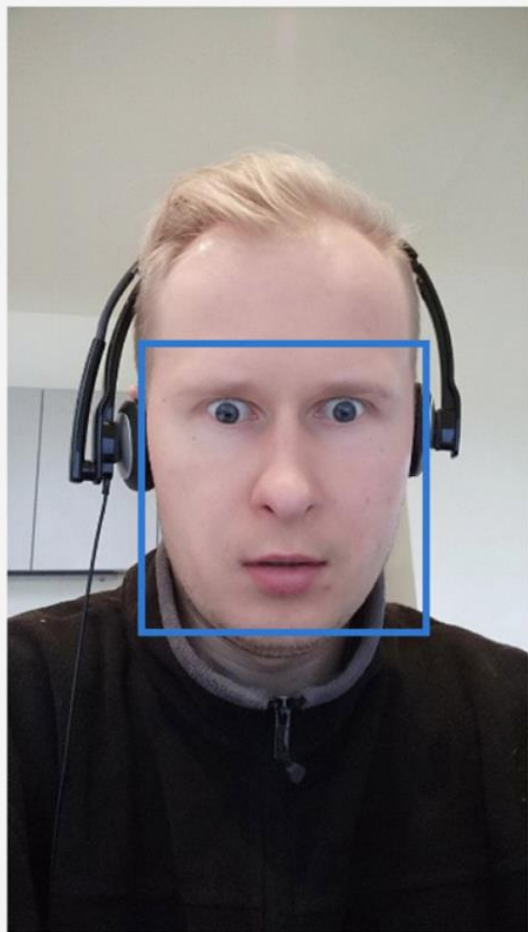


Lessons Learned on Building a Data Science Environment on a Public Cloud

`valdas@maksimavicius.eu`

Speaker – Valdas Maksimavičius

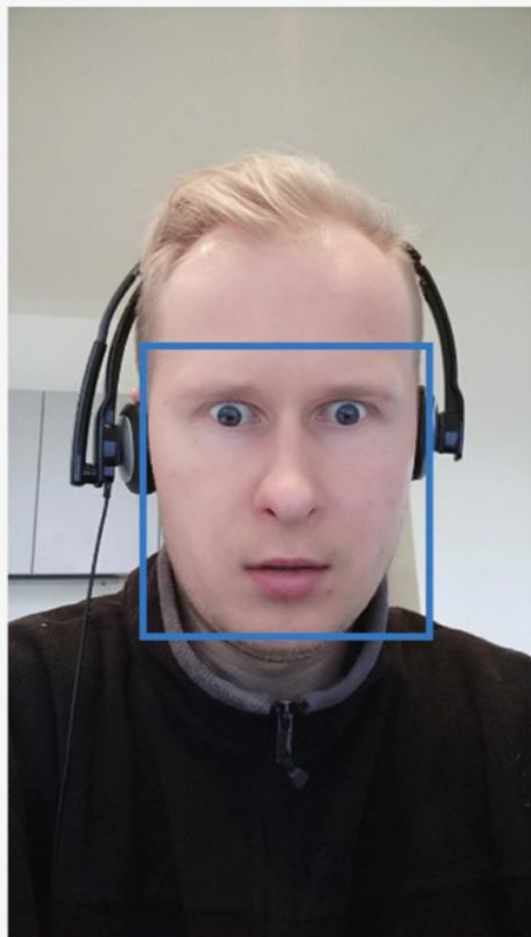
- Big Data Architect at Cognizant
- Insurance, Pensions & Manufacturing industries
- Vilnius Microsoft Data Platform Meetup /
Hack4Vilnius Hackathon
- Valdas.blog



Detection result:

JSON:

```
[
  {
    "faceId": "f783a705-c1c9-4cf1-bb24-064f951f4e52",
    "faceRectangle": {
      "top": 415,
      "left": 163,
      "width": 366,
      "height": 366
    },
    "faceAttributes": {
      "hair": {
        "bald": 0.13,
        "invisible": false,
        "hairColor": [
          {
            "color": "brown",
            "confidence": 0.91
          },
          {
            "color": "red",
            "confidence": 0.9
          },
          {
            "color": "blond",
            "confidence": 0.58
          }
        ]
      }
    }
  }
]
```

