

# Building Efficient ML Pipelines and Responsible AI Solutions

Adi Polak  
Microsoft

@adipolak





Trust



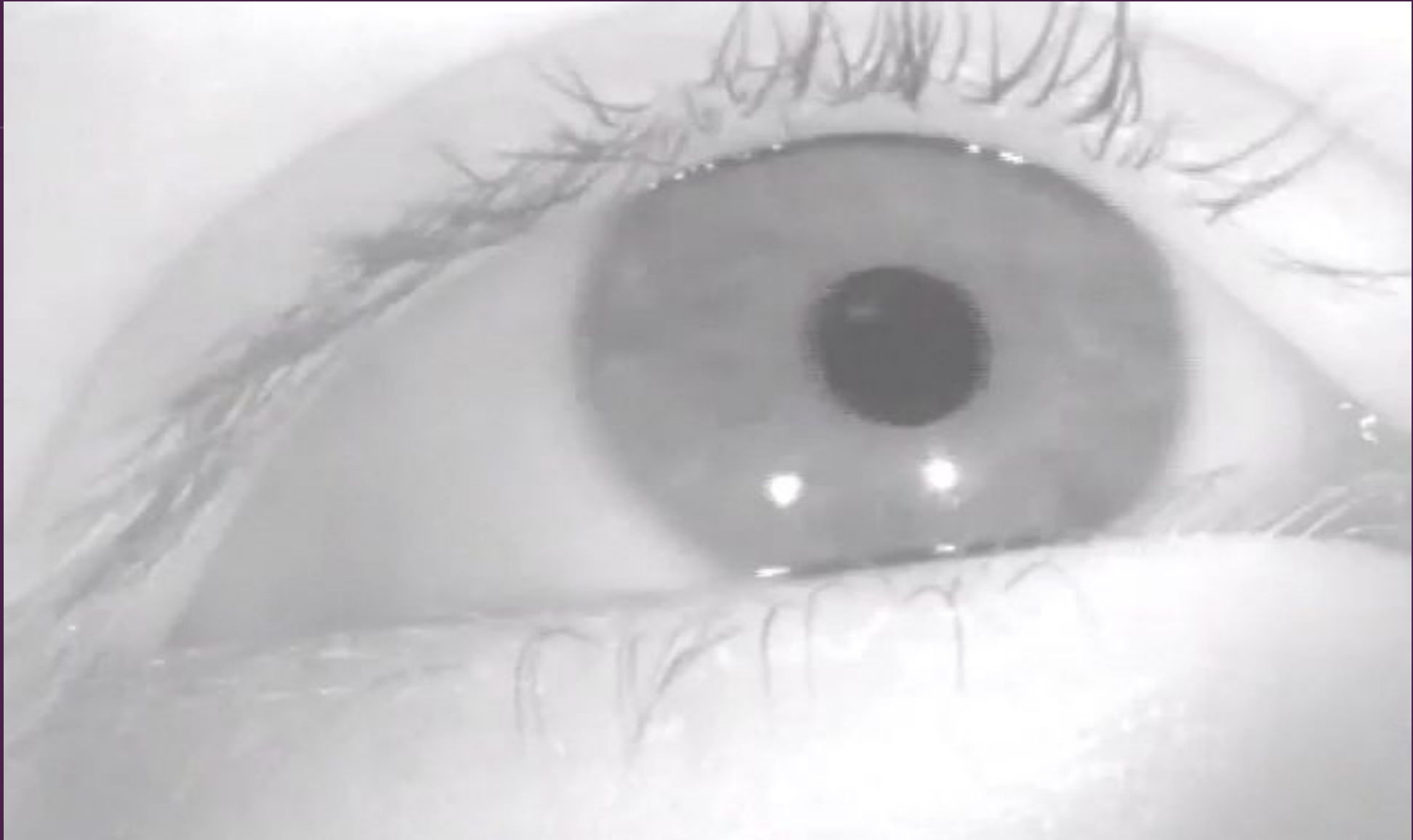
- LET'S START FROM THE BEGINNING.

What happens when  
we get raw data?



@adipolak

Movie copyright (C) 2008, Blender Foundation



@adipolak

# ML Process / Life Cycle

- 1 Gather Data
- 2 Feature Extract, Clean and Normalize
- 3 Select algorithm
- 4 Evaluate model
- 5 Data/Insights visualization

 *Repeat!*



@adipolak



But in real life:

Accuracy < 0.5  
ROC curve 😞

---







Aim for high  
Accuracy



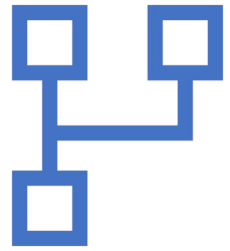
What can  
you do?

---

Automate!

