



UNIVERSITAT DE  
BARCELONA

# Master in Fundamental Principles of Data Science

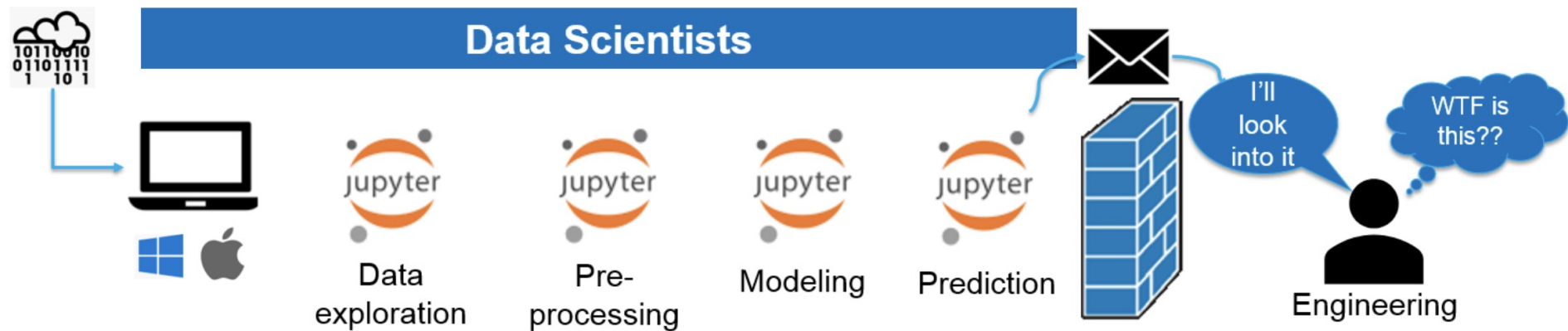
Dr Rohit Kumar



UNIVERSITAT DE  
BARCELONA

# ML Models in Production

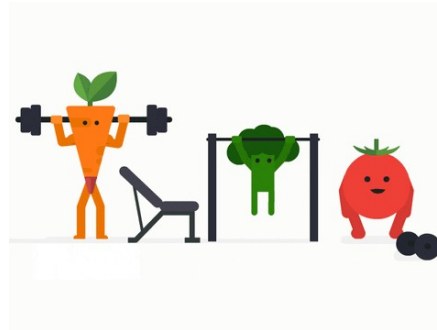
# Typical flow



# Data Science Life Cycle



Data wrangling



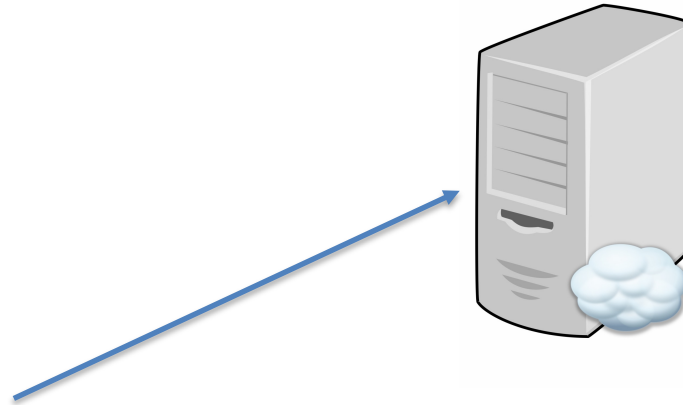
Training



Prediction

# Training

## 1. One Off training



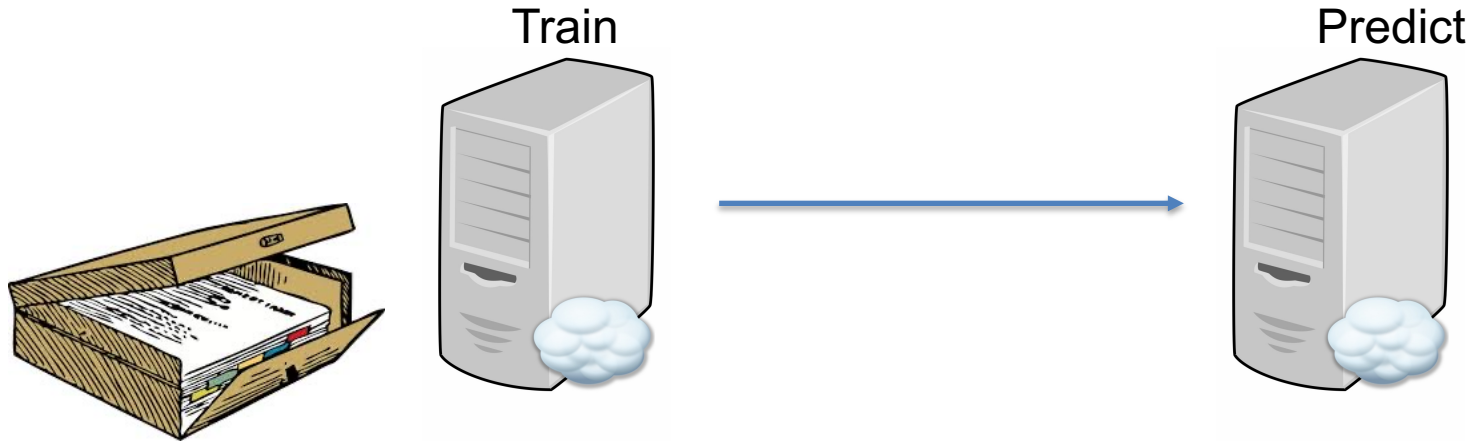
Where do you save your notebook/code?  
Where and how do you save your models?