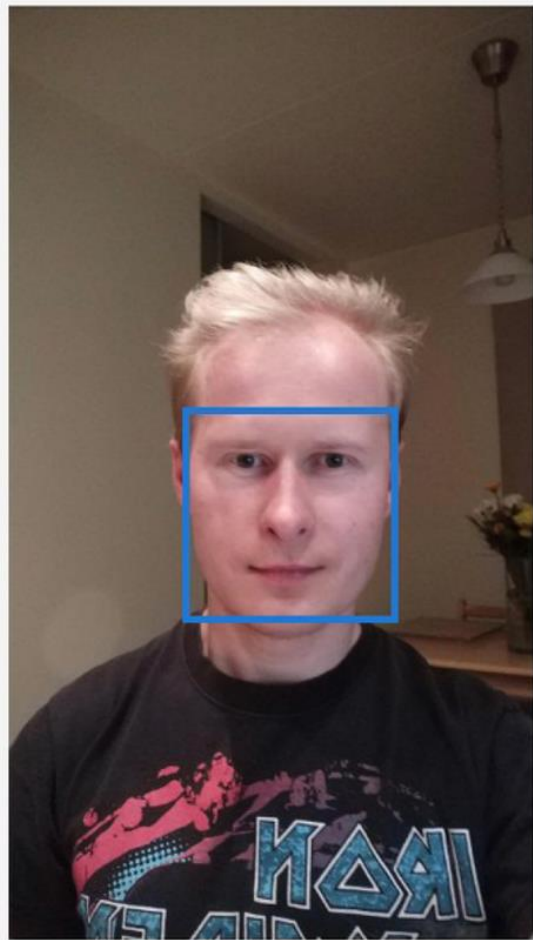


Detection result:

JSON:

```
[
  {
    "faceId": "f783a705-c1c9-4cf1-bb24-064f951f4e52",
    "faceRectangle": {
      "top": 415,
      "left": 163,
      "width": 366,
      "height": 366
    },
    "faceAttributes": {
      "hair": {
        "bald": 0.13,
        "invisible": false,
        "hairColor": [
          {
            "color": "brown",
            "confidence": 0.91
          },
          {
            "color": "red",
            "confidence": 0.9
          },
          {
            "color": "blond",
            "confidence": 0.58
          }
        ]
      }
    }
  }
]
```

"bald": 0.13



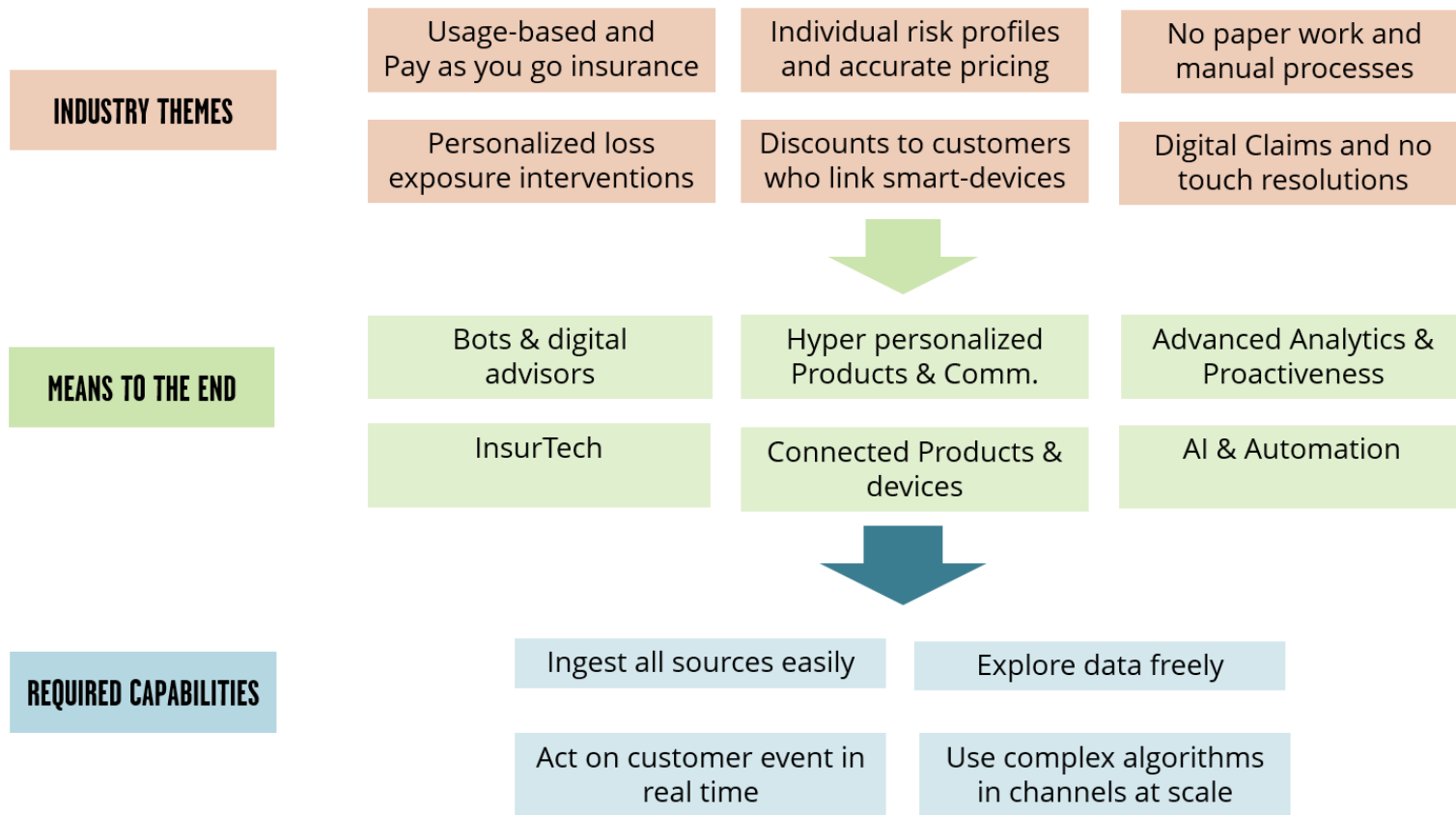
Detection result:

