



Machine Learning with the Apache Kafka Open Source Ecosystem

How to build a Scalable, Mission-Critical Machine Learning Infrastructure?

Kai Waehner

Technology Evangelist
kontakt@kai-waehner.de
LinkedIn
@KaiWaehner
www.confluent.io
www.kai-waehner.de



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

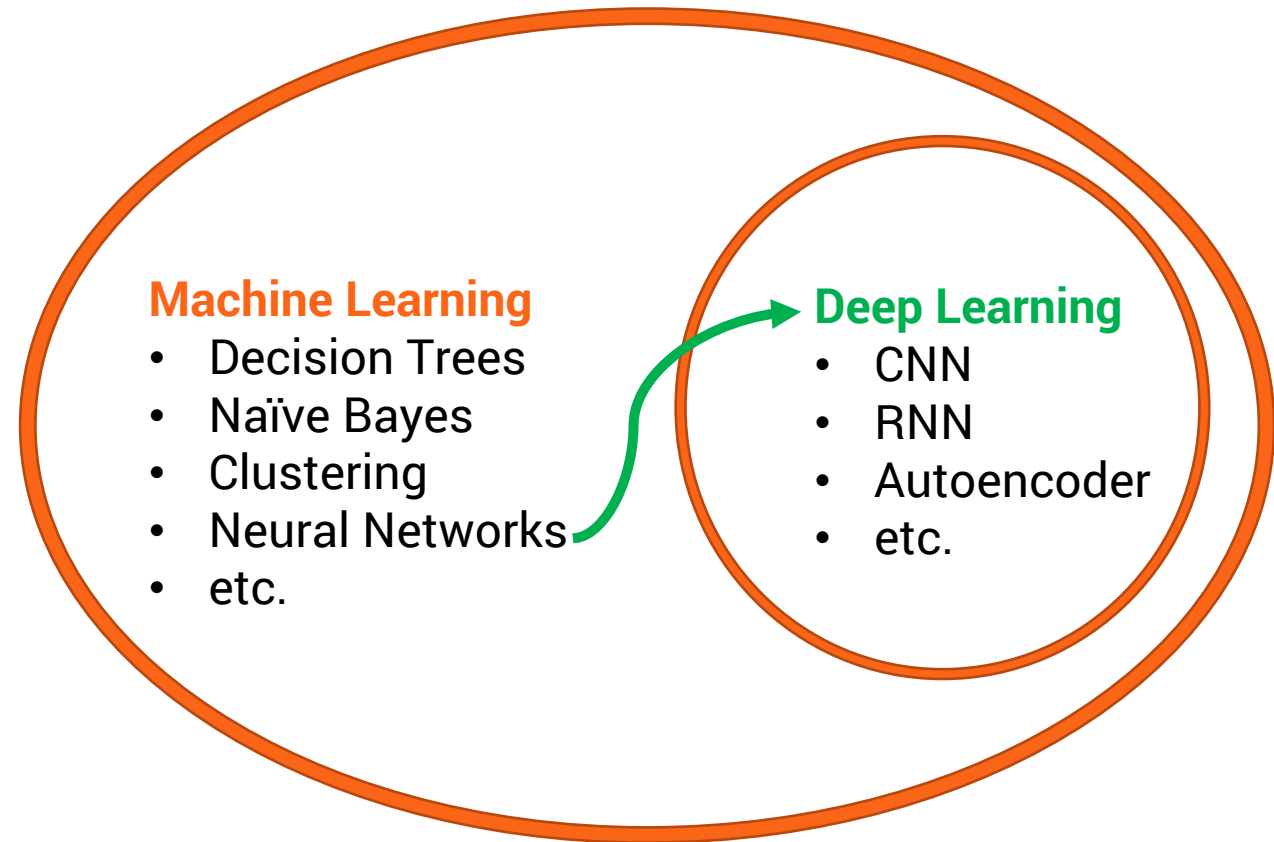
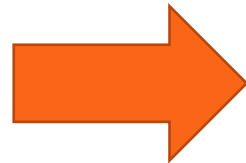
- 1) Added Business Value via Machine Learning
- 2) Apache Kafka Ecosystem as Infrastructure for Machine Learning
- 3) Data Ingestion and Preprocessing with Apache Kafka for Model Training
- 4) Predictions in Real Time with Kafka Streams and KSQL
- 5) DevOps and Monitoring of a Machine Learning Infrastructure

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 for Model Training
- 4) Predictions in Real Time with Kafka Streams and KSQL
- 5) DevOps and Monitoring of a Machine Learning Infrastructure

Machine Learning

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



Real World Examples of Machine Learning



Spam Detection



Search Results +
Product Recommendation



Picture Detection
(Friends, Locations, Products)

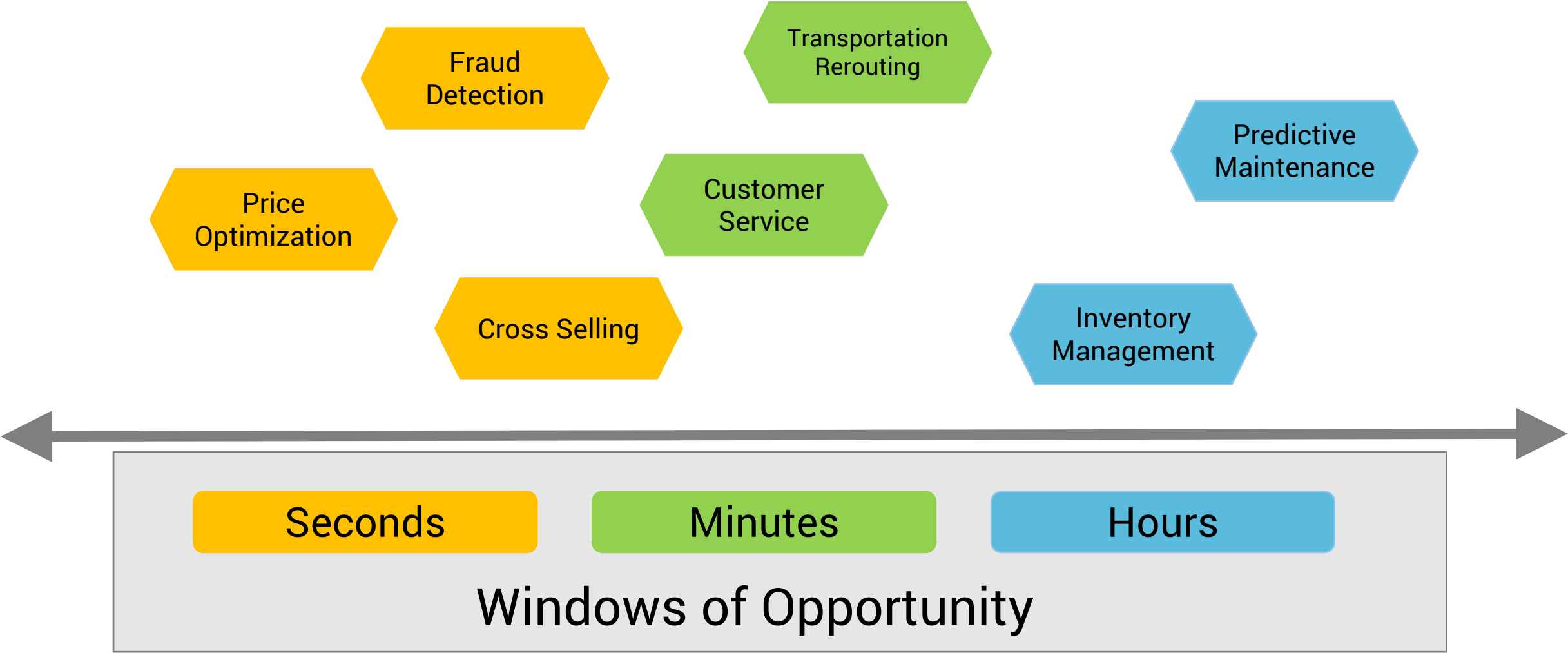


The Next Disruption:
Google Beats Go Champion



Your Company

Leverage Machine Learning to Analyze and Act on Critical Business Moments

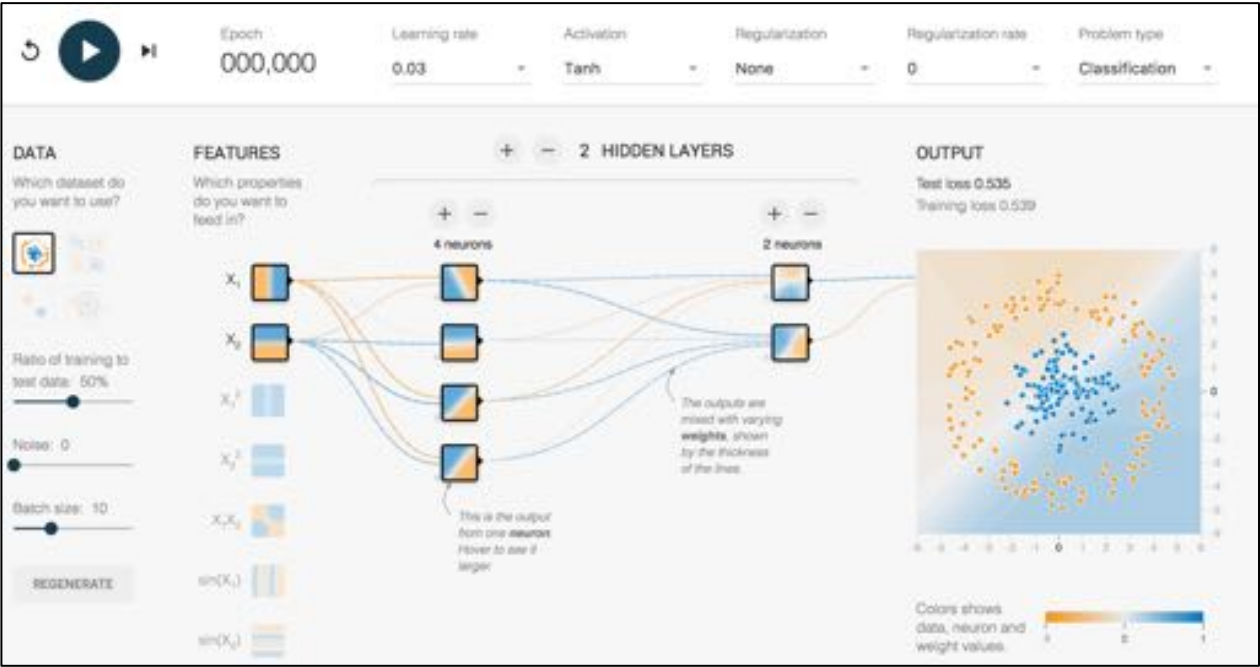


Live Demo – Building an Analytic Model



Neural Networks in Action

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



Languages, Frameworks and Tools for Machine Learning



theano

dmlc
mxnet

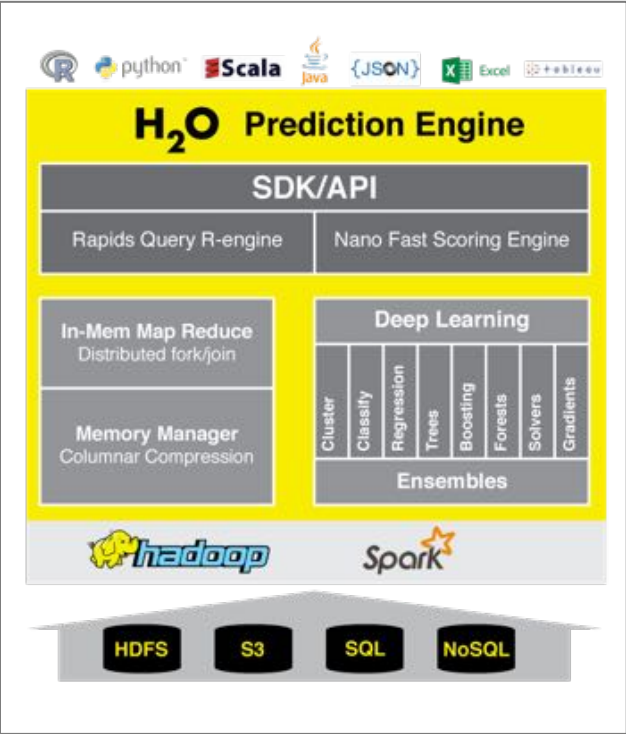


DL4J
DEEPLEARNING4J



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 in H2O. The 'Select an algorithm' dropdown is set to 'Deep Learning'. The 'PARAMETERS' section includes fields for 'model_id', 'training_frame', 'validation_frame', 'nfold', 'response_column', and 'ignored_columns'. The 'ignored_columns' section shows a list of columns with checkboxes for selection. The 'ignore_const_cols' checkbox is checked. The 'activation' dropdown is set to 'Rectifier'. The 'hidden' field is set to '200, 200'. The 'epochs' field is set to '10'. The 'variable_importances' checkbox is checked.