Models don't necessarily need to be continuously trained in order to be pushed to production. Quite often a model can be just trained ad-hoc by a data-scientist, and pushed to production until its performance deteriorates enough that they are called upon to refresh it.

Example: Image recognition, self driving cars

# Training

## 2. Batch training (continuous learning system)



Not simply retraining the model weights !!

But also feature processing, feature selection, model selections and parameter optimization.

# Training

## 3. Real time training

Train



Incremental Training using “Online Machine Learning Algorithms”

# AutoML

**AutoWEKA** is an approach for the simultaneous selection of a machine learning algorithm and its hyperparameters; combined with the WEKA package it automatically yields good models for a wide variety of data sets.

**TPOT** is a data-science assistant which optimizes machine learning pipelines using genetic programming. (Good for Python users)

**H2O AutoML** provides automated model selection and ensembling for the H2O machine learning and data analytics platform. (Nice for Java users)

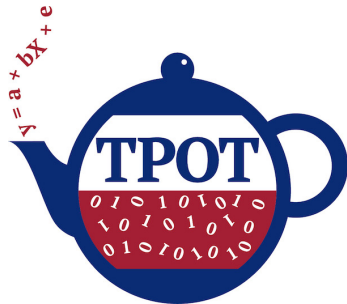
**TransmogrifAI** is an AutoML library running on top of Spark.





# Airflow

To manage different data workflows

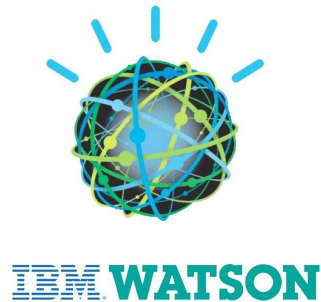


To automate Machine Learning

Ofcourse there are easy to use fancy platforms !!  
(if you have money and do not want to get your hands too dirty !!)



Azure ML



# Saving ML Models

In Python world:

- 1) Pickle:
- 2) Joblib: significantly faster on large numpy arrays

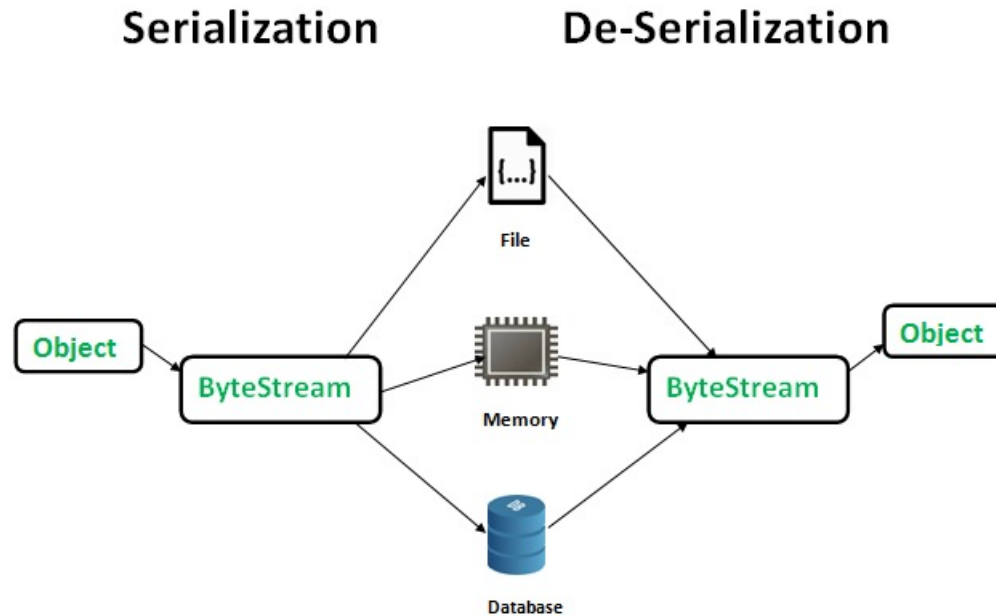
In R world:

- 1) saveRDS/readRDS – can be assigned to a new model variable
- 2) save/load – returns the same model variable

In Java (H2O world):