JSON:

```
[
  {
    "faceId": "588778e2-bf43-48a2-9dbd-181fb26aaf41",
    "faceRectangle": {
      "top": 1644,
      "left": 704,
      "width": 883,
      "height": 883
    },
    "faceAttributes": {
      "hair": {
        "bald": 0.17,
        "invisible": false,
        "hairColor": [
          {
            "color": "blond",
            "confidence": 0.98
          },
          {
            "color": "brown",
            "confidence": 0.73
          },
          {
            "color": "gray",
            "confidence": 0.52
          }
        ]
      }
    }
  }
]
```

"bald": 0.17

Let's start with industry trends (Insurance)



Select relevant end-to-end use cases

Retail

CONSUMER ENGAGEMENT



**Real-time Pricing
Optimization**

Financial

RISK AND REVENUE MANAGEMENT



**Risk and Fraud, Threat
Detection**

Oil/Gas & Energy

GRID OPS, ASSET OPTIMIZATION



Industrial IoT

Security

ACTIONABLE THREAT INTELLIGENCE



Security Intelligence

Healthcare

SENSOR DATA



**IoT DEVICE
ANALYTICS**

Advertising

RECOMMENDATION ENGINE



**Next Best and
Personalized Offers**

Media Entertainment

CONSUMER ENGAGEMENT
ANALYSIS



Sentiment Analysis

Perform gap analysis

[illegible]

