IT Home

IT HOME
CONFERENCE

*Homemade cookies are always better than Bought*

# Gregory Kulga
## Data Quality Engineer

- Core skills: DWBI, data quality, data analysis
- 3-year experience in Big data testing, Automation testing, Data analysis
- Hobbies: quizzes, camping, snowboard
That's not me. Not yet...

Agenda

Big Data paradigm

Project structure and Targets

DQ Goals and opportunities

60 mins

Overall Estimation

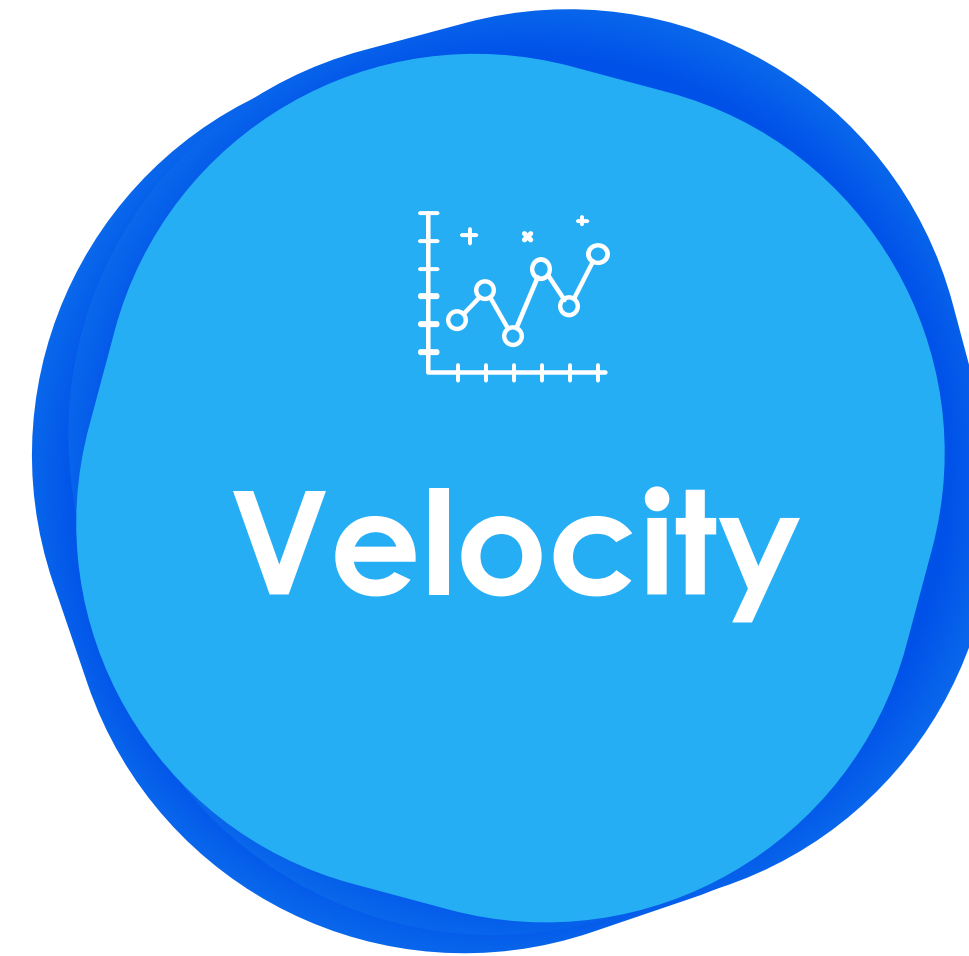AMAZON Deequ As a rapid tool

Apache Griffin as a COMPLEX tool

IT Home

# Big Data paradigm

# Big Data paradigm
## On what it is based

IT Home

## Volume

The quantity of generated and stored data

## Velocity

The speed at which the data is generated and processed

## Variety

Type and nature of the data.
Semi-structured and unstructured

# Project Structure

# Our customer – Lead American health insurance provider

Business Objective is to predict and assess the possible increase or decrease diseases among specific groups of people
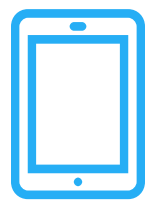
These Groups of people are divided by
- age
- gender
- race
- genetic diseases
- features of growth and maturation
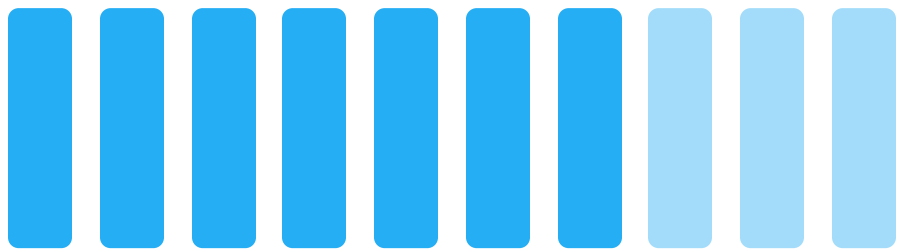- features of the working environment