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] = 0.18894640014126883;
            sa[3] = 0.5257610635956896;
            sa[4] = 0.00209932098808304;
```



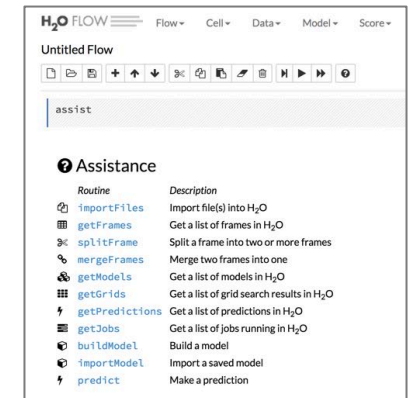
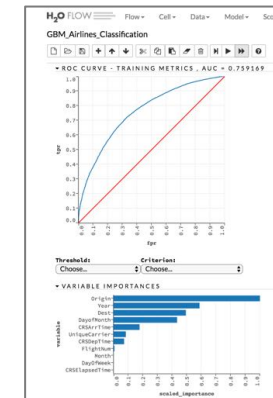
H2O Engine

Live Demo – Building an Analytic Model

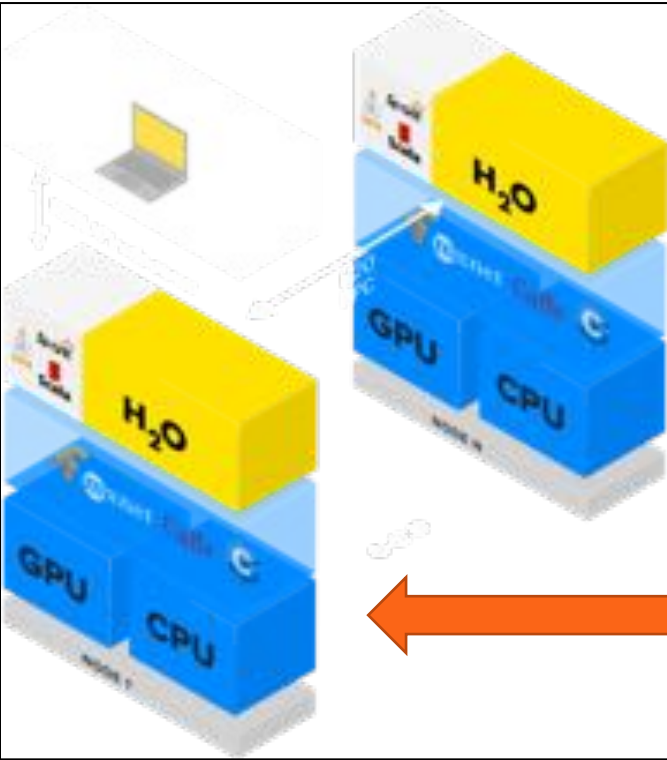
Use Case:
Airline Flight Delay Prediction

Machine Learning Algorithm:
Deep Learning
using Neural Networks

Technology:
H2O.ai, TensorFlow



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



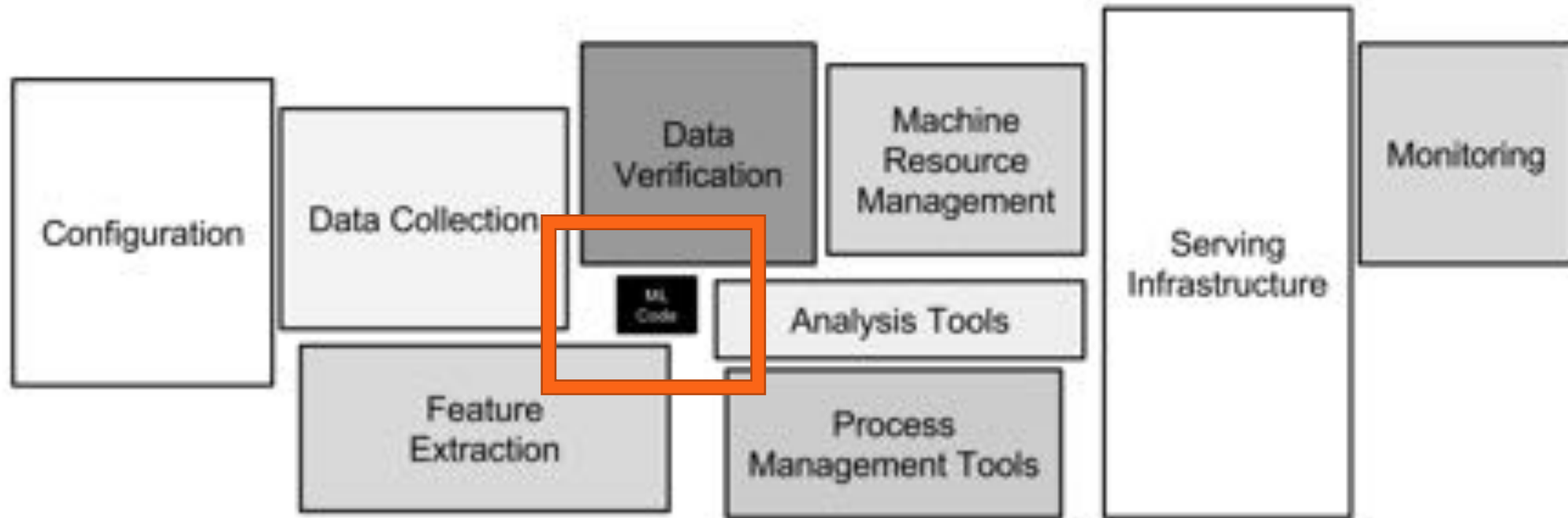
max_runtime_secs	~ (Choose...)	Maximum allowed runtime in seconds for model training. Use 0 to disable.
backend	mxnet	Deep Learning Backend.
image_shape	caffe	Width and height of image.
channels	3	Number of (color) channels.
network_definition_file		Path of file containing network definition (graph, architecture).
network_parameters_file		Path of file containing network (initial) parameters (weights, biases).

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

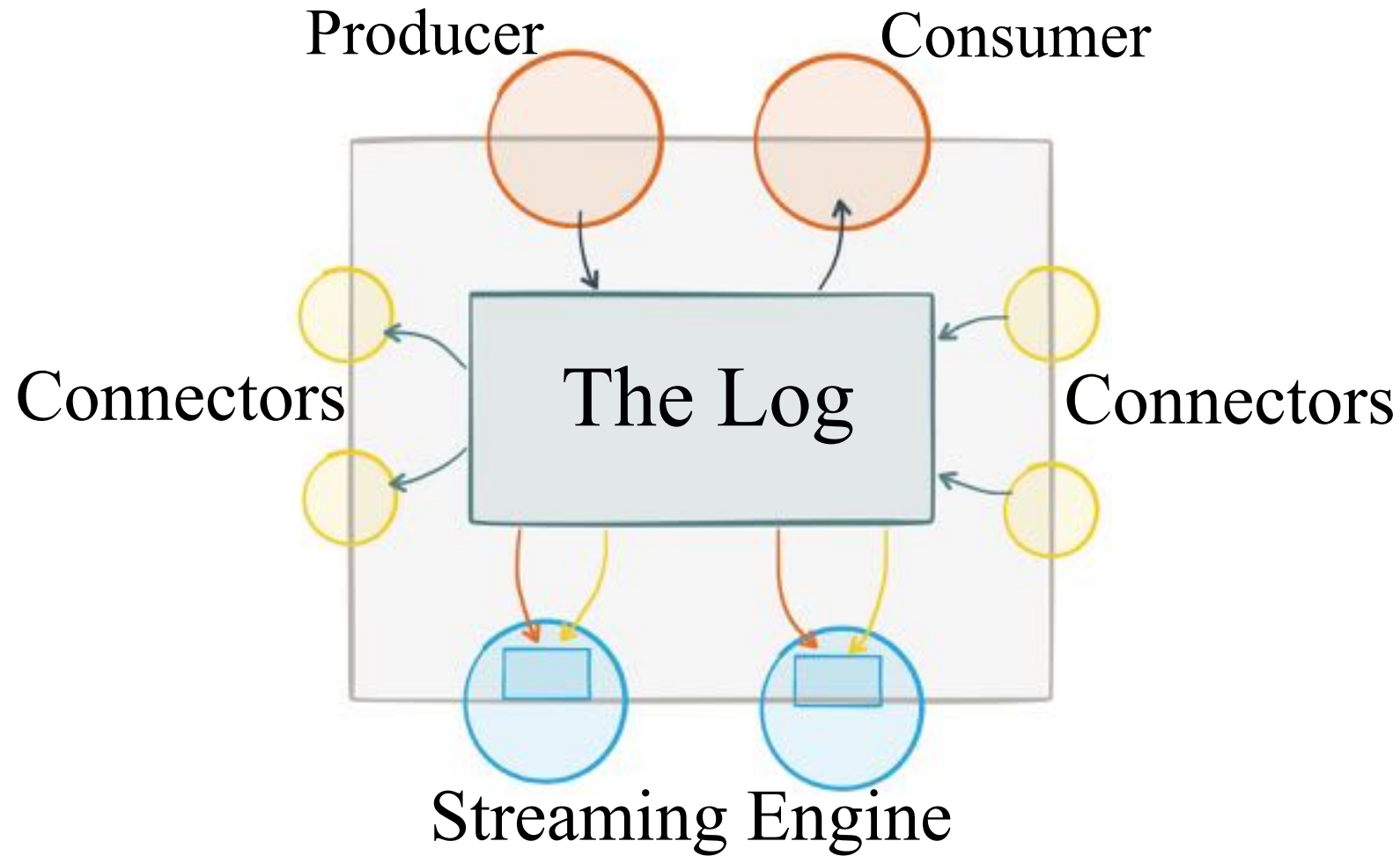


<https://papers.nips.cc/paper/5656-hidden-technical-debt-in-machine-learning-systems.pdf>

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 for Model Training
- 4) Predictions in Real Time with Kafka Streams and KSQL
- 5) DevOps and Monitoring of a Machine Learning Infrastructure

Apache Kafka – The Rise of a Streaming Platform



1 PUB/SUB
2 STORE
3 PROCESS

Apache Kafka at Scale



(2018)

Operation Challenges

- The scale of Kafka deployment @LinkedIn
 - 2,100+ brokers
 - ~ 60,000 topics
 - ~ 4.2 million partitions
 - > 4.5 trillion messages / day

> 4.5 trillion messages / day

Strata
DATA CONFERENCE



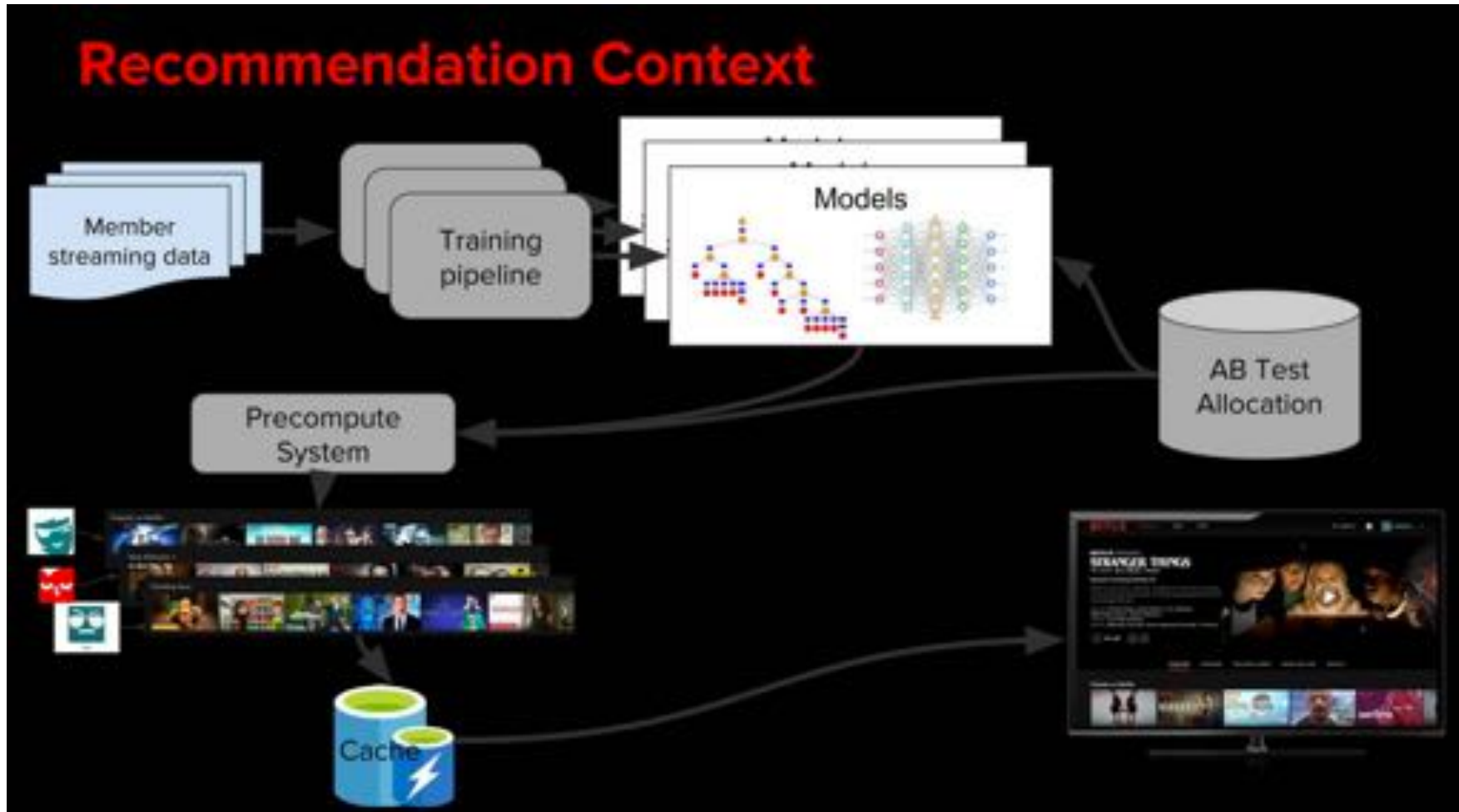
(2018)