The Azure data landscape



Azure Data Factory



Azure Import/Export service



Azure CLI



Azure SDK



Azure IoT Hub



Azure event hubs



Kafka on Azure HDInsight



Azure SQL DB



Azure Cosmos DB



Azure SQL data warehouse



Azure Analysis Services



Power BI



Azure Blob Storage



Azure Data Lake Store



Azure HDInsight



Azure Databricks



Azure ML



ML Server



Azure Databricks



Azure Search



Azure Data Catalog



Azure Stream Analytics



Azure HDInsight



Azure Databricks



Bot service



Cognitive services



Azure ExpressRoute



Azure Active Directory



Azure network security groups



Azure key management service



Operations Management Suite



Azure Functions

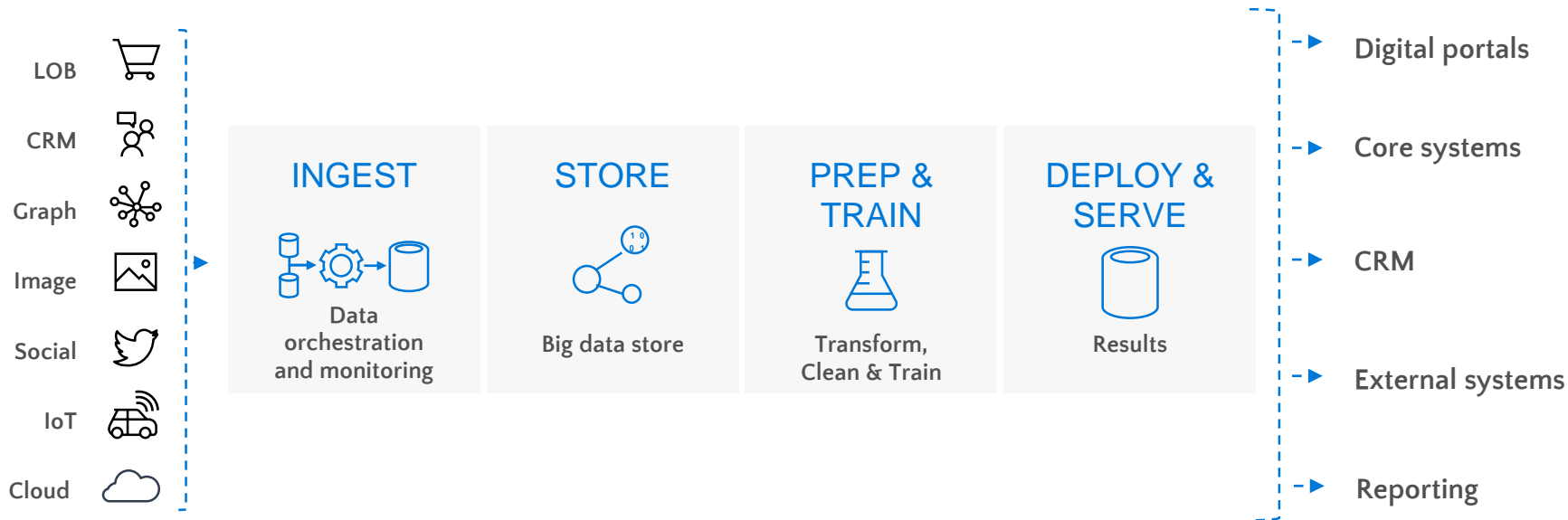


Visual Studio

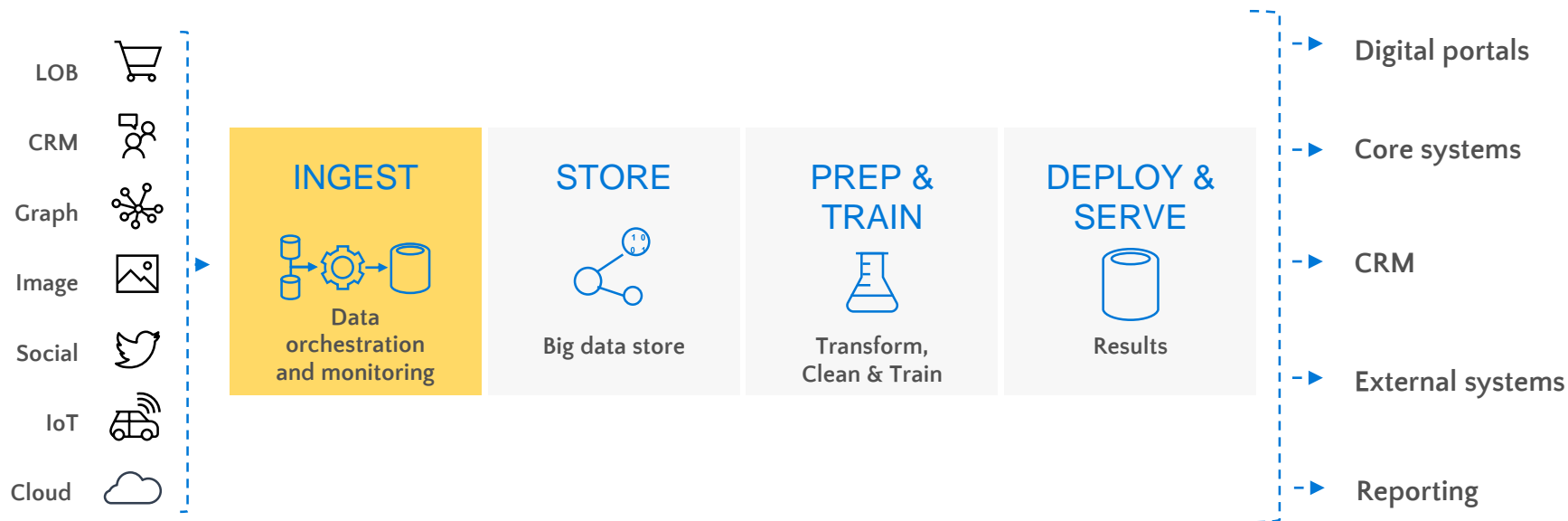
Our job, pretty much



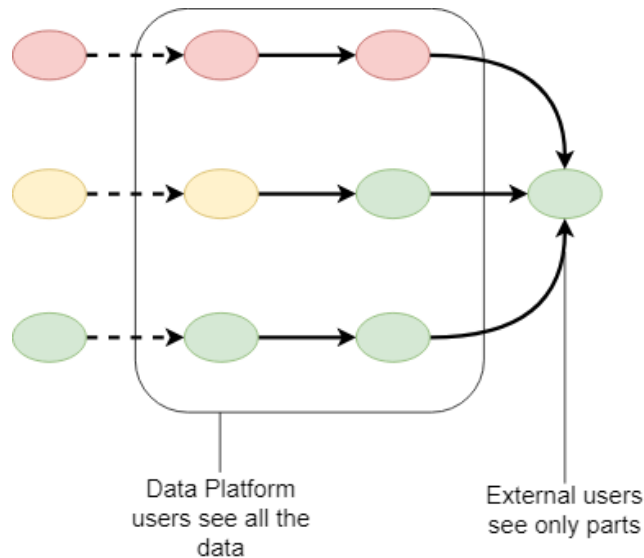
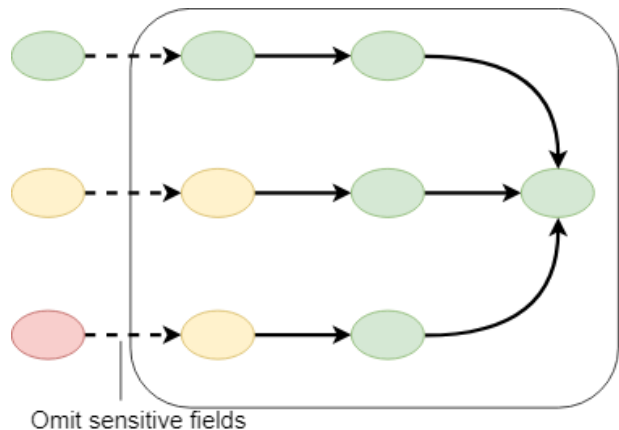
Getting things done



