



# Deep Learning at Extreme Scale (in the Cloud) with the Apache Kafka Open Source Ecosystem

Apache Kafka, Connect, Streams, KSQL, etc...

Kai Waehner

Technology Evangelist

kontakt@kai-waehner.de

LinkedIn

@KaiWaehner

www.confluent.io

www.kai-waehner.de



DL4J  
DEEPLearning4J



Google Cloud Platform



kubernetes

# What is eXtreme Scale?

---

- High Volume of Events (millions, billions, trillions)
- Big Data Sets for Analytics (GB, TB, PB)
- Dynamic Scalability for Training (minutes, hours, days)
- Real Time Prediction Process for Deployment (ms)
- Hybrid Deployments (different frameworks and clouds)



EXTREME

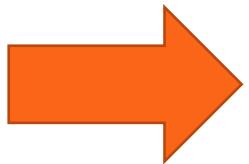
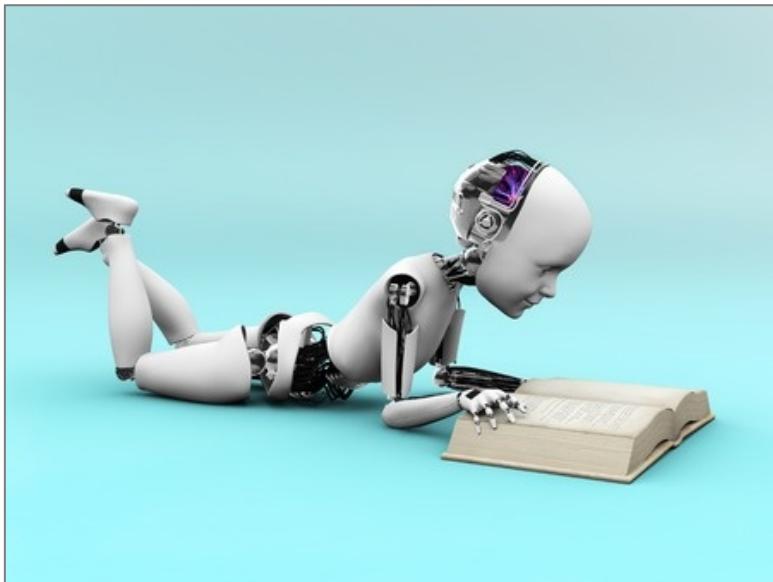
# Agenda

---

- 1) Added Business Value via Machine Learning
- 2) Apache Kafka Ecosystem as Infrastructure for Machine Learning
- 3) Data Ingestion and Preprocessing with Apache Kafka
- 4) Predictions in Real Time with Kafka Streams and KSQL
- 5) Automation and DevOps of a Machine Learning Infrastructure

- 1) Added Business Value via Machine Learning**
- 2) Apache Kafka Ecosystem as Infrastructure for Machine Learning
- 3) Data Ingestion and Preprocessing with Apache Kafka
- 4) Predictions in Real Time with Kafka Streams and KSQL
- 5) Automation and DevOps of a Machine Learning Infrastructure

... allows computers **to find hidden insights without being explicitly programmed where to look.**



## Machine Learning

- Decision Trees
- Naïve Bayes
- Clustering
- Neural Networks
- etc.

## Deep Learning

- CNN
- RNN
- Autoencoder
- etc.

# Real World Examples of Machine Learning



Spam Detection



Search Results +  
Product Recommendation



Picture Detection  
(Friends, Locations, Products)

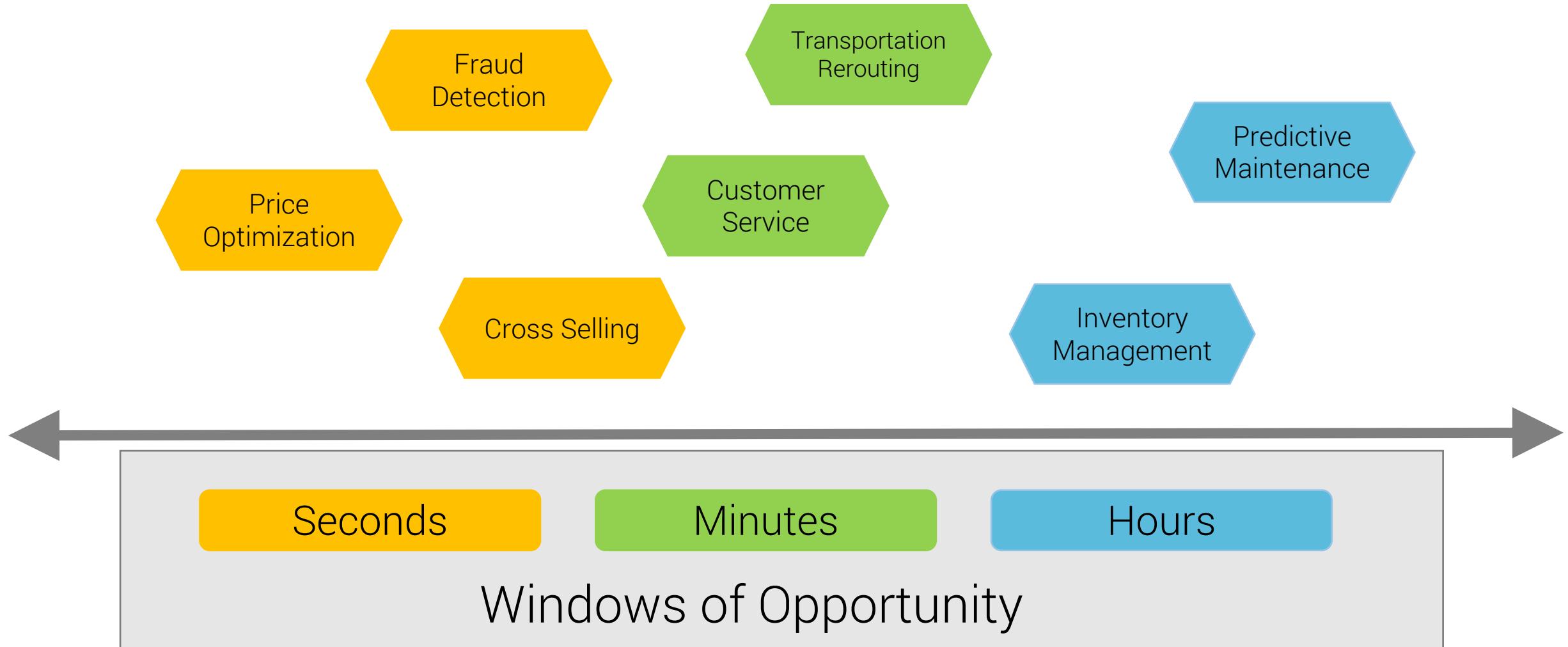


The Next Disruption:  
Google Beats Go Champion



Your Company

# Leverage Machine Learning to Analyze and Act on Critical Business Moments

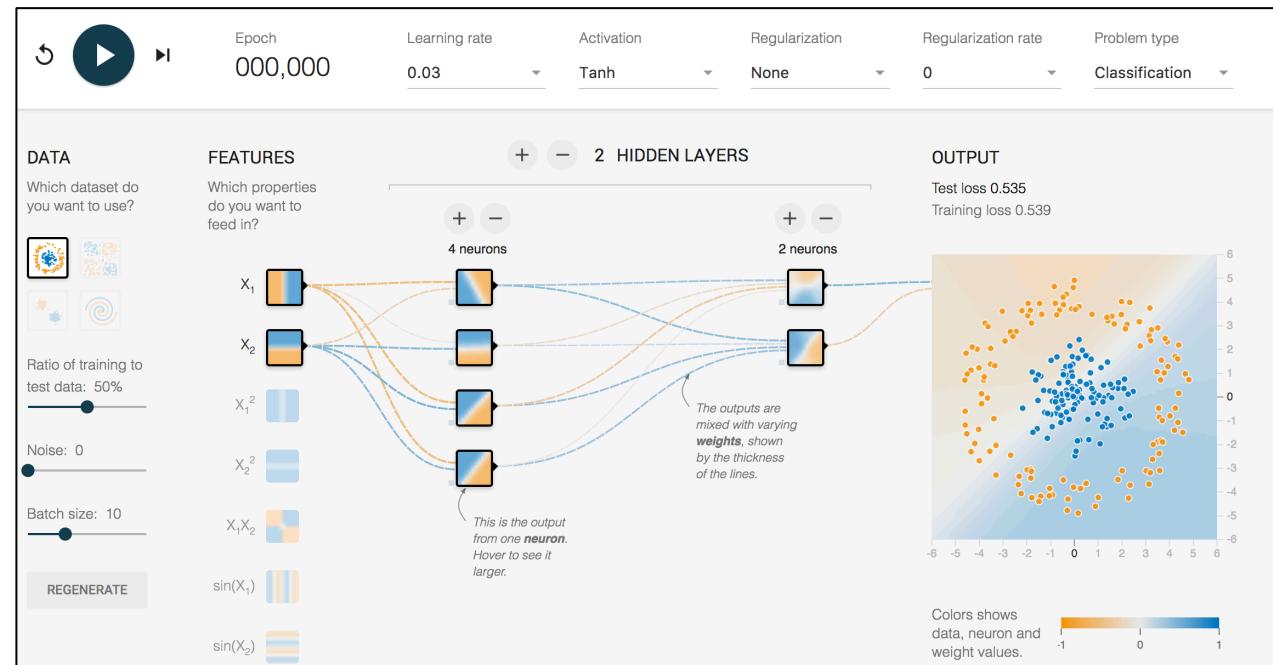


# Live Demo – Building an Analytic Model



## Neural Networks in Action

<http://playground.tensorflow.org/>



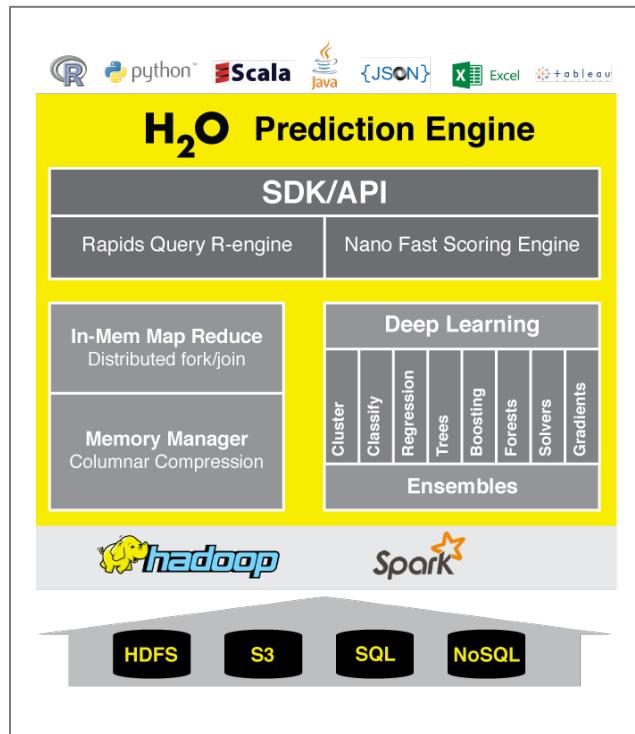
# Languages, Frameworks and Tools for Machine Learning

---



There is no Allrounder → ML-independent infrastructure needed!

# Machine Learning with H2O.ai



R / Python /  
Scala / Flow UI



The screenshot shows the "Build a Model" interface for "Deep Learning". It includes fields for "model\_id", "training\_frame", "validation\_frame", "nfolds", "response\_column", and "ignored\_columns". A preview table shows data for "Cancelled", "CancellationCode", "Diverted", "CarrierDelay", "WeatherDelay", and "NASDelay". Parameters for "activation", "hidden", and "epochs" are also visible.

Java Code

```
@ModelPojo(name="deeplearning_fe7c1f02_08ec_4070_b784_c2531147e451", algorithm="deeplearning")
public class deeplearning_fe7c1f02_08ec_4070_b784_c2531147e451 extends GenModel {
    public hex.ModelCategory getModelCategory() { return hex.ModelCategory.Binomial; }
    public boolean isSupervised() { return true; }
    public int nfeatures() { return 12; }
    public int nclasses() { return 2; }
    // Thread-local storage for input neuron activation values.
    final double[] NUMS = new double[10];
    static class NORMMUL implements java.io.Serializable {
        public static final double[] VALUES = new double[10];
        static {
            NORMMUL_0.fill(VALUES);
        }
        static final class NORMMUL_0 implements java.io.Serializable {
            static final void fill(double[] sa) {
                sa[0] = 0.1573591362493411;
                sa[1] = 0.5316756588306932;
                sa[2] = 0.10894640014128883;
                sa[3] = 0.5257616635956896;
                sa[4] = 0.00209932098808304;
            }
        }
    }
}
```

H2O Engine

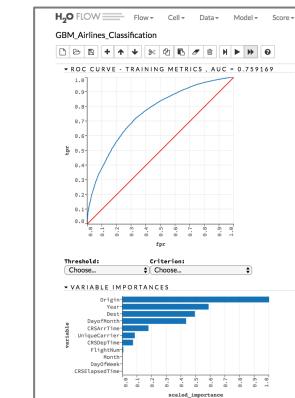
confluent

# Live Demo

Use Case:  
Airline Flight Delay Prediction

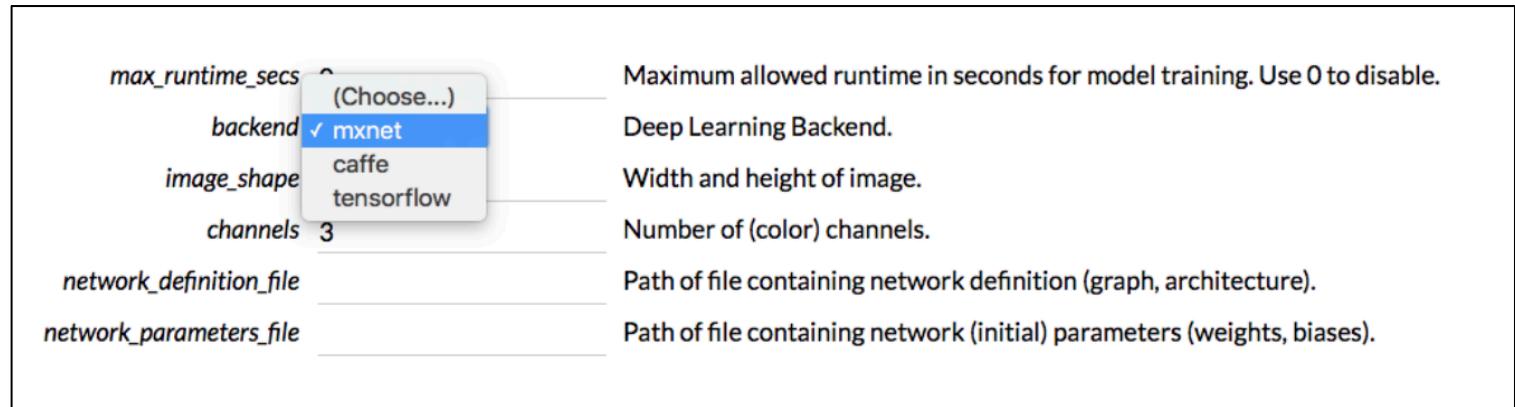
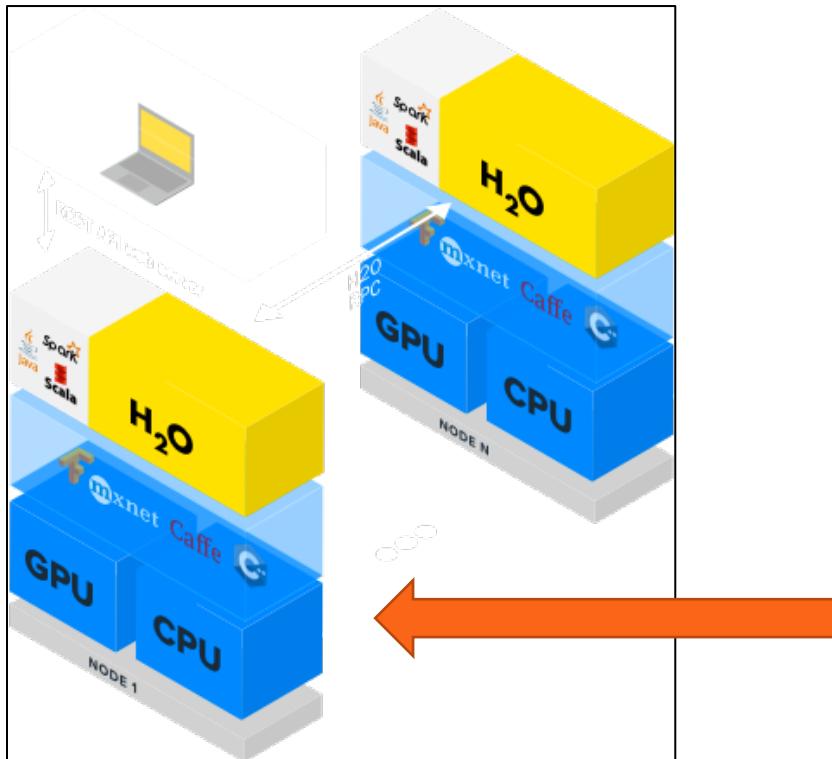
Machine Learning Algorithm:  
Deep Learning  
using Neural Networks

Technology:  
H2O.ai, TensorFlow



A screenshot of the H2O Flow interface showing the "Assistance" menu. It lists various H2O routines such as importFiles, getFrames, splitFrame, mergeFrames, getModels, getGrids, getPredictions, getJobs, buildModel, importModel, and predict, along with their descriptions.