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

6 Petabytes

QCon



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

Meet Michelangelo: Uber's Machine Learning Platform

By Jeremy Hermann & Mike Del Balso

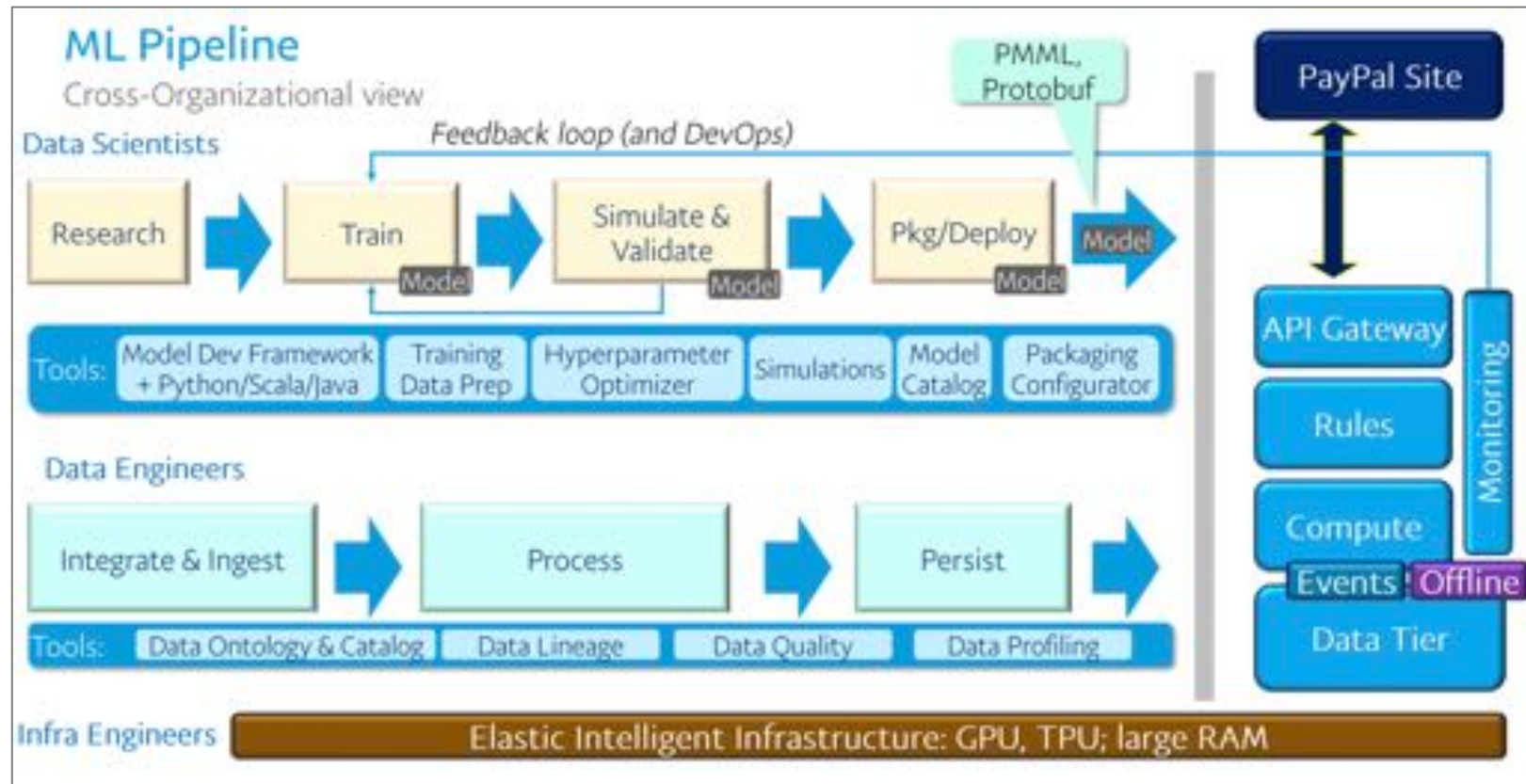
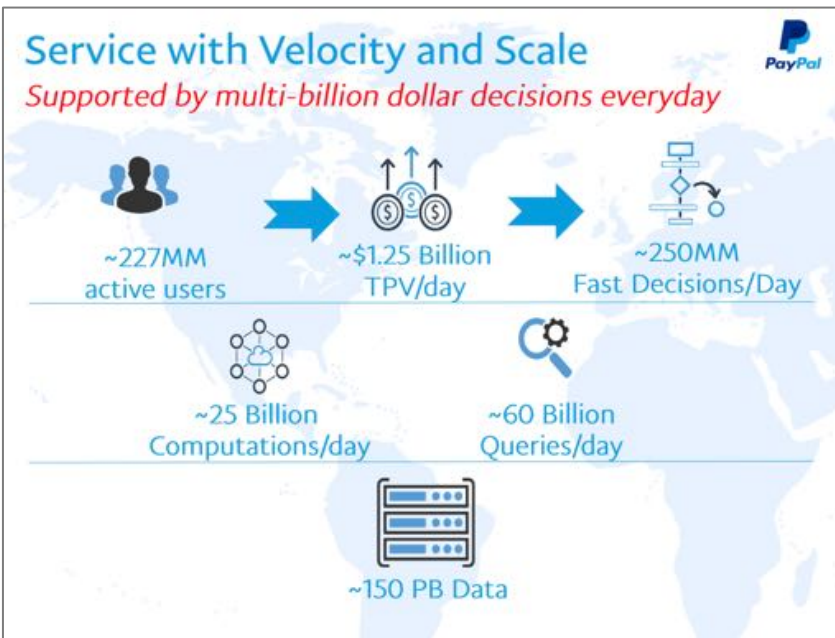
September 5, 2017



- 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

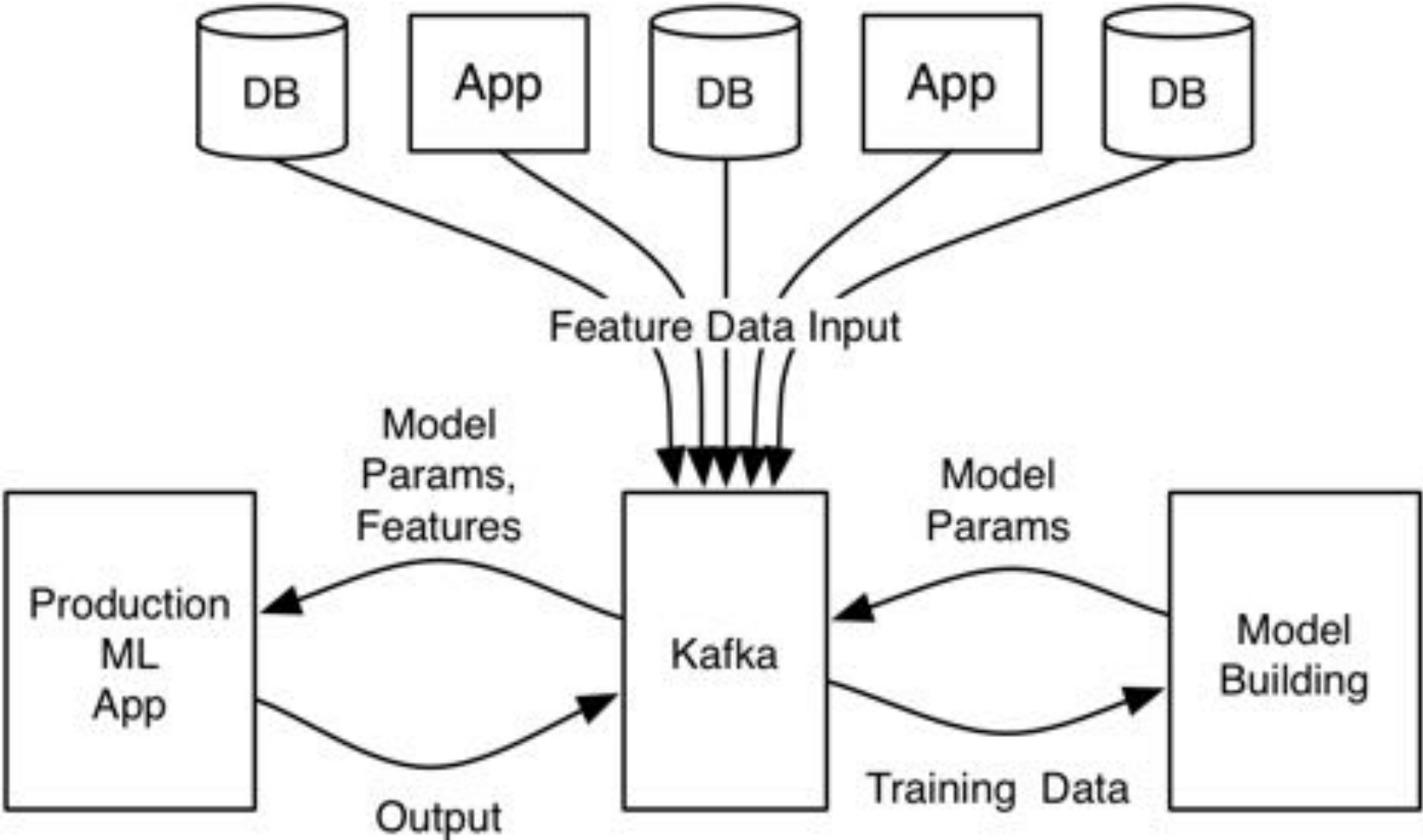
<https://eng.uber.com/michelangelo>

Paypal: Real Time Fraud Detection at Scale

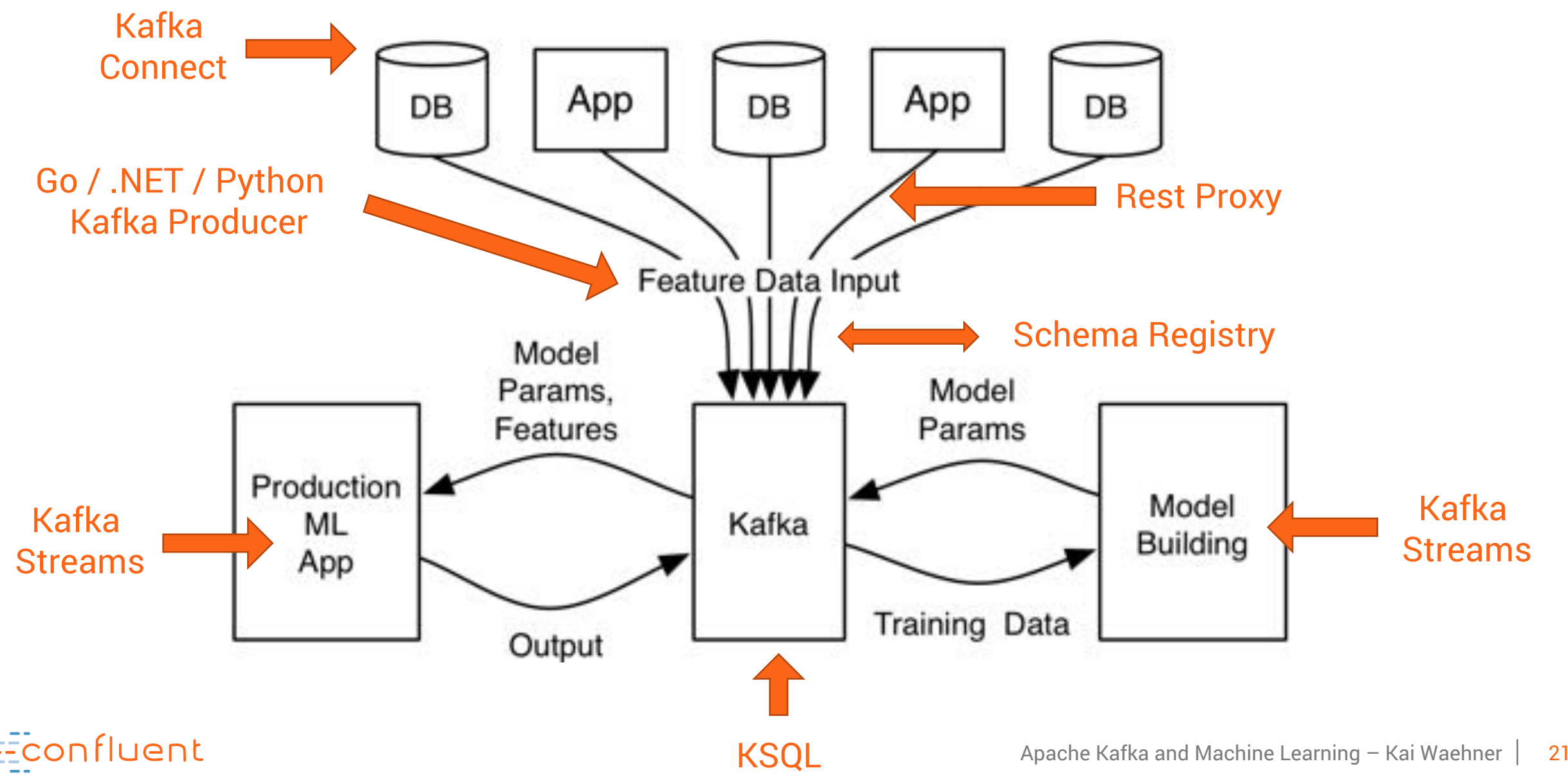


“Scalable model / framework independent infrastructure for fraud detection”

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

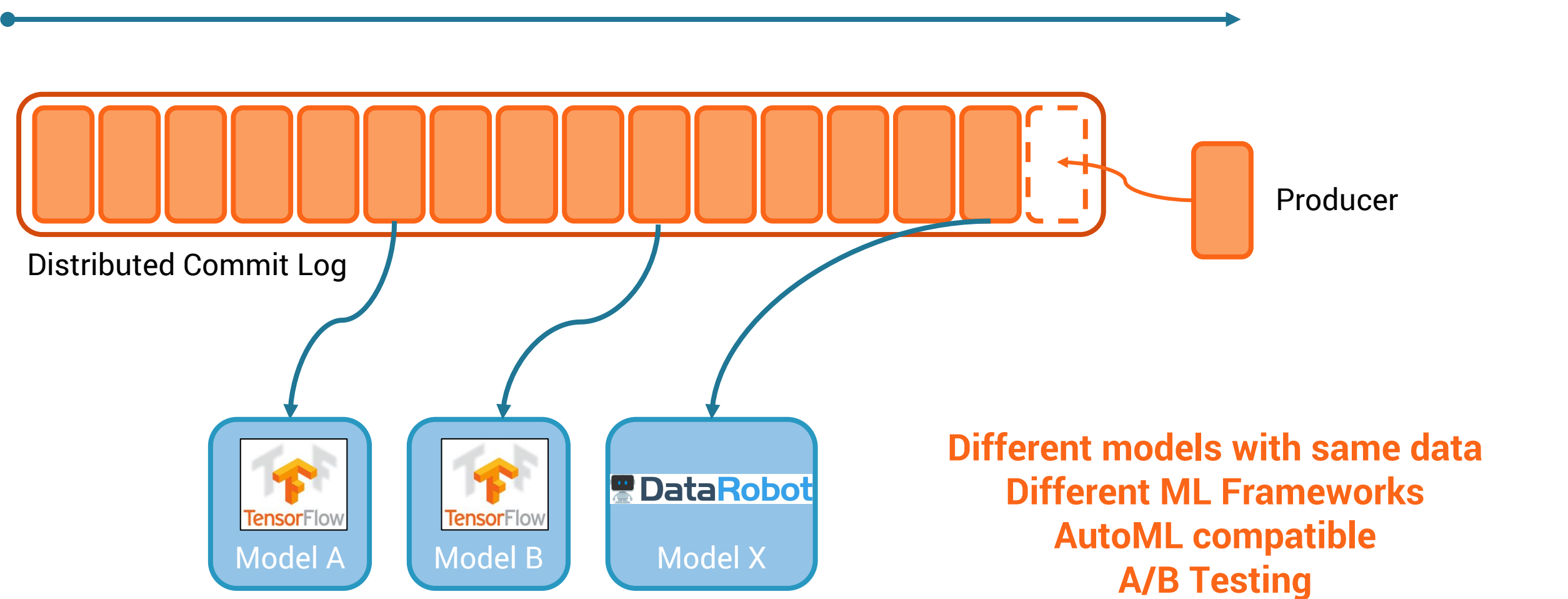


Apache Kafka's **Open Source Ecosystem** as Infrastructure for Machine Learning



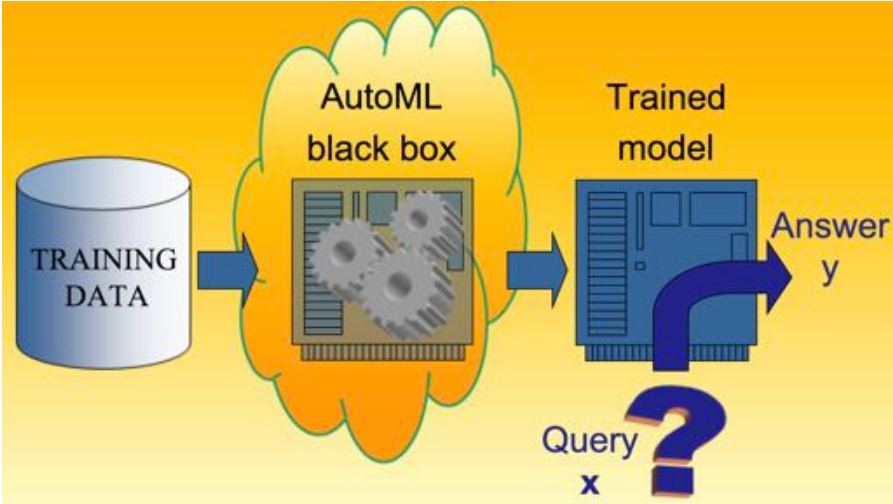
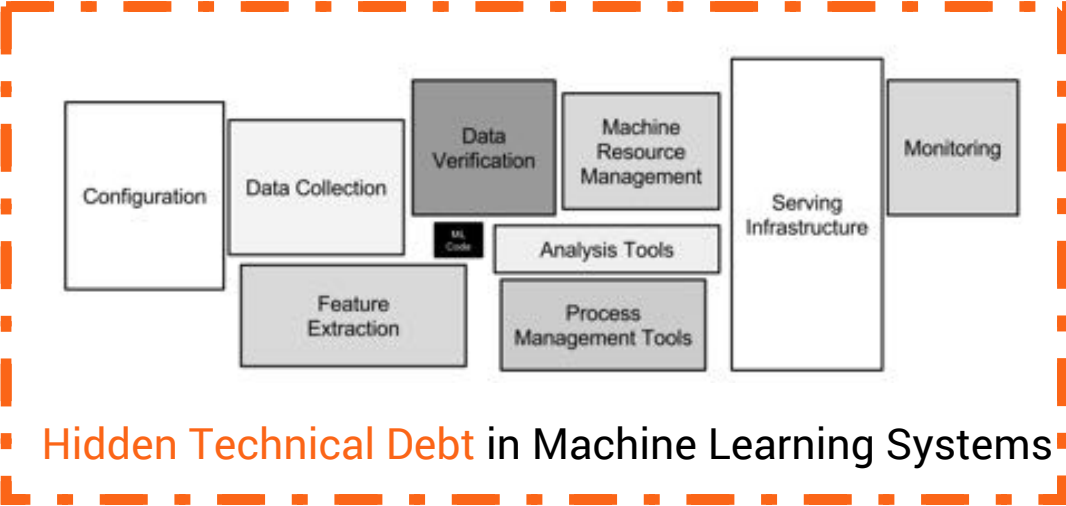
Replay-ability – A log never forgets!

Time



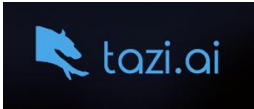
Different models with same data
Different ML Frameworks
AutoML compatible
A/B Testing

AutoML → No Data Scientist available for the ML Tasks?



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

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

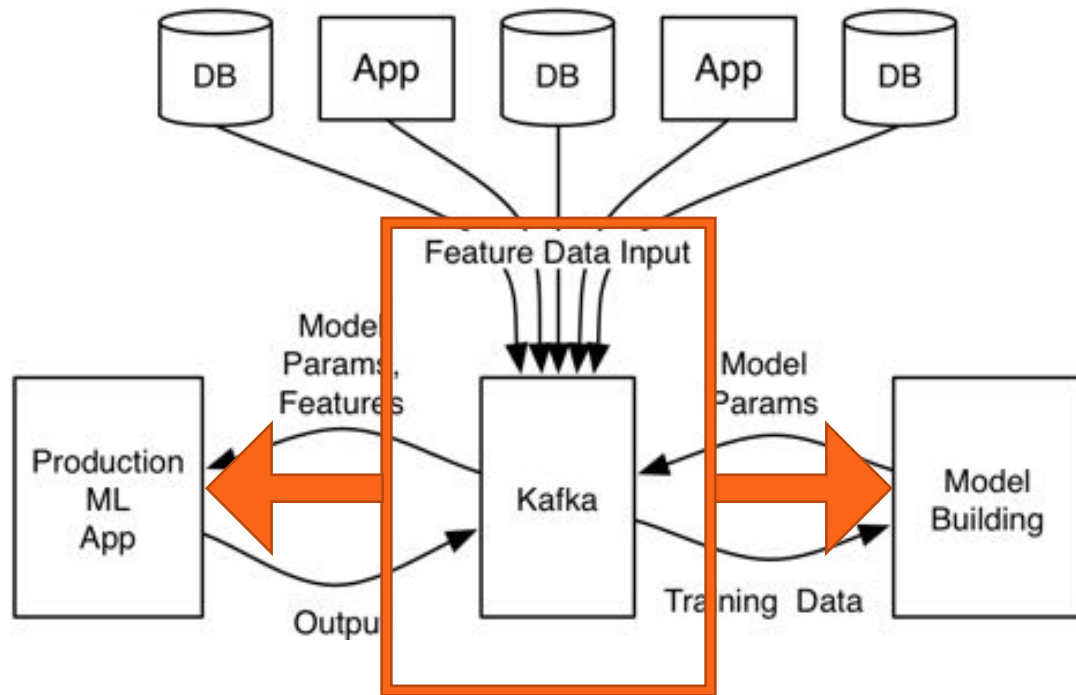


This talk does not focus on building models, but scalable infrastructure for ML

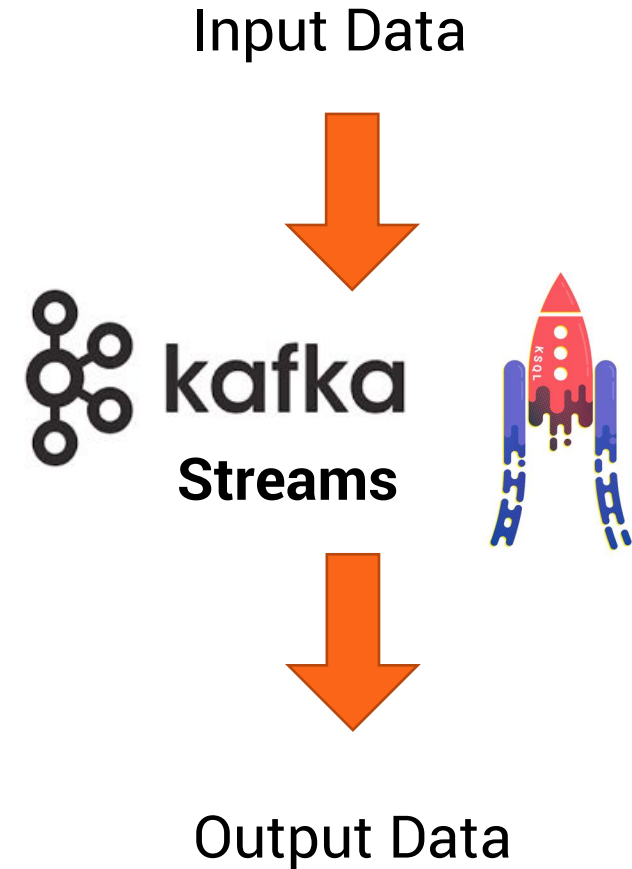
Agenda

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

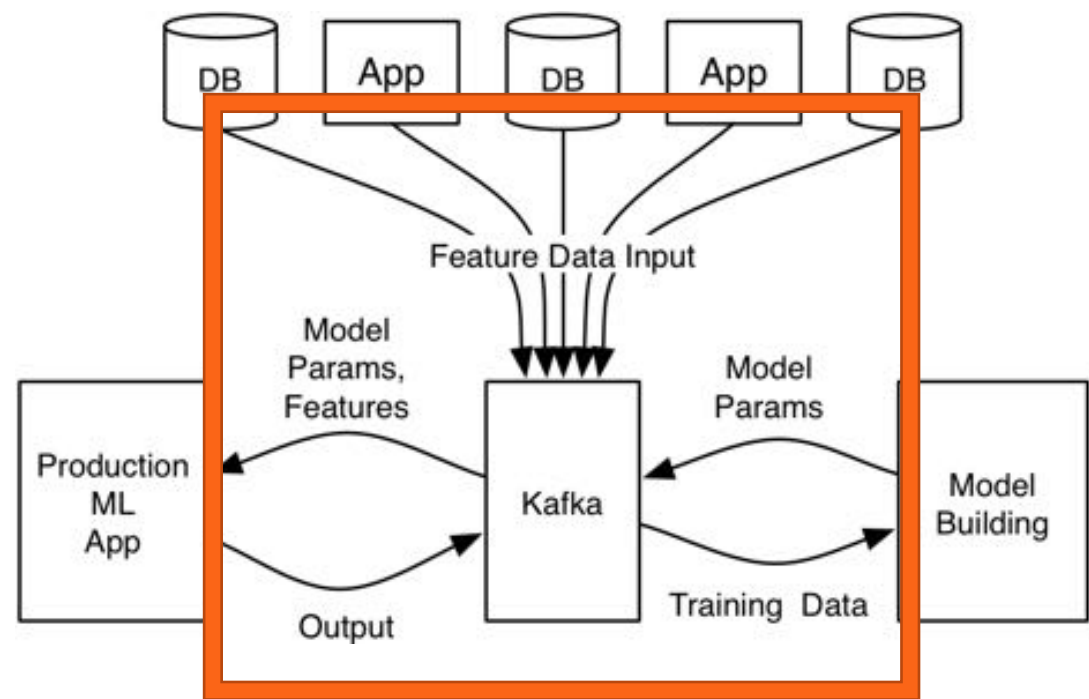
Kafka Streams / KSQL for Data Preprocessing



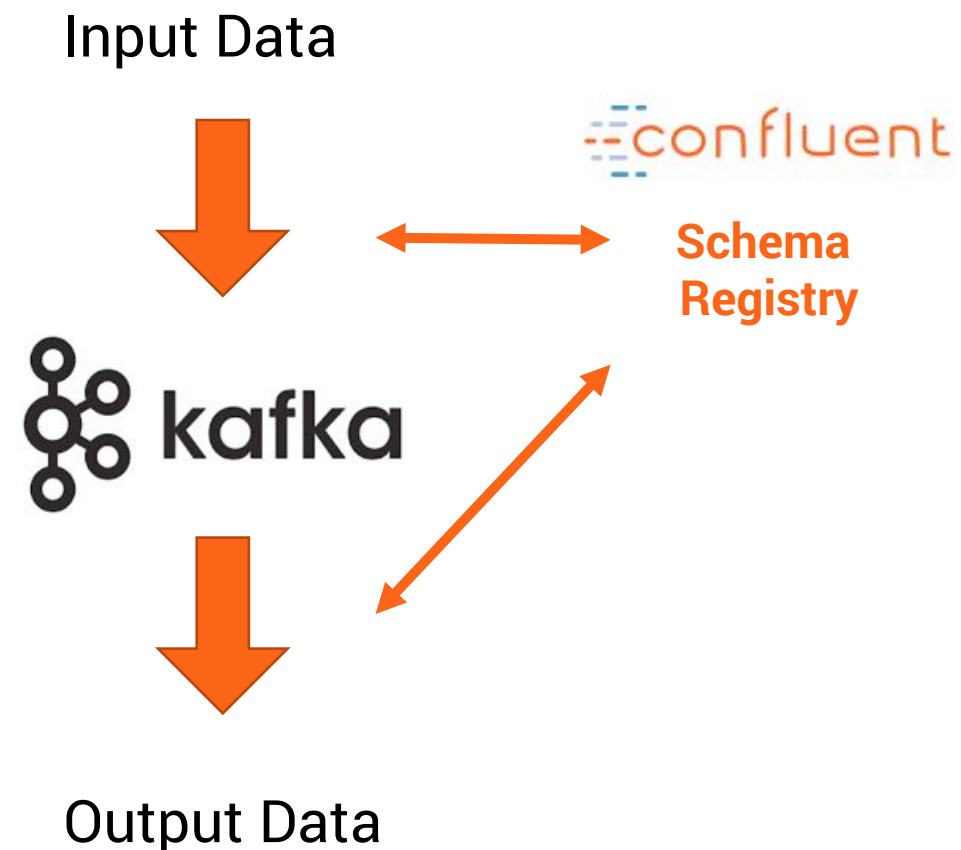
“Kafka benefits under the hood”
Streaming ETL
Same Pipeline for Training and Serving



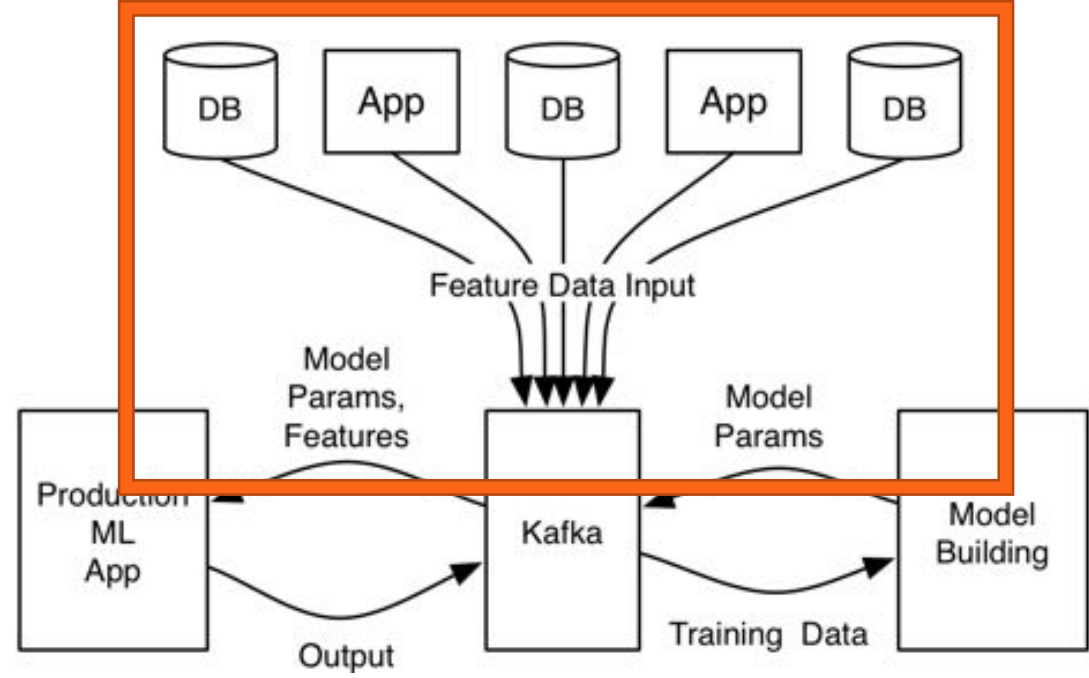
Confluent Schema Registry for Message Validation



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