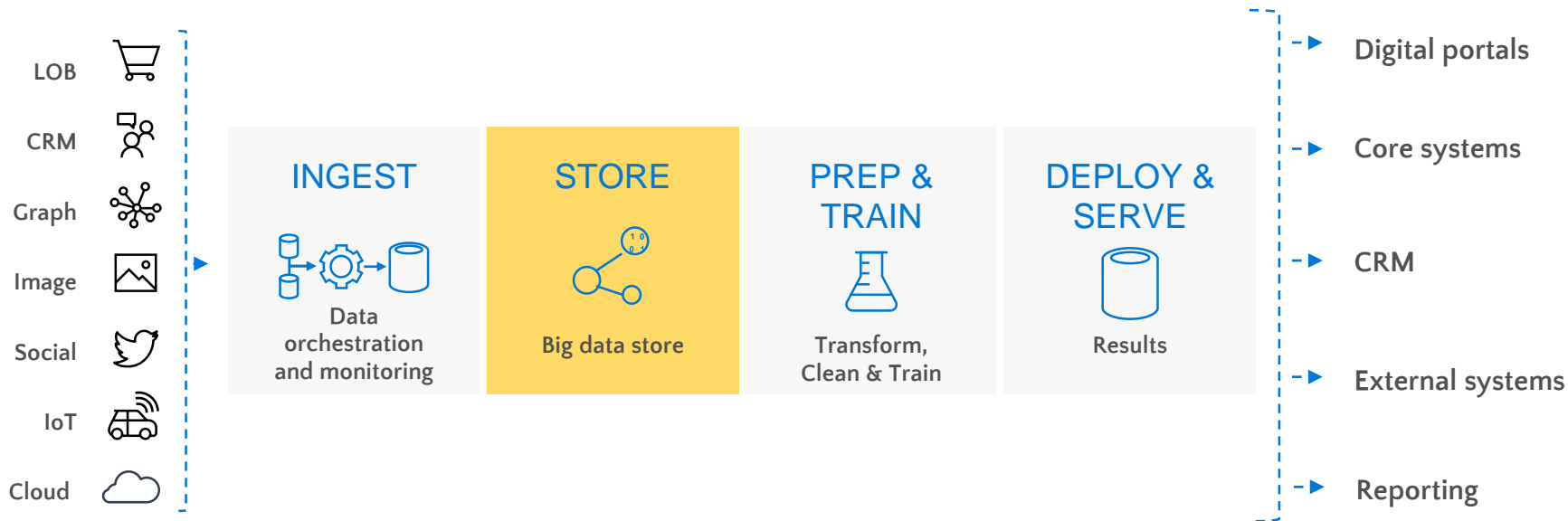
Ingest data from source systems to a data lake



Lesson 1: Build privacy protection patterns



Store data in its native format



Lesson 2: Use cloud storage offerings instead of Hadoop

Machine Learning & Big Data Blog

Is Hadoop Dead? How Kubernetes and Cloud-Native Could Displace Hadoop



The Death of Hadoop?

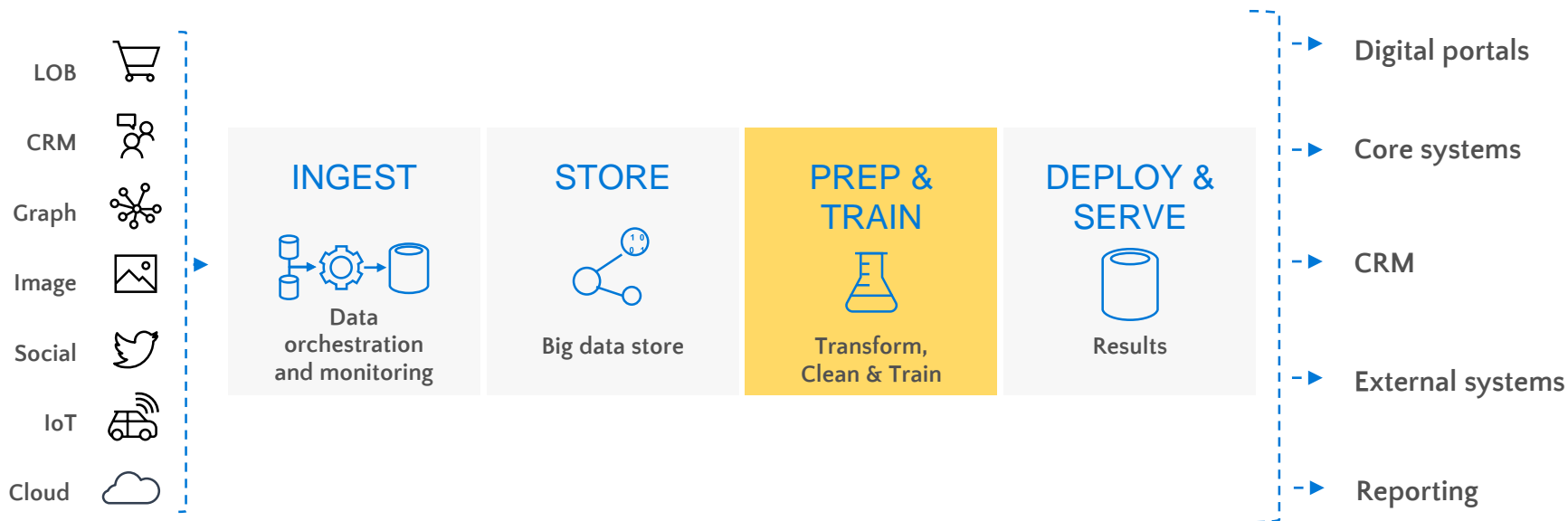
Is Hadoop dead? Not so fast. Plan on supporting multiple environments for some time to come.