- 1) Pojo : Plain old java object
- 2) Mojo: Model Object, Optimized

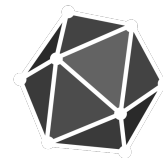
# Saving ML Models



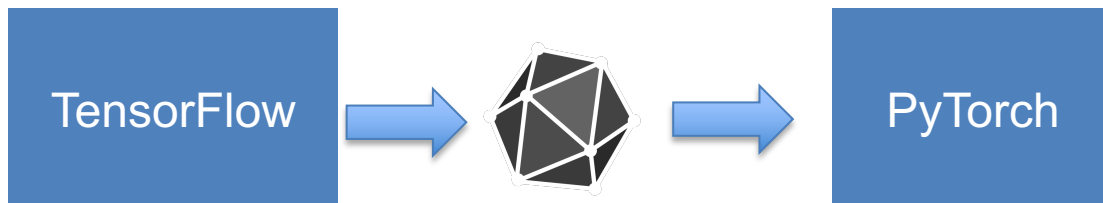
**“Change**, which is the only constant in life, can be **secured and facilitated** by standardization.”

# Saving ML Models

**ONNX** (Open Neural Network Exchange) format: developed by Facebook and Microsoft in collaboration it is an open-source standard for representing deep learning models that will let you transfer them between CNTK, Caffe2 and PyTorch.



Link: <https://github.com/onnx/tutorials>



# Saving ML Models

Framework / Tool	Installation	Tutorial
Caffe	<a href="#">apple/coremltools</a> and <a href="#">onnx/onnxmltools</a>	<a href="#">Example</a>
Caffe2	<a href="#">part of caffe2 package</a>	<a href="#">Example</a>
Chainer	<a href="#">chainer/onnx-chainer</a>	<a href="#">Example</a>
Cognitive Toolkit (CNTK)	<a href="#">built-in</a>	<a href="#">Example</a>
CoreML (Apple)	<a href="#">onnx/onnxmltools</a>	<a href="#">Example</a>
Keras	<a href="#">onnx/keras-onnx</a>	<a href="#">Example</a>
LibSVM	<a href="#">onnx/onnxmltools</a>	<a href="#">Example</a>
LightGBM	<a href="#">onnx/onnxmltools</a>	<a href="#">Example</a>
MATLAB	<a href="#">Deep Learning Toolbox</a>	<a href="#">Example</a>
ML.NET	<a href="#">built-in</a>	<a href="#">Example</a>
MXNet (Apache)	<a href="#">part of mxnet package</a> <a href="#">docs</a> <a href="#">github</a>	<a href="#">Example</a>
PyTorch	<a href="#">part of pytorch package</a>	<a href="#">Example1</a> , <a href="#">Example2</a> , <a href="#">export for Windows ML</a> , <a href="#">Extending support</a>
SciKit-Learn	<a href="#">onnx/sklearn-onnx</a>	<a href="#">Example</a>
SINGA (Apache) - Github (experimental)	<a href="#">built-in</a>	<a href="#">Example</a>
TensorFlow	<a href="#">onnx/tensorflow-onnx</a>	<a href="#">Examples</a>

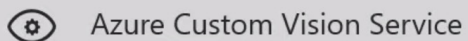


## Create

### Frameworks



### Services



ONNX Model

## Deploy

### Azure

Azure Machine Learning services

Ubuntu VM

Windows Server 2019 VM

### Devices

Windows devices

IoT Devices

Other devices (iOS, etc.)

# Saving ML Models

**PMML** (Predictive model markup language): Developed by Data Mining Group (DMG) a consortium of commercial and open-source data mining companies and supported by IBM, AWS and Google.



**PFA** (Portable Format for Analytics): it is more flexible (using JSON instead of XML) and also handling data preparation along with the models themselves.





# Saving ML Models



MLWriter/MLReader: Inherent to Spark

Supports not only saving a Model but a ML Pipeline.

The result of **save** for pipeline model is a JSON file for metadata while Parquet for model data, e.g. coefficients.

## **Is a model or Pipeline saved using Apache Spark ML persistence in Spark version X loadable by Spark version Y?**

- Major versions: No guarantees, but best-effort.
- Minor and patch versions: Yes; (backwards compatible)
- Note about the format: There are no guarantees for a stable persistence format, but model loading itself is designed to be backwards compatible.

## **Does a model or Pipeline in Spark version X behave identically in Spark version Y?**

- Major versions: No guarantees, but best-effort.
- Minor and patch versions: Identical behavior, except for bug fixes.