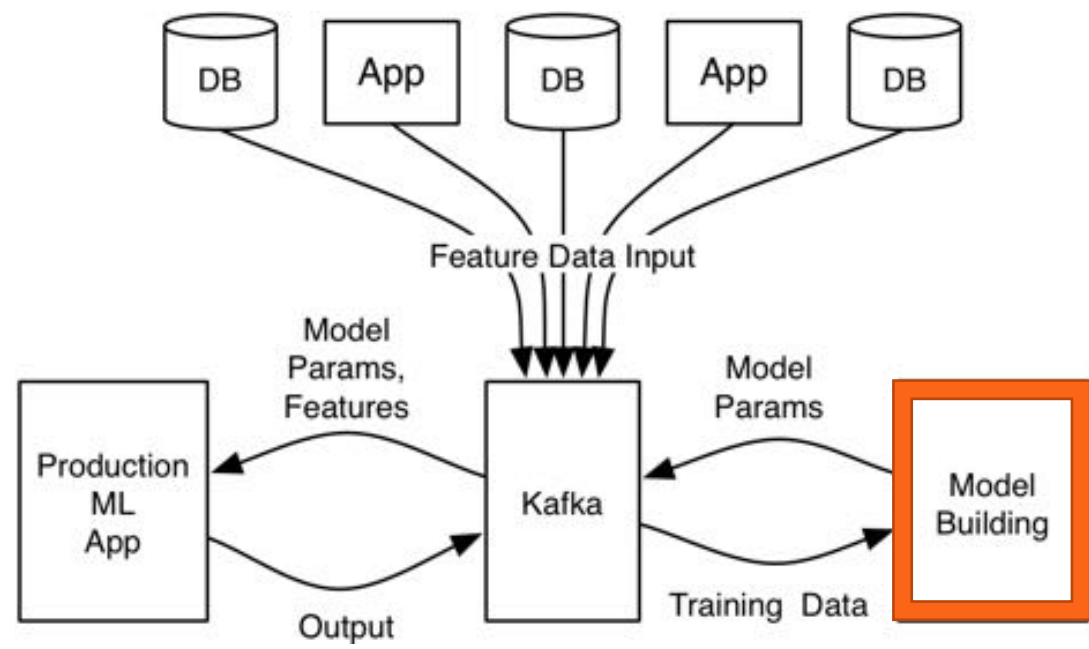
Kafka Connect for Data Ingestion



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

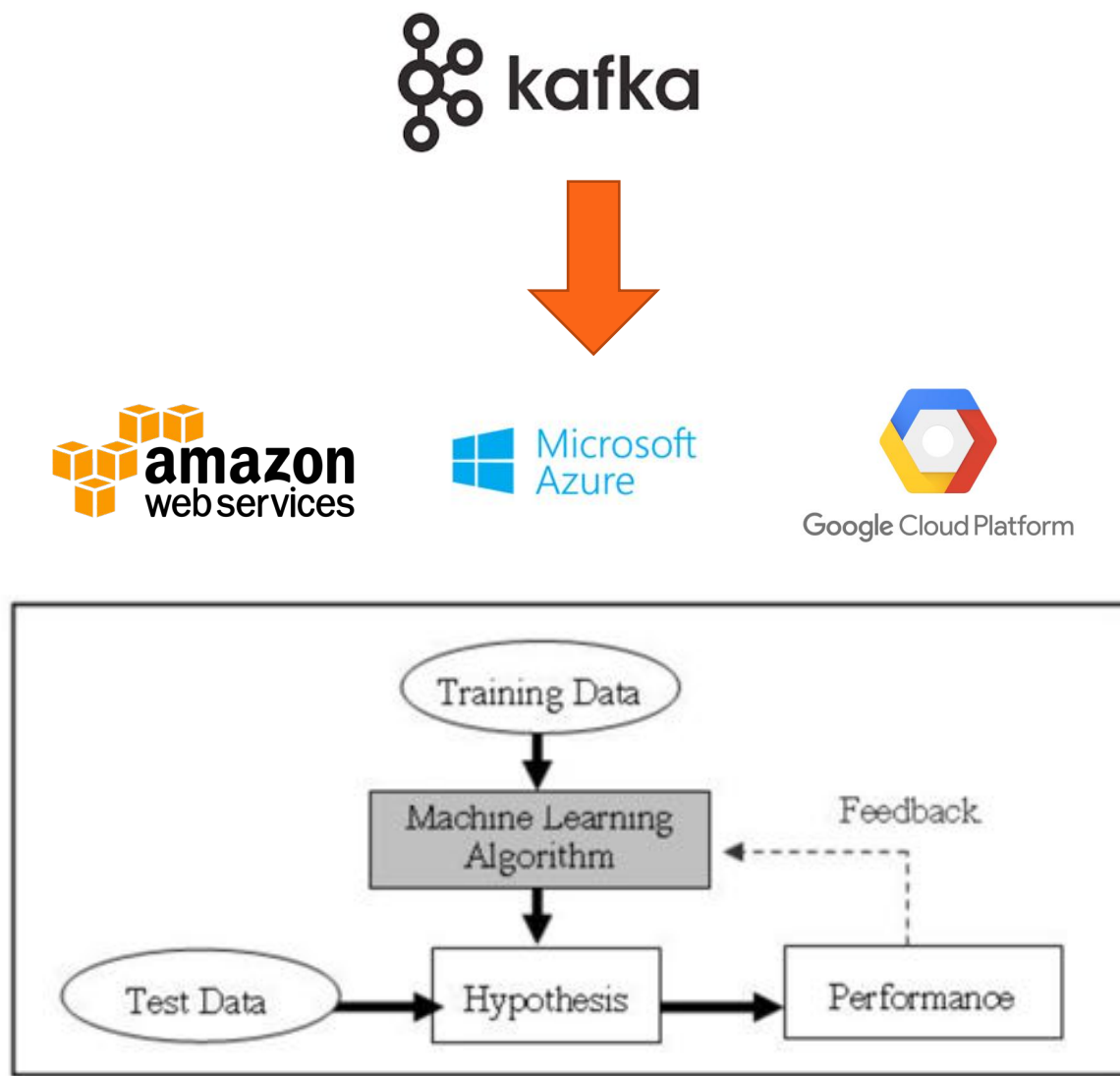


Model Training



Benefits of Public Cloud

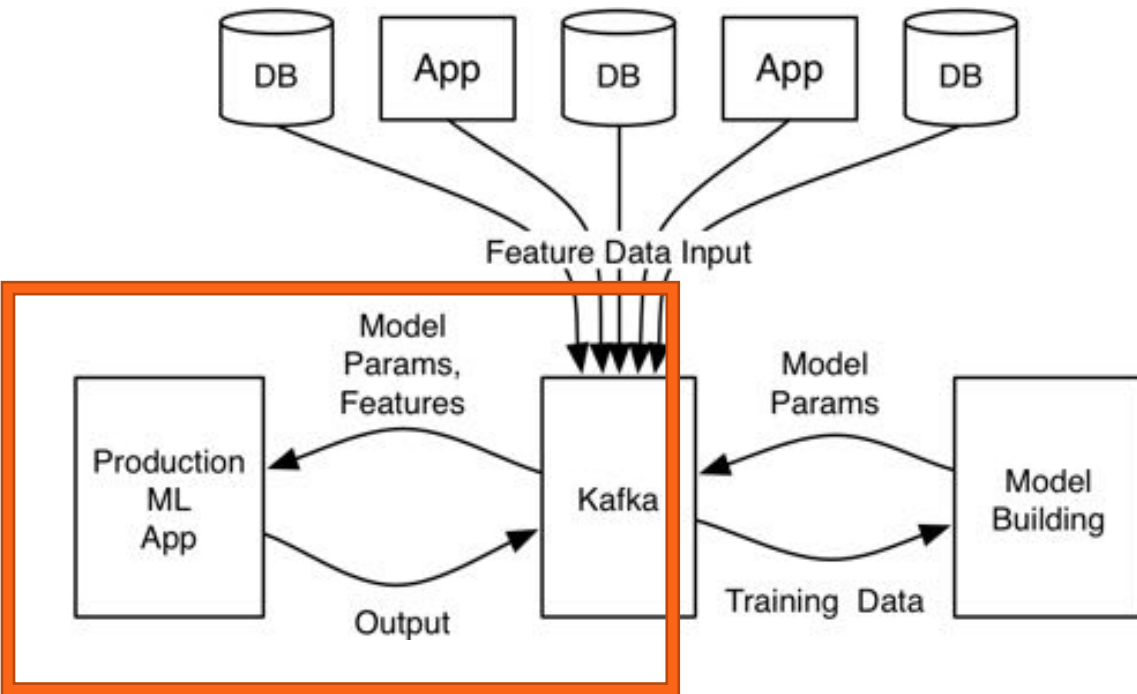
- Extreme Scale
- Dynamic Instances
- Special Hardware



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 for Model Training
- 4) Predictions in Real Time with Kafka Streams and KSQL**
- 5) DevOps and Monitoring of a Machine Learning Infrastructure

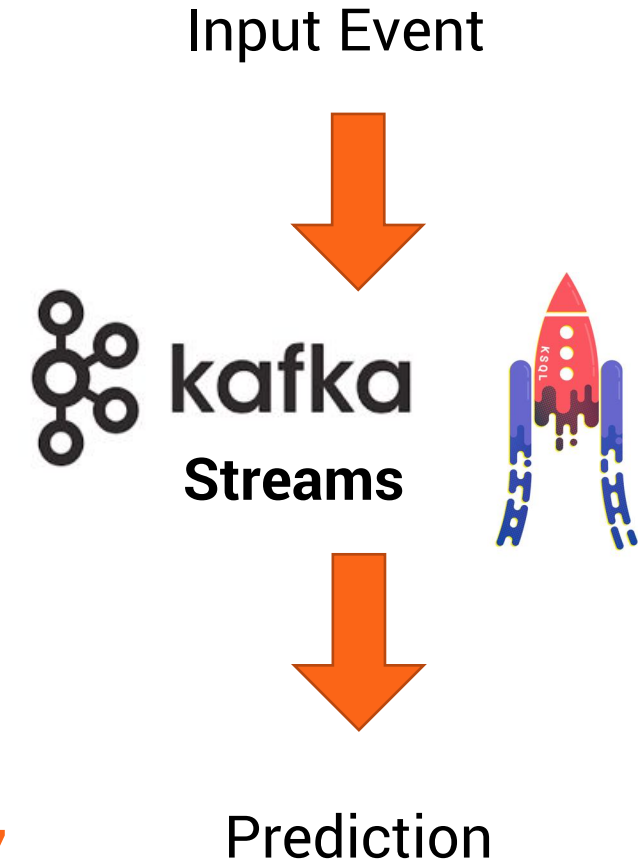
Model Serving / Inference / Deployment / Scoring



Kafka Streams
KSQL

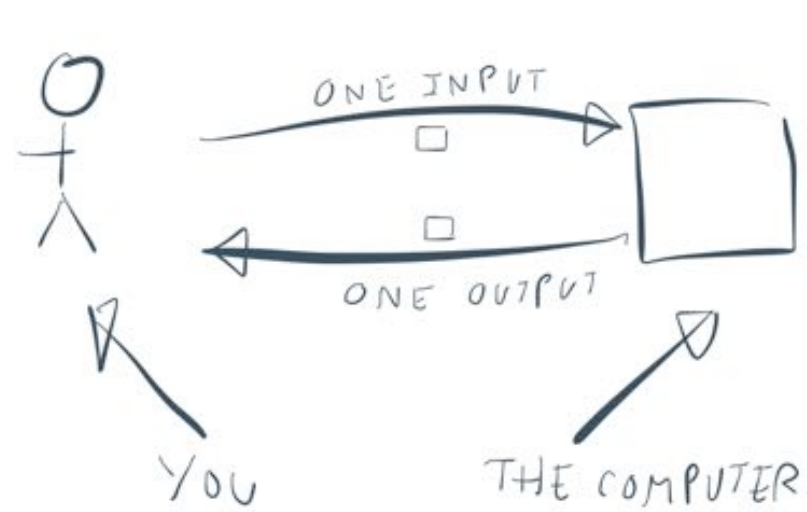
**“Kafka benefits under the hood”
Continuous Stream Processing**

**Reuse Preprocessing Logic from Ingestion Pipeline
Serving within the application (not via REST interface)
Predictions in real time**



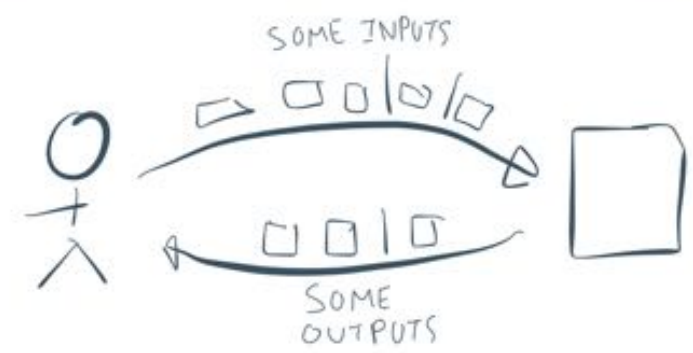
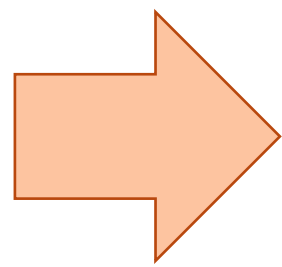
Stream Processing

REQUEST/RESPONSE



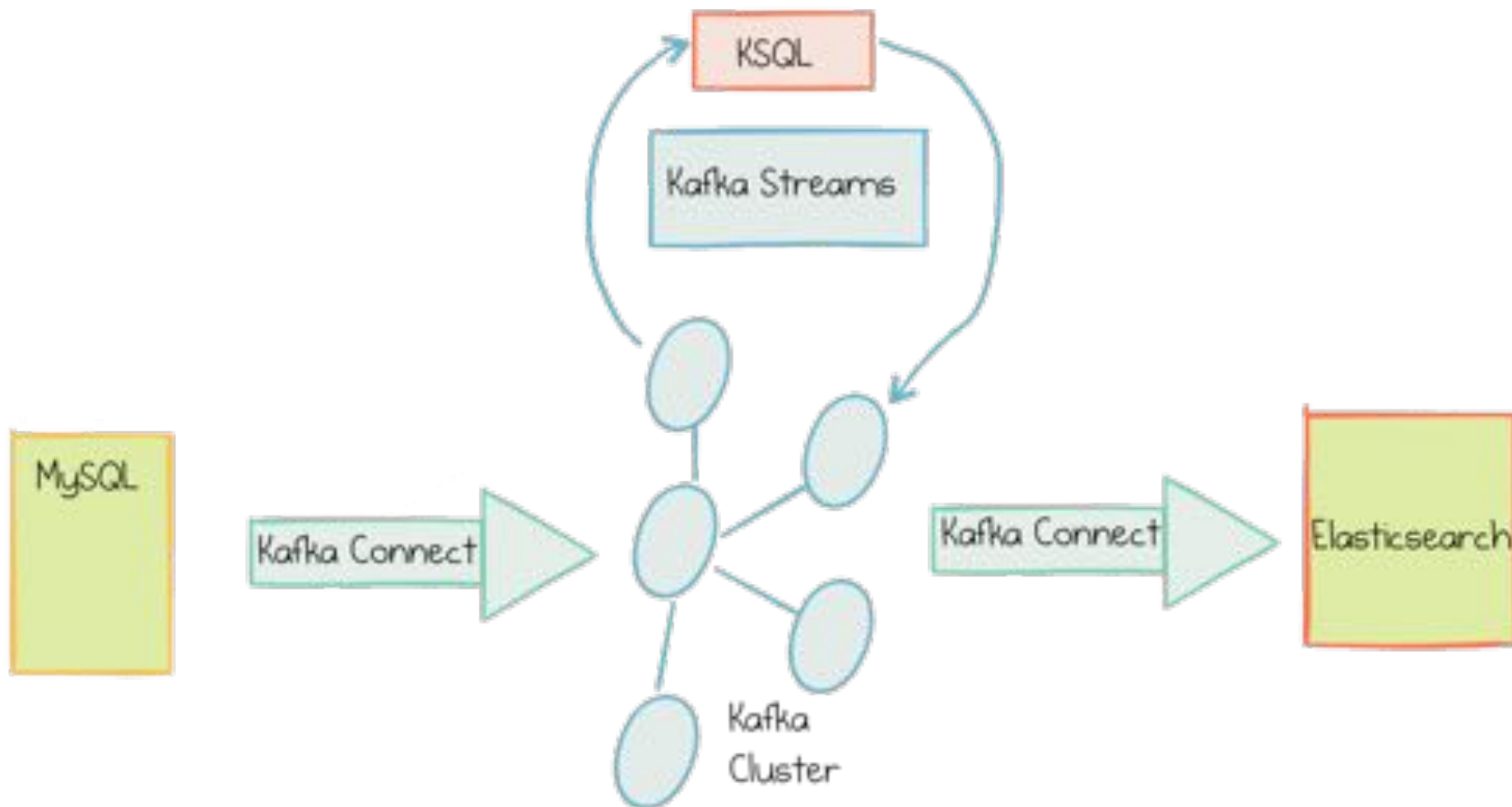
Data at Rest

STREAM PROCESSING



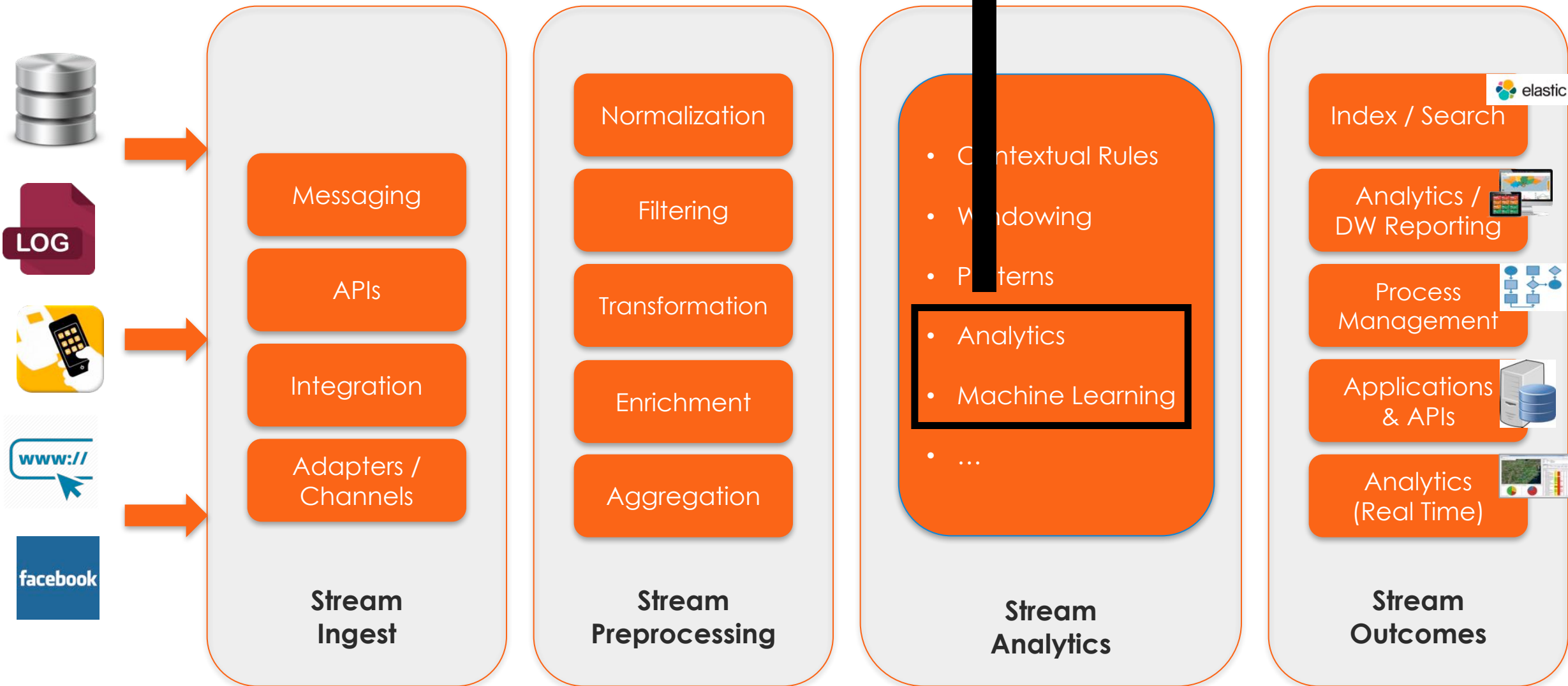
Data in Motion

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

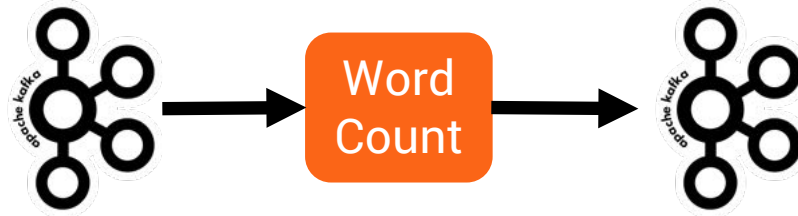


Stream Processing Pipeline

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)

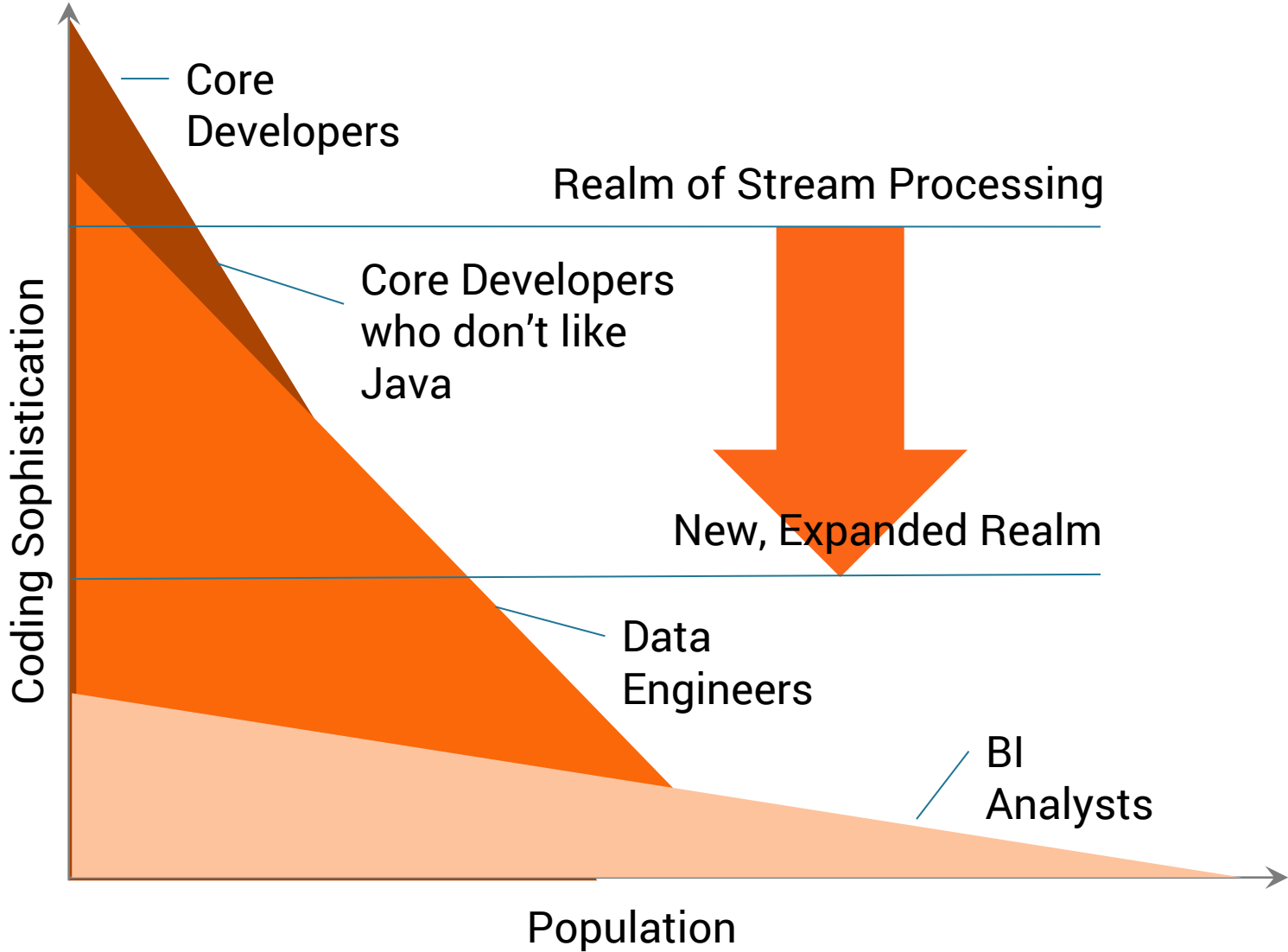
Start processing

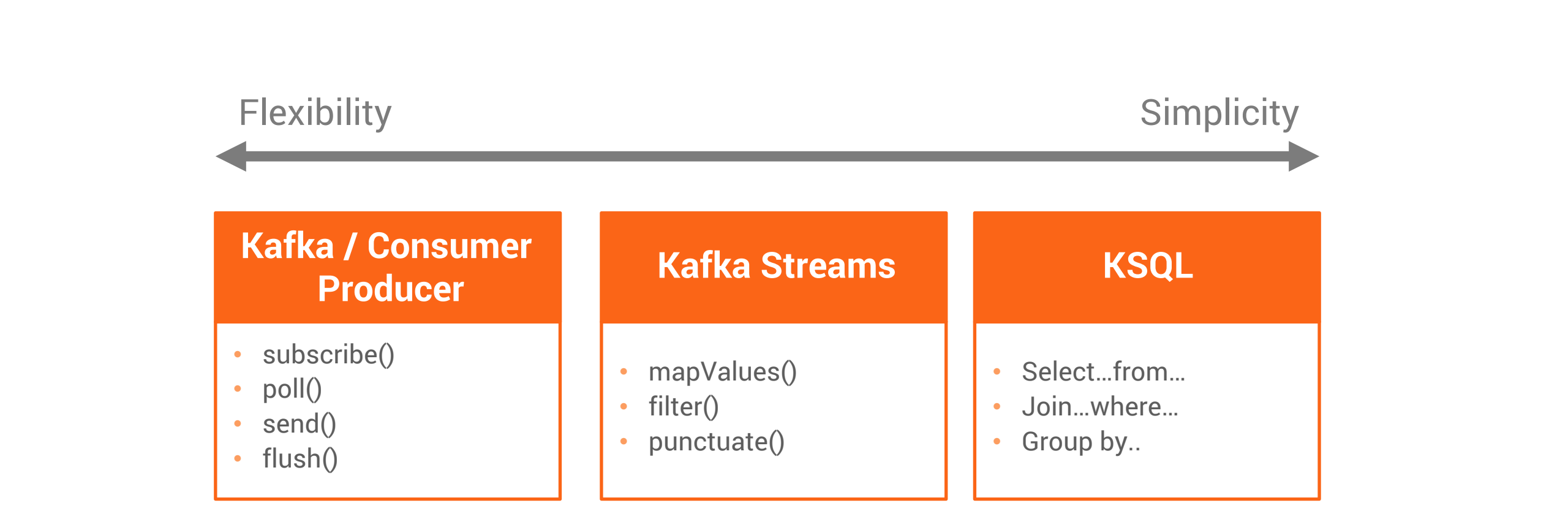
Why KSQL?

**Kafka
Streams**

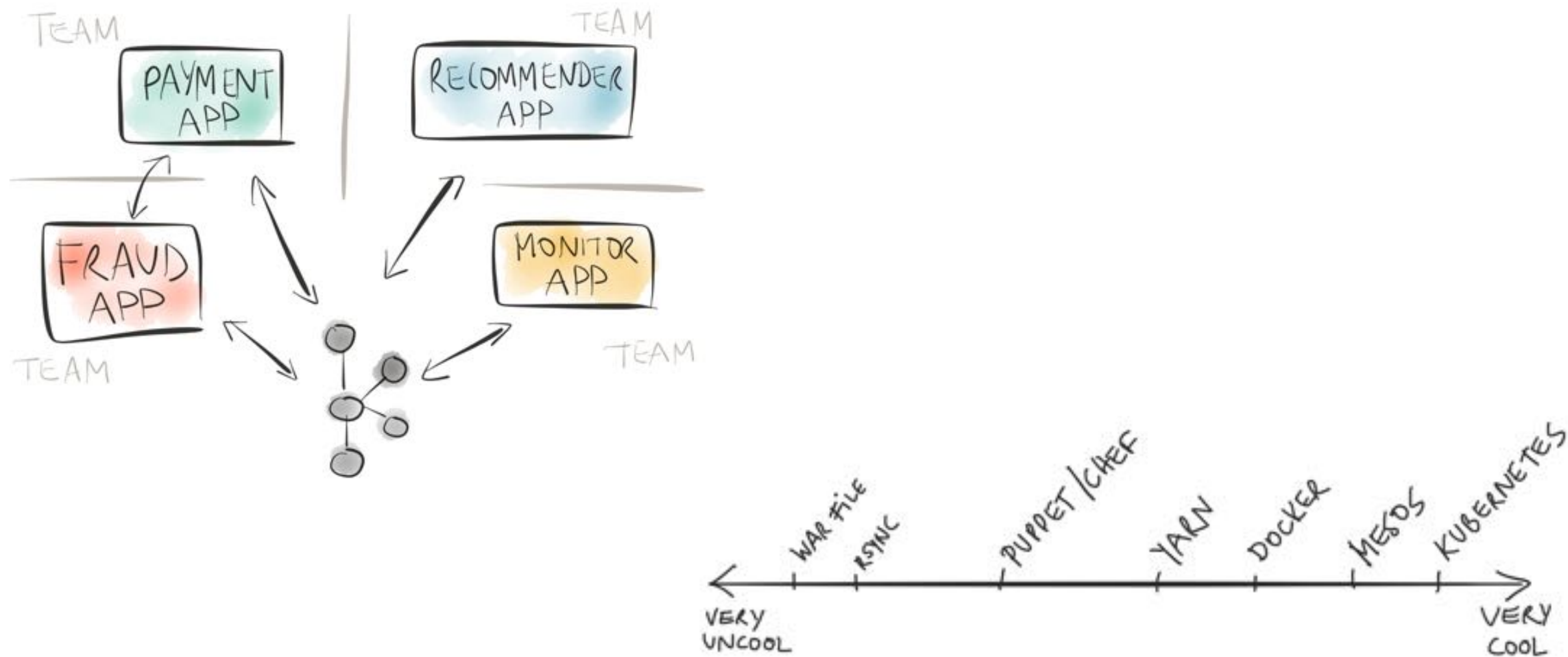


KSQL

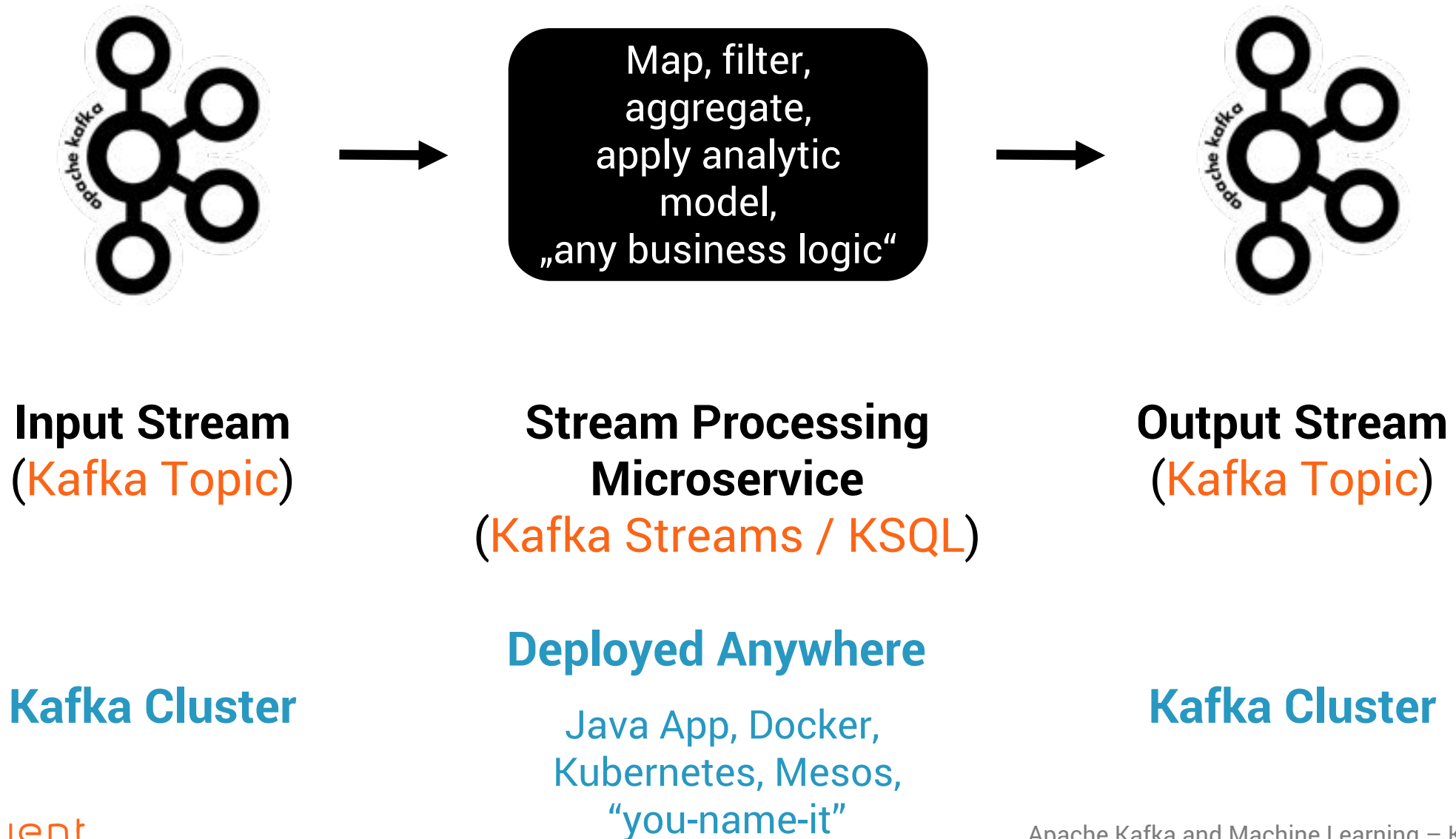




When to use Kafka Streams or KSQL for Stream Processing?



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.

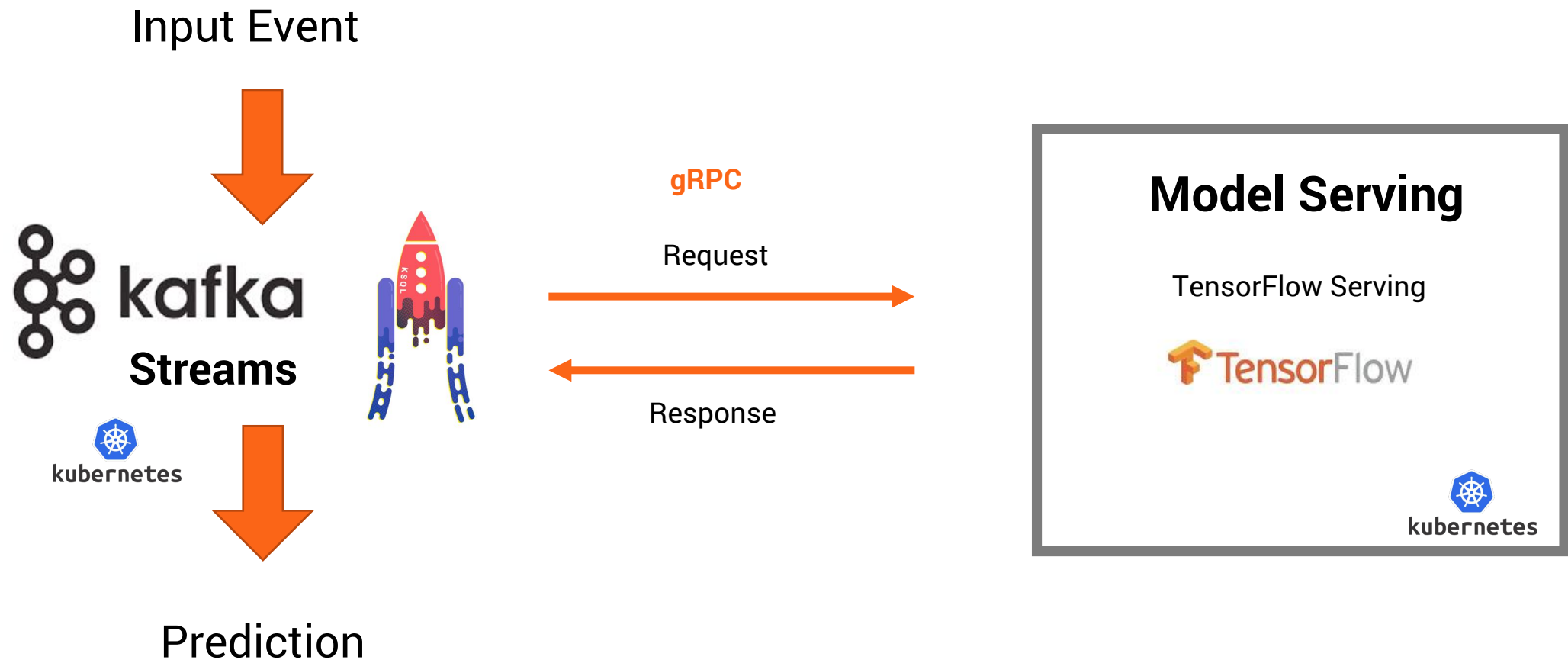


Ok.

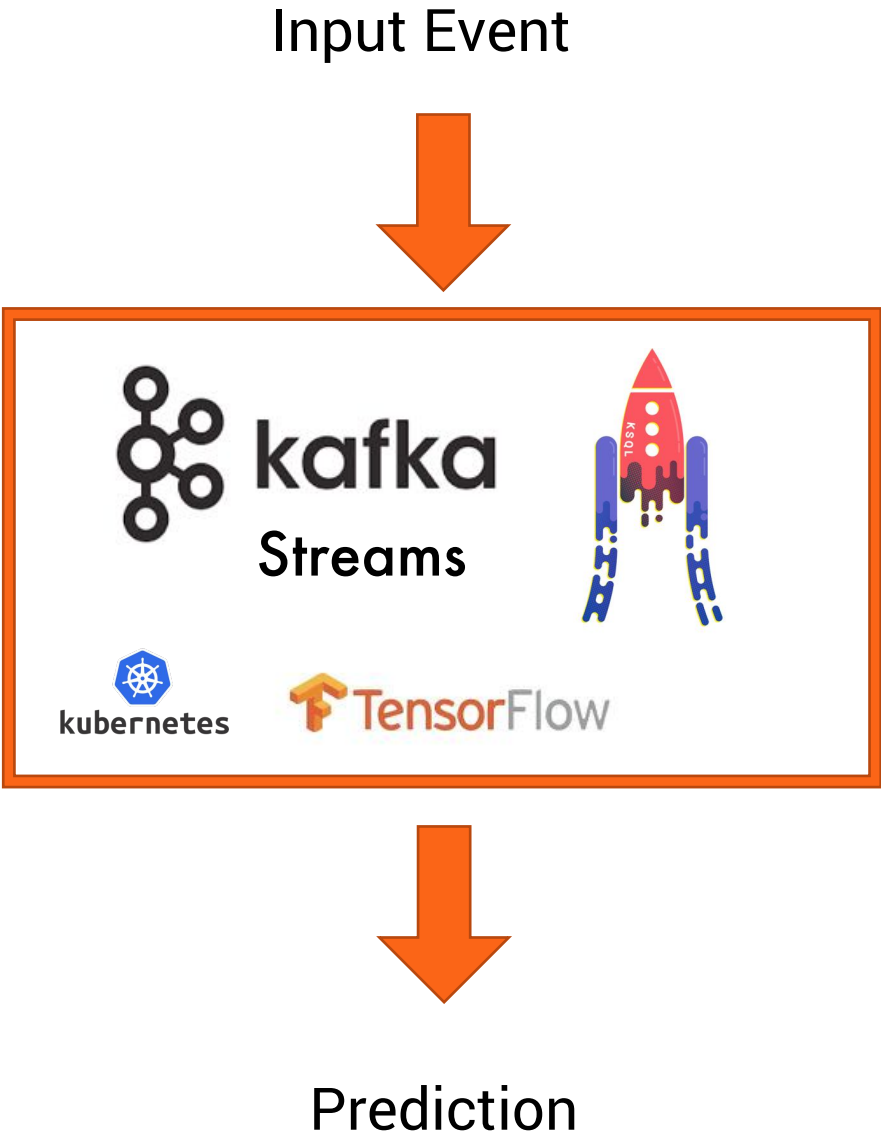


Ok.

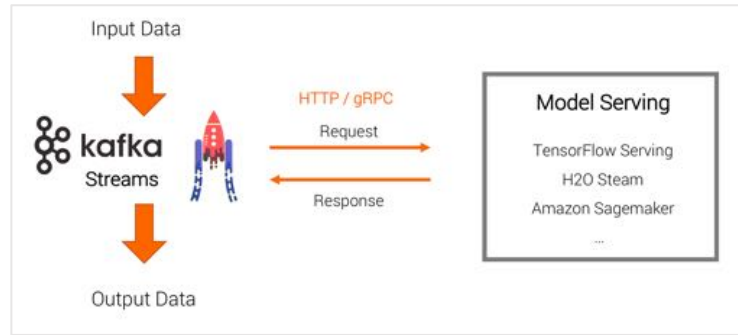