# H2O Deep Water (TensorFlow, MXNet, ...)

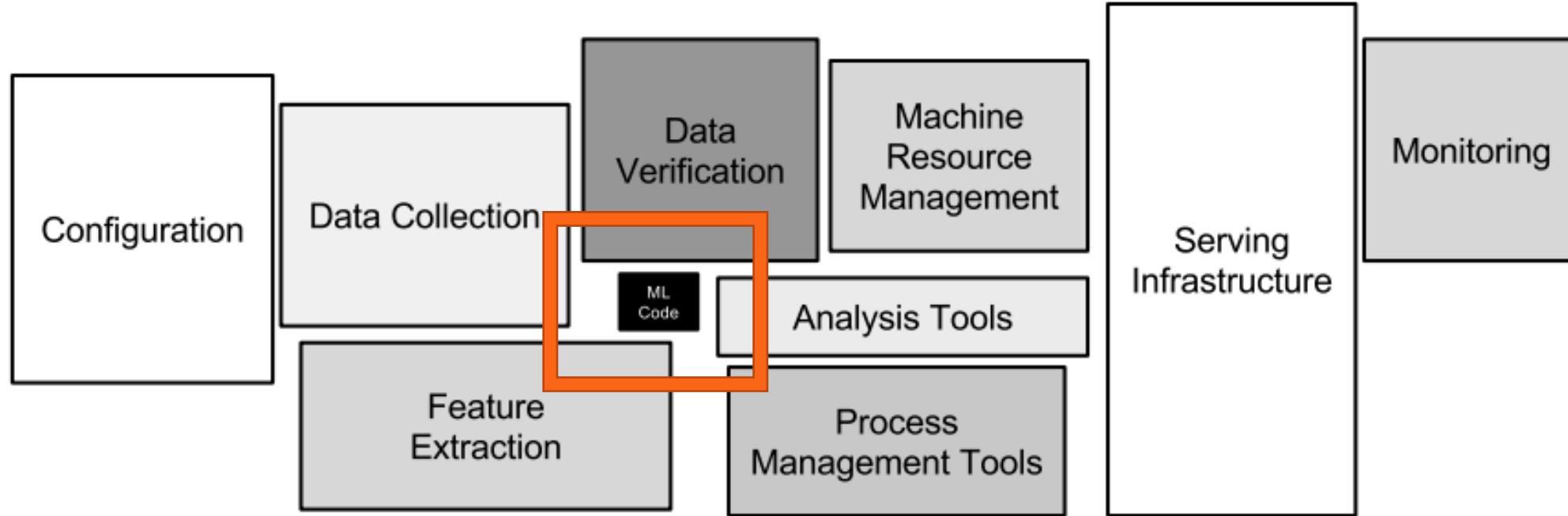


Deep Water  
(H2O + TensorFlow)

Pre-Defined Networks  
+  
User-Defined Networks

<https://h2o-release.s3.amazonaws.com/h2o/rel-vapnik/1/docs-website/h2o-docs/booklets/DeepWaterBooklet.pdf>

# Hidden Technical Debt in Machine Learning Systems



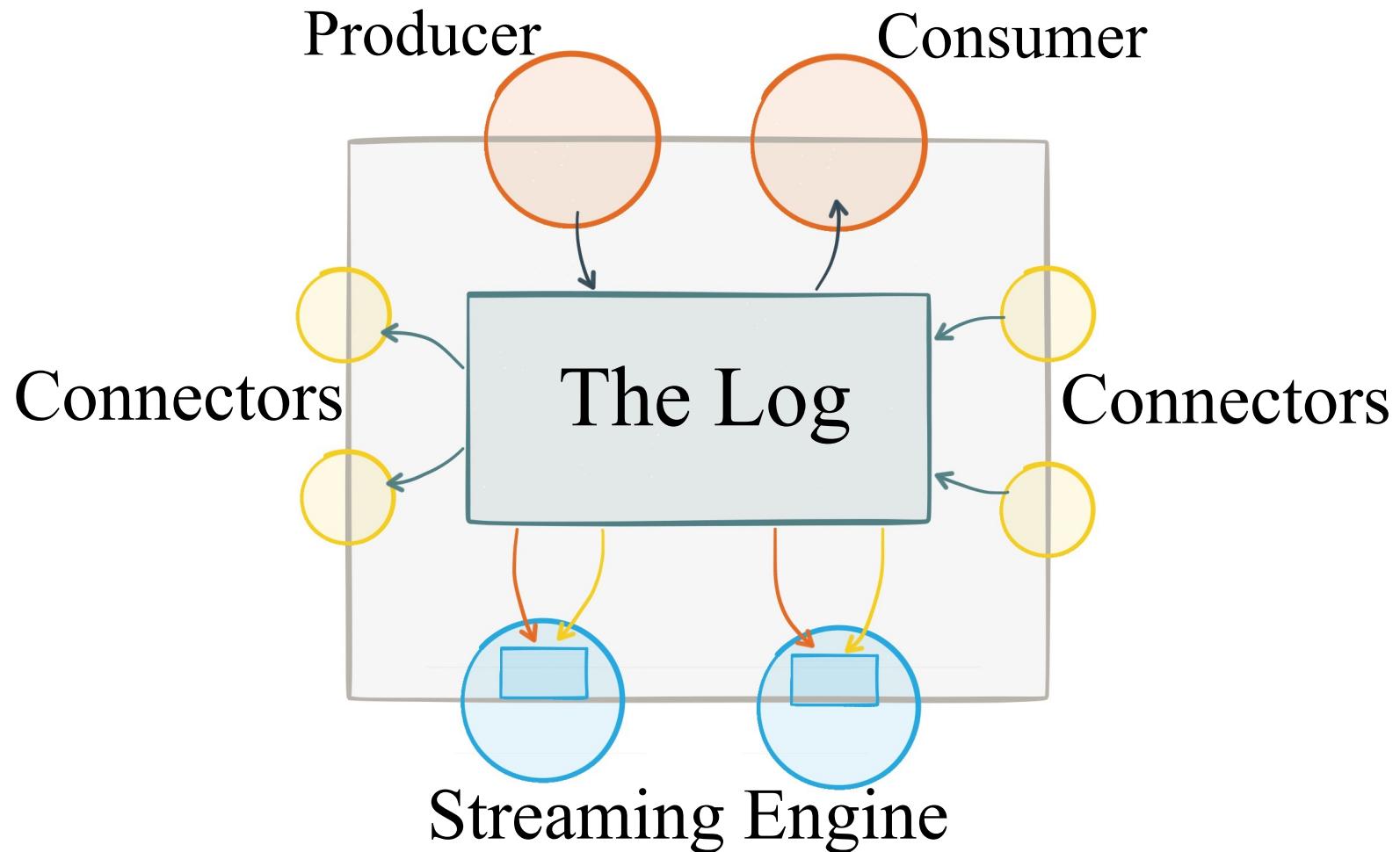
<https://papers.nips.cc/paper/5656-hidden-technical-debt-in-machine-learning-systems.pdf>

# Agenda

---

- 1) Added Business Value via Machine Learning
- 2) Apache Kafka Ecosystem as Infrastructure for Machine Learning**
- 3) Data Ingestion and Preprocessing with Apache Kafka
- 4) Predictions in Real Time with Kafka Streams and KSQL
- 5) Automation and DevOps of a Machine Learning Infrastructure

# Apache Kafka – The Rise of a Streaming Platform



# Apache Kafka at Large Scale → No need to do a POC



## Operation Challenges

- The scale of Kafka deployment @LinkedIn
  - 2,100+ brokers
  - ~ 60,000 topics
  - 1.2 million partitions
  - > 4.5 trillion messages / day
- Huge operation overhead
  - Hardware failures are norm
  - Workload skews



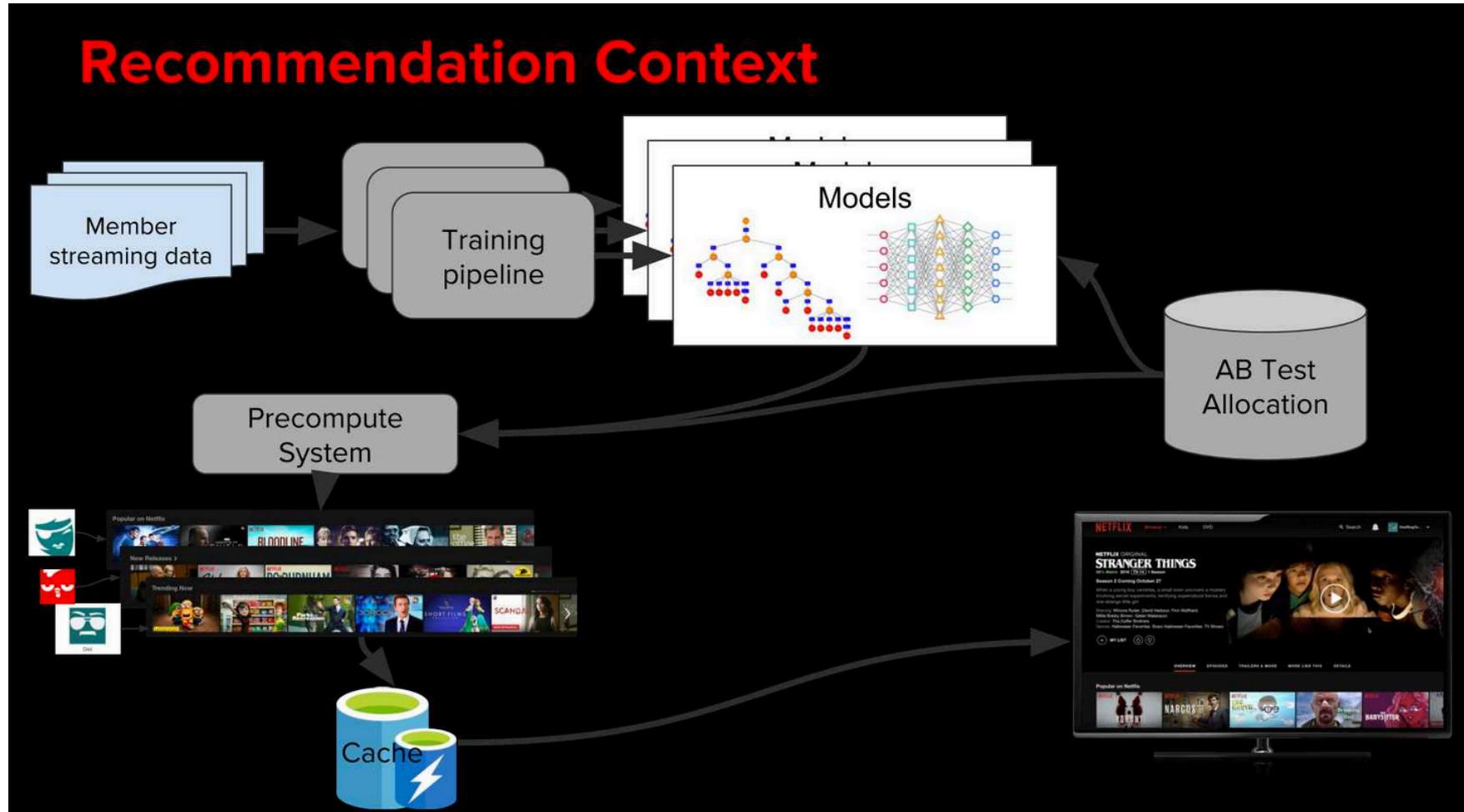
## Kafka @ Netflix Scale

- 4,000+ brokers and ~50 clusters in 3 AWS regions
- > 1 Trillion messages per day
- At peak (New Years Day 2018)
  - 2.2 trillion messages (1.3 trillion unique)
  - 6 Petabytes

**Strata**  
DATA CONFERENCE

**QCon**

<https://conferences.oreilly.com/strata/strata-ca/public/schedule/detail/63921>  
<https://qconlondon.com/london2018/presentation/cloud-native-and-scalable-kafka-architecture>