# Machine Learning model stages

The based data split up for 3 stages

IT Home
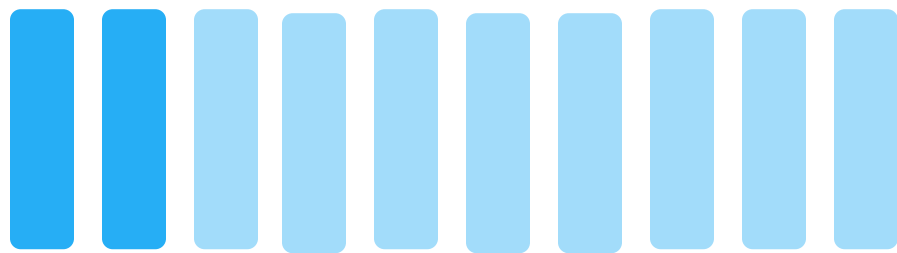
### Training Set
Prediction is created

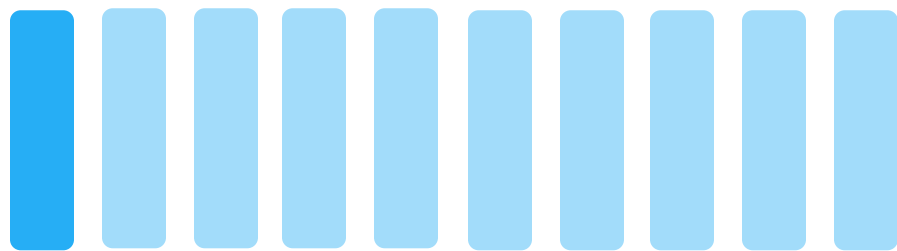**70%**

### Validation Set
Model check for stability

**20%**

### Test Set
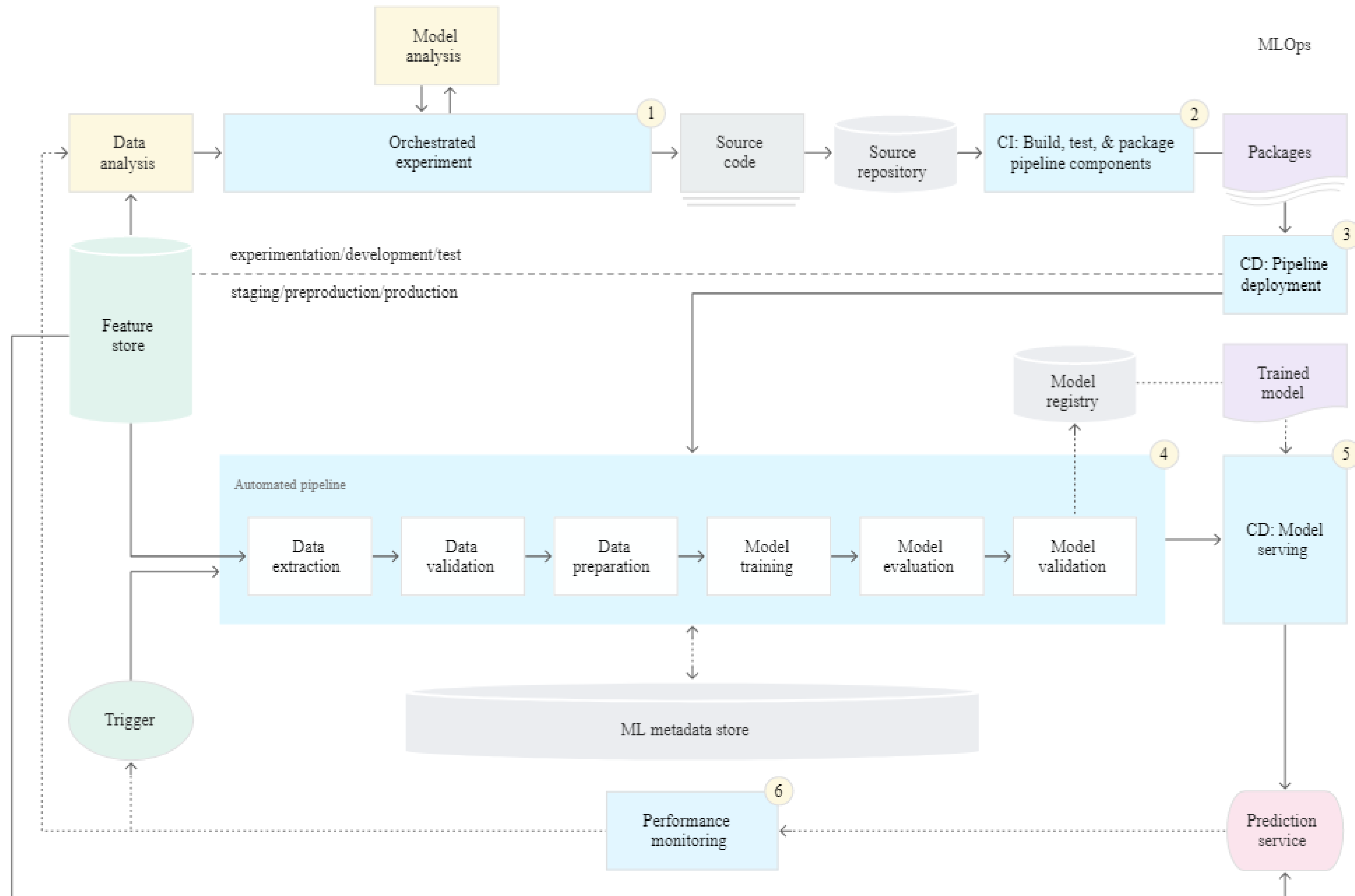Model check as prediction and compared with previous data

**10%**

# Data train model scheme
## If Something goes wrong Let's retrain!

IT Home

MLOps

Model analysis

Data analysis

Orchestrated experiment  ①

Source code

Source repository

CI: Build, test, & package pipeline components  ②

Packages

experimentation/development/test

staging/preproduction/production

CD: Pipeline deployment  ③

Feature store

Model registry

Trained model

**Automated pipeline**  ④

Data extraction → Data validation → Data preparation → Model training → Model evaluation → Model validation