# Batch vs. Real-time Prediction



**Load implications**



**Cost Implications**

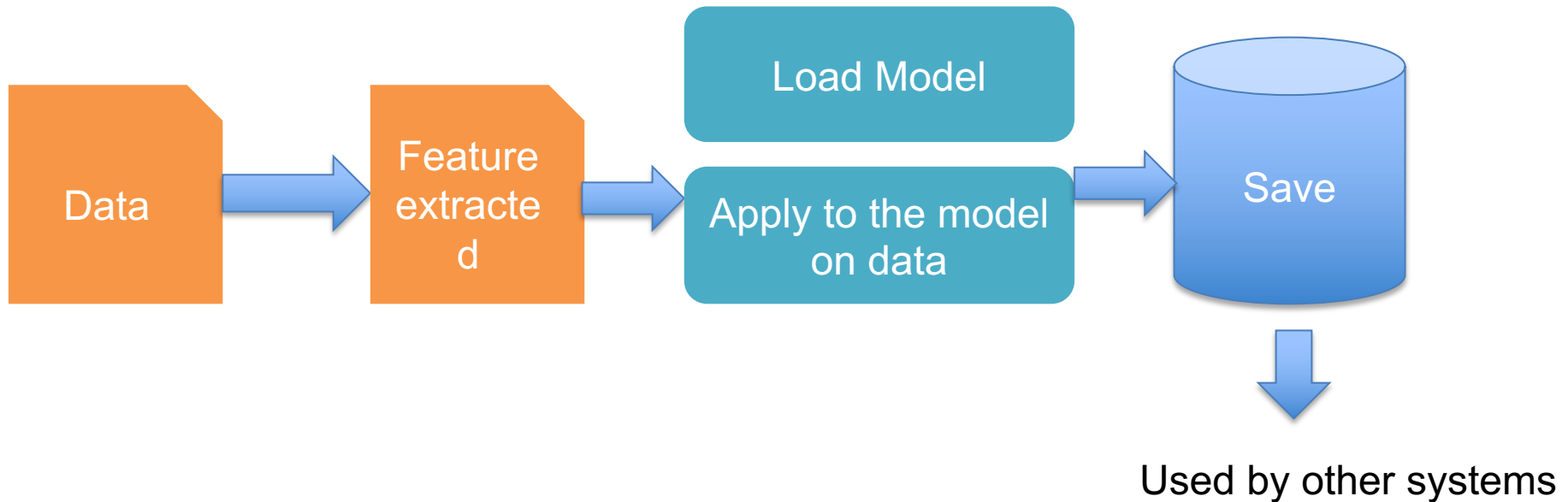


**Infrastructure Implications**



**Evaluation Implication**

# Batch Prediction



Use workflow managers like Airflow, or as simple as crontab etc.

# Real Time Prediction

## 1. Database Trigger based prediction:

### **When?**

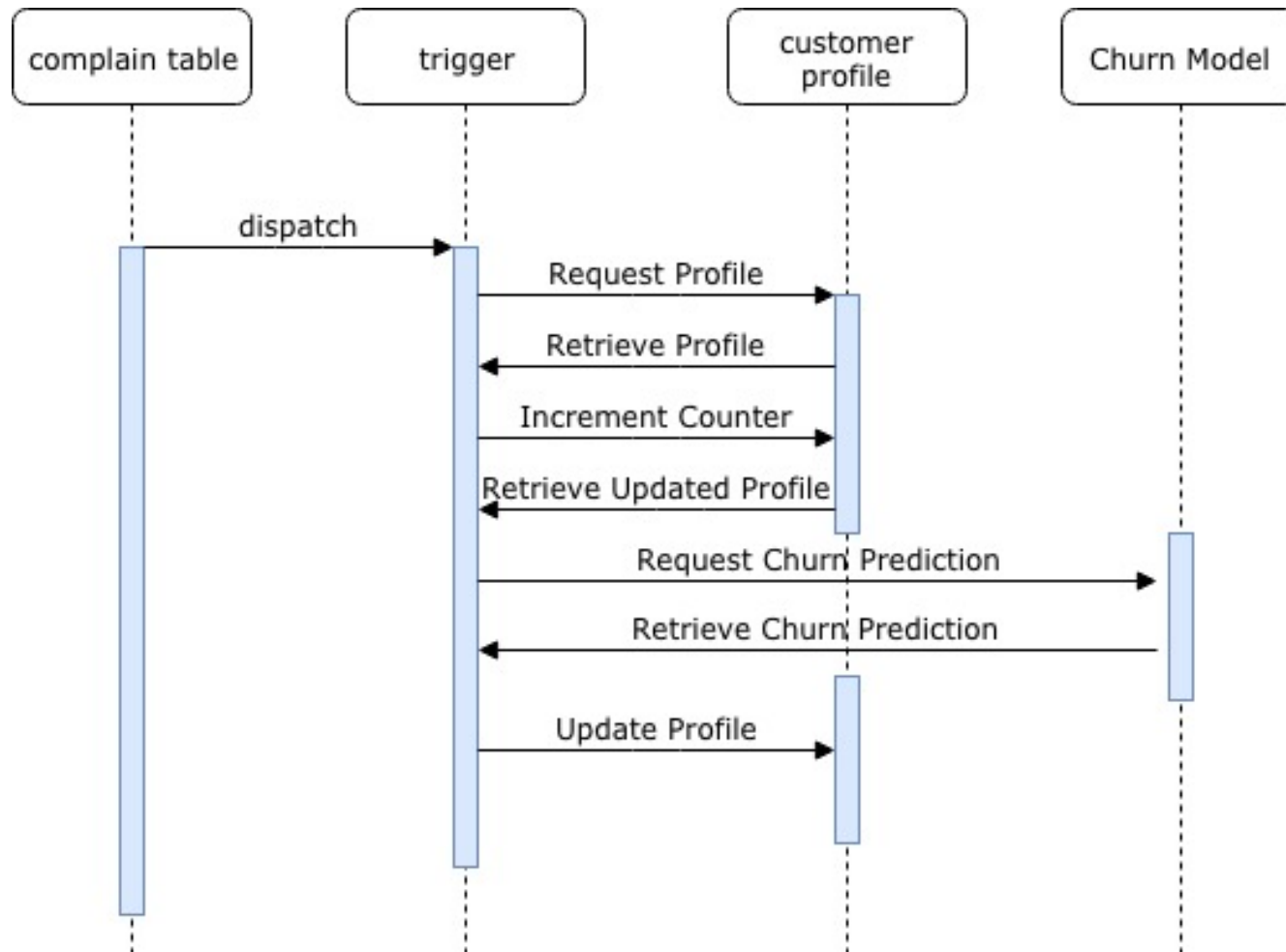
If the overall size of your database is fairly small ( $< 1\text{M}$  user profile) and the update frequency is occasional it can make sense to integrate some of the real-time update process directly within the database.