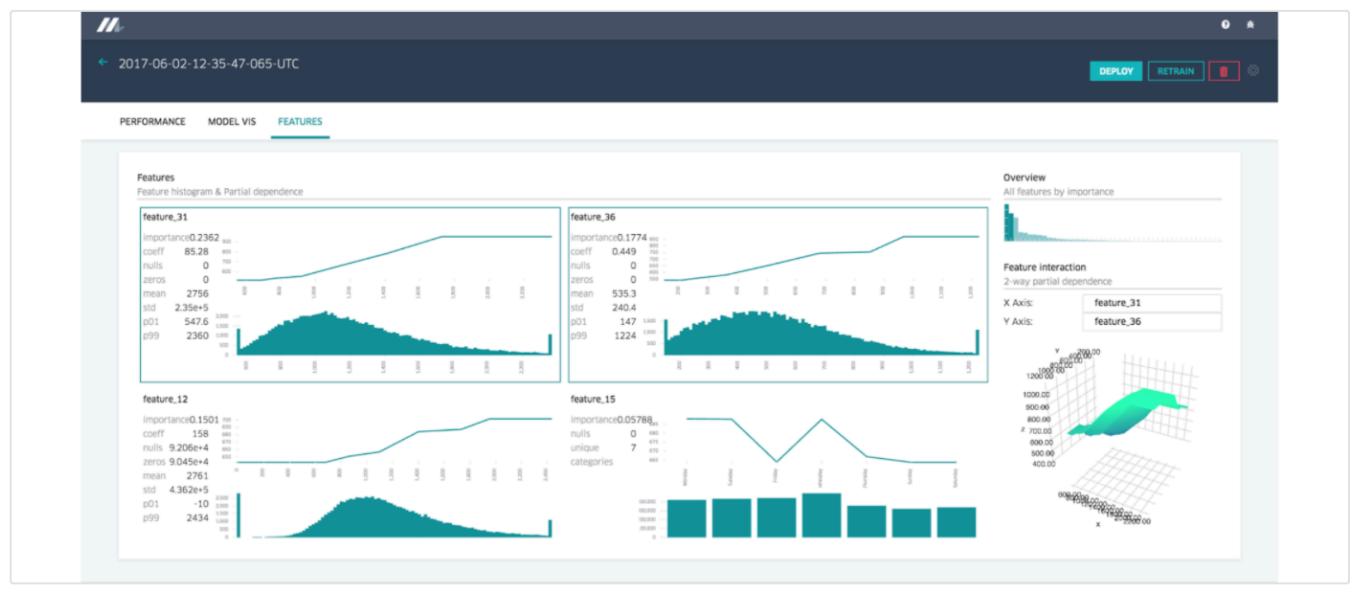


<https://www.infoq.com/presentations/netflix-ml-meson>

## Meet Michelangelo: Uber's Machine Learning Platform

By Jeremy Hermann & Mike Del Balso

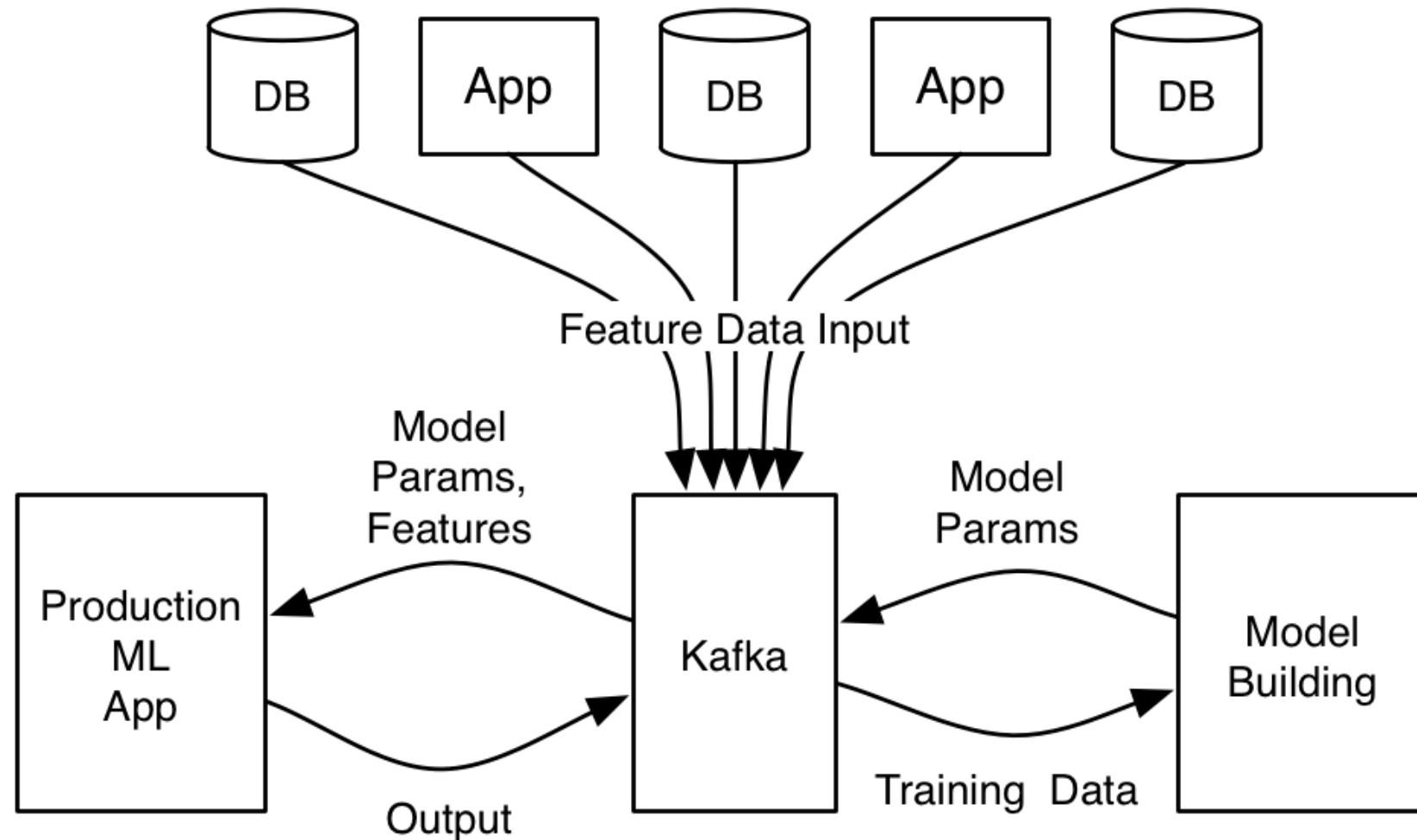
September 5, 2017



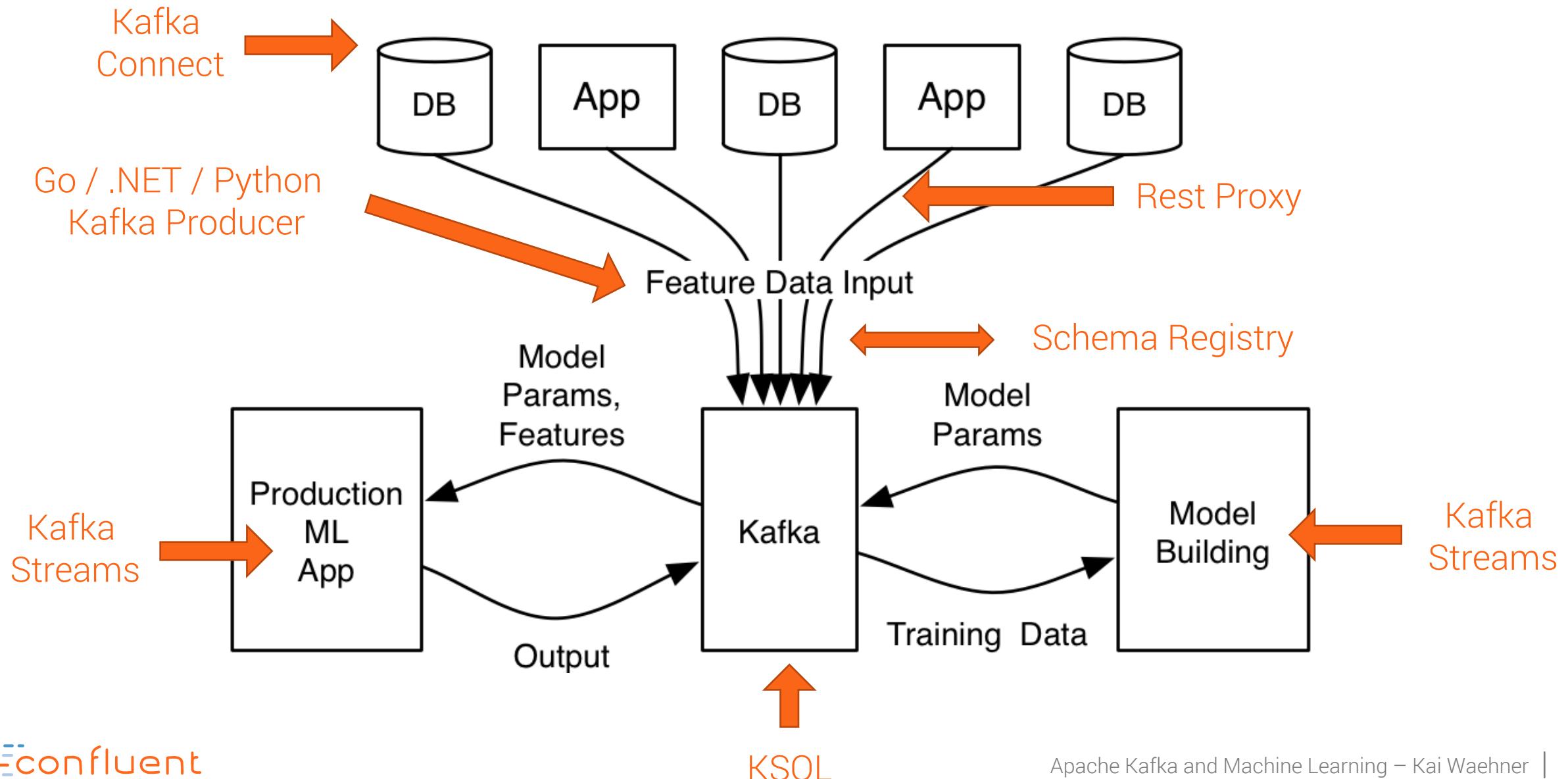
<https://eng.uber.com/michelangelo>

- Cover the **end-to-end ML workflow**: manage data, train, evaluate, and deploy models, make predictions, and monitor predictions
- **Supports various AI technologies**: Traditional ML models, time series forecasting, and deep learning

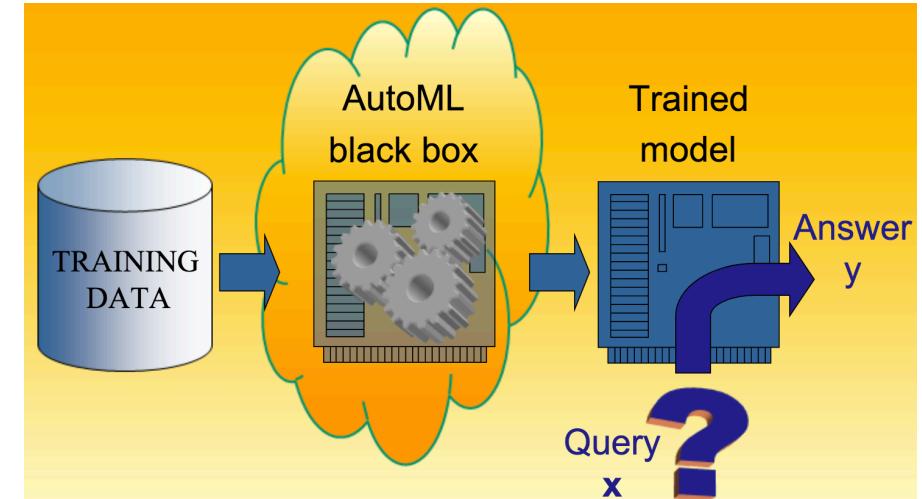
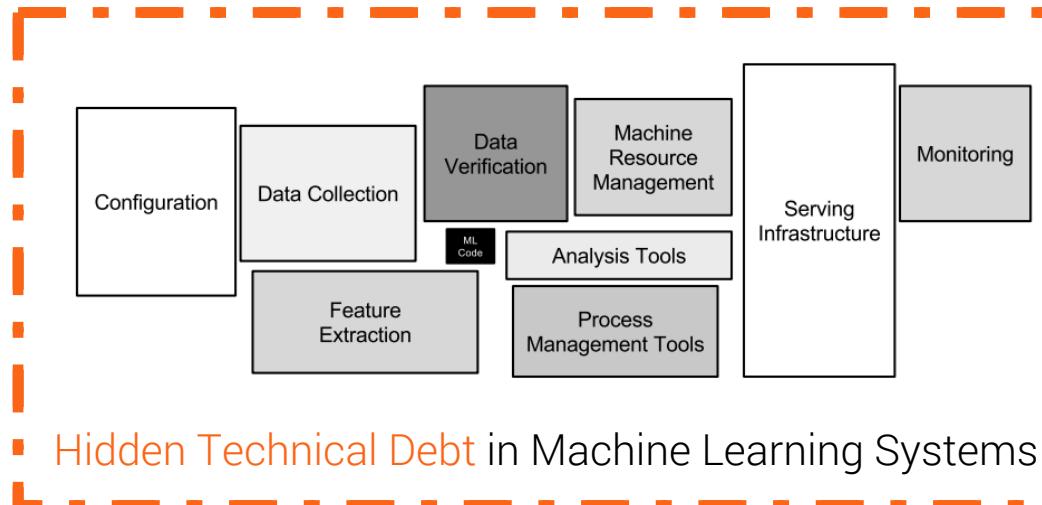
# Apache Kafka's Open Source Ecosystem as Infrastructure for Machine Learning



# Apache Kafka's Open Source Ecosystem as Infrastructure for Machine Learning



# AutoML → No Data Scientist available for the ML Tasks?



<http://slideplayer.com/slide/10575150/>

“One-Click Data-In  
Model-Out simplicity”

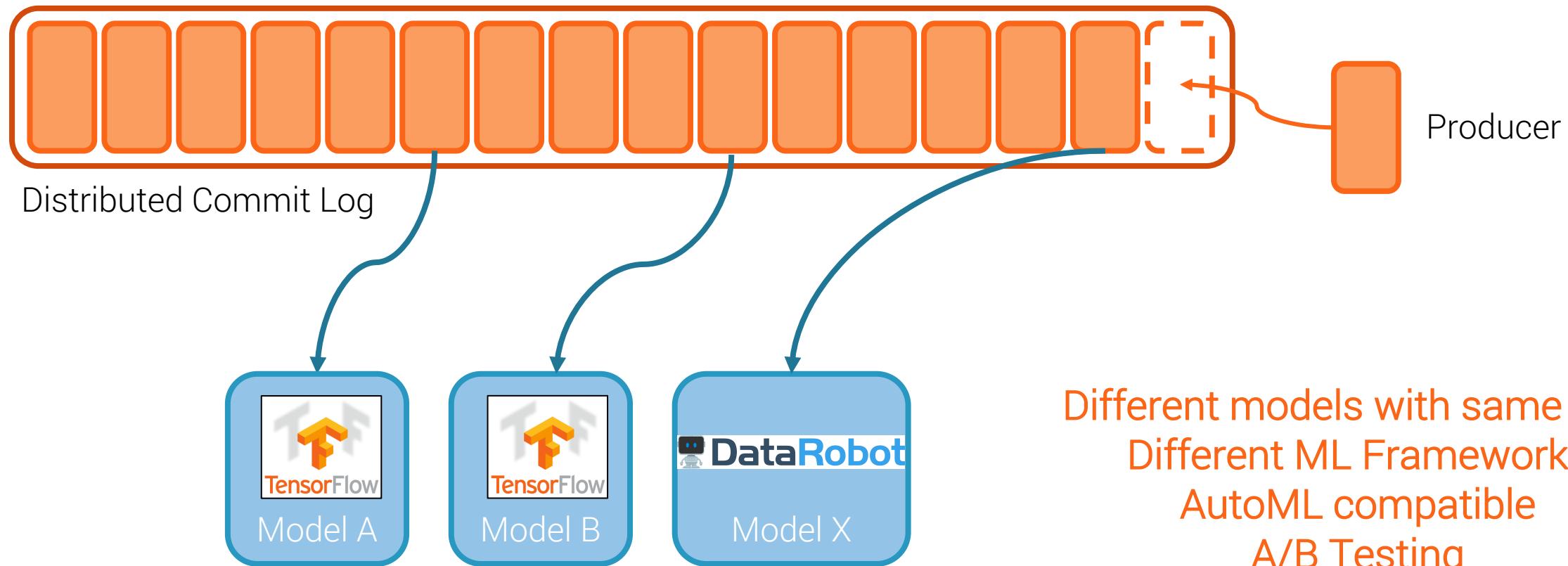


Driverless AI: Your Expert System for AI



# Replay-ability – A log never forgets!

Time

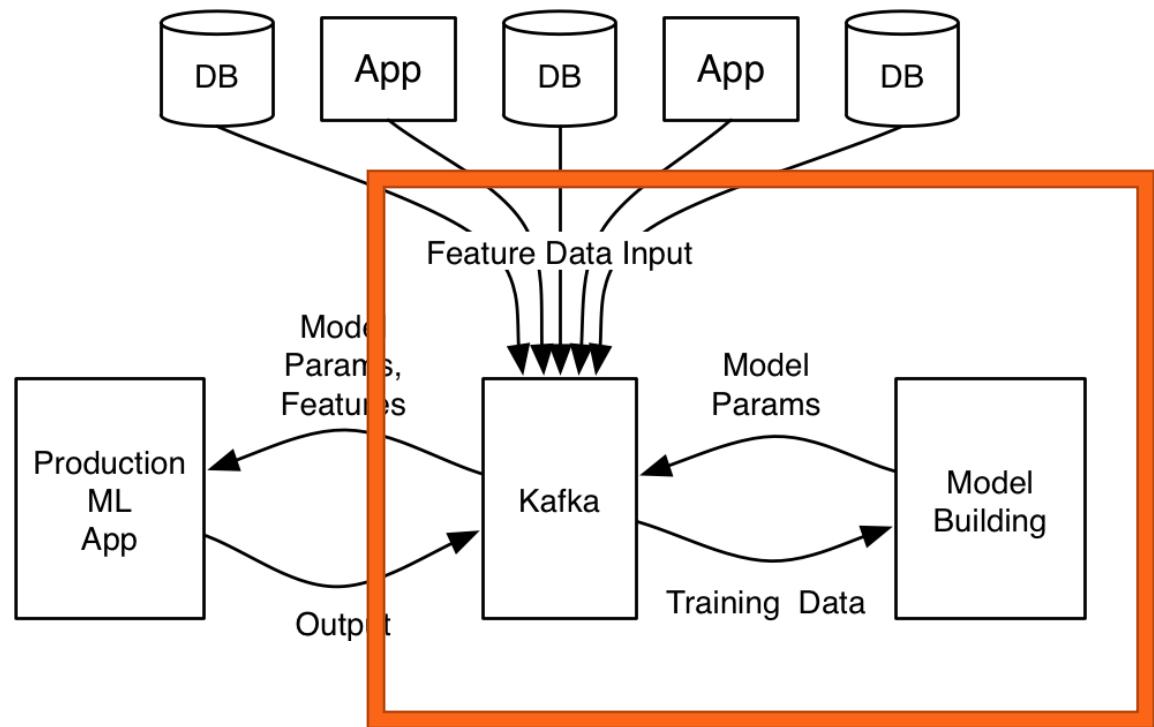


# Agenda

---

- 1) Added Business Value via Machine Learning
- 2) Apache Kafka Ecosystem as Infrastructure for Machine Learning
- 3) Data Ingestion and Preprocessing with Apache Kafka**
- 4) Predictions in Real Time with Kafka Streams and KSQL
- 5) Automation and DevOps of a Machine Learning Infrastructure

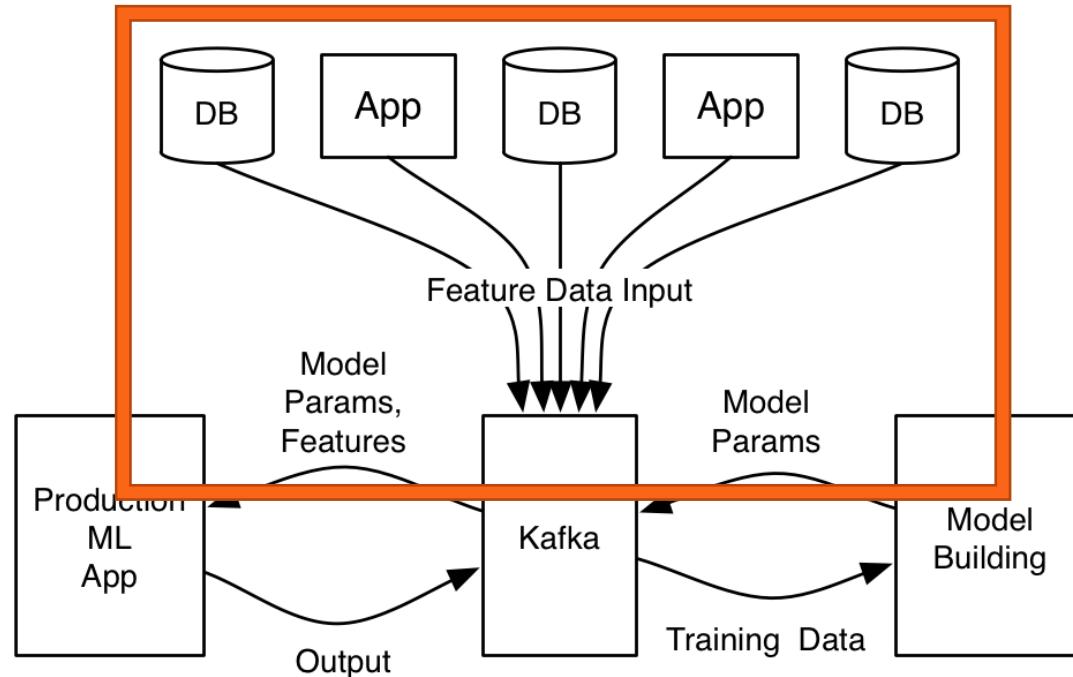
# Apache Kafka for ML Pipelines at Large Scale