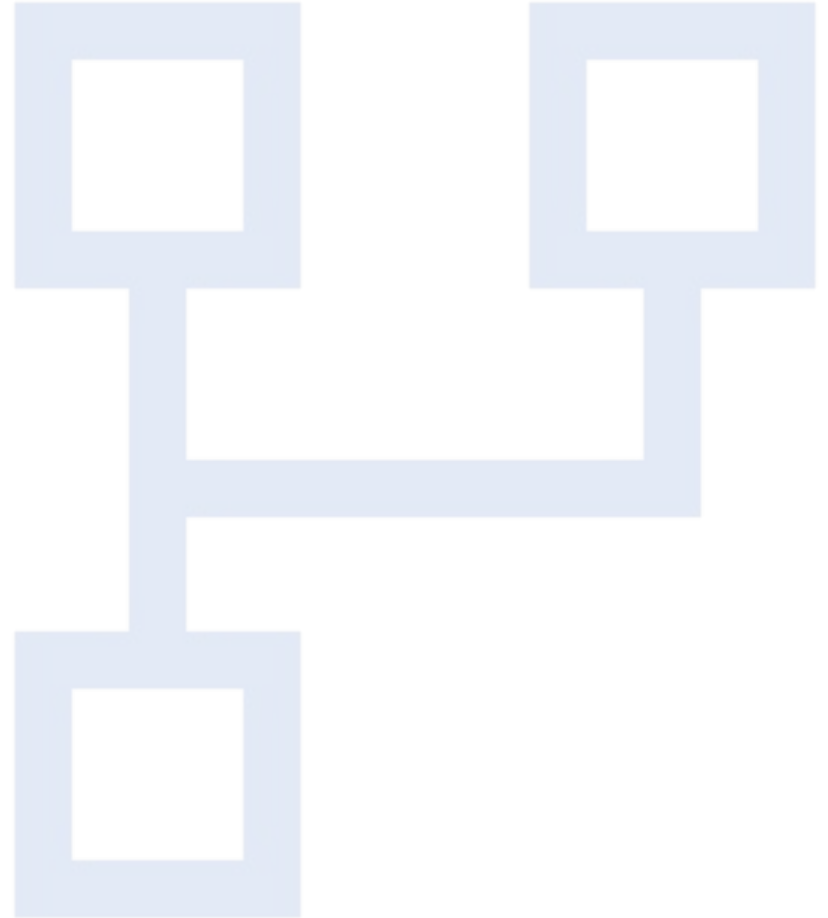
@adipolak





HOW?

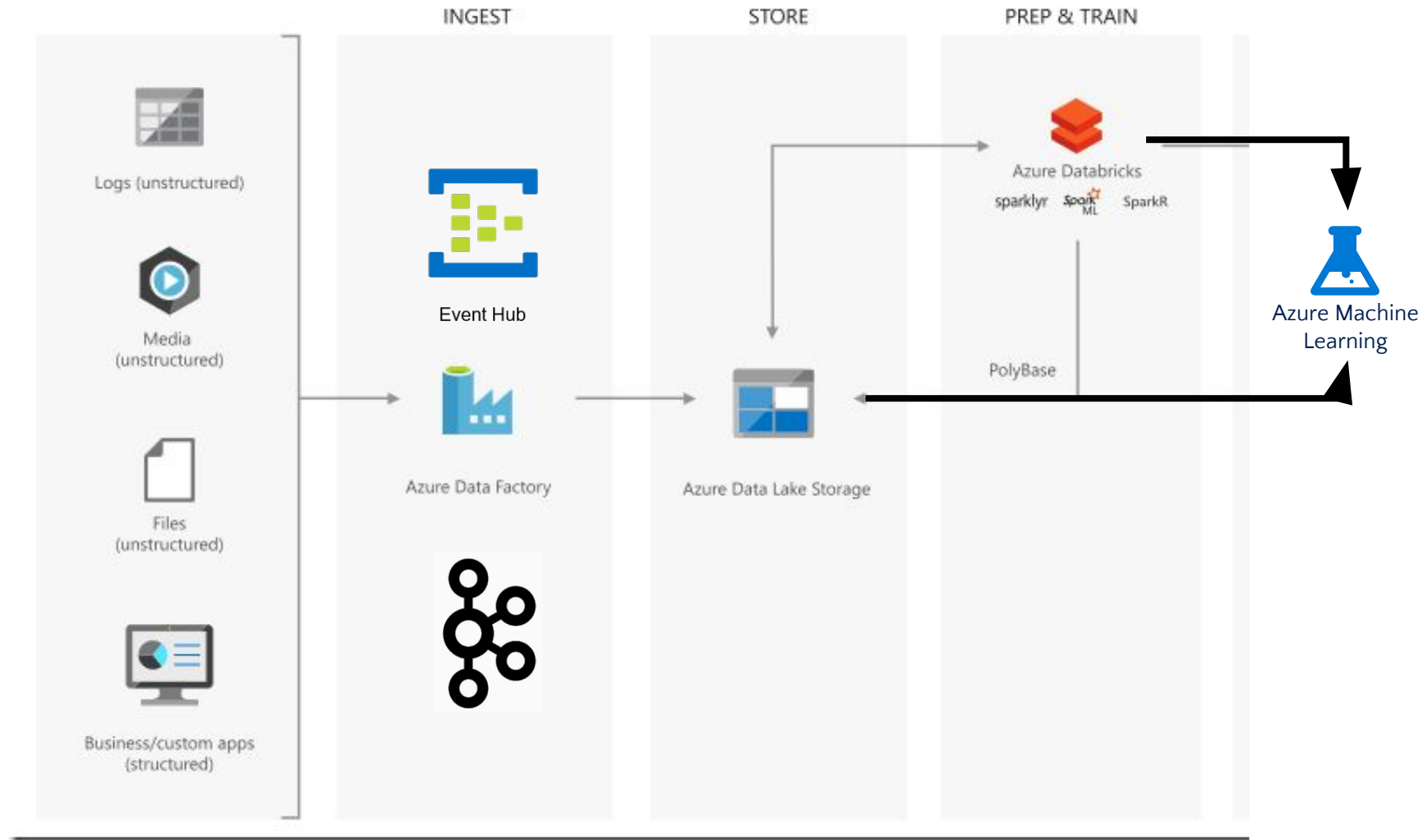
Pipelines!