Option 1: gRPC communication to do model inference



Option 2: Model inference natively integrated into the App



Stream Processing vs. Request-Response for Model Serving



Pros of a Model Server:

- **Simple integration** with existing technologies and organizational processes
- **Easier to understand** if you come from non-streaming world
- **Later migration** to real streaming is also **possible**
- **Model management built-in** for different models, versioning and A/B testing

Cons (== Pros of Deployment in the Streaming App):

- **Worse latency** as remote call instead of local inference
- **No offline inference** (devices, edge processing, etc.)
- **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)

Live Demo – Deployment of a Trained Model

Use Case:

Airline Flight Delay Prediction

Machine Learning Algorithm:

Neural Network

built with H2O and TensorFlow

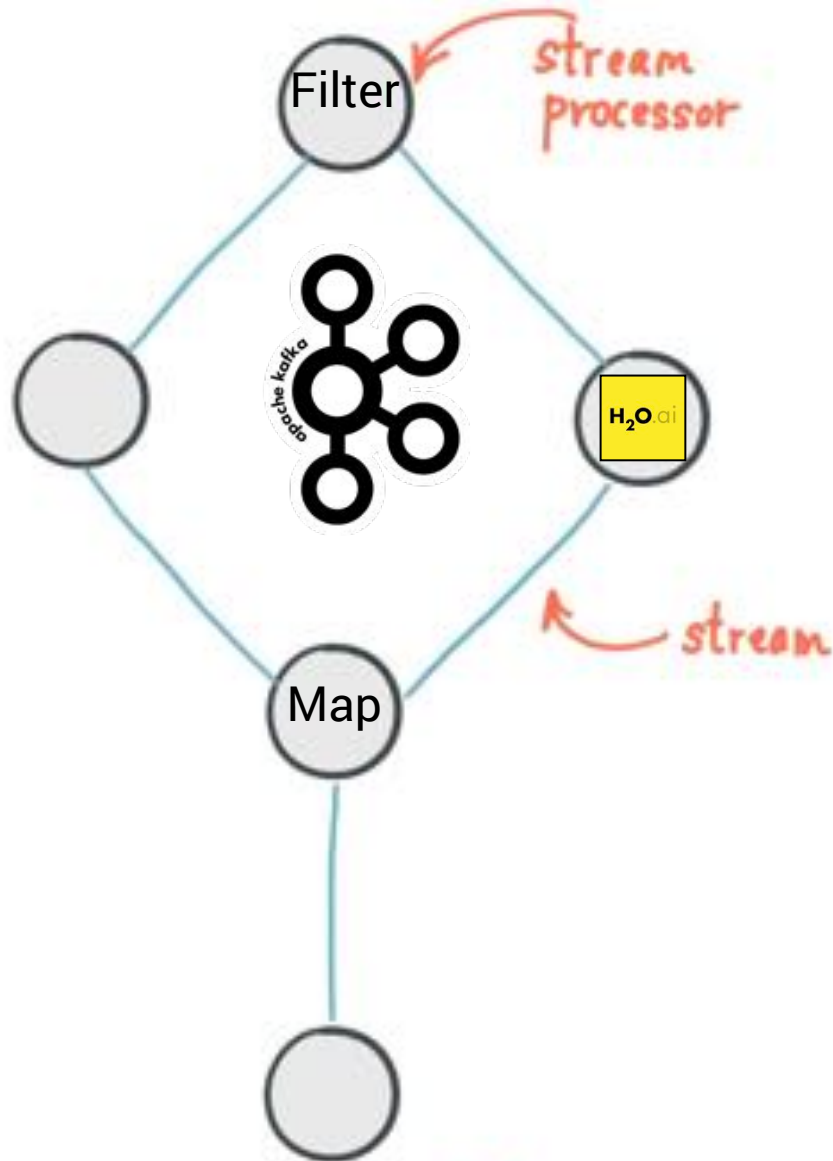


Streaming Platform:

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,ActualTime
        // 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,YES,YES

        if (value != null && !value.equals("")) {
            System.out.println("*****");
            System.out.println("Flight Input: " + value);

            String[] valuesArray = value.split(",");

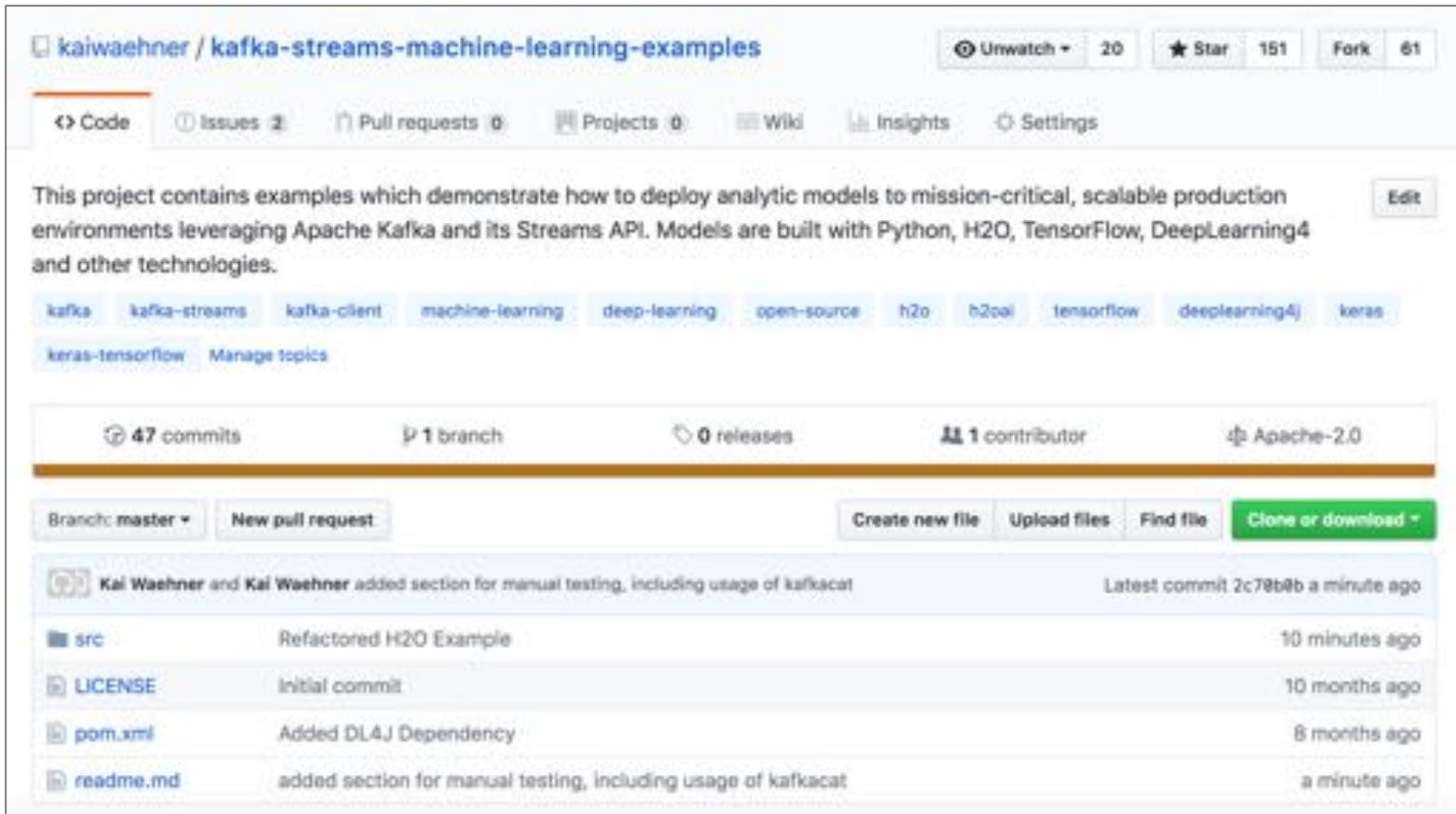
            RowData row = new RowData();
            row.put("Year", valuesArray[0]);
            row.put("Month", valuesArray[1]);
            row.put("DayOfMonth", valuesArray[2]);
            row.put("DayOfWeek", valuesArray[3]);
            row.put("CRSDepTime", valuesArray[4]);
            row.put("UniqueCarrier", 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.println("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();
```