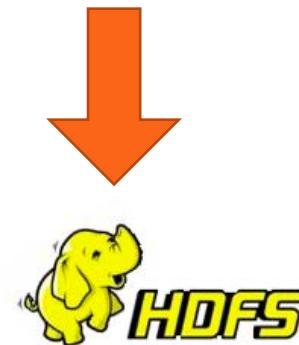
Google Cloud Platform

Extreme Scale  
Dynamic Instances  
Special Hardware

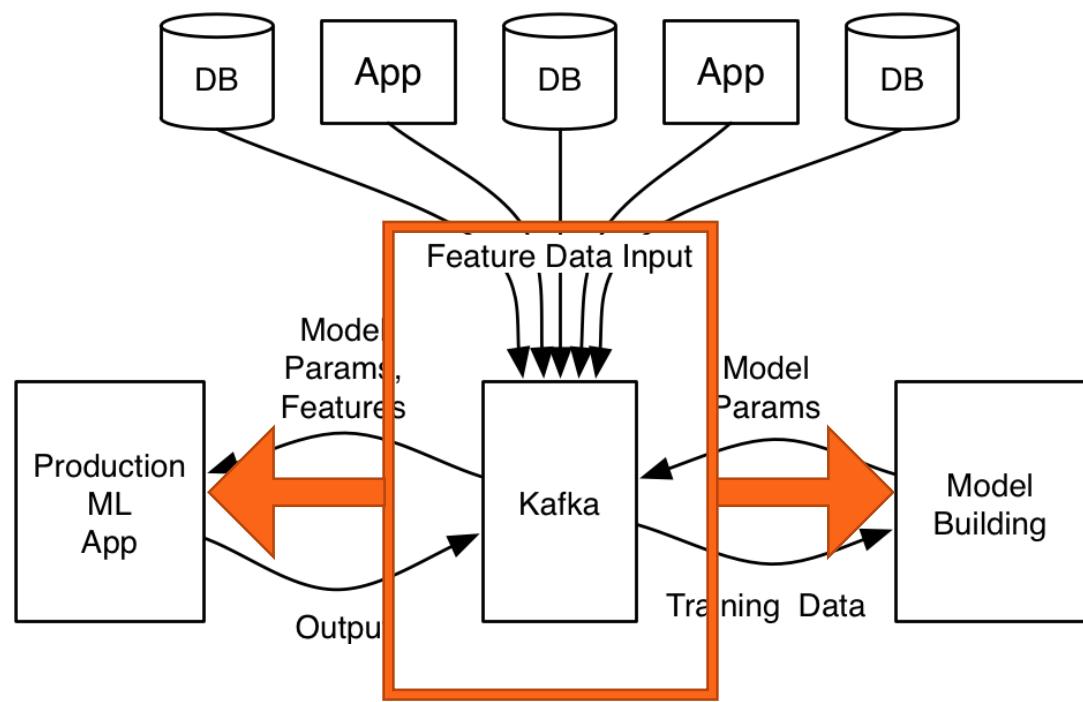
# Kafka Connect for Data Ingestion



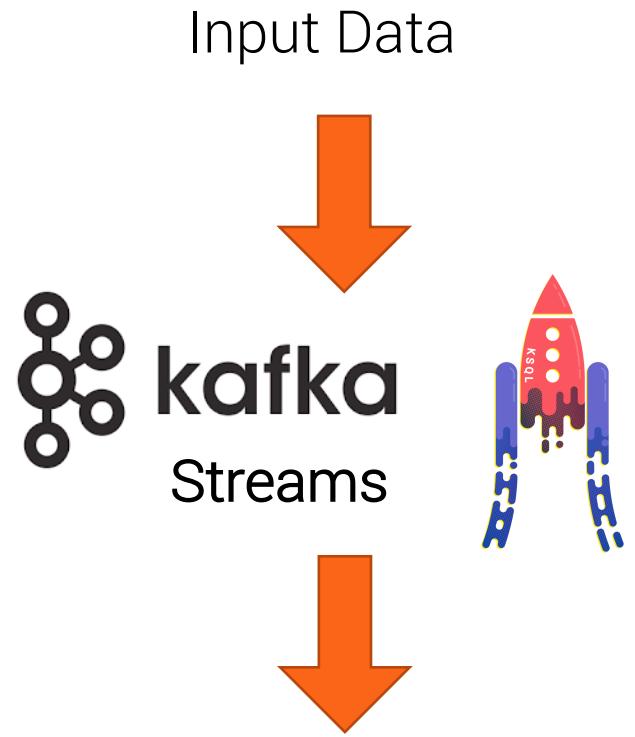
“Kafka benefits under the hood”  
Out-of-the-Box Connectivity  
Data Format Conversion  
Simple Message Transformation



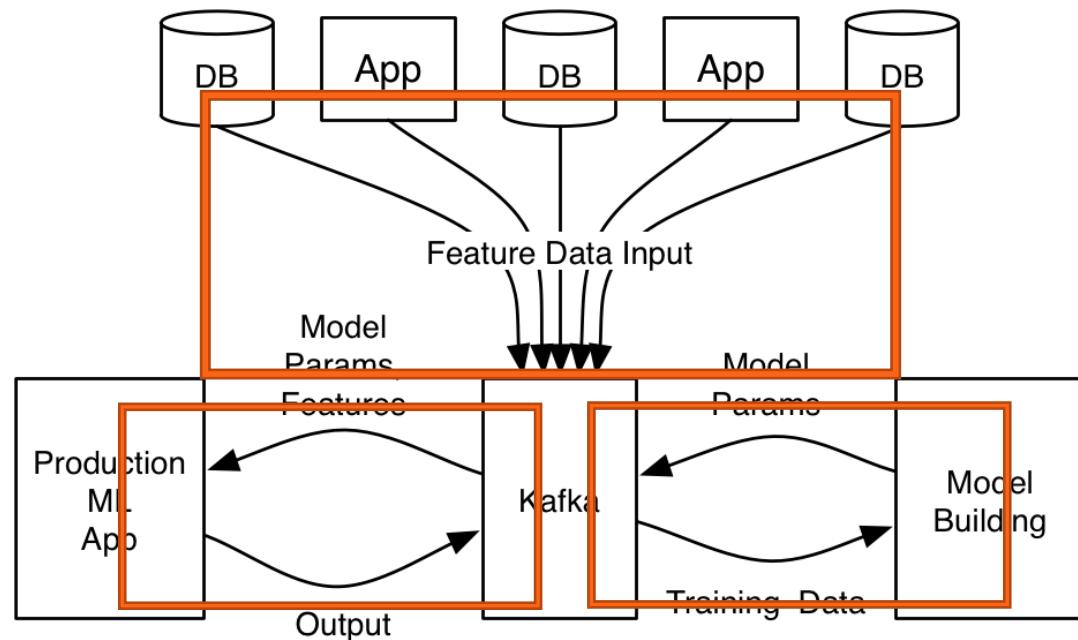
# Kafka Streams / KSQL for Data Preprocessing



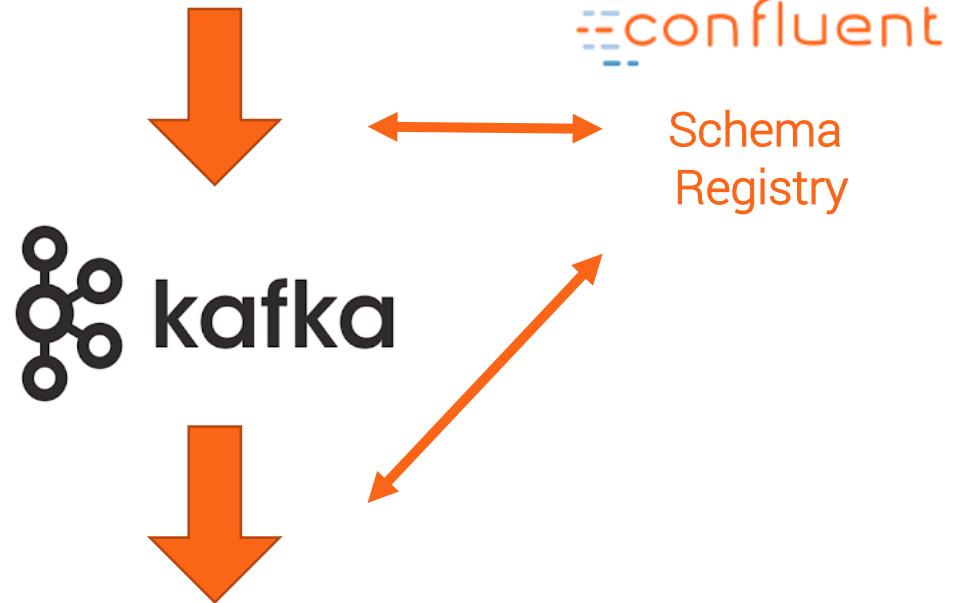
“Kafka benefits under the hood”  
Streaming ETL  
Same Pipeline for Training and Serving



# Confluent Schema Registry for Message Validation



Input Data



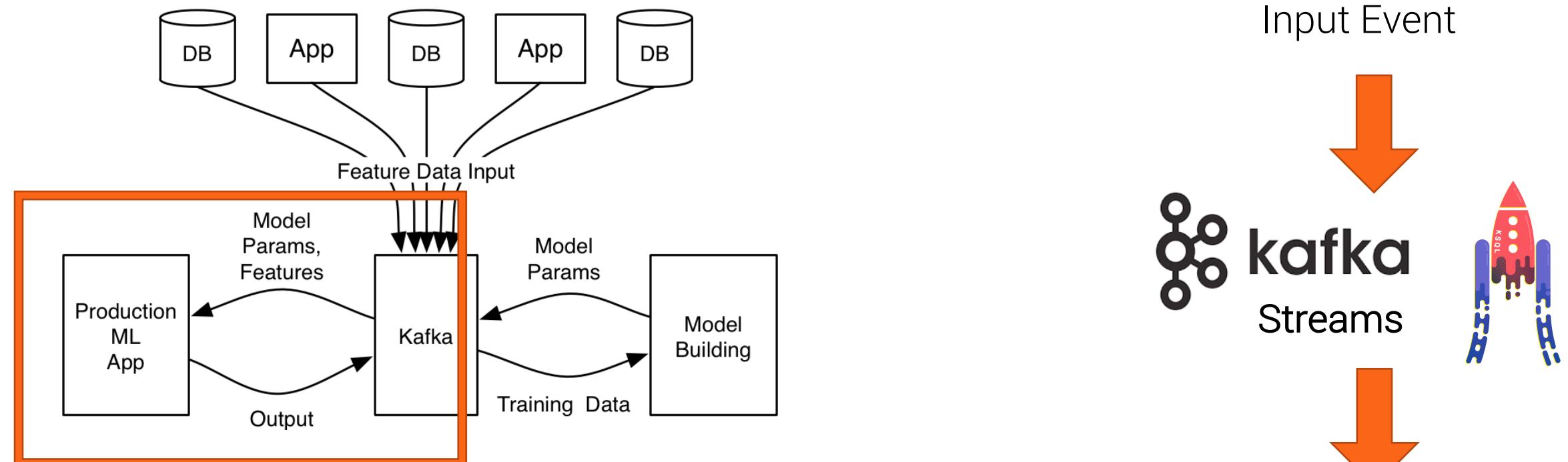
"Kafka benefits under the hood"  
Schema Definition + Evolution  
Forward and Backward Compatibility

# Agenda

---

- 1) Added Business Value via Machine Learning
- 2) Apache Kafka Ecosystem as Infrastructure for Machine Learning
- 3) Data Ingestion and Preprocessing with Apache Kafka
- 4) Predictions in Real Time with Kafka Streams and KSQL**
- 5) Automation and DevOps of a Machine Learning Infrastructure

# Model Serving / Inference / Deployment / Scoring

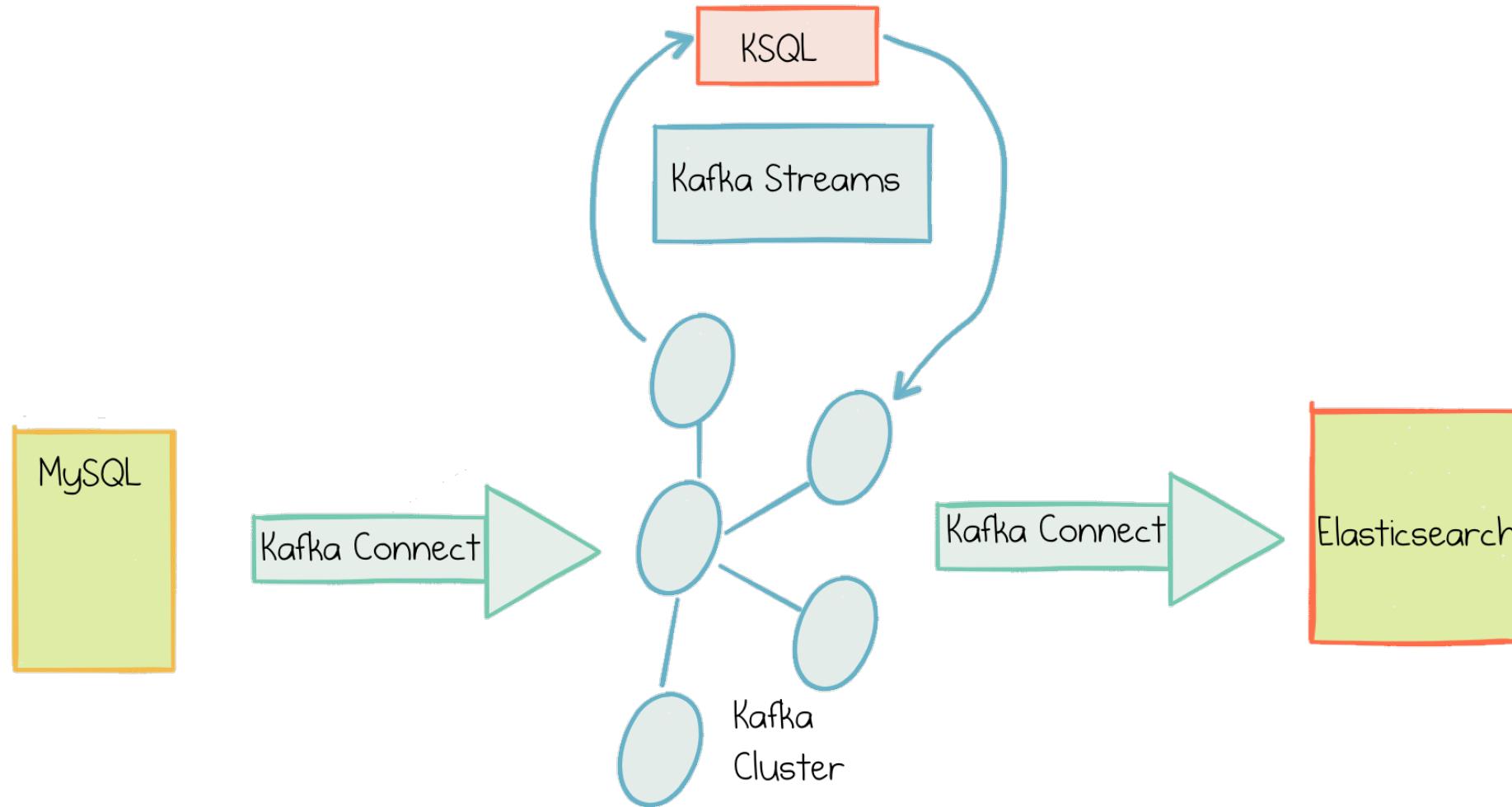


Kafka Streams  
KSQL

“Kafka benefits under the hood”  
Continuous Stream Processing  
Predictions in real time

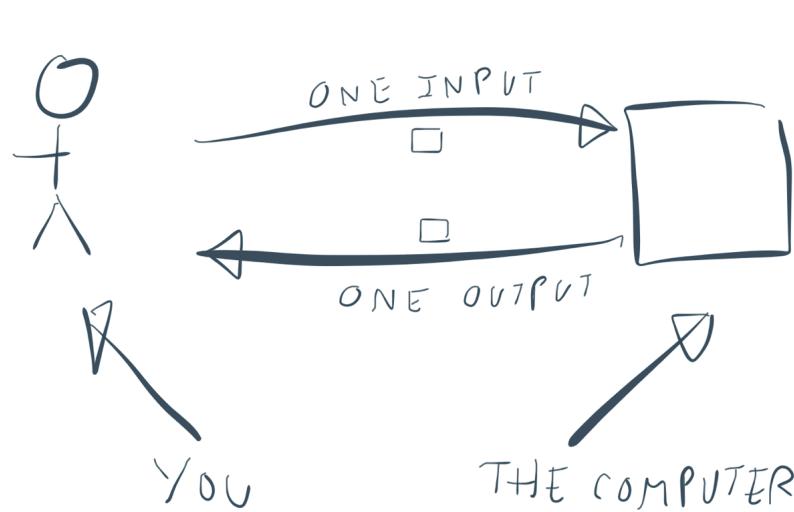
Scalability like any other Kafka Microservice