### **How?**

Postgres possess an integration that allows to run Python code as functions or stored procedure called PL/Python. This implementation has access to all the libraries that are part of the PYTHONPATH, and as such are able to use libraries such as Pandas and SKlearn to run some operations.

# Real Time Prediction

## Example Case



# Real Time Prediction

## 1. Database Trigger based prediction

- PostGres supports native Python integration and can run Python script on trigger.
- Ms SQL Server has **Machine Learning Services (In-Database)**
- Teradata can run R/Python script through an external script command.
- Oracle supports PMML model through its data mining extension.

# Real Time Prediction

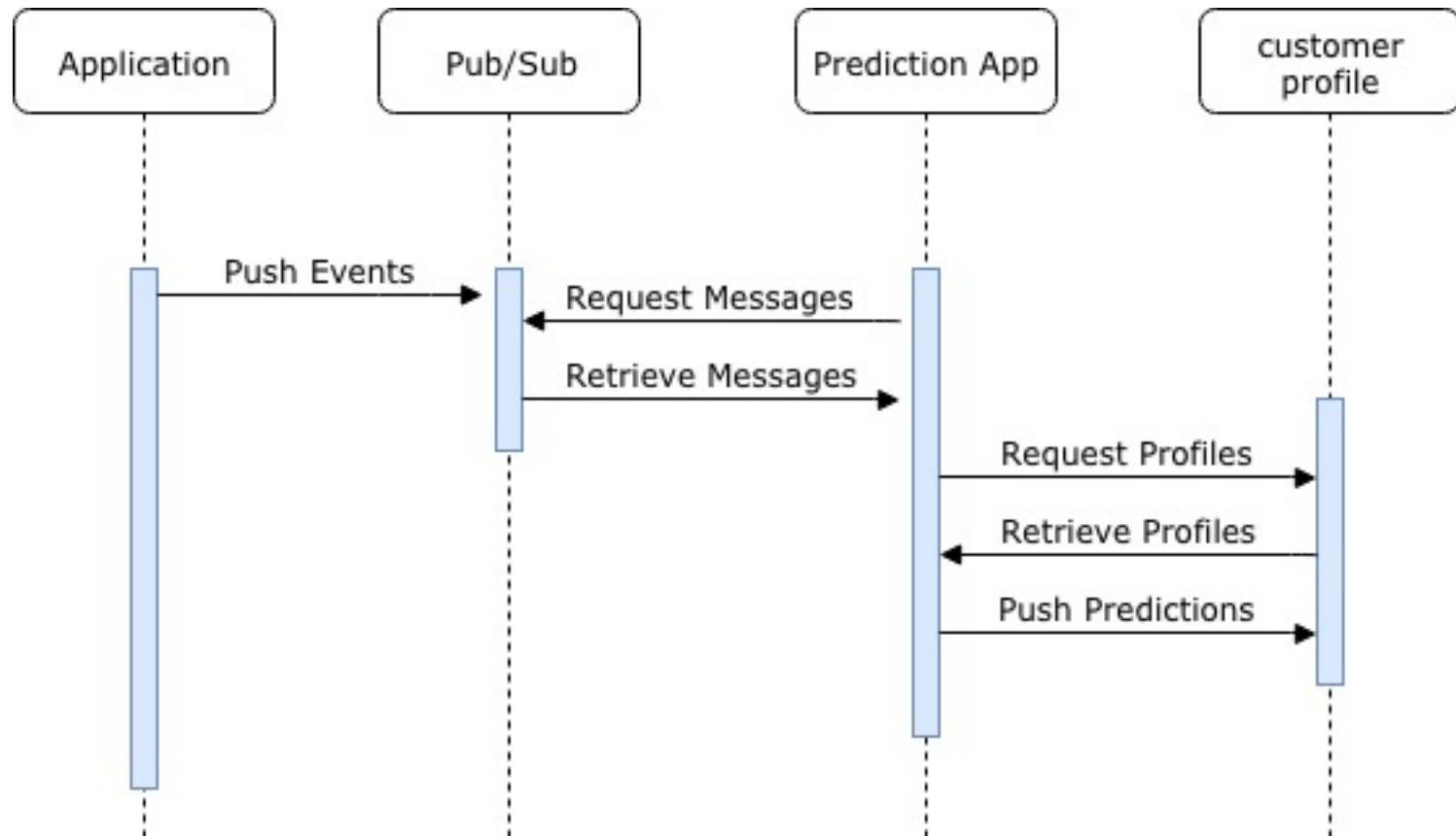
## 2. Message queue based prediction service