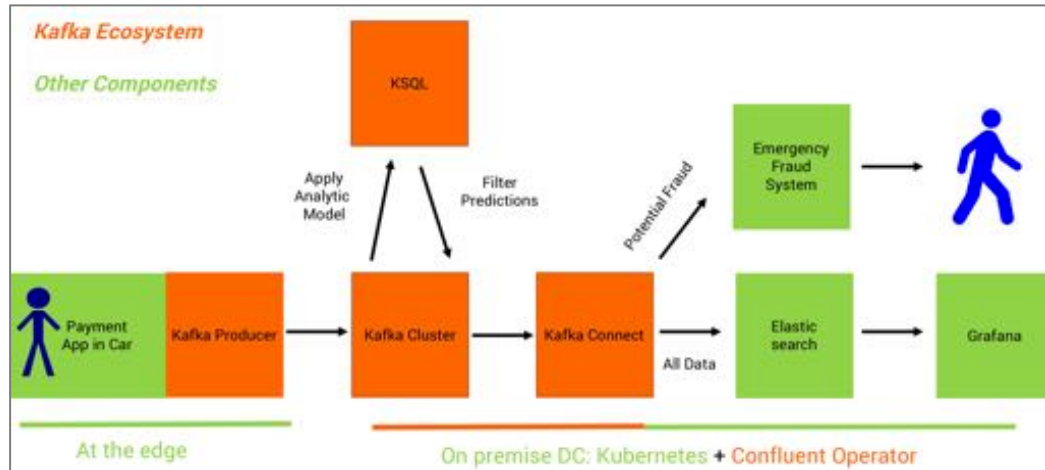
Github Examples: Kafka + Machine Learning