# Kafka Streams (shipped with Apache Kafka) / KSQL (Confluent Open Source)



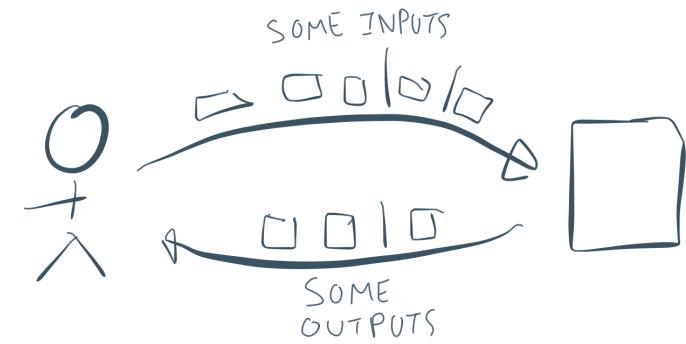
# Stream Processing – The only option for extreme scale

REQUEST / RESPONSE



Data at Rest

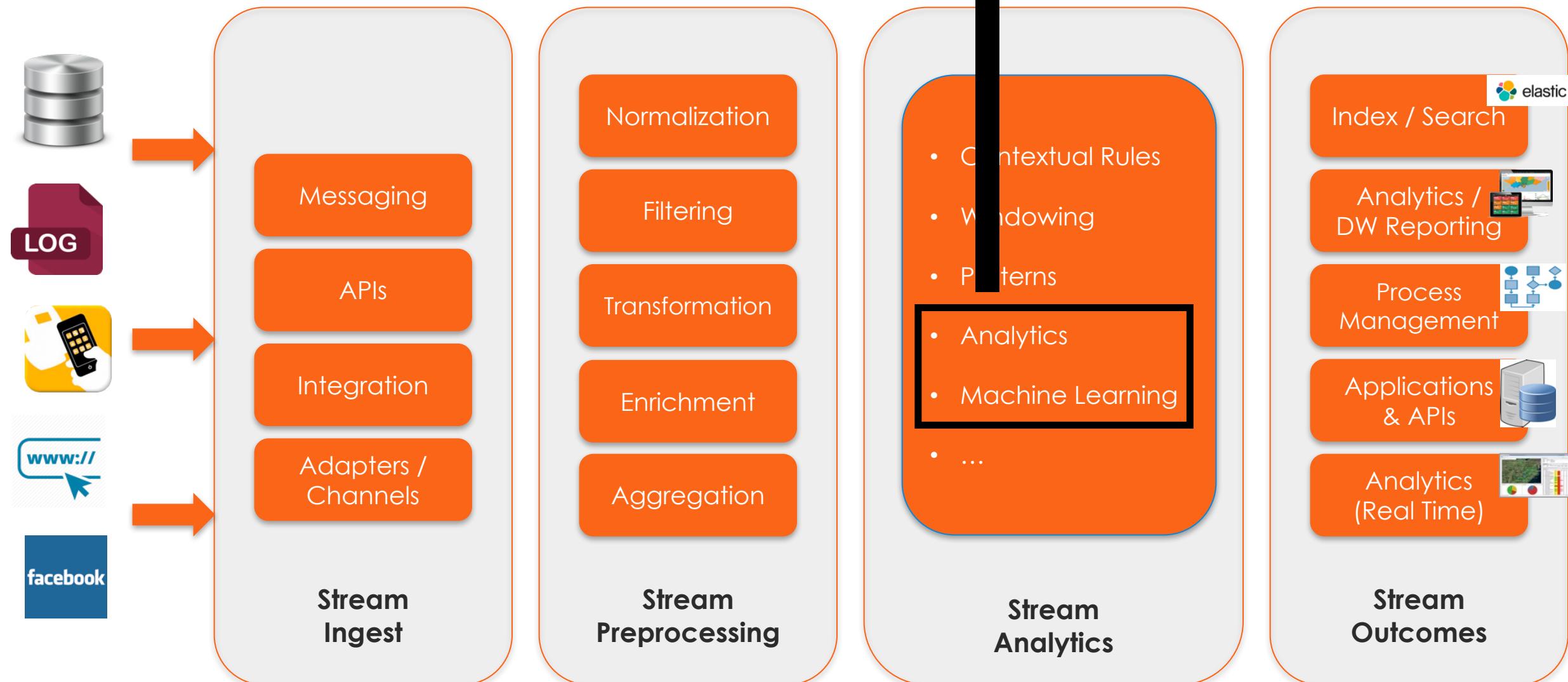
STREAM PROCESSING



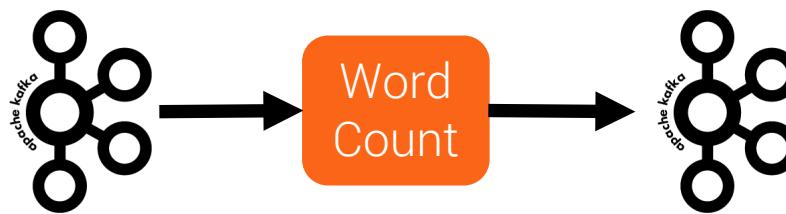
Data in Motion

# Stream Processing Pipeline

Applying an Analytic Model  
is just a piece of the puzzle!



# A complete streaming microservices, ready for production at large-scale



```
1 public static void main(final String[] args) throws Exception {  
2     Properties config = new Properties();  
3     config.put(StreamsConfig.APPLICATION_ID_CONFIG, "wordcount-example");  
4     config.put(StreamsConfig.BOOTSTRAP_SERVERS_CONFIG, "kafka-broker1:9092");  
5     config.put(StreamsConfig.KEY_SERDE_CLASS_CONFIG, Serdes.String().getClass().getName());  
6     config.put(StreamsConfig.VALUE_SERDE_CLASS_CONFIG, Serdes.String().getClass().getName());  
7  
8     KStreamBuilder builder = new KStreamBuilder();  
9     KStream<String, String> textLines = builder.stream("TextLinesTopic");  
10    KStream<String, Long> wordCounts = textLines  
11        .flatMapValues(value -> Arrays.asList(value.toLowerCase().split("\\W+")))  
12        .groupBy((key, word) -> word)  
13        .count("Counts")  
14        .toStream();  
15    wordCounts.to(Serdes.String(), Serdes.Long(), "WordsWithCountsTopic");  
16  
17    KafkaStreams streams = new KafkaStreams(builder, config);  
18    streams.start();  
19 }
```

App configuration

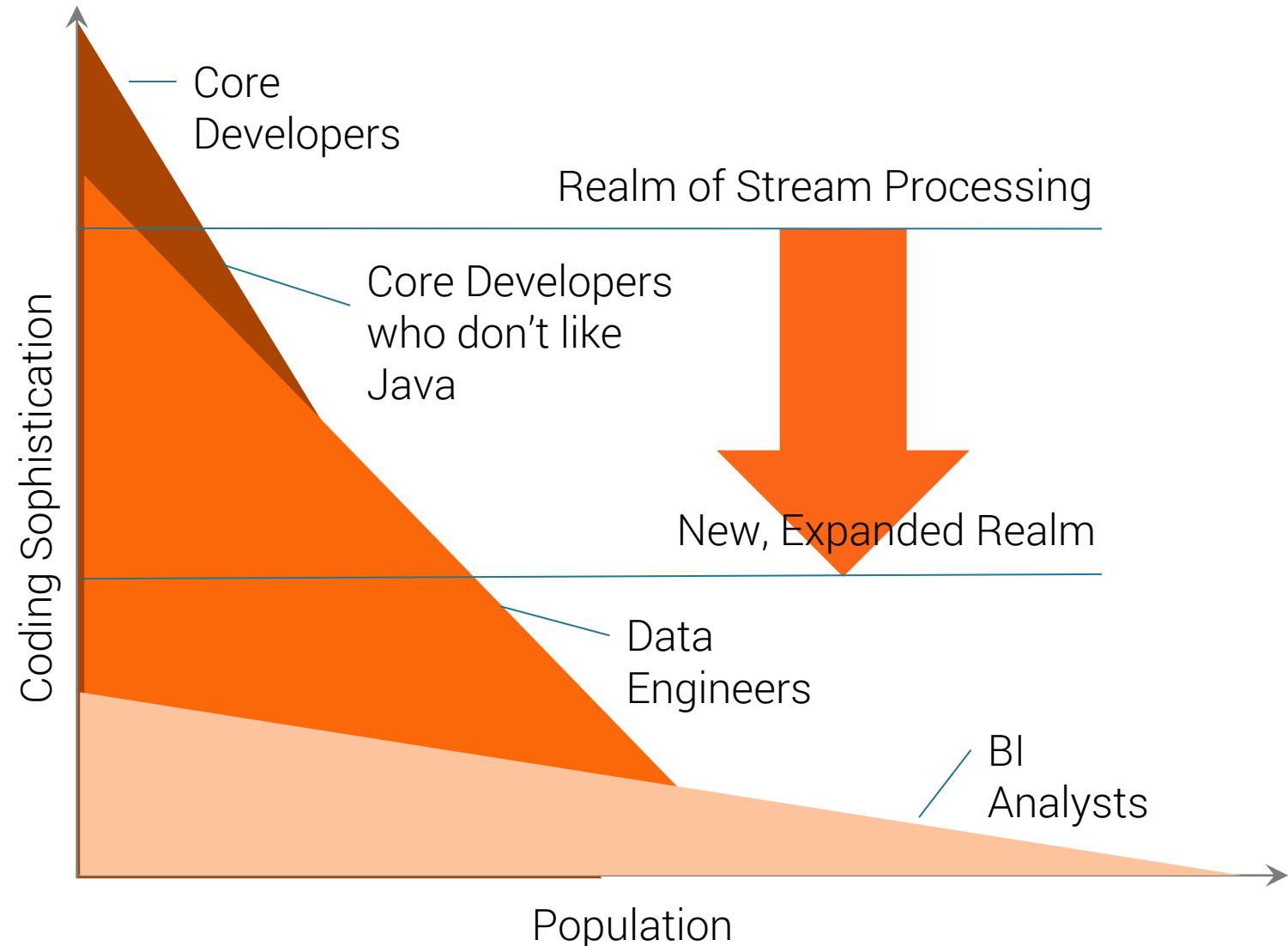
Define processing  
(here: WordCount)

Start processing

# Why KSQL?

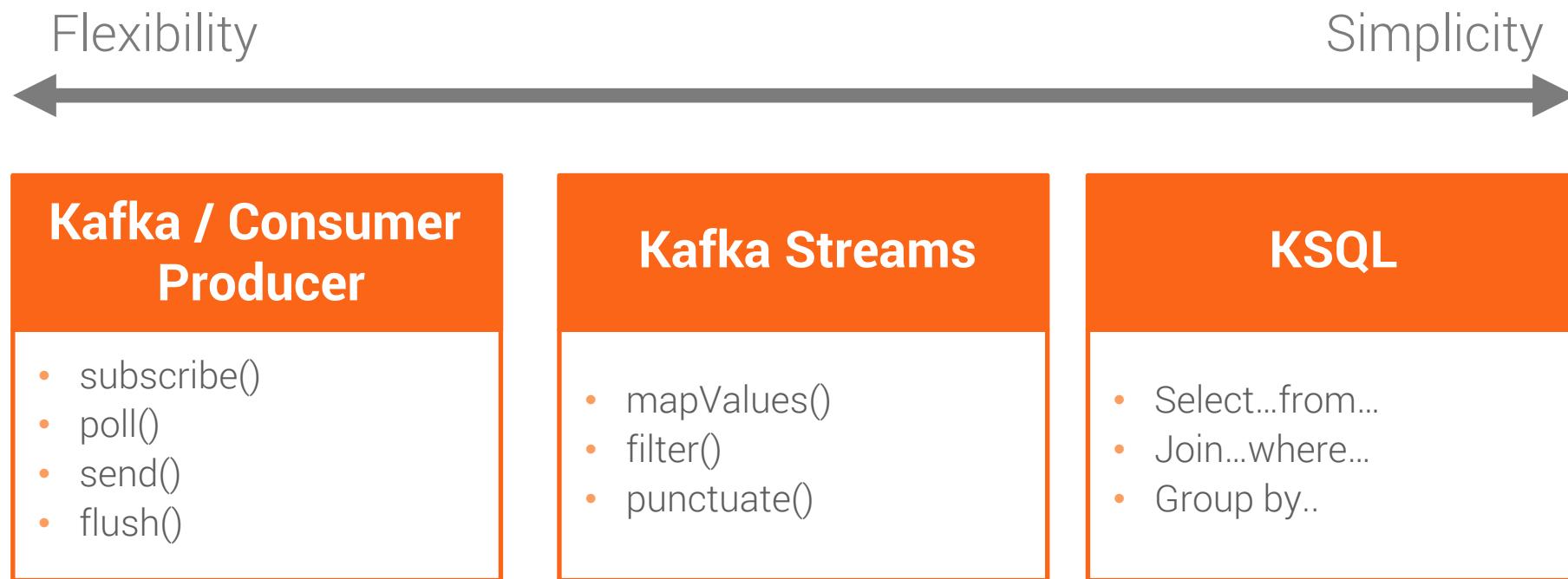
Kafka Streams

KSQL

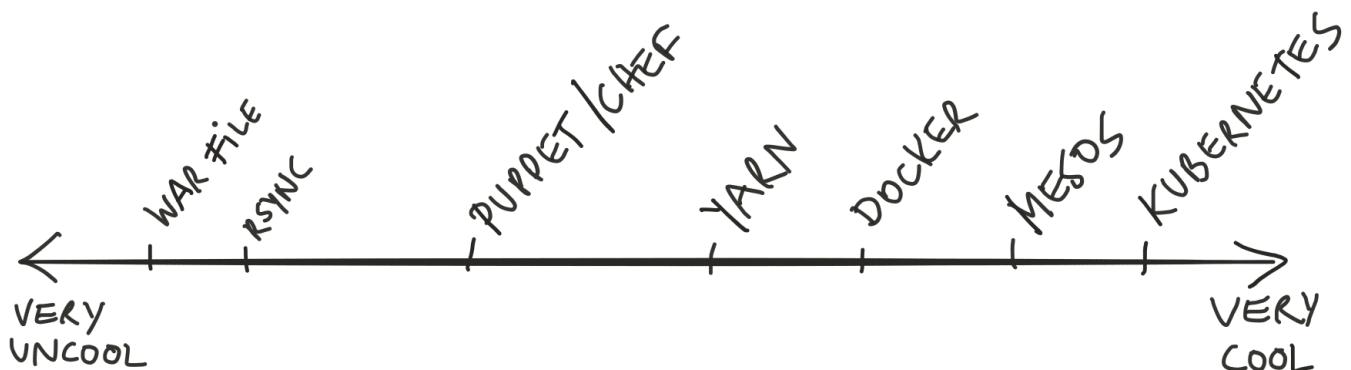
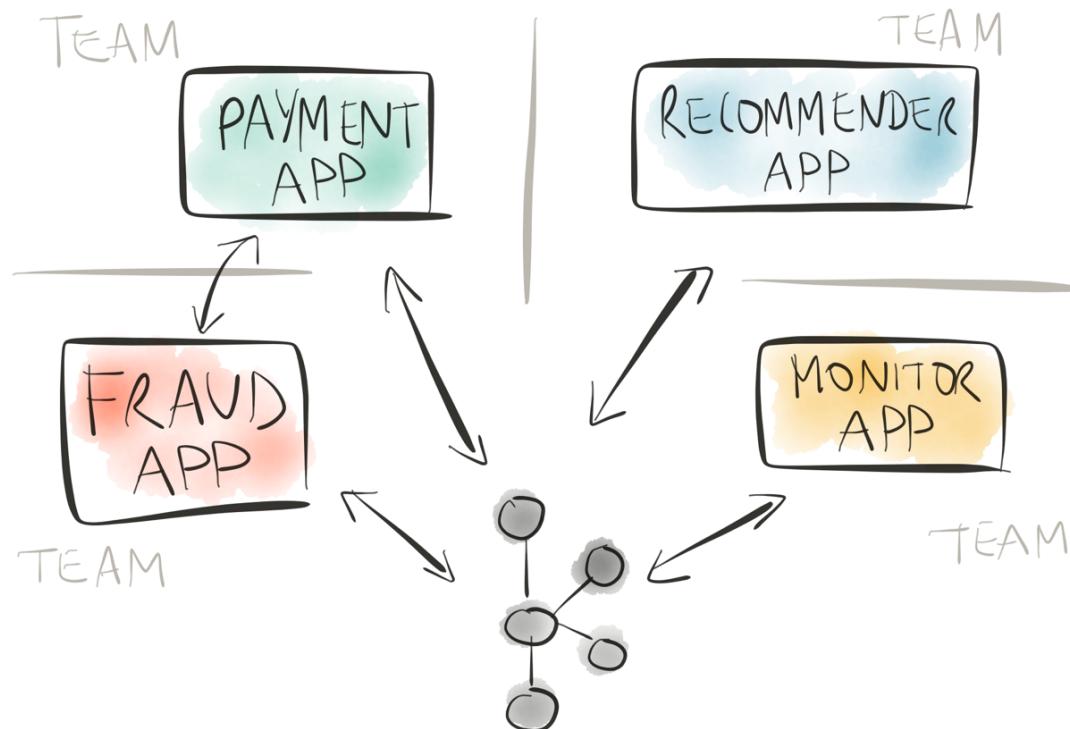