- A combination of Kafka for message queue and Spark Streaming for processing.
- Azure eventhub with Azure functions
- Google Pub/Sub with Apache Beam.



# Real Time Prediction

## Example Case



# Real Time Prediction

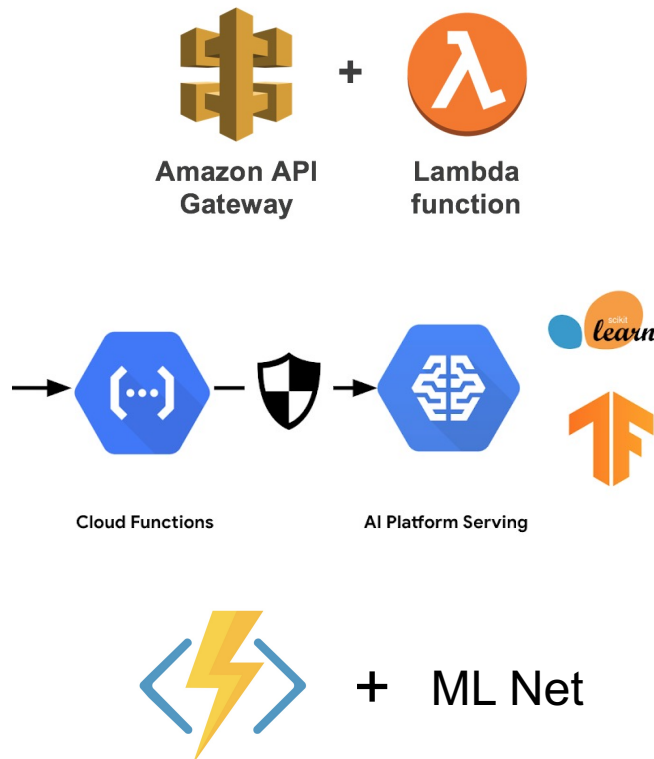
## 3. Web services based deployment

- Get unique ids and let web service pull related data from backend system and make prediction.
- Get the complete payload with data and make the prediction.
- A combination of both.

The web service or rest api based deployment models are more popular as they are very easy to scale and also can be used to be integrated in different frontend applications like mobile, web or desktop.

# Real Time Prediction

## 3.1 Using Functions!!





```
def main(req: func.HttpRequest) -> func.HttpResponse:
    logging.info('Python HTTP trigger function processed a request.')
    if req.body:
        try:
            logging.info("Converting Request to DataFrame")
            req_body = req.get_json()
            df_body = pd.DataFrame([req_body])

            logging.info("Loading the Prediction Model")
            filename = "model.pkl"
            loaded_model = joblib.load(filename)
            # Features names need to have been added to the pickled model
            feature_names = loaded_model.feature_names
            # subselect only the feature names

            logging.info("Subselecting the dataframe")
            df_subselect = df_body[feature_names]

            logging.info("Predicting the Probability")
            result = loaded_model.predict_proba(df_subselect)
            # We are looking at the probba prediction for class 1
            prediction = result[0][1]

            return func.HttpResponse("{prediction}".format(prediction=predicti

        except ValueError:
            pass
    else:
        return func.HttpResponse(
            "Please pass a name on the query string or in the request body",
            status_code=400
        )
```

---

Can any one tell me what is the issue with this approach?

# Real Time Prediction

## 3.2 Using Container



- Ununiform environments across models.
- Ununiform library requirements across models.
- Ununiform resource requirements across models.
- Scaling at model level

# Real Time Prediction

## 3.3 Directly In-App



Google ML Kit or the likes of Caffe2 or Apple's core ML allows to leverage models within native applications



Tensorflow.js or ONNXJs allow running models directly on browser

# Real Time Prediction

## 3.4 Using Notebooks!!!!



# Why ML Life Cycle is more complicated than SDLC!

- Variety of tools and packages at disposal.
- Hard to track useful experiments.
- Reproducibility is hard.
- Deploying is hard.



# ML Life Cycle Management Platforms

Azure ML services (launched in 2015)

## What's possible with AML?

Using Azure Machine Learning service, you can track your models as you build, train, deploy, and manage them at cloud scale.



### Run & Monitor Experiments

Submit Experiments for training and automatically track their progress and view logs.



### Register Models

Save scoring logic operations into models to create Docker Images and Deployments.



### Build Images

Quickly create Docker images that encapsulate models, scripts, and any associated files.



### Deploy Models

Send scoring requests to web services in Azure Container Instances, Azure Kubernetes Service, or field programmable gate arrays (FPGA).



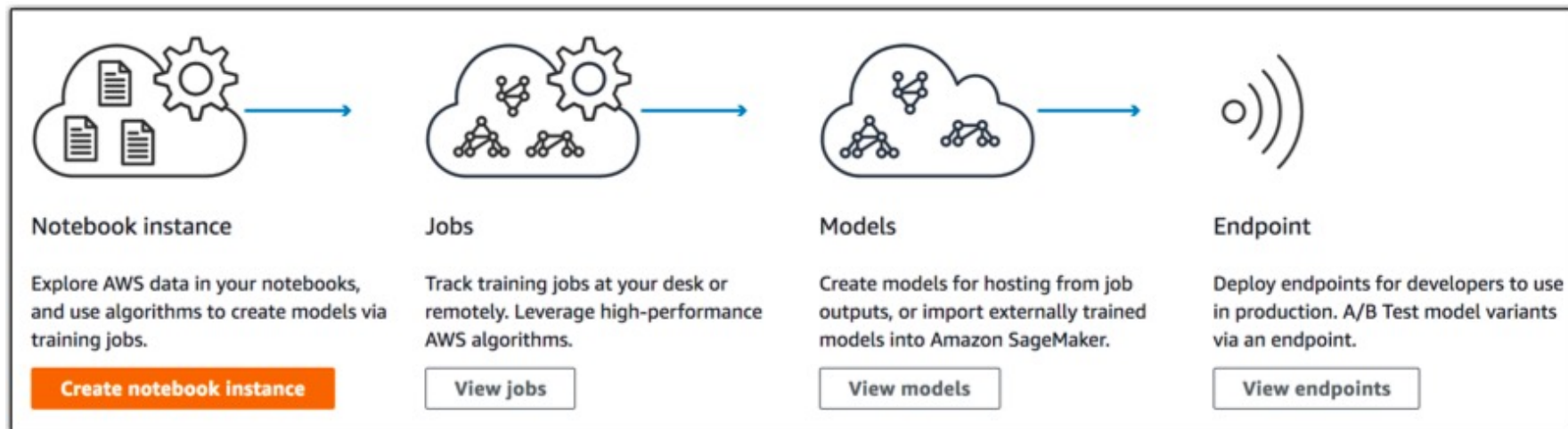
### Create Pipeline

Pipelines are used to build, optimize, and manage machine learning workflows.



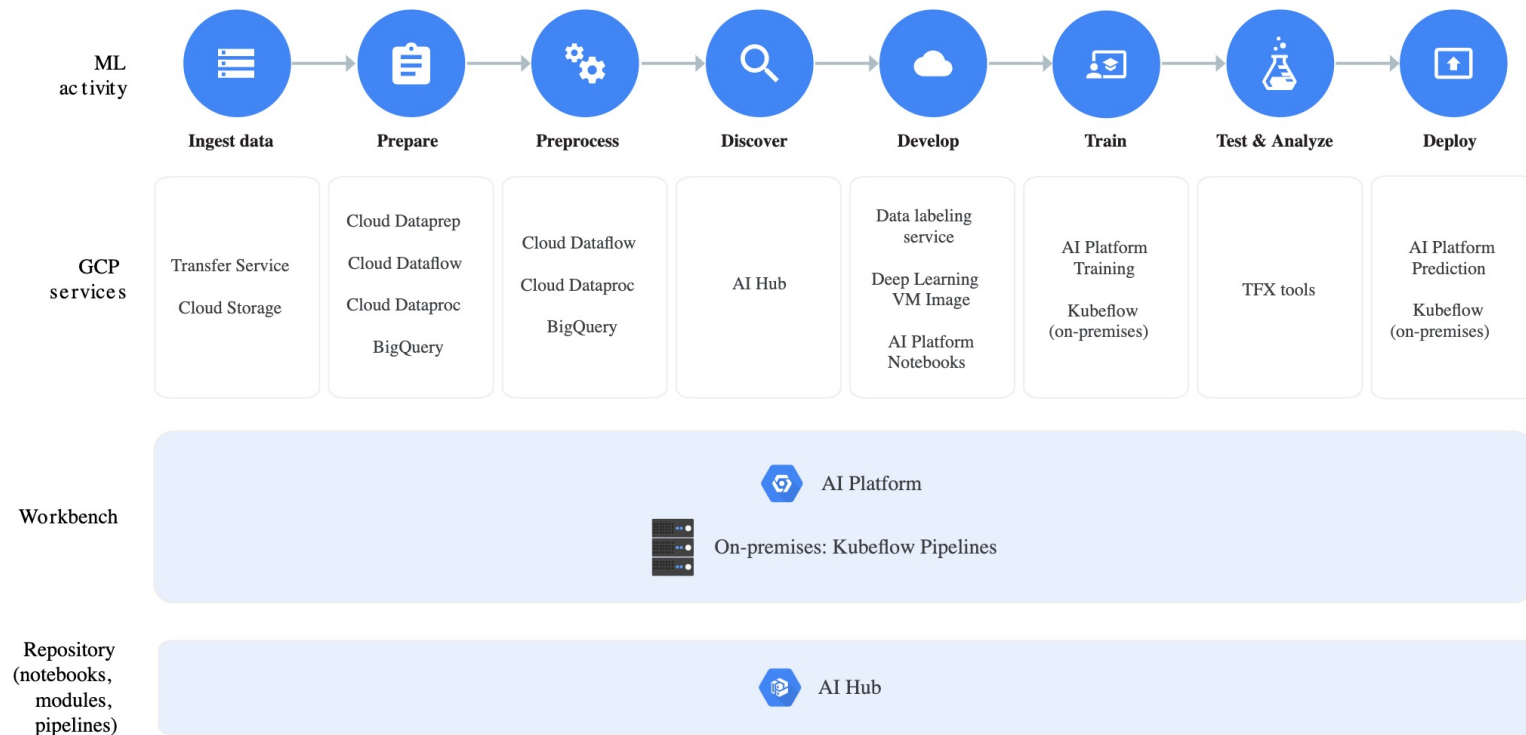
# ML Life Cycle Management Platforms

Amazon SageMaker (launched in 2017)



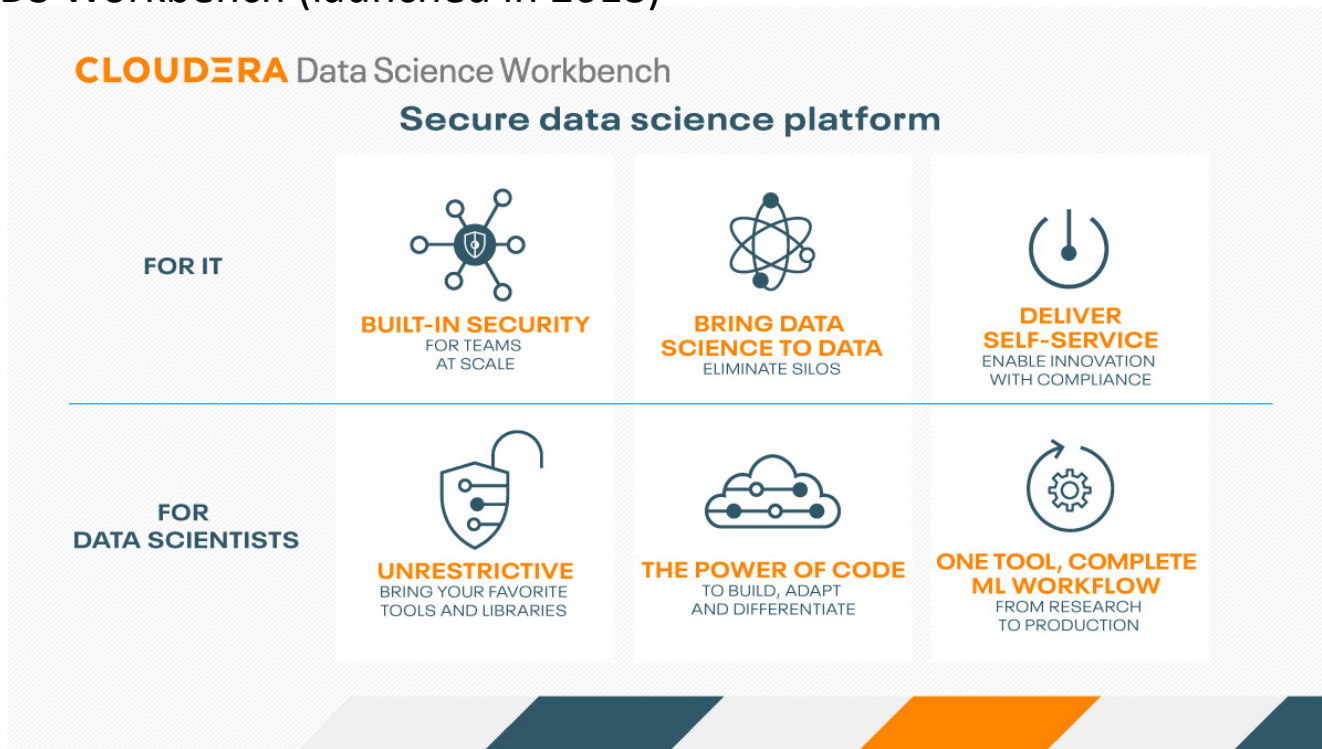
# ML Life Cycle Management Platforms

Google AI (launched in 2018)



# ML Life Cycle Management Platforms

Cloudera DS Workbench (launched in 2018)



# ML Life Cycle Management Platforms

Databricks Managed ML Flow (launched in 2019) : ML Flow in itself is a library and not a hosted platform.

