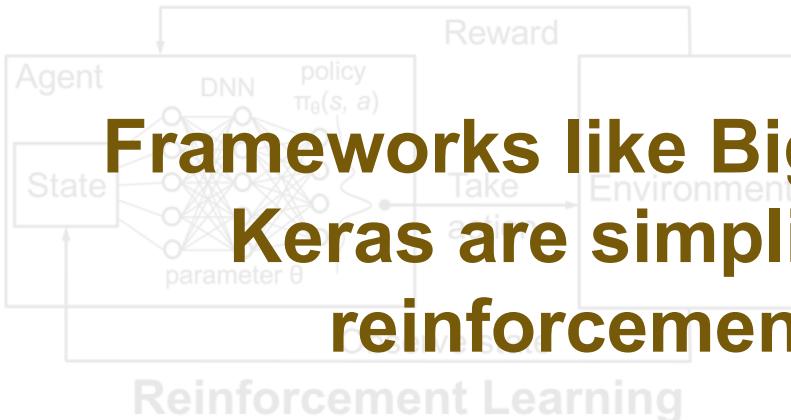
Xplus: A SQL-Tuning-Aware Query Optimizer

Herodotos Herodotou and Shivnath Babu*
 Department of Computer Science
 Duke University
 {hero,shivnath}@cs.duke.edu

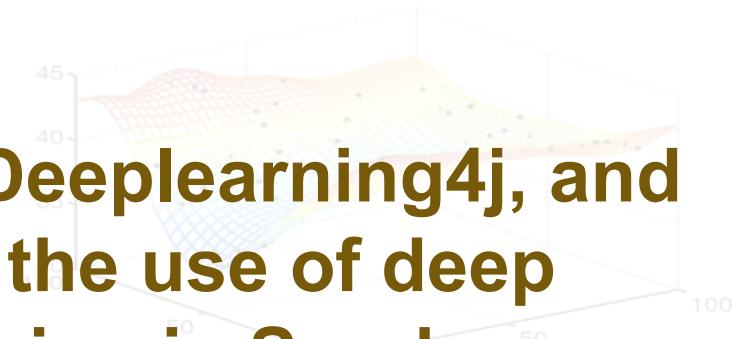
ABSTRACT

The need to improve a suboptimal execution plan picked by the query optimizer for a repeatedly run SQL query arises routinely. Complex expressions, skewed or correlated data, and changing con-

step in to lead the optimizer towards a good plan [6]. This process of improving the performance of a “problem query” is referred to in the database industry as *SQL tuning*. Tuning a problem query is critical in two settings:



Frameworks like BigDL, DeepLearning4j, and Keras are simplifying the use of deep reinforcement learning in Spark



Tuning Database Configuration Parameters with iTuned

See our talk on Unravel Sessions at Spark+AI Summit 2018 on this problem

Songyun Duan, Vamsi Malladi, Shivnath Babu

Department of Computer Science

Duke

Durham, North Carolina 27708

{syduan,vamsi,shivnath}@cs.duke.edu

Xplus: A SQL-Tuning-Aware Query Optimizer

Herodotus Herodotou and Shivnath Babu

Department of Computer Science

Duke University

durham.cs.duke.edu

ABSTRACT

Database systems have a large number of configuration parameters that control memory distribution, I/O optimization, costing of query plans, and many aspects of logging, recovery, and

Amy recalls that the database has *configuration parameters*. For lack of better understanding, she had set them to default values during installation. The parameters may need tuning, so Amy pulls out the 1000+ page database tuning manual. One finds many sections on how to tune various parameters for different workloads. Amy

ABSTRACT

The need to improve a suboptimal execution plan picked by the query optimizer for a repeatedly run SQL query arises routinely. Complex expressions, skewed or correlated data, and changing con-

text in the query can make this task challenging. In this paper, we propose Xplus, a SQL-tuning-aware query optimizer that uses machine learning to predict the performance of a query under different parameter settings. Xplus can automatically tune the database configuration parameters to lead the optimizer towards a good plan [6]. This process of improving the performance of a “problem query” is referred to in the database industry as *SQL tuning*. Tuning a problem query is

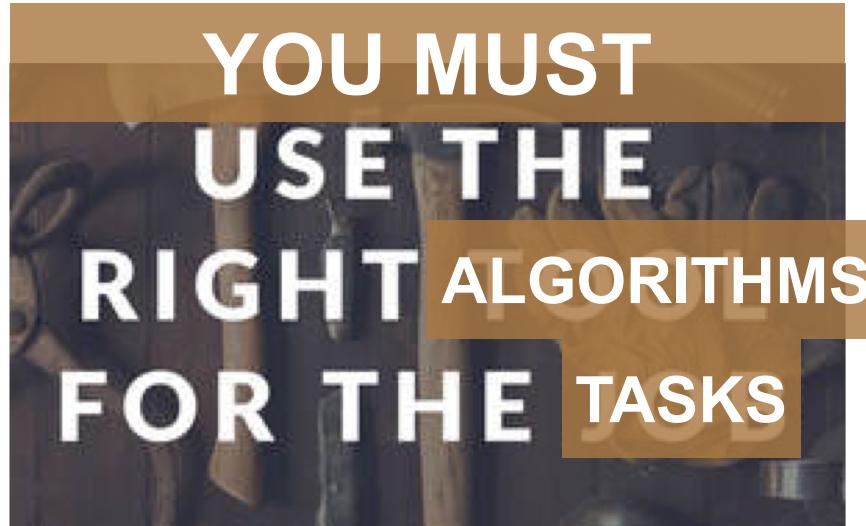
In Summary

- Spark Performance Management can be addressed using Spark itself



In Summary

- Spark Performance Management can be addressed using Spark itself



Sign up for a free trial. We value your feedback!
<https://unraveldata.com/free-trial>

Join the Unravel team!
adrian & shivnath [@unraveldata.com]