# Trade-Offs

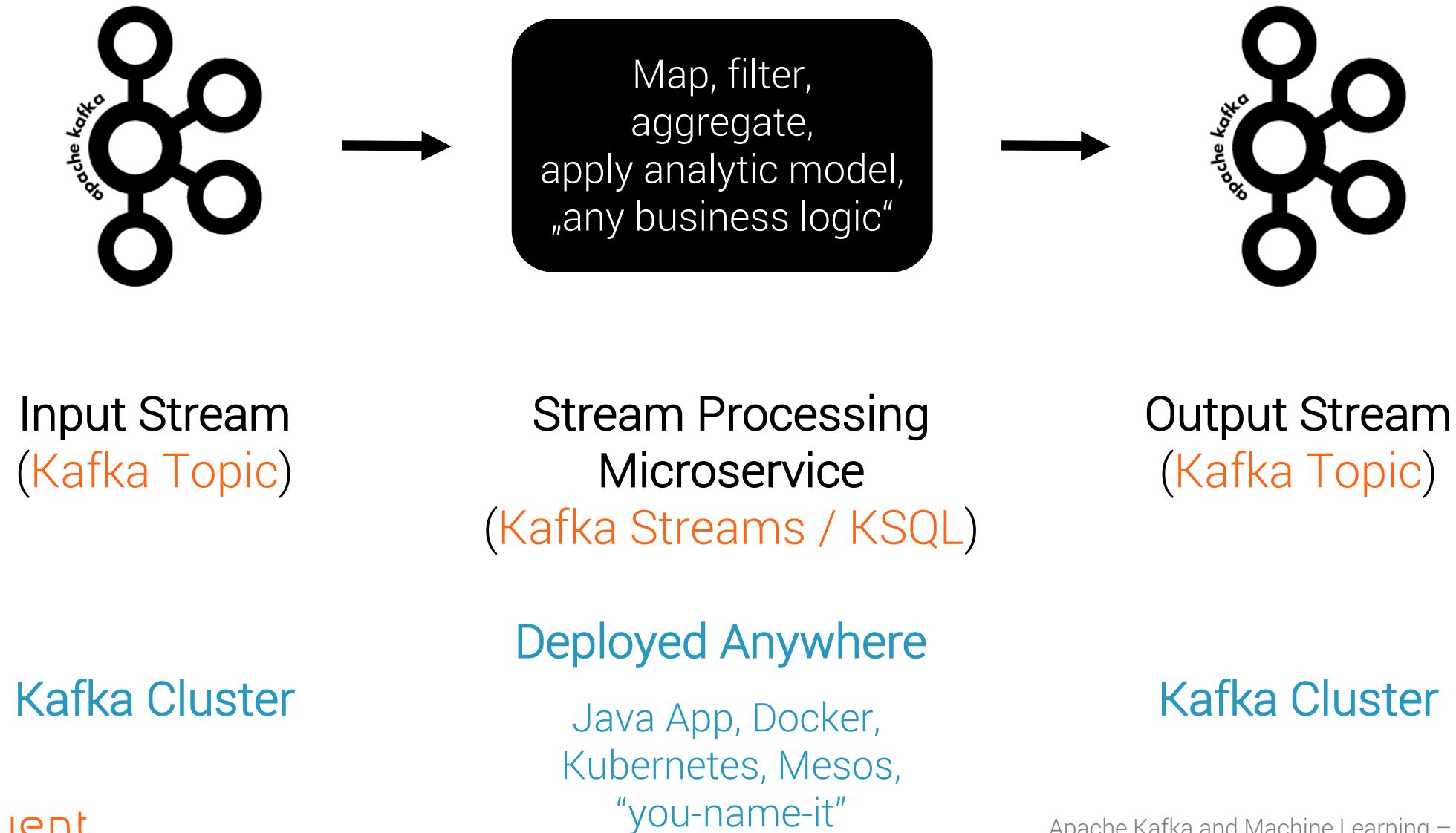
---



# When to use Kafka Streams or KSQL for Stream Processing?



# Kafka Streams (shipped with Apache Kafka) / KSQL (Confluent Open Source)



# Kafka Streams and KSQL

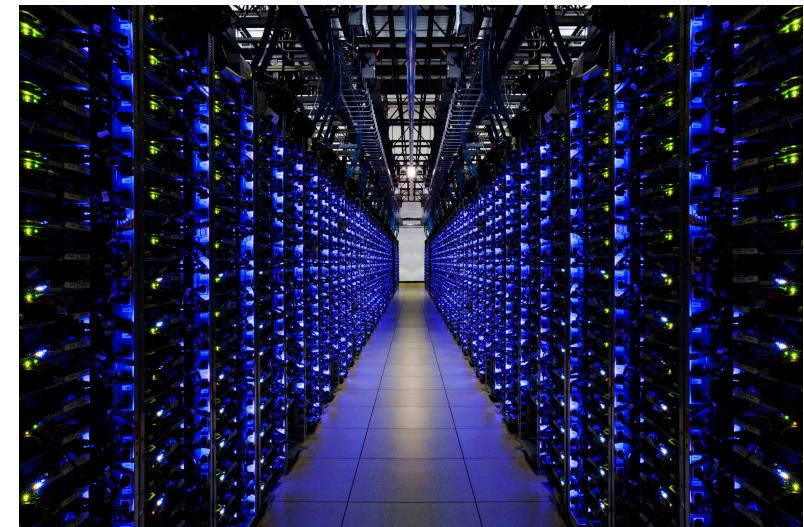
are viable for S / M / L / XL / XXL use cases



Ok.



Ok.

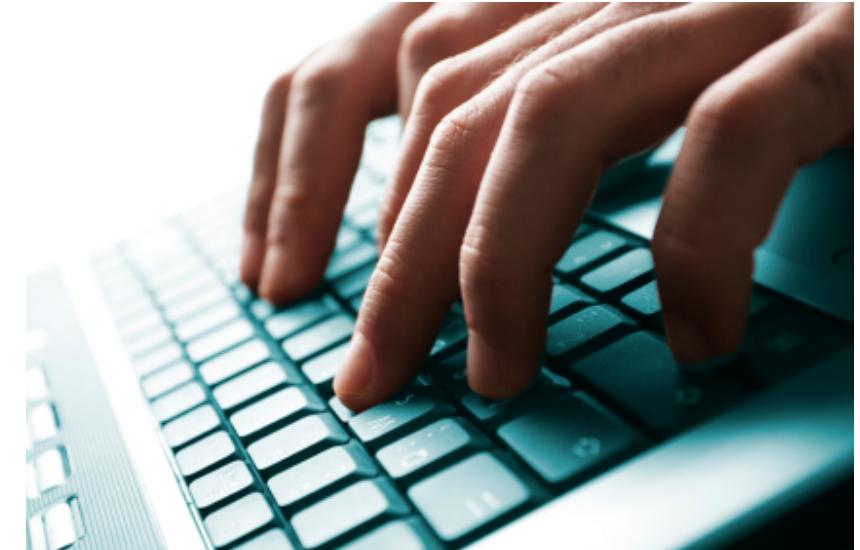


Ok.

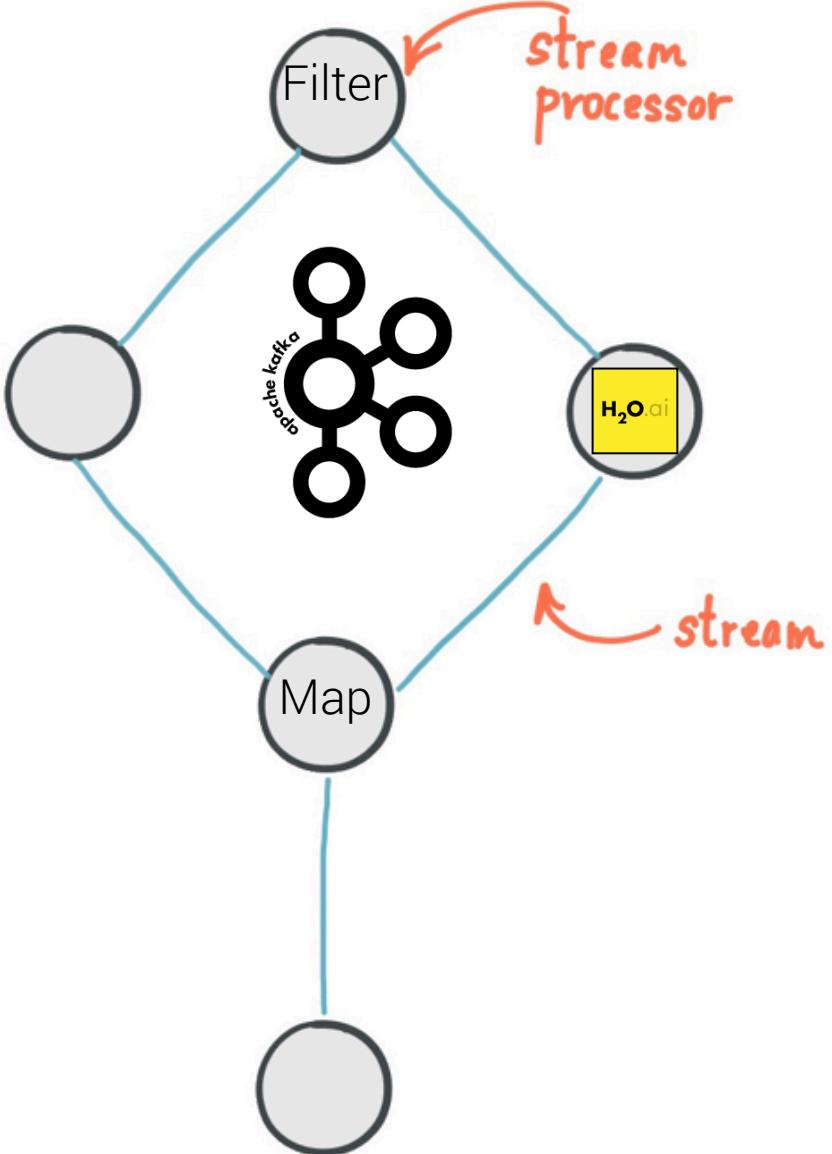
Use Case:  
Airline Flight Delay Prediction

Machine Learning Algorithm:  
Neural Network  
built with H2O and TensorFlow

Streaming Platform:  
Apache Kafka and Kafka Streams



# H2O.ai Model + Kafka Streams



## 1) Create H2O DL model

```
// Create H2O object (see deeplearning_fe7c1f02_08ec_4070_b784_c2531147e451.java)
hex.genmodel.GenModel rawModel;
rawModel = (hex.genmodel.GenModel) Class.forName(modelClassName).newInstance();
EasyPredictModelWrapper model = new EasyPredictModelWrapper(rawModel);
```

## 2) Configure Kafka Streams Application

```
// Configure Kafka Streams Application
final String bootstrapServers = args.length > 0 ? args[0] : "localhost:9092";
final Properties streamsConfiguration = new Properties();
// Give the Streams application a unique name. The name must be unique
// in the Kafka cluster
// against which the application is run.
streamsConfiguration.put(StreamsConfig.APPLICATION_ID_CONFIG, "machine-learning-example");
// Where to find Kafka broker(s).
streamsConfiguration.put(StreamsConfig.BOOTSTRAP_SERVERS_CONFIG, bootstrapServers);
```

### 3) Apply H2O DL model to Streaming Data

```
airlineInputLines.foreach(new ForeachAction<String, String>() {
    public void apply(String key, String value) {
        // Year,Month,DayofMonth,DayOfWeek,DepTime,CRSDepTime,ArrTime,CRSArrTime,UniqueCarrier,FlightNum,TailNum,ActualElapsedTime,CRSElapsedTime,OneWay,Origin,Airline,Dest,Distance,FlightDelayMin
        // value:
        // 1987,10,14,3,741,730,912,849,PS,1451,NA,91,79,NA,23,11,SAN,SFO,447,NA,NA,0,NA,0,NA,NA,NA,NA,NA,YES,YES
        if (value != null && !value.equals("")) {
            System.out.println("#####");
            System.out.println("Flight Input:" + value);
            String[] valuesArray = value.split(",");
            RowData row = new RowData();
            row.put("Year", valuesArray[0]);
            row.put("Month", valuesArray[1]);
            row.put("DayofMonth", valuesArray[2]);
            row.put("DayOfWeek", valuesArray[3]);
            row.put("CRSDepTime", valuesArray[5]);
            row.put("UniqueCarrier", valuesArray[8]);
            row.put("Origin", valuesArray[16]);
            row.put("Dest", valuesArray[17]);
            BinomialModelPrediction p = null;
            try {
                p = model.predictBinomial(row);
            } catch (PredictException e) {
                e.printStackTrace();
            }
            System.out.println("Label (aka prediction) is flight departure delayed: " + p.label);
            System.out.print("Class probabilities: ");
            for (int i = 0; i < p.classProbabilities.length; i++) {
                if (i > 0) {
                    System.out.print(",");
                }
                System.out.print(p.classProbabilities[i]);
            }
        }
    }
})
```

#### 4) Start Kafka Streams App

```
// Start Kafka Streams Application to process new incoming messages from Input Topic
final KafkaStreams streams = new KafkaStreams(builder, streamsConfiguration);
streams.cleanUp();
streams.start();
```

# Github Examples: Kafka + Machine Learning

[kaiwaechner / kafka-streams-machine-learning-examples](#)

Code Issues 0 Pull requests 0 Projects 0 Wiki Settings Insights ▾

This project contains examples which demonstrate how to deploy analytic models to mission-critical, scalable production environments leveraging Apache Kafka and its Streams API. Models are built with Python, H2O, TensorFlow, DeepLearning4 and other technologies.

Edit

kafka kafka-streams kafka-client machine-learning deep-learning open-source h2o h2oai tensorflow deeplearning4j keras  
keras-tensorflow Manage topics

30 commits 1 branch 0 releases 1 contributor Apache-2.0

Branch: master New pull request Create new file Upload files Find file Clone or download ▾

Kai Waehner Added TensorFlow CNN Example for image recognition Latest commit 293fd23 10 days ago

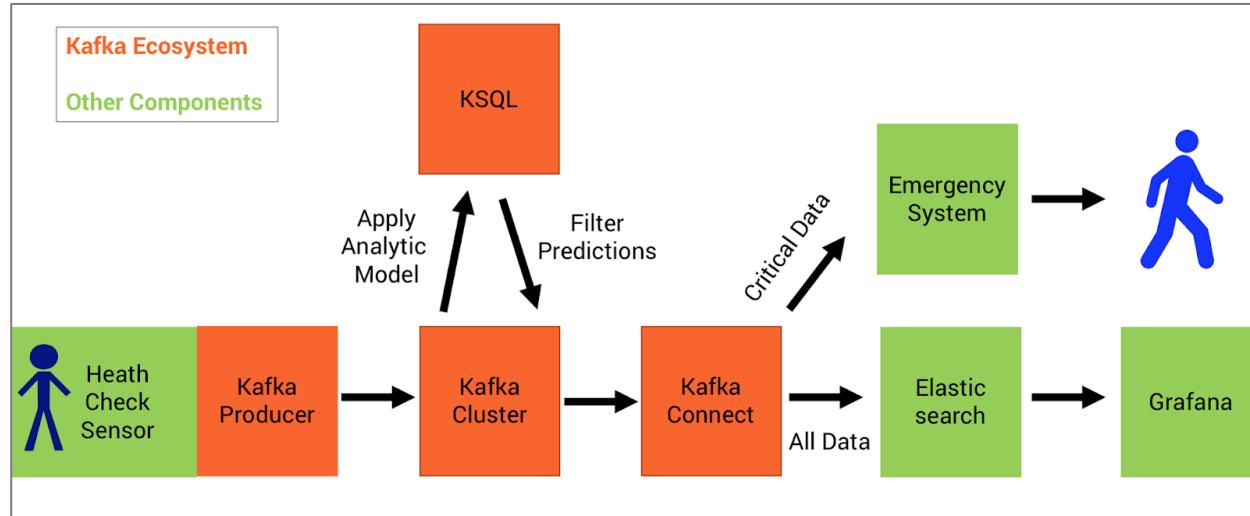
src	Added TensorFlow CNN Example for image recognition	10 days ago
LICENSE	Initial commit	19 days ago
pom.xml	Added TensorFlow CNN Example for image recognition	10 days ago
readme.md	Added TensorFlow CNN Example for image recognition	10 days ago

<https://www.confluent.io/blog/build-deploy-scalable-machine-learning-production-apache-kafka/>  
<https://github.com/kaiwaechner/kafka-streams-machine-learning-examples>

1) git clone → 2) mvn clean package → 3) look at implementations and unit tests



# KSQL and Deep Learning (Autoencoder) for IoT Sensor Analytics



“SELECT event\_id, anomaly(SENSORINPUT) FROM health\_sensor;”



KSQL UDF using an analytic model under the hood  
→ Write once, use in any KSQL statement

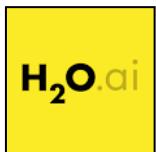
# Live Demo

---

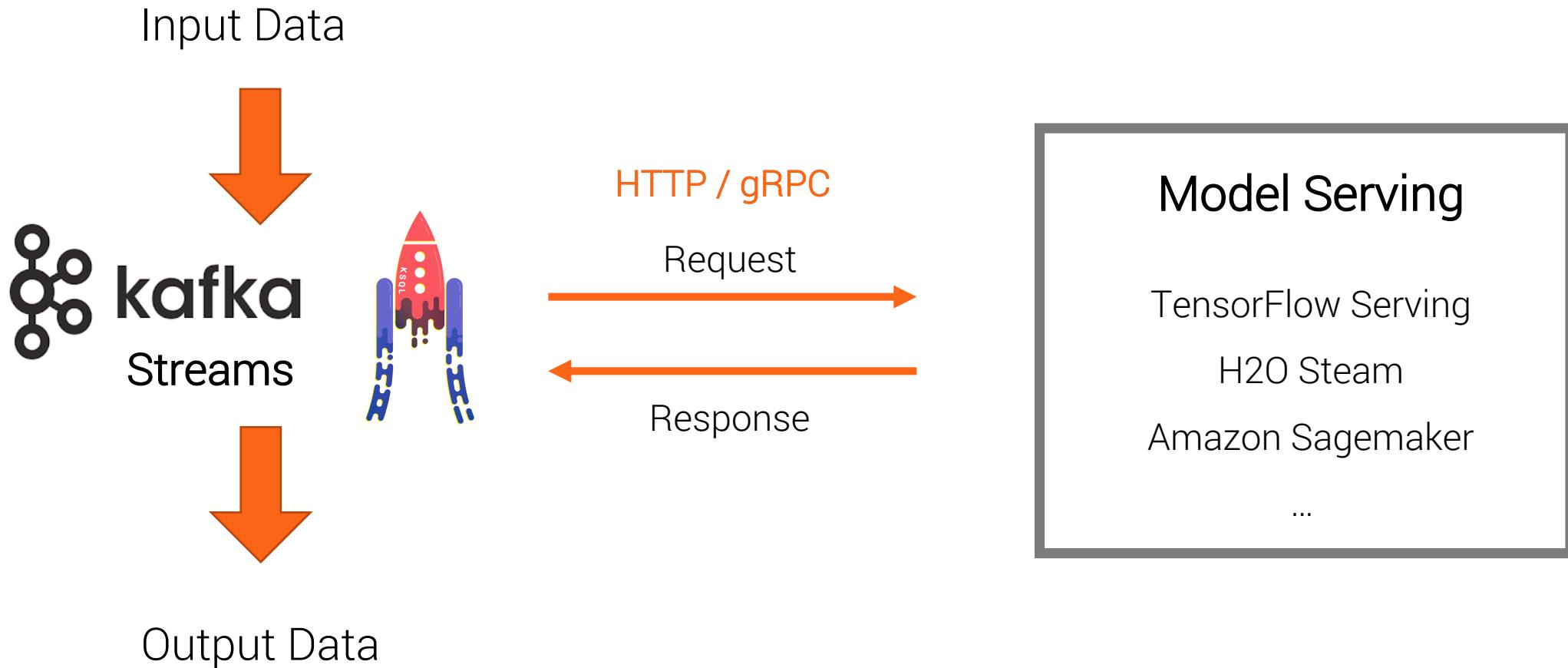
Use Case:  
Anomaly Detection  
(Sensor Healthcheck)

Machine Learning Algorithm:  
Autoencoder built with H2O

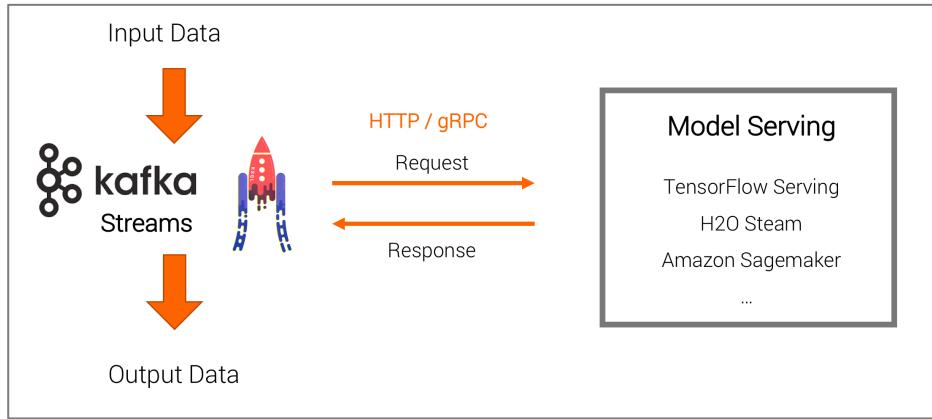
Streaming Platform:  
Apache Kafka and KSQL



# Stream Processing vs. Request-Response for Model Serving



# Stream Processing vs. Request-Response for Model Serving



## Pros:

- Simple integration with existing systems and technologies
- Easier to understand if you come from non-streaming world
- Later migration to real streaming is also possible

## Cons:

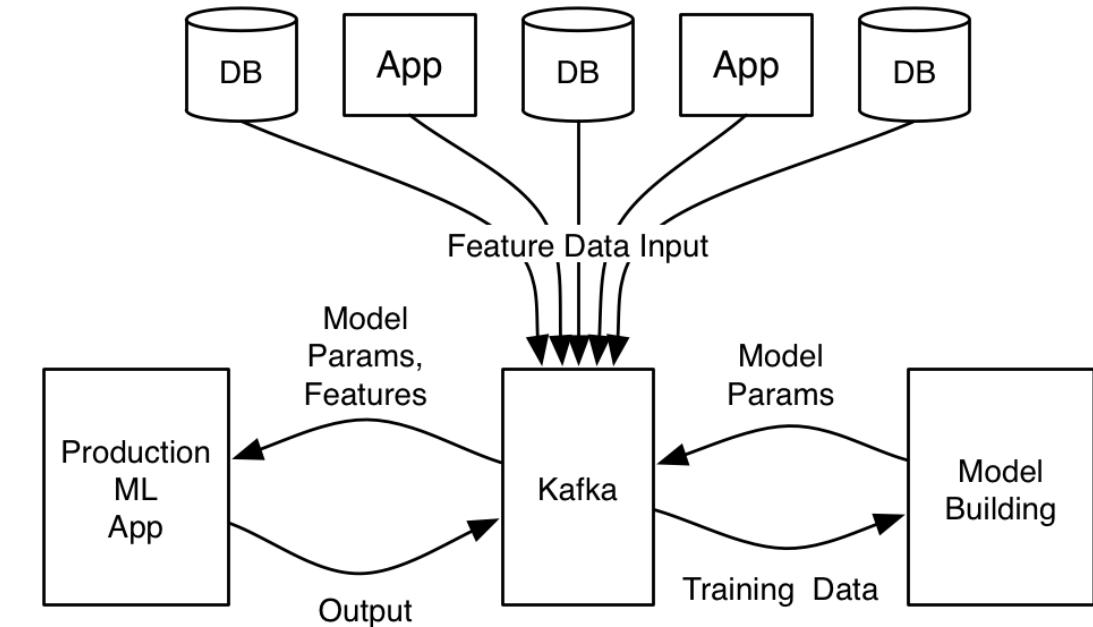
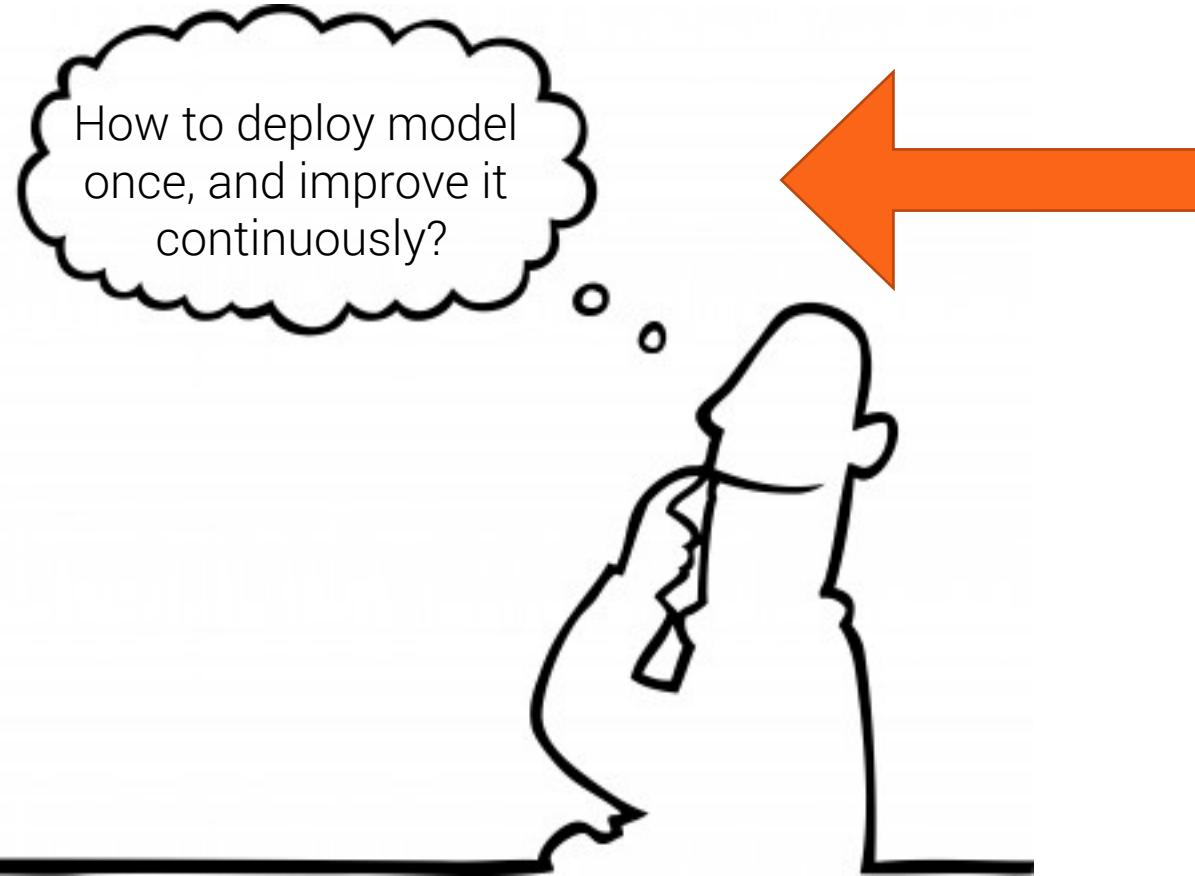
- Coupling the availability, scalability, and latency/throughput of your Kafka Streams application with the SLAs of the RPC interface
- Side-effects (e.g. in case of failure) not covered by Kafka processing (e.g. Exactly Once)

# Agenda

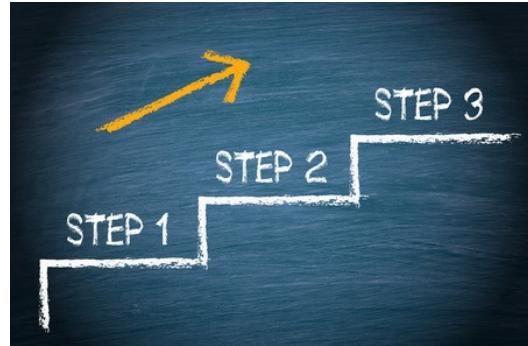
---

- 1) Added Business Value via Machine Learning
- 2) Apache Kafka Ecosystem as Infrastructure for Machine Learning
- 3) Data Ingestion and Preprocessing with Apache Kafka
- 4) Predictions in Real Time with Kafka Streams and KSQL
- 5) Automation and DevOps of a Machine Learning Infrastructure**

# Automated Model Improvement with Apache Kafka and Kafka Streams



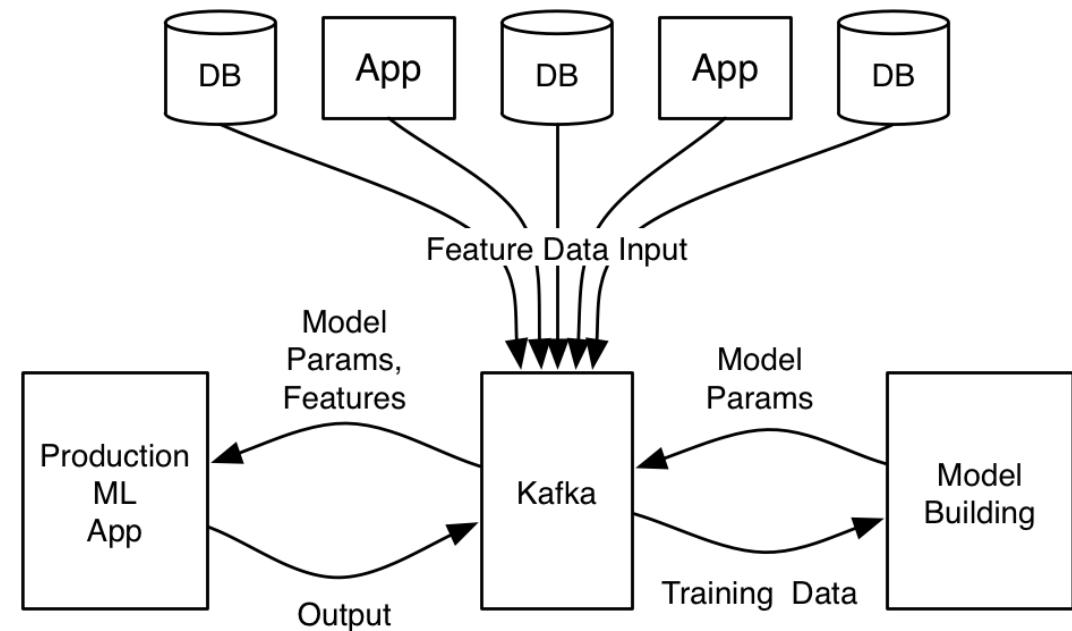
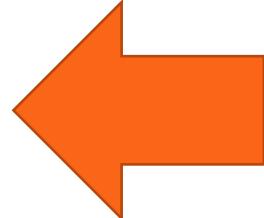
# Automated Model Improvement with Apache Kafka and Kafka Streams



How to improve models?

1. Manual Update
2. Continuous Batch Updating
3. Real Time → Online Model Training

Your choice... All possible with Kafka!



# Caveats for Online Model Training

---

- Processes and infrastructure not ready
- Validation needed before production
- Slows down the system
- Only a few ML implementations → Build your own!
- Only possible for unsupervised ML (e.g. clustering)
- Many use cases do not need it

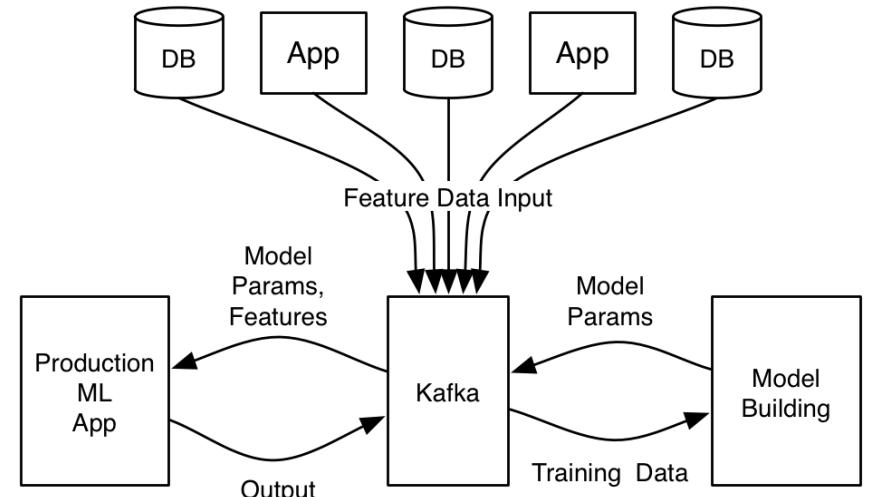
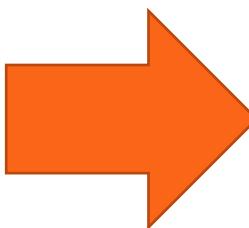
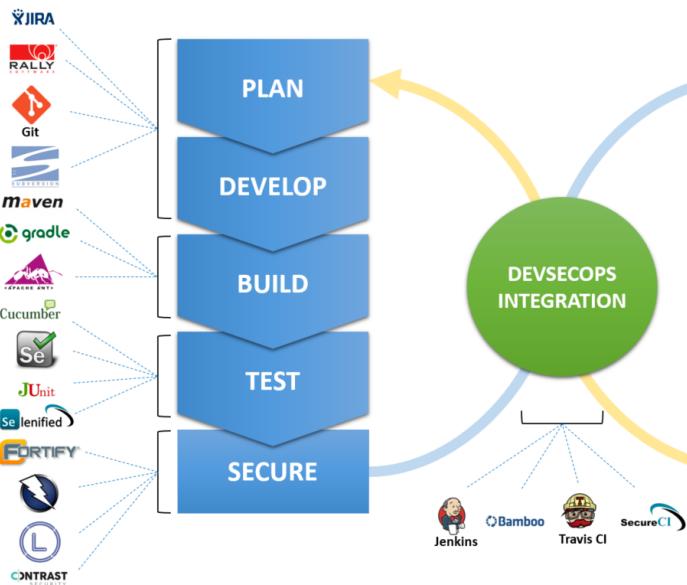
→ Do it only when feasible!