<https://www.confluent.io/blog/build-deploy-scalable-machine-learning-production-apache-kafka/>
<https://github.com/kaiwaehner/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 Fraud Detection



**“CREATE STREAM FraudDetection AS
SELECT payment_id, **applyFraudModel(payment_input)**
FROM payment_engine;”**



User Defined Function (UDF)



<https://www.confluent.io/blog/build-udf-udaf-ksql-5-0>
<https://github.com/kaiwaehner/ksql-udf-deep-learning-mqtt-iot>

Live Demo – Prebuilt Model Embedded in KSQL Function

Use Case:

Anomaly Detection
(Payment Fraud Detection)

Machine Learning Algorithm:
Autoencoder built with H2O

Streaming Platform:
Apache Kafka and KSQL



Github Examples: KSQL + Deep Learning

kaiwaehner / ksql-udf-deep-learning-mqtt-iot

Unwatch

8

Star

25

Fork

8

Code

Issues

Pull requests

Projects

Wiki

Insights

Settings

Deep Learning UDF for KSQL for Streaming Anomaly Detection of MQTT IoT Sensor Data

kafka

kafka-connect

kafka-client

confluent

confluent-platform

open-source

deep-learning

machine-learning

tensorflow

h2oai

java

ksql

mqtt

Manage topics

30 commits

1 branch

0 releases

2 contributors

Apache-2.0

Branch: master

New pull request

Create new file

Upload files

Find file

Clone or download

kaiwaehner

Updated requirements: Confluent Platform needs to be at least version...

Latest commit 6a86378 5 days ago

pictures

Added picture with MQTT / Kafka architecture

27 days ago

src/main

Removed Kafka Streams code (not needed for the UDF)

27 days ago

LICENSE

Initial commit

a month ago

README.md

Updated requirements: Confluent Platform needs to be at least version...

5 days ago

live-demo.adoc

Updated README to clarify steps for running demo

6 days ago

pom.xml

Removed Kafka Streams dependency (not needed for the UDF)

27 days ago

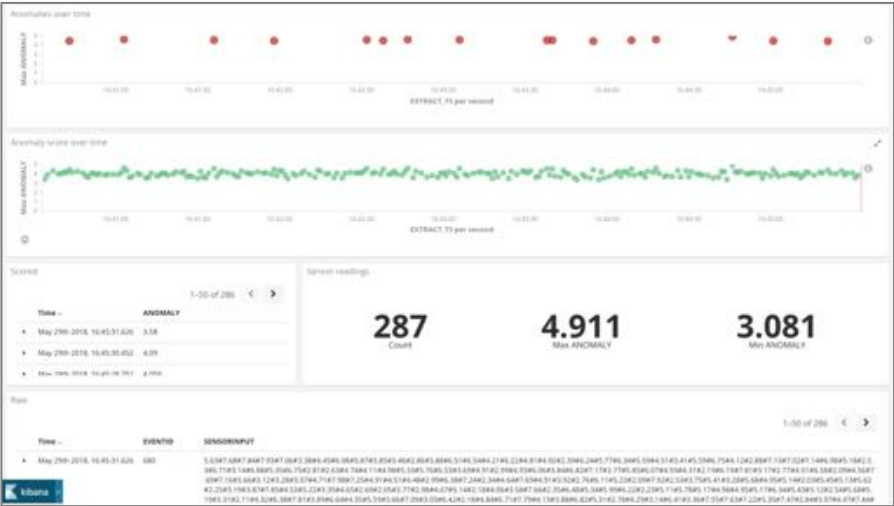
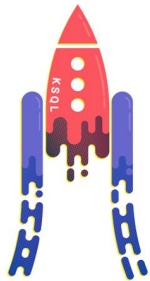
sensor_generator.sh