CD: Model serving  ⑤

ML metadata store

Trigger

Performance monitoring  ⑥

Prediction service

Manual Stack

How should it work?

Automation Workflow

# ML Pipeline Architecture

NOTEBOOK REQUEST

CloudFormation

Source Data (Lake)

Prediction Consumers

(or Jenkins)

AWS CodePipeline

S3

## Experiment

SageMaker Notebook

New Data Sets

.ipynb

Bitbucket

## Model Packager

Lambda (serverless) OR Fargate (containers)

Model Package / Artifacts

## Model Trainer

Training Job

Data: Prep — Extract — Validate

Model: Train — Evaluate — Validate

Data Catalog

TBD

Testing Stack

Trained Model

Model Mart

## Model Hosting

Hot Endpoint — Cold Endpoint — Batch Transform

## Model Monitor

TBD

Model Retrain

---

**Legend**

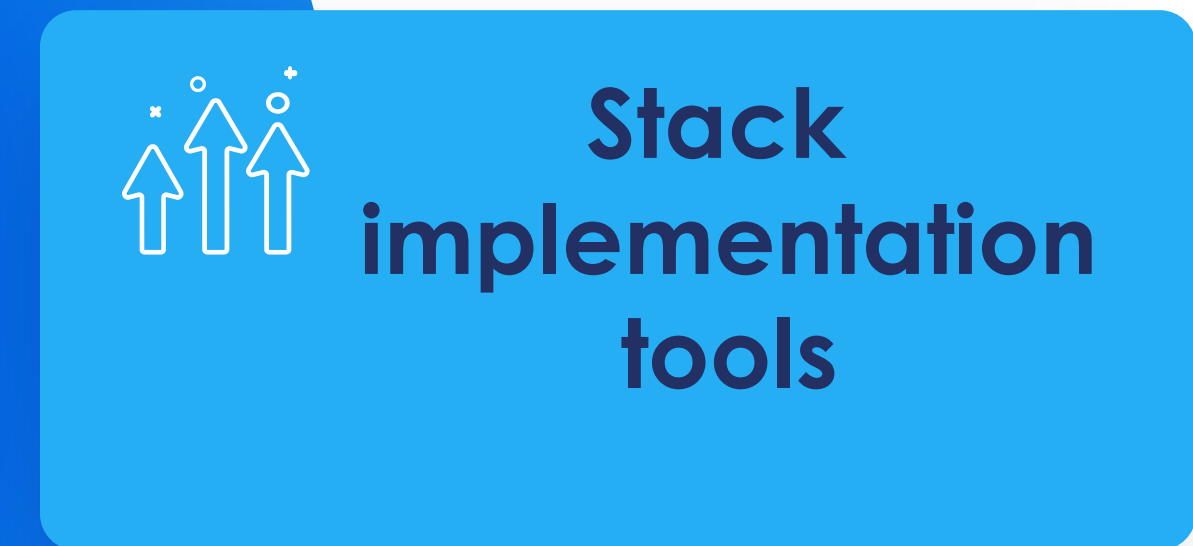