# Continuous Batch Updating as “Best Feasible Option”

## DevOps Pipeline

1. Apply the model online to make predictions
2. Collect data and train a new model
3. Automated Re-Deployment



<https://www.confluent.io/blog/predicting-flight-arrivals-with-the-apache-kafka-streams-api/>  
<https://www.coveros.com/services/devops/>

# Kubernetes – The Winner of the Container and DevOps Wars!



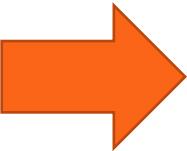
INFOWORLD TECH WATCH  
By Matt Asay, InfoWorld | SEP 9, 2016

About

Informed news analysis every weekday

## Why Kubernetes is winning the container war

It's all about knowing how to build an open source community -- plus experience running applications in Linux containers, which Google invented



# kubernetes



Google Cloud Platform



Amazon EKS



CLOUD FOUNDRY



docker

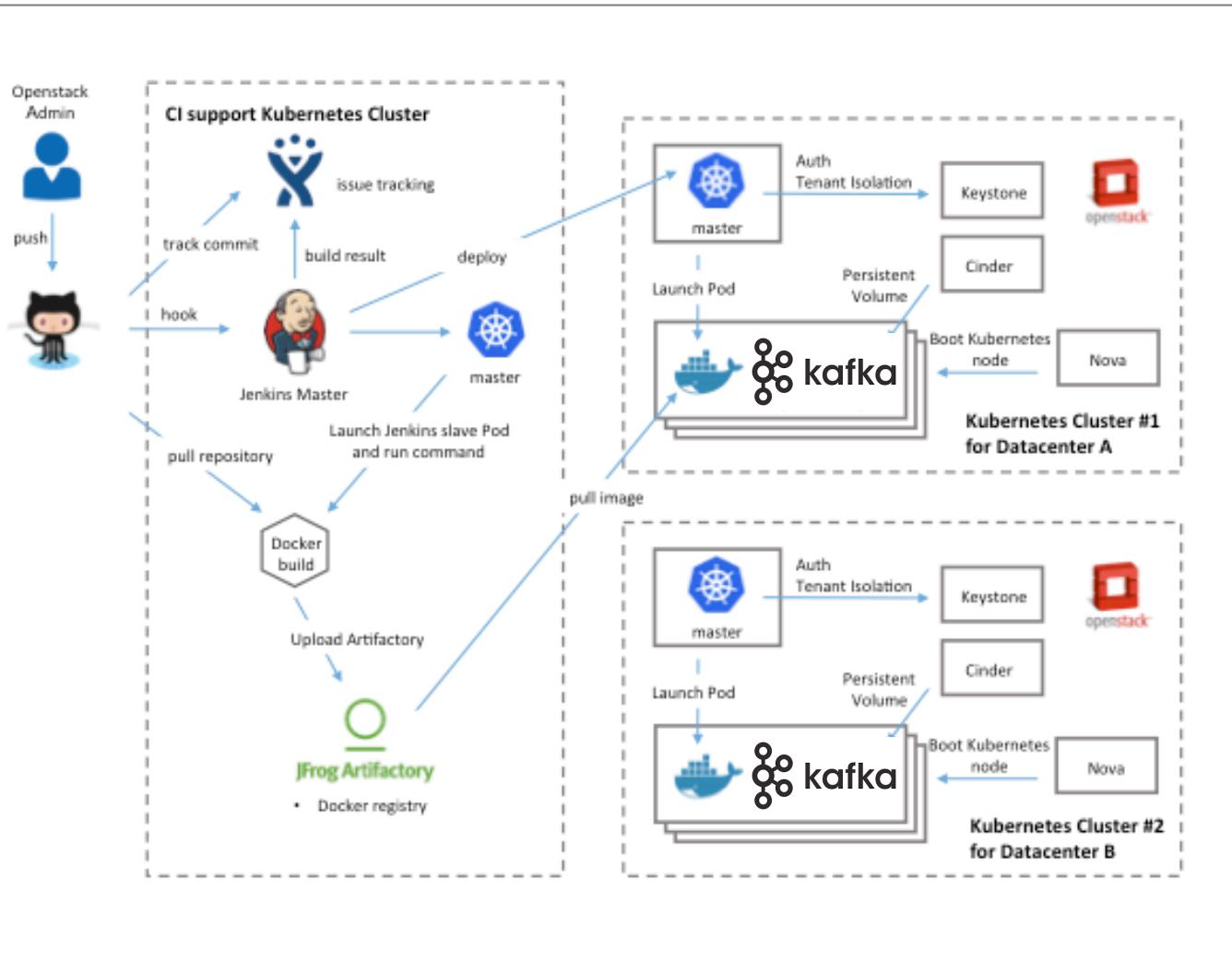
Docker, Inc



MESOSPHERE

<https://www.infoworld.com/article/3118345/cloud-computing/why-kubernetes-is-winning-the-container-war.html>  
<http://techgenix.com/year-of-kubernetes/>

# Kubernetes for Infrastructure Deployment



## Backend

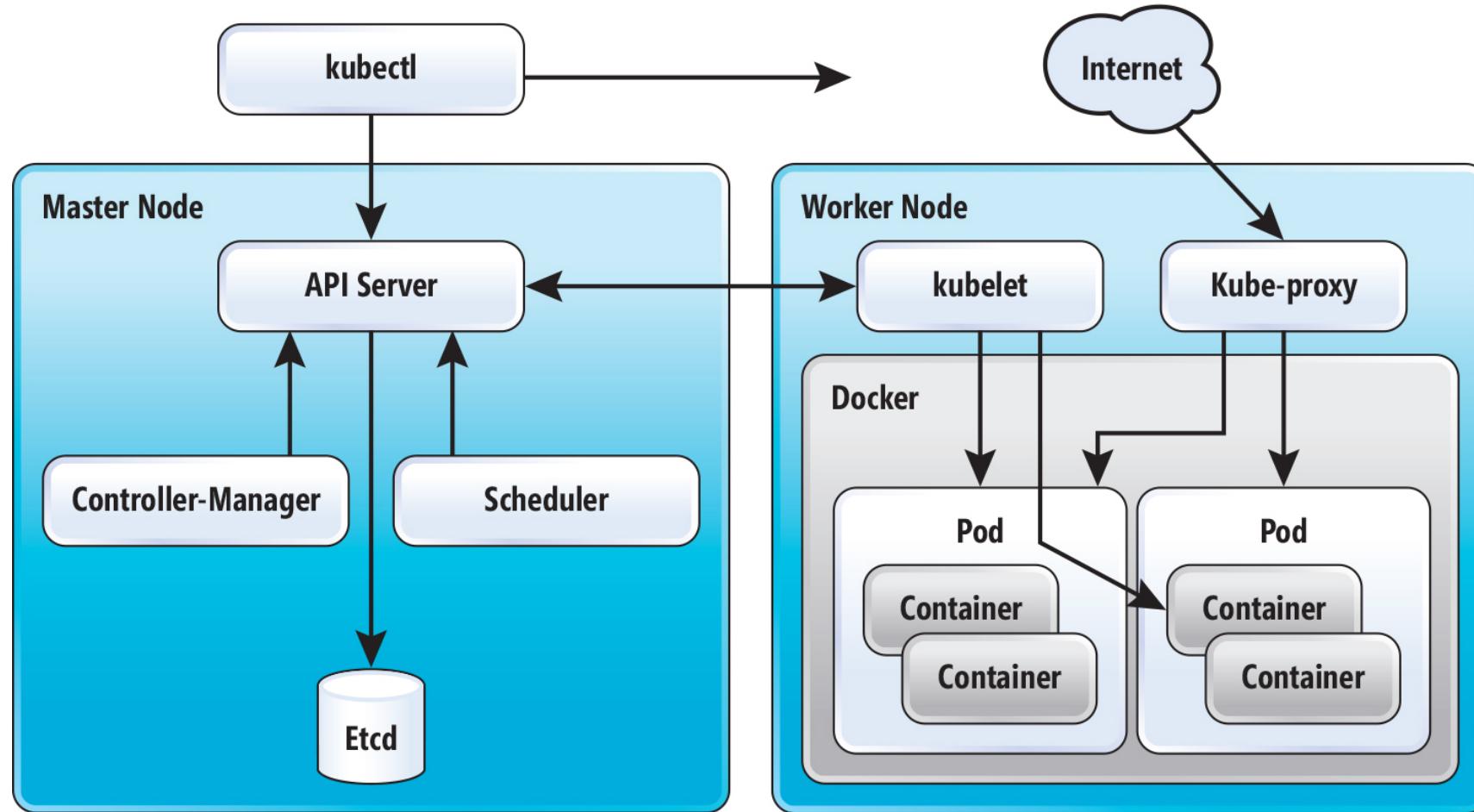
- Zookeeper and Kafka Broker Pods
- REST Proxy, Schema Registry
- Persistent Volumes
- Kubernetes Operator

## Clients

- Java / .Net / Go / Python Kafka Clients
- Kafka Streams / KSQL Apps
- Scalability and Elasticity

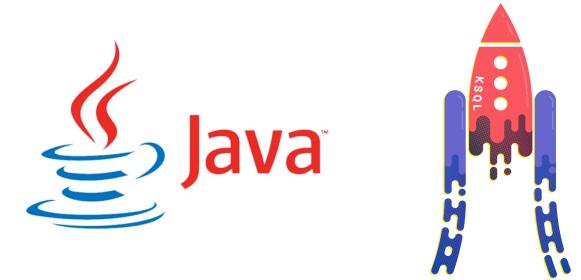
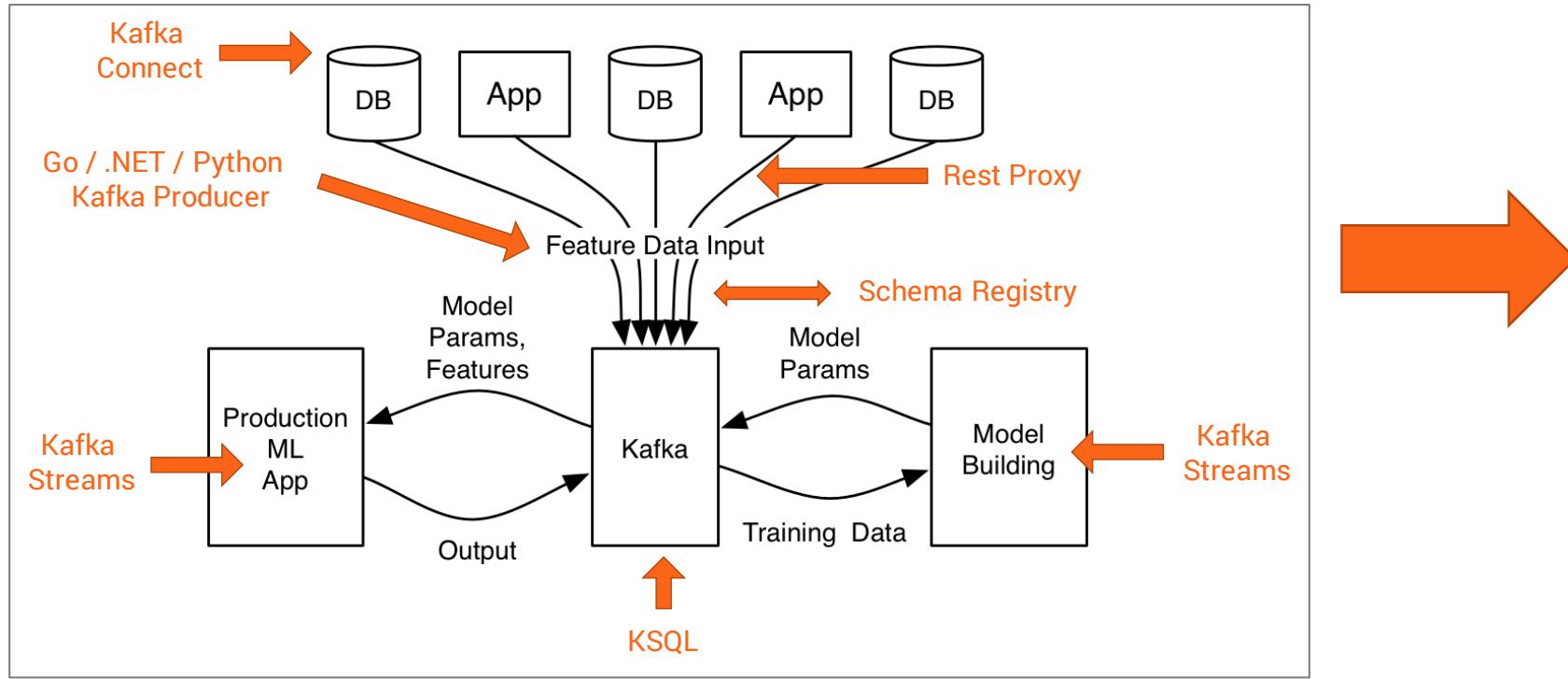
<https://kubernetes.io/blog/2016/10/kubernetes-and-openstack-at-yahoo-japan>

# Kubernetes for Kafka Deployment



[https://redmondmag.com/articles/2017/08/01/~media/ECG/redmondmag/Images/2017/08/0817red\\_Kubernetes\\_Figure1\\_hires.ashx](https://redmondmag.com/articles/2017/08/01/~media/ECG/redmondmag/Images/2017/08/0817red_Kubernetes_Figure1_hires.ashx)

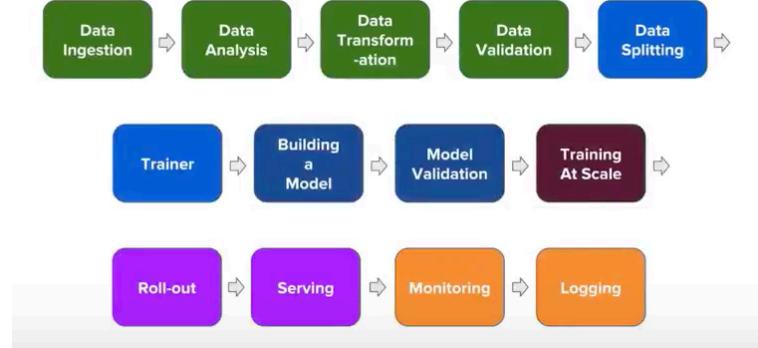
# Monitoring the Infrastructure for Machine Learning



Build vs. Buy  
Hosted vs. Managed  
Basic vs. Advanced

# Kubernetes Deployment of ML Workflows

The screenshot shows the GitHub repository page for 'kubeflow / kubeflow'. The repository has 171 pull requests, 3,011 stars, and 354 forks. It includes sections for Code, Issues (161), Pull requests (20), Projects (0), Wiki, and Insights. Tags at the bottom include ml, kubernetes, minikube, tensorflow, notebook, jupyterhub, and google-kubernetes-engine.



## Kubeflow

The Kubeflow project is dedicated to making **deployments** of machine learning (ML) workflows on **Kubernetes** simple, portable and scalable. Our goal is **not** to recreate other services, but to provide a straightforward way to deploy best-of-breed open-source systems **for ML** to diverse infrastructures. Anywhere you are running Kubernetes, you should be able to run Kubeflow.

Warning:

Early Stage with focus on TensorFlow Training, TensorFlow Serving, Jupyter...  
Bigger ecosystem expected soon... Including Kafka components for ingestion, serving, monitoring...

<https://github.com/kubeflow/kubeflow>

# Key Take-Aways

---



- Data Scientist and Developers have to work together continuously (org + tech!)
- Mission critical, scalable production infrastructure is key for success of Machine Learning projects
- Apache Kafka Ecosystem + Cloud = Machine Learning at Extreme Scale (Ingestion, Processing, Training, Inference, Monitoring)



DL4J  
DEELEARNING4J



Questions? Feedback?  
Please contact me!

## Kai Waehner

Technology Evangelist

[kontakt@kai-waehner.de](mailto:kontakt@kai-waehner.de)  
[@KaiWaehner](https://twitter.com/KaiWaehner)  
[www.kai-waehner.de](http://www.kai-waehner.de)  
[www.confluent.io](http://www.confluent.io)  
[LinkedIn](#)