Updated README to clarify steps for running demo

6 days ago



+ Kafka Connect
+ Elasticsearch

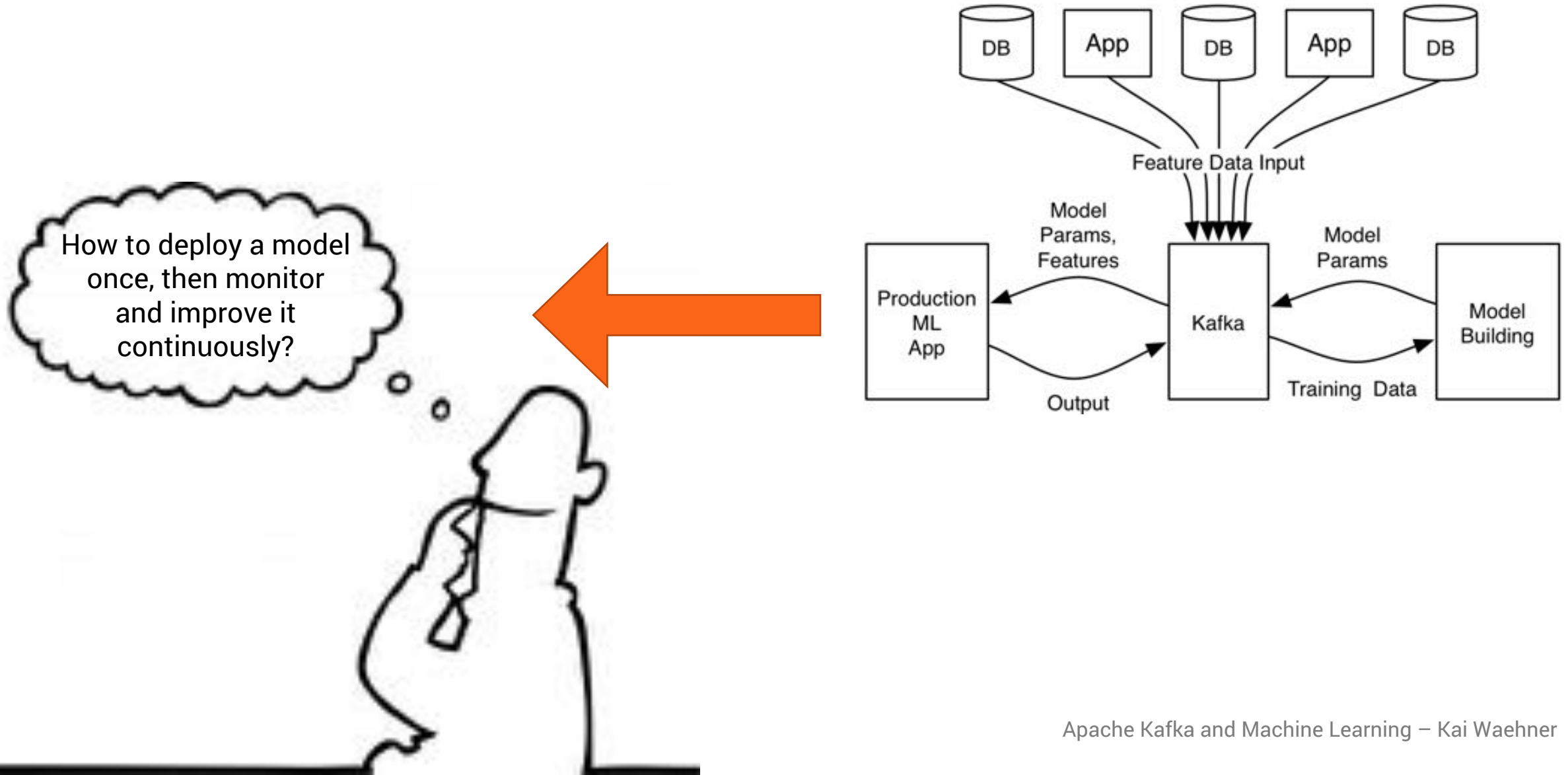


<https://www.confluent.io/blog/build-deploy-scalable-machine-learning-production-apache-kafka/>
<https://github.com/kaiwaehner/ksql-udf-deep-learning-mqtt-iot>
<https://github.com/kaiwaehner/ksql-fork-with-deep-learning-function>

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 for Model Training
- 4) Predictions in Real Time with Kafka Streams and KSQL
- 5) DevOps and Monitoring of a Machine Learning Infrastructure**

Automated Model Improvement with Apache Kafka and Kafka Streams



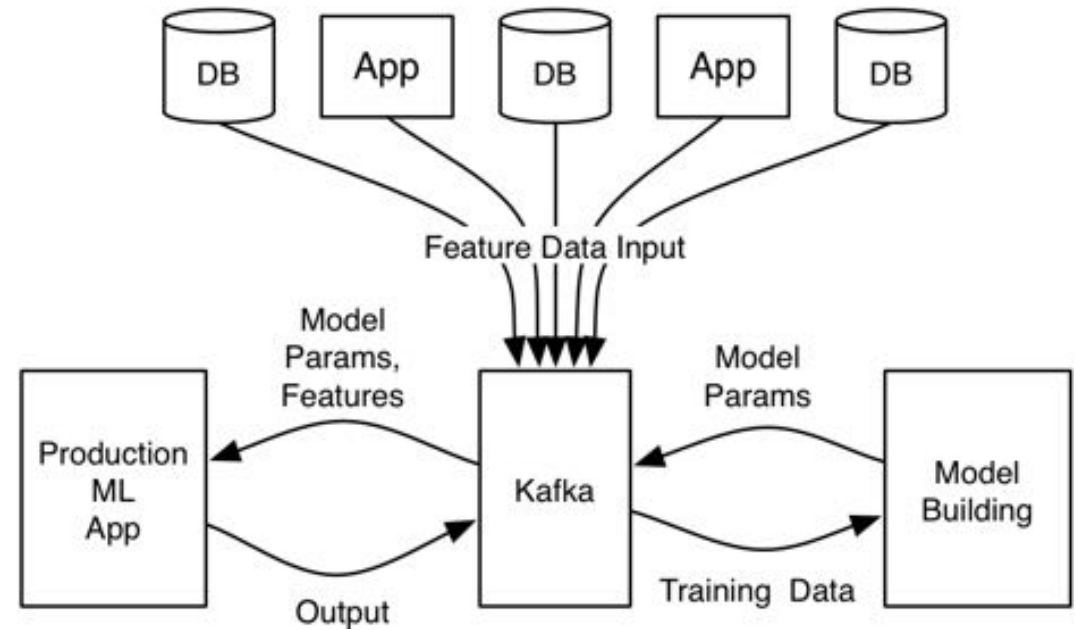
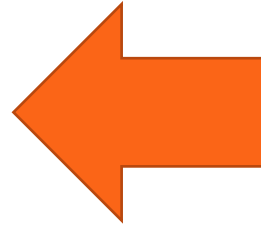
Automated Model Improvement with Apache Kafka and Kafka Streams



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

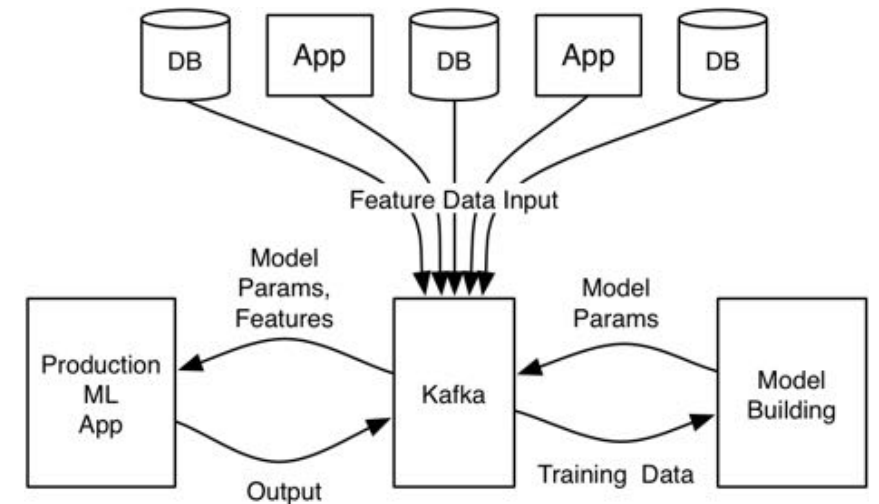
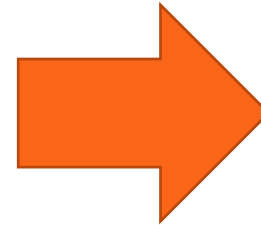
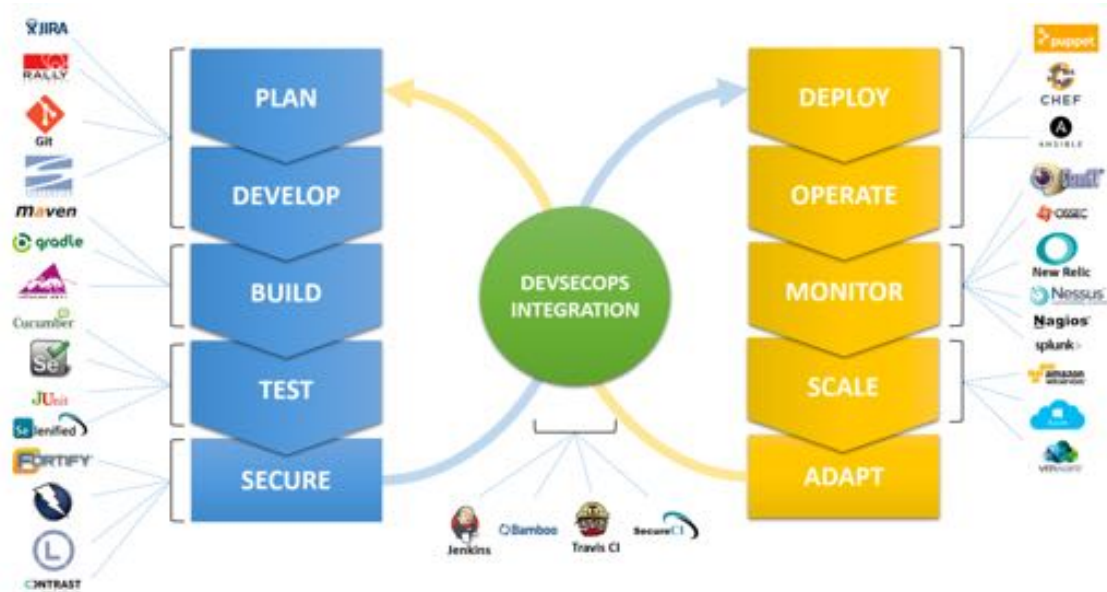
→ Do it only when feasible!



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 (e.g. via a Kafka Topic)



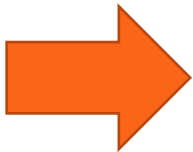
Kubernetes – The Winner of the Container and DevOps Wars!

 **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


RED HAT[®] OPENSIFT


Microsoft Azure

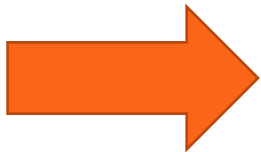
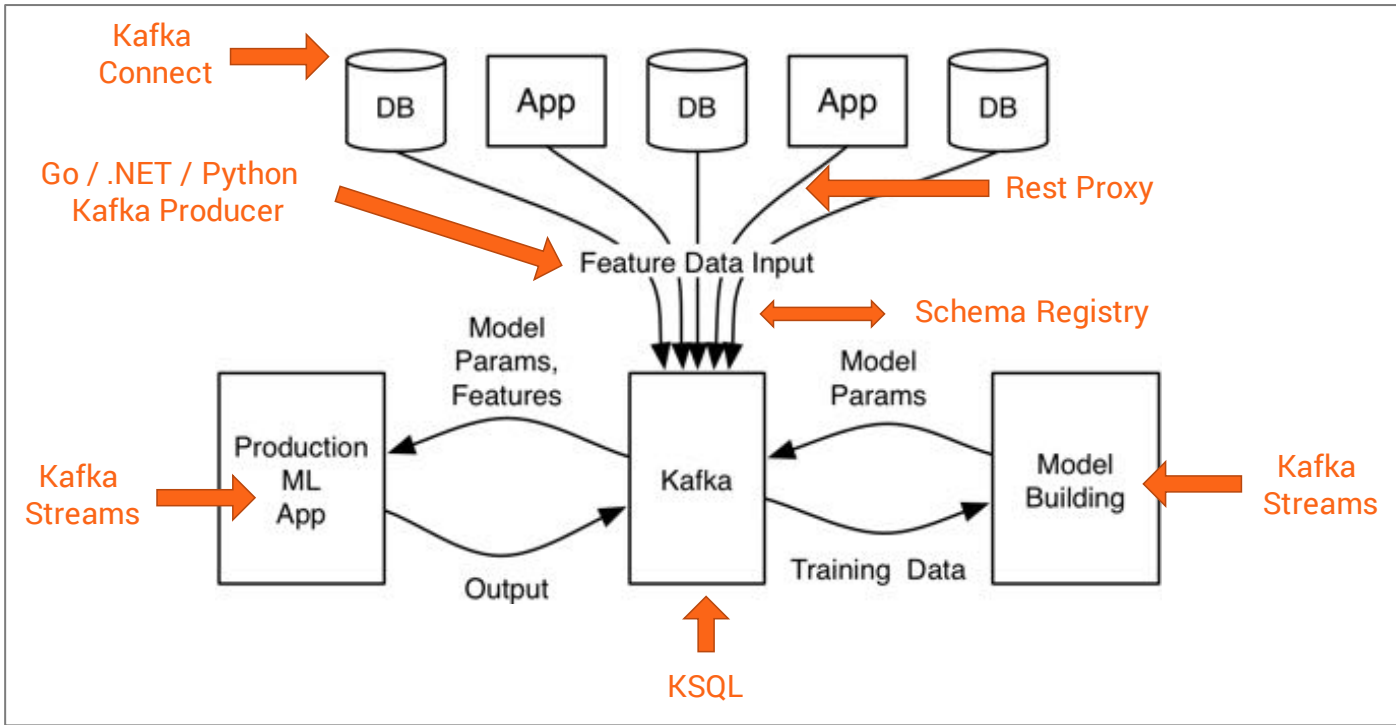

CLOUDFOUNDRY


docker
Docker, Inc


MESOSPHERE

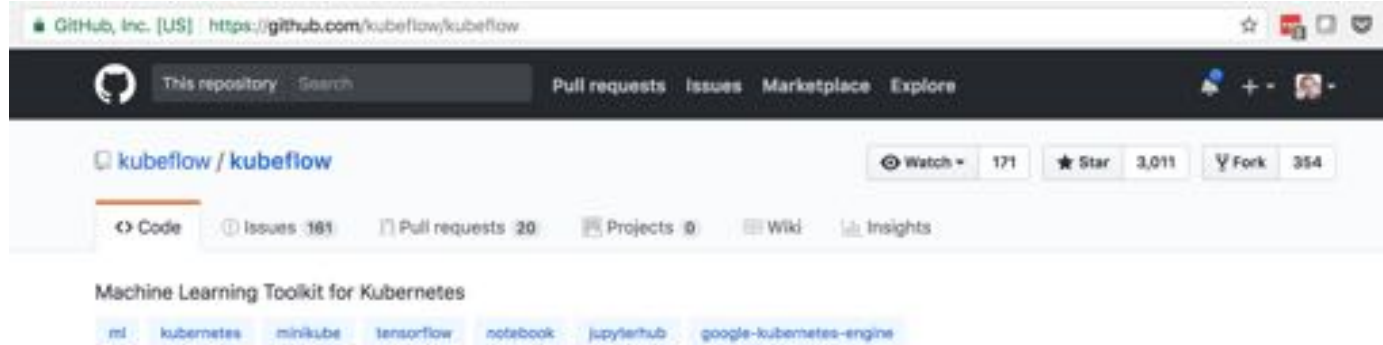
<https://www.infoworld.com/article/3118345/cloud-computing/why-kubernetes-is-winning-the-container-war.html>
<http://techgenix.com/year-of-kubernetes/>

Monitoring the Infrastructure for Machine Learning



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

Kubernetes Deployment of ML Workflows

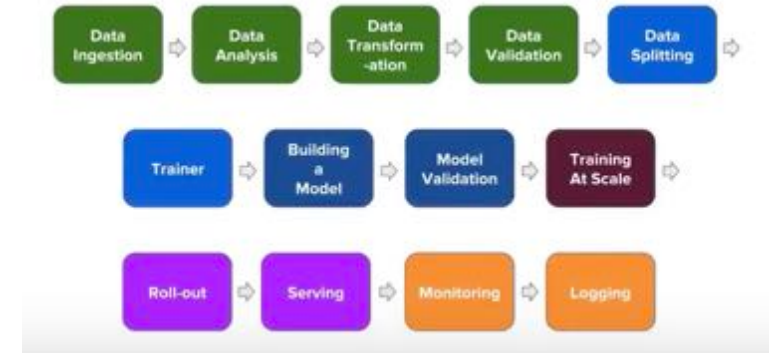


Kubeflow

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...



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)



Questions? Feedback? Please contact me!

Kai Waehner

Technology Evangelist

kontakt@kai-waehner.de

[@KaiWaehner](#)

www.kai-waehner.de

www.confluent.io

[LinkedIn](#)