By [Barry Devlin](#)

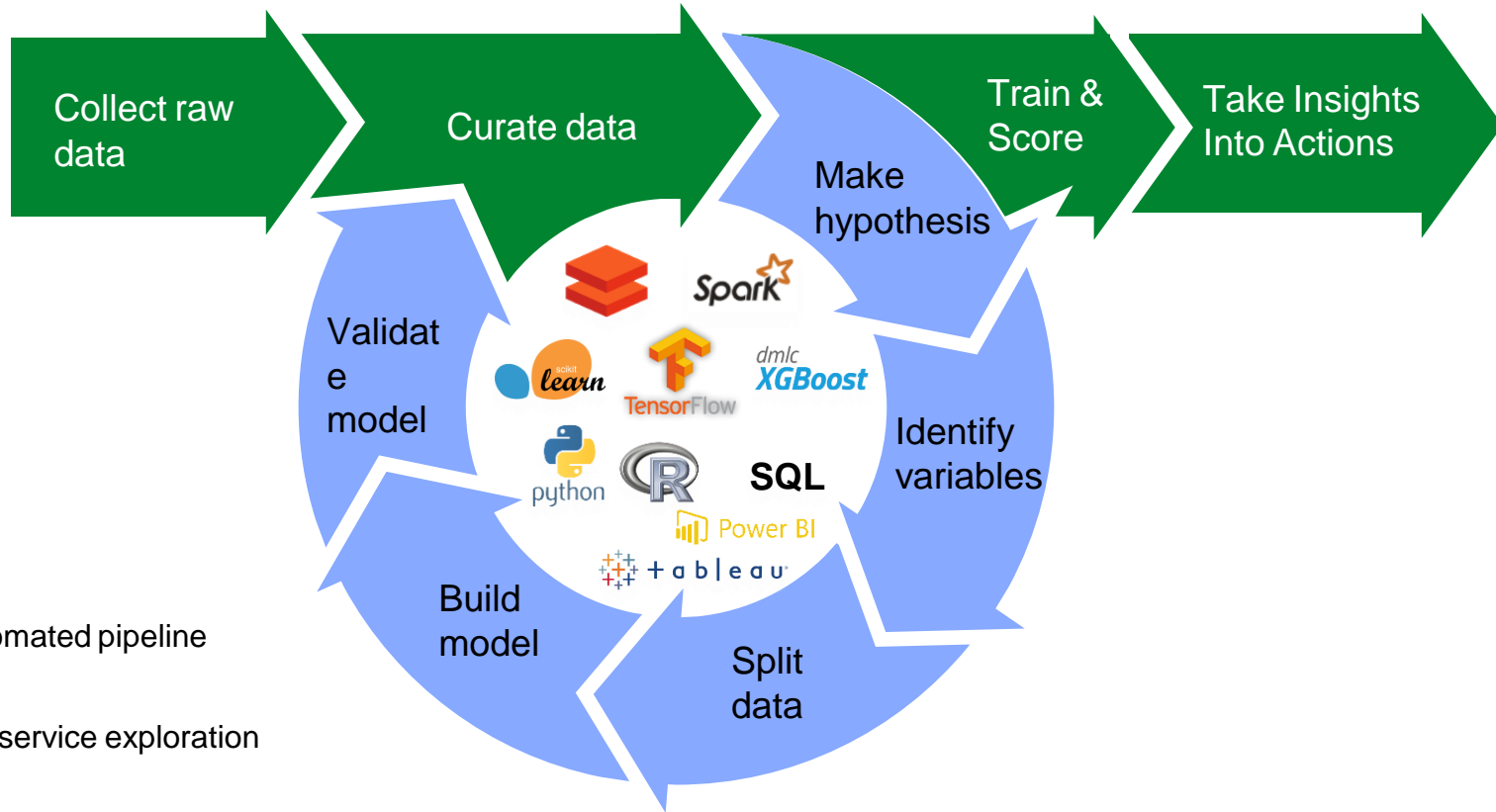
February 26, 2019

The recent "merger of equals" between Cloudera and Hortonworks has triggered speculation about the possible imminent demise of Hadoop. Market observers question if the merger indicates a shrinking Hadoop ecosystem market that can no longer support its two largest competing beasts.








Provide tools and resources to prep & train




Lesson 3: Offer self-service tools



Azure Machine Learning Studio



Search experiment items 

▶ Saved Datasets

▶ Trained Models

▶ Transforms

▶ Data Format Conversions

▶ Data Input and Output

▶ Data Transformation

▶ Feature Selection

▶ Machine Learning

▶ OpenCV Library Modules

▶ Python Language Modules

▶ R Language Modules


▶ Statistical Functions


▶ Text Analytics


▶ Time Series


▶ Web Service


Binary Classification: Direct marketing


Finished running 


Draft saved at 12:38:31 


Import Data 


Edit Metadata  ▼


Select Columns in Dataset  ▼


Split Data 


Two-Class Boosted Decision... 