A photograph of a large industrial facility, likely a refinery or chemical plant, at night. The scene is illuminated by numerous yellow lights from within the structures, creating a stark contrast against the dark blue twilight sky. In the foreground, there is a dark, silhouetted line of trees and a grassy field. The background features several tall, cylindrical distillation columns and a complex network of pipes and structural steel. The overall atmosphere is industrial and somewhat somber due to the low light.

What are pipelines?

# Big Data/ ML Pipelines



# Demo

## Apache Spark ML Pipelines



## Configure ML pipeline stages

```
1
2 // Configure an ML pipeline, which consists of three stages: tokenizer, hashingTF, and lr.
3 val tokenizer = new Tokenizer()
4   .setInputCol("text")
5   .setOutputCol("words")
6
7
8 val hashingTF = new HashingTF()
9   .setNumFeatures(1000)
10  .setInputCol(tokenizer.getOutputCol)
11  .setOutputCol("features")
12
13
14 val lr = new LogisticRegression()
15   .setMaxIter(10)
16   .setRegParam(0.001)
17
18 run.logParam("lr_maxIter", "10")
19 run.logParam("lr_RegParam", "0.001")
20
```

[Show result](#)

## Set the pipeline

```
1  val pipeline = new Pipeline()  
2    .setStages(Array(tokenizer, hashingTF, lr))
```

[Show result](#)

## Fit the data

```
1 // Fit the pipeline to training documents.  
2 val model : PipelineModel = pipeline.fit(training)
```

## Save and load model

```
1 // Now we can optionally save the fitted pipeline to disk
2 model.write.overwrite().save("/tmp/spark-logistic-regression-model")
3
4 // We can also save this unfit pipeline to disk
5 pipeline.write.overwrite().save("/tmp/unfit-lr-model")
6
7 // And load it back in during production
8 val sameModel = PipelineModel.load("/tmp/spark-logistic-regression-model")
```

## Create Mock test set

```
1 // Prepare test documents, which are unlabeled (id, text) tuples.
2 val test = spark.createDataFrame(Seq(
3     (4L, "spark i j k"),
4     (5L, "l m n"),
5     (6L, "spark hadoop spark"),
6     (7L, "apache hadoop")
7 ))).toDF("id", "text")
8
```

cmd 3

# Prepare file printer for MLflow artifact file

```
1 val filePrinter = new PrintWriter("/tmp/output.txt")
2 filePrinter.write("LR run")
```



## Make prediction using the ML model

```
1 // Make predictions on test documents.
2 model.transform(test)
3   .select("id", "text", "probability", "prediction")
4   .collect()
5   .foreach { case Row(id: Long, text: String, prob: Vector, prediction: Double) =>
6     //println(s"($id, $text) --> prob=$prob, prediction=$prediction")
7     filePrinter.write(s"($id, $text) --> prob=$prob, prediction=$prediction")
8   }
```

## Close and save

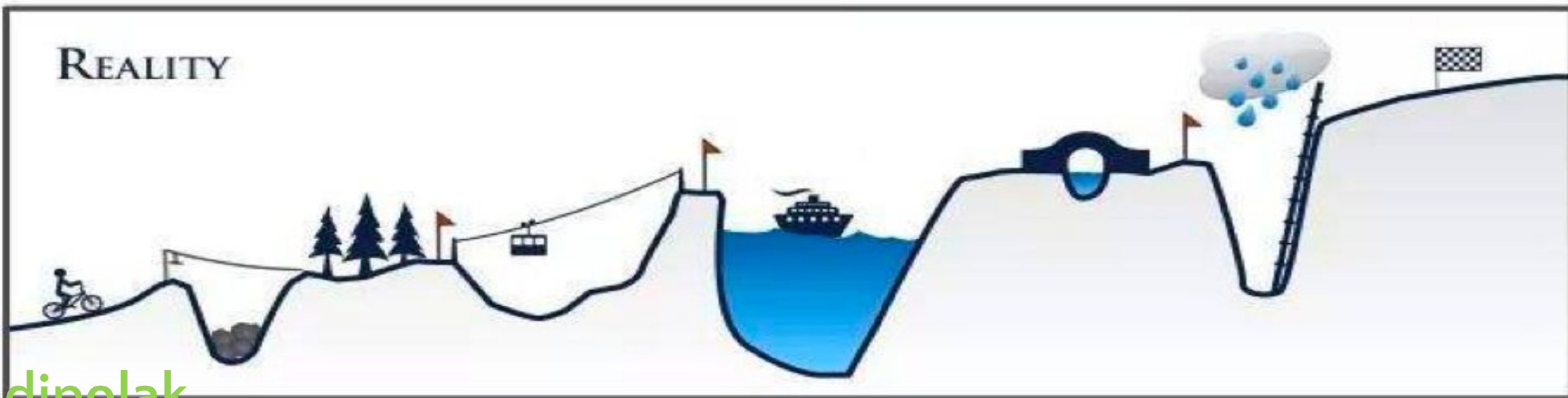
```
1 run.logArtifact(Paths.get("/tmp/output.txt"))  
2 run.endRun()  
3 filePrinter.close()
```

[Show result](#)

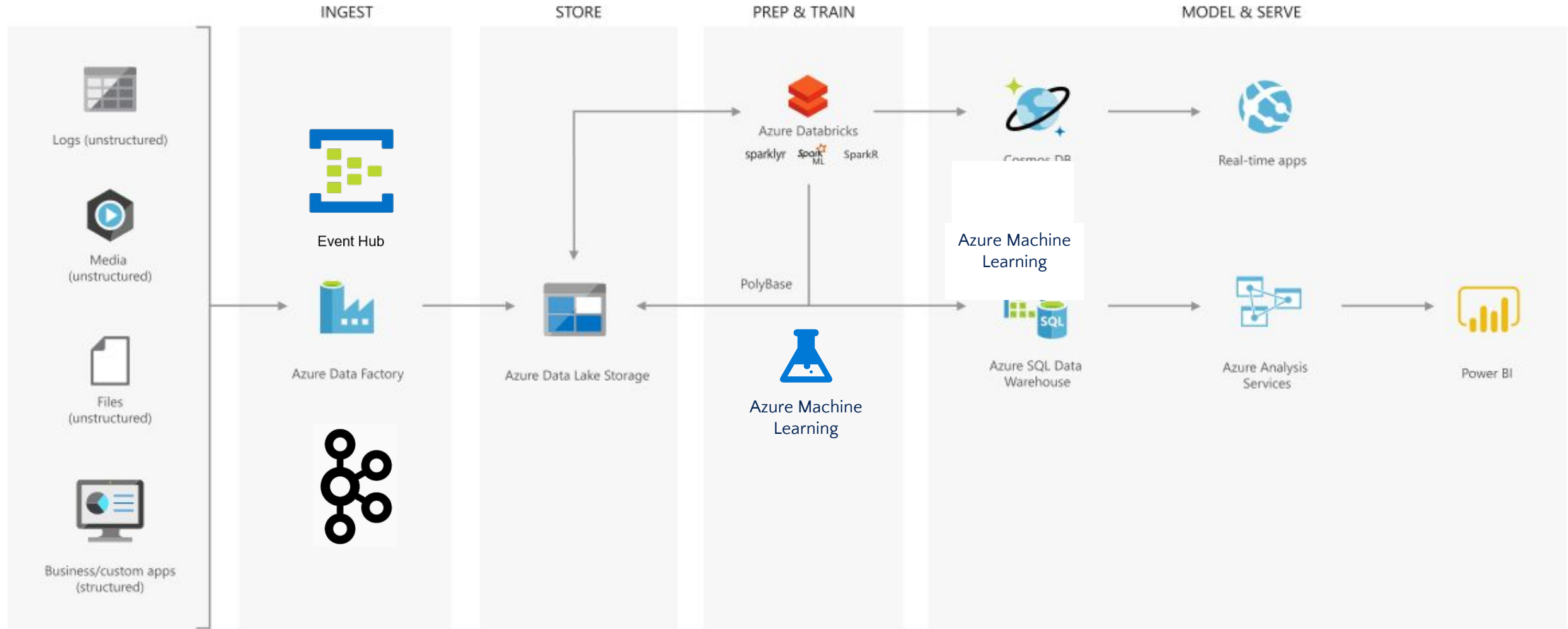
## YOUR PLAN



## REALITY



# Big Data/ ML Pipelines





High  
accuracy!

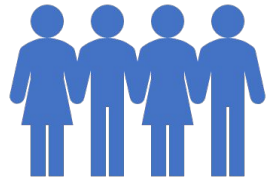
But, at what  
cost?

---



false positives





# Human centric

## Responsible AI







Pixabay

# Our updated goals:



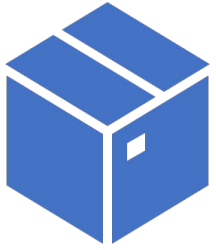
Lawful



Ethical



Robust

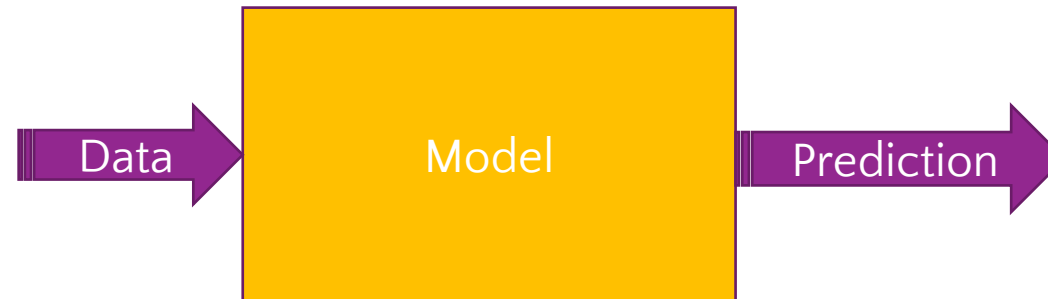


# ML is a black box

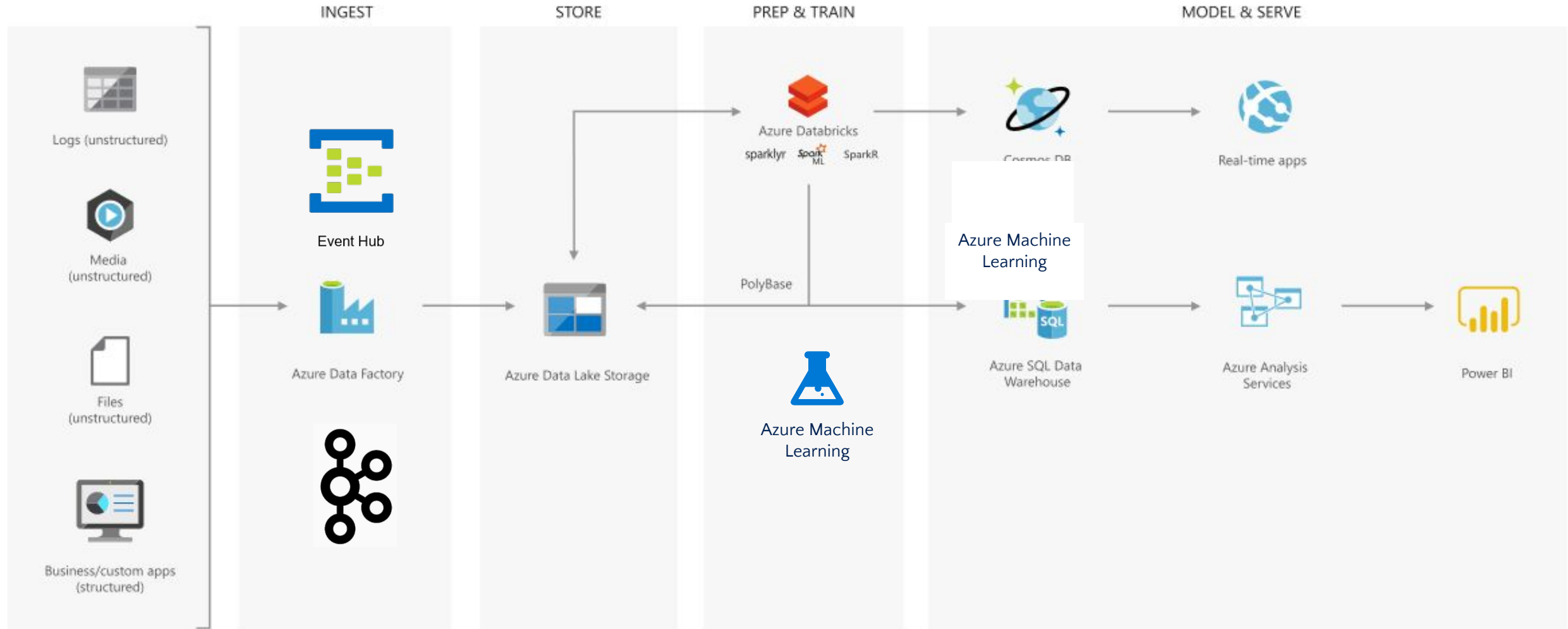
Training:



Testing/Prediction:

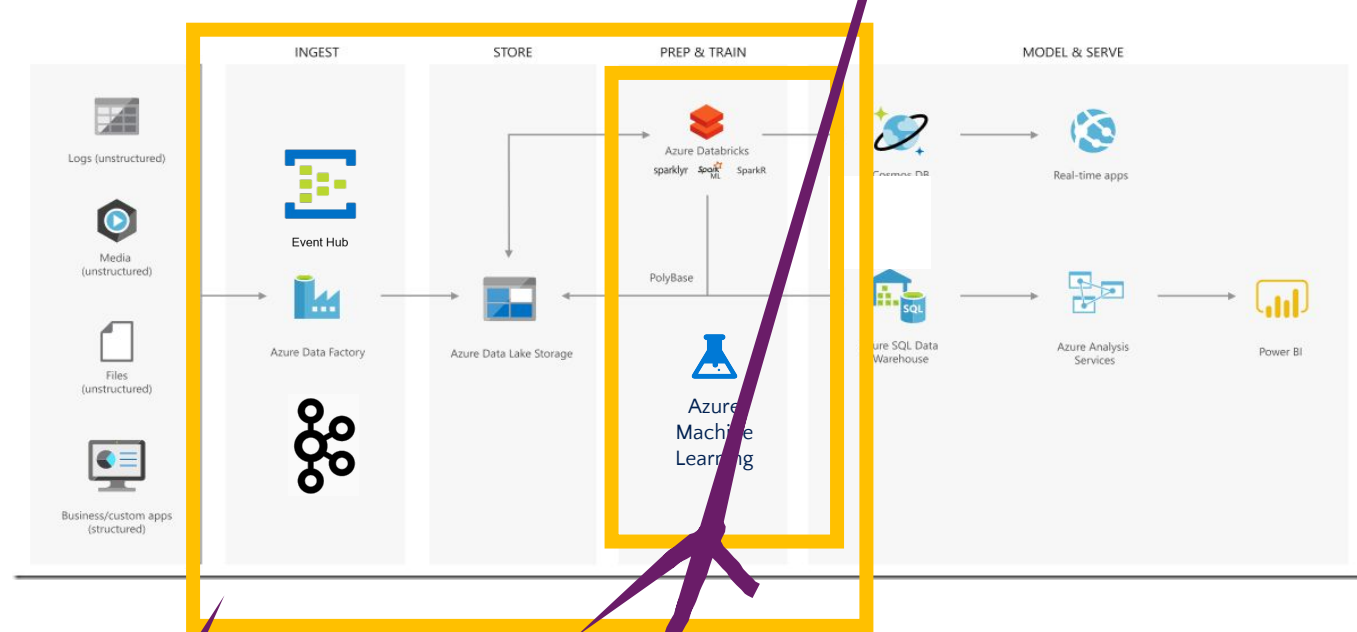


# Big Data/ ML Pipelines





# Big Data/ ML Pipelines



Check and transform the data

Visualize the model

Balance the data

Explainers

# How Microsoft support Responsible AI



ONLINE **FREE**  
HIGH QUALITY  
COURSES



INVESTED **1B\$** IN  
OPEN AI



**115M\$** GRANT FOR  
AI FOR GOOD

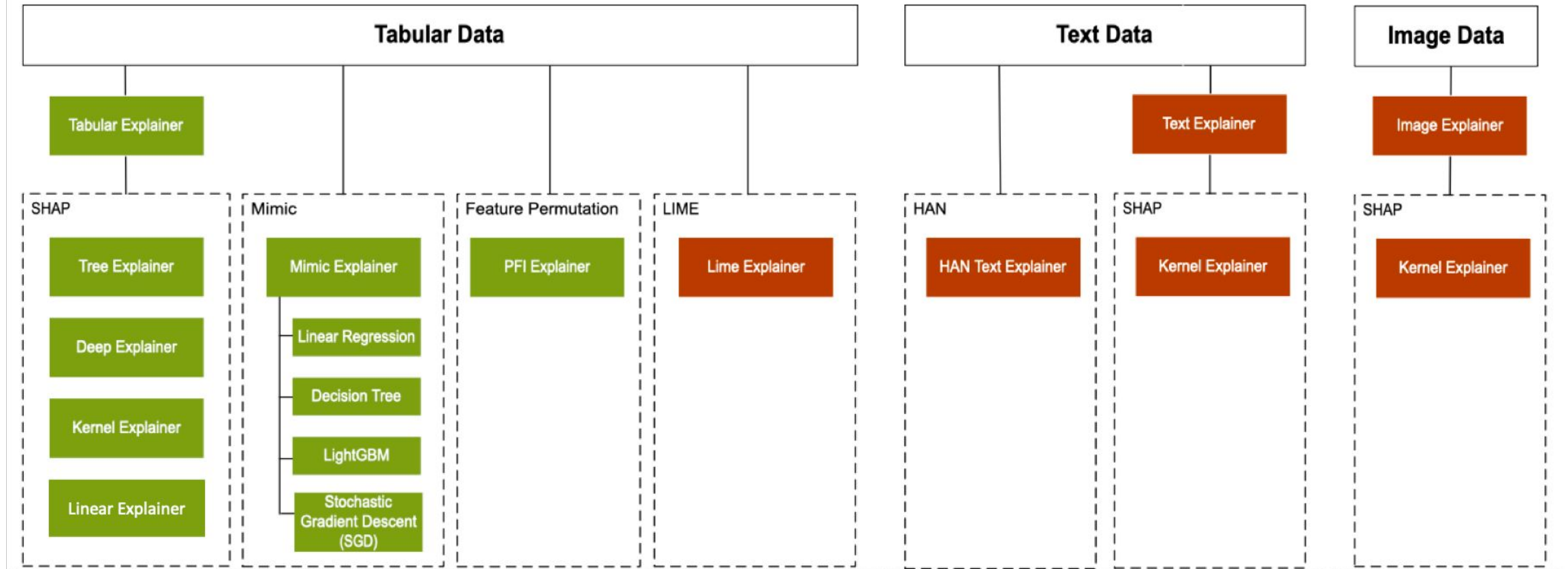


**OPEN SOURCE**

# Machine Learning Interpretability

Legend:

■ Main Package  
■ Contrib Package



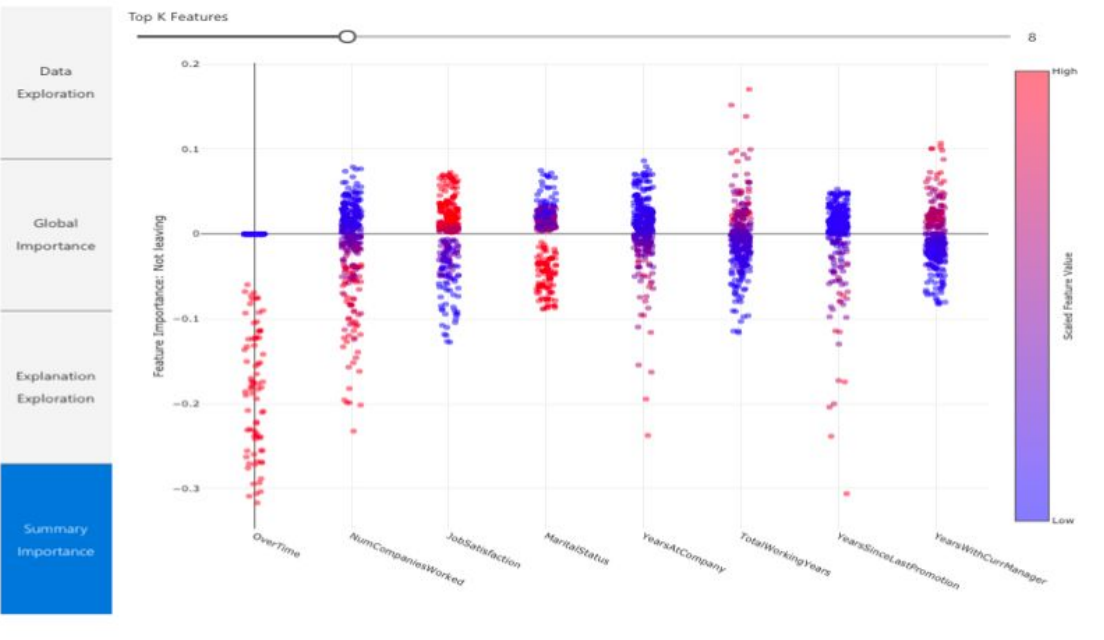
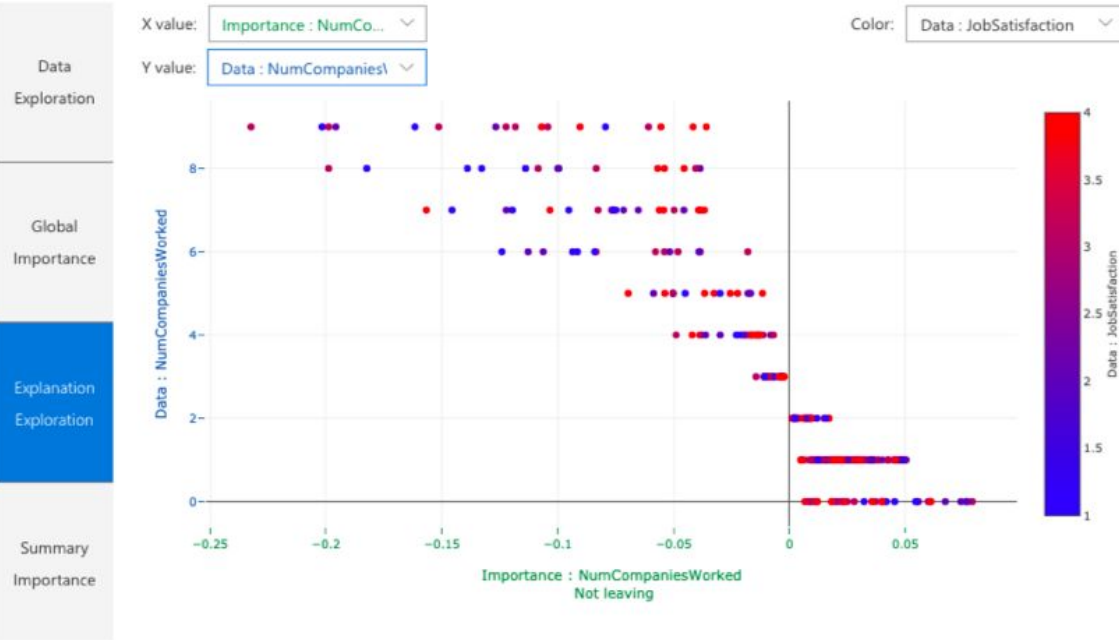
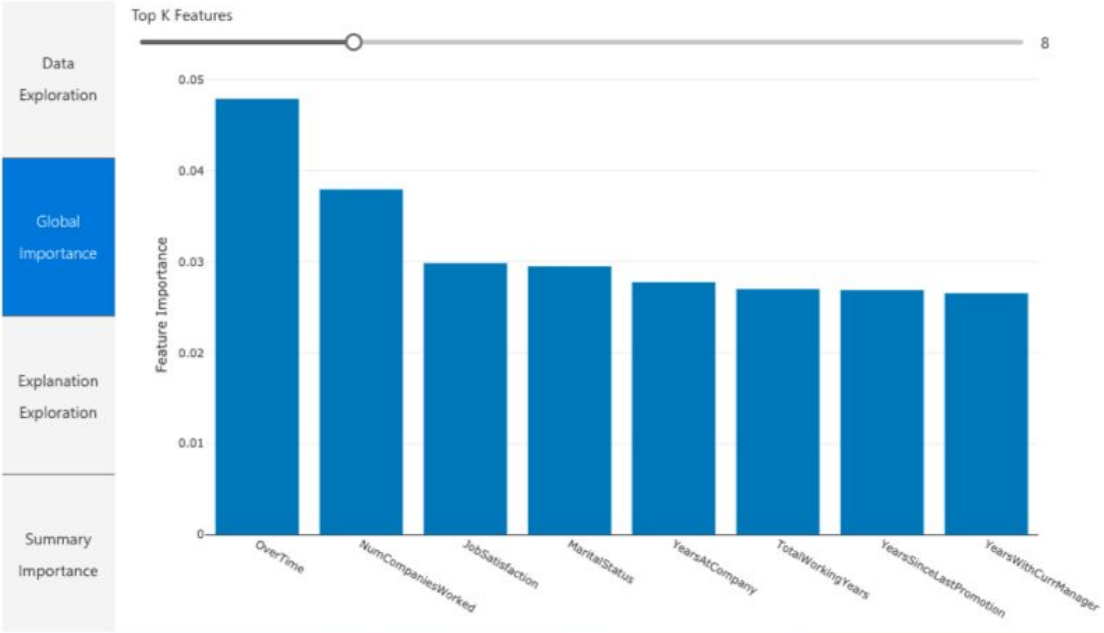
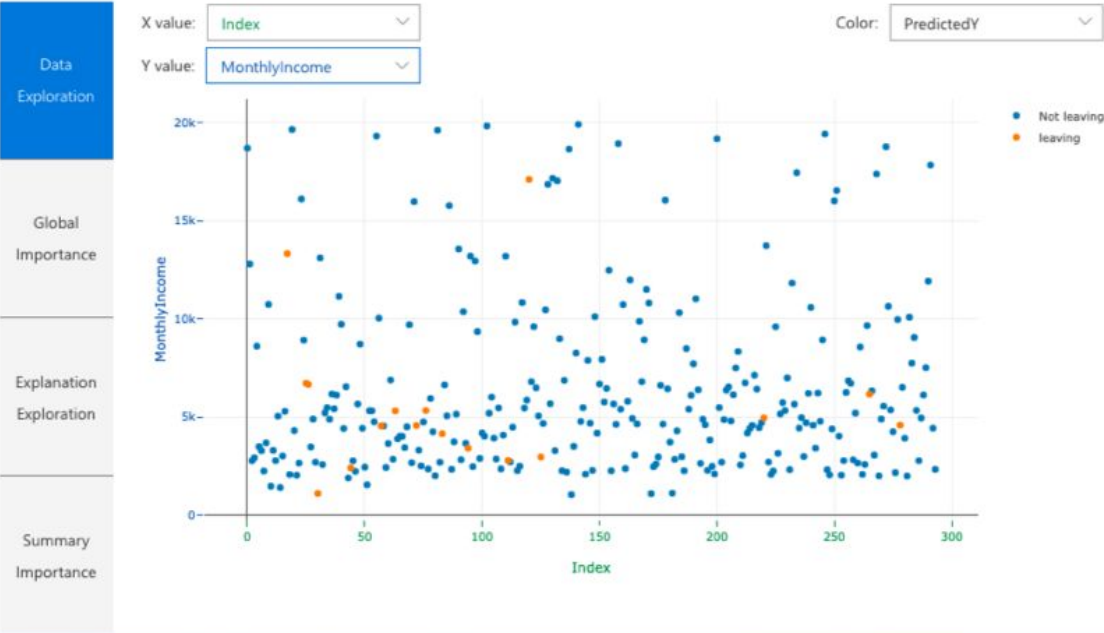
```
# load people dataset
from sklearn.datasets import load_breast_cancer
from sklearn import svm
from sklearn.model_selection import train_test_split
people_data = load_people()
classes = people_data.target_names.tolist()

# split data into train and test
from sklearn.model_selection import train_test_split
x_train, x_test, y_train, y_test = train_test_split(people_data.data,
                                                    people_data.target,
                                                    test_size=0.2,
                                                    random_state=0)

clf = svm.SVC(gamma=0.001, C=100., probability=True)
model = clf.fit(x_train, y_train)
```

```
from interpret.ext.blackbox import TabularExplainer

# "features" and "classes" fields are optional
explainer = TabularExplainer(model,
                              x_train,
                              features=people_data.feature_names,
                              classes=classes)
```



aka.ms/ml-interpretability-to  
ol

# Azure Cognitive Services

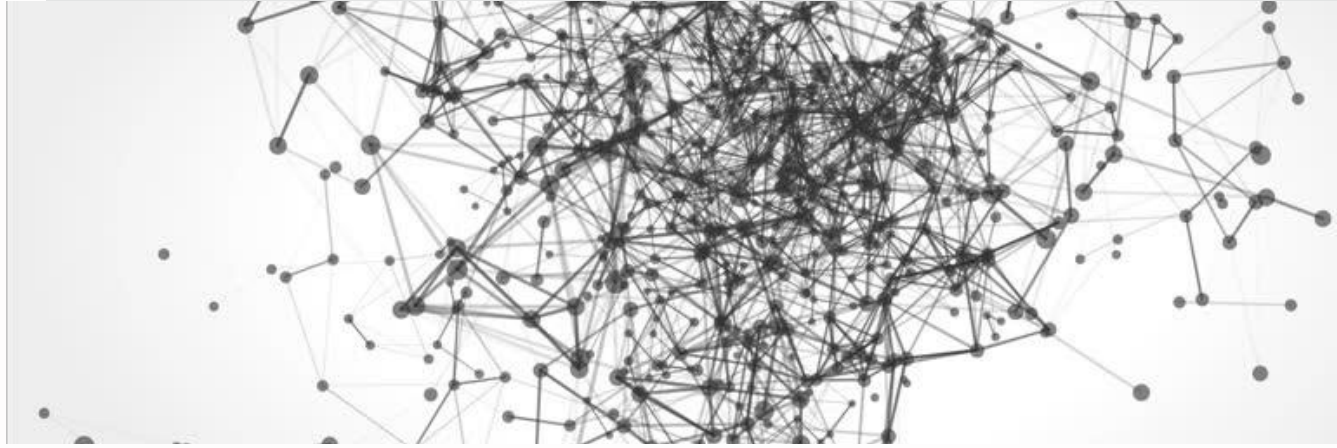
---

[aka.ms/AA6kex](https://aka.ms/AA6kex)  
s





# Tools



**Spark  
Streaming**



**Spark ML**



**Spark SQL**



**MLflow**







Demo

Apache Spark ML Pipelines with Cognitive Services

## Define MLflow experiment

```
1  val mlflowContext = new MlflowContext()
2  val experimentName = "/Shared/PipelinesExp"
3  val client = mlflowContext.getClient()
4  val experimentOpt = client.getExperimentByName(experimentName);
5  if (!experimentOpt.isPresent()) {
6      client.createExperiment(experimentName)
7  }
8  mlflowContext.setExperimentName(experimentName)
9
10 val run = mlflowContext.startRun("run")
```

[Show result](#)

# Demo



# *You are only Good as your Data is*

Use explainers

Understand your data



# Learn more !

[aka.ms/free-responsible-ai-course](https://aka.ms/free-responsible-ai-course)

[aka.ms/twitter\\_sentiment\\_analysis](https://aka.ms/twitter_sentiment_analysis)

[aka.ms/ml-interpretability-tool](https://aka.ms/ml-interpretability-tool)

[aka.ms/ai-for-good-grant](https://aka.ms/ai-for-good-grant)

*Thank you !*



@adipolak

# What is Machine learning

---

- Lifecycle:
- Gather data
- Data preparation – clean it
- Data wrangling
- Data analysis
- Feature extraction
- Train model
- Test model
- Deployment