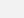
Two-Class Support Vector ... 

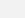
Split Data 


Split Data 



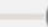
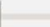
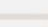


Tune Model Hyperparameters 

Tune Model Hyperparameters 

Score Model 

Score Model 

Evaluate Model 



Azure Databricks

Azure Databricks

Home

Workspace

Recents

Data

Clusters

Jobs

Search

CarrierDelay (Python)

7

Attached: thecluster File View: Code Permissions Stop Execution Clear

Schedule Comments Runs Revision history

```
1 from pyspark.ml.feature import OneHotEncoderEstimator, StringIndexer
2 from pyspark.ml.feature import VectorAssembler
3 from pyspark.ml import Pipeline
4 from pyspark.ml.regression import LinearRegression
5 from pyspark.ml.evaluation import RegressionEvaluator
6 from pyspark.ml.tuning import CrossValidator, ParamGridBuilder
7
8 # Delete (if exists) parquet folder
9 dbutils.fs.rm("/mnt/parquet/carrierdelay", True)
10
11 carrierDF = (spark.read.format('csv')
12             .options(header='true', inferschema='true')
13             .load("/mnt/flightdelay3/*.csv"))
14
15 carrierDF = (carrierDF.withColumn("QUARTER", carrierDF["QUARTER"].cast("double"))
16             .withColumn("MONTH", carrierDF["MONTH"].cast("double"))
17             .withColumn("DAY_OF_MONTH", carrierDF["DAY_OF_MONTH"].cast("double"))
18             .withColumn("DAY_OF_WEEK", carrierDF["DAY_OF_WEEK"].cast("double"))
19             .withColumn("DISTANCE_GROUP", carrierDF["DISTANCE_GROUP"].cast("double")))
20
21 carrierDF = carrierDF.drop('DEST_AIRPORT_ID', '_c17', 'YEAR', 'DEP_DELAY').na.drop()
22
23 # carrierDF.printSchema()
24
25 carrierDF.write.parquet("/mnt/parquet/carrierdelay")
26
27 # carrierDF = spark.read.format('parquet').options(header='true', inferschema='true').load("/mnt/parquet/carrierdelay")
28 print("Number of carrier flights in the database: ", carrierDF.count())
```

Lesson 4: Use on-demand resources



Storage



10 TB
+ 24/7
= ~300 Eur/month

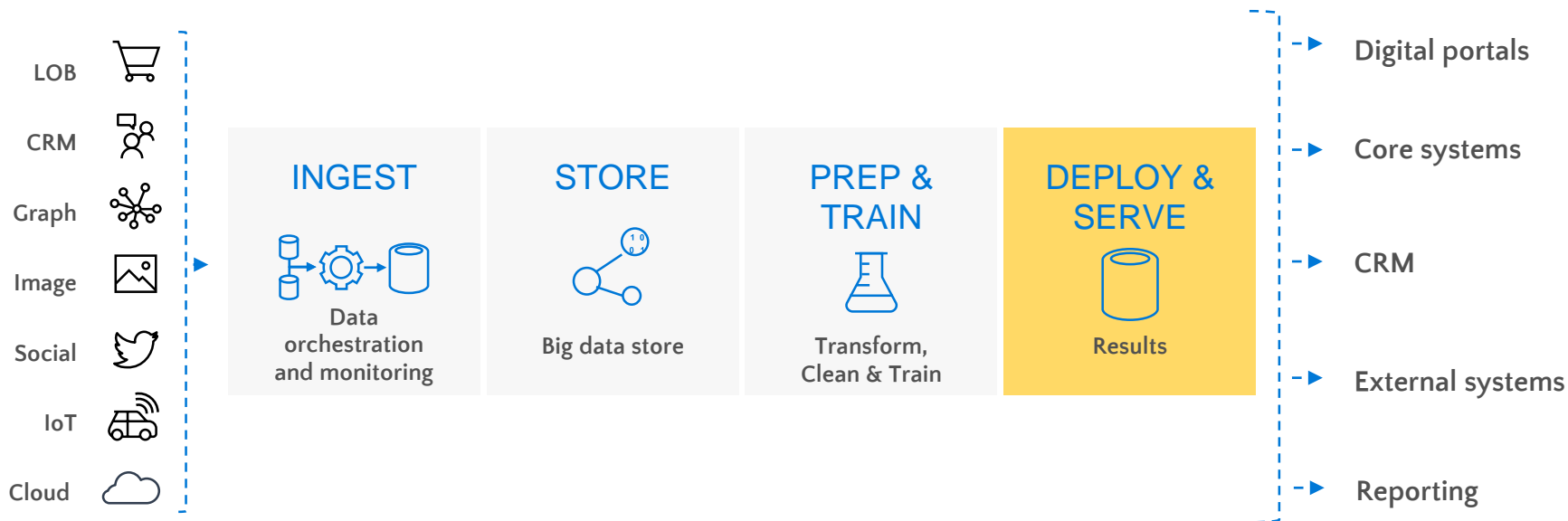


databricks

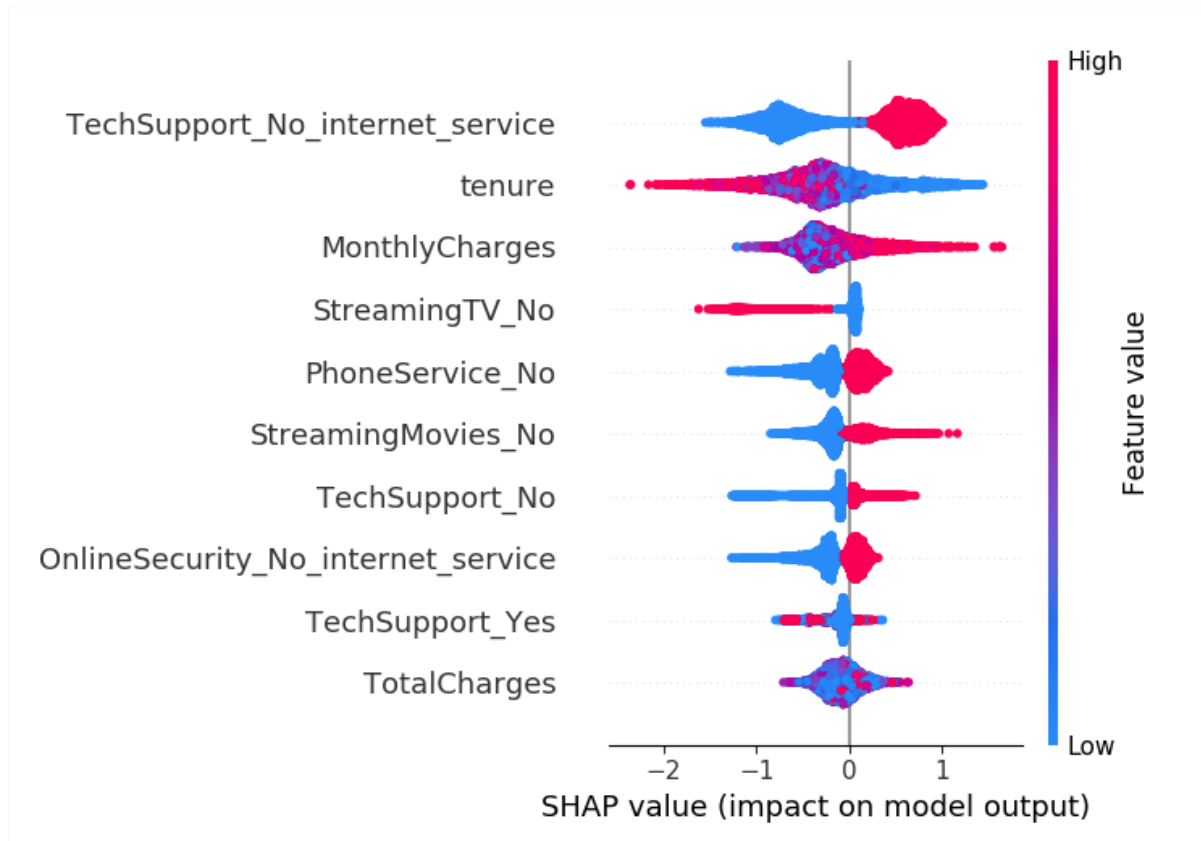


32 cores / 128 GB RAM
+ 160 hours
= ~700 Eur/month

Serve results to end consumers



Lesson 5: Explain your models to business users (e.g. SHAP)



Key Takeaways

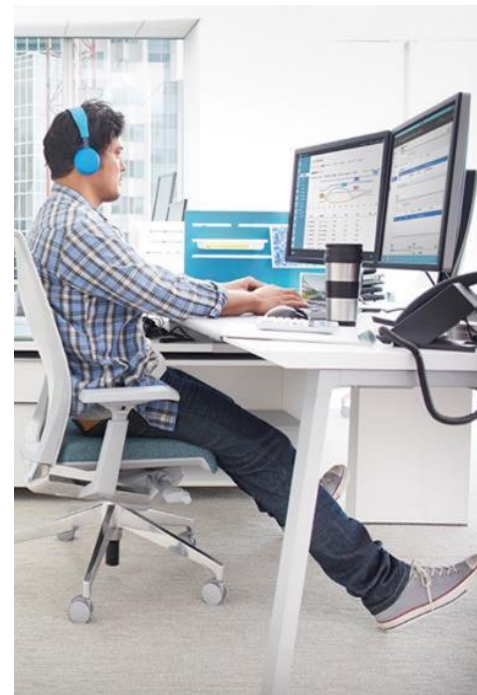
Lesson 1 Build privacy protection patterns

Lesson 2 Use cloud storage instead of Hadoop

Lesson 3 Offer self-service tools

Lesson 4 Use on-demand resources

Lesson 5 Explain ML models to business



Questions?

Valdas Maksimavičius

valdas@maksimavicius.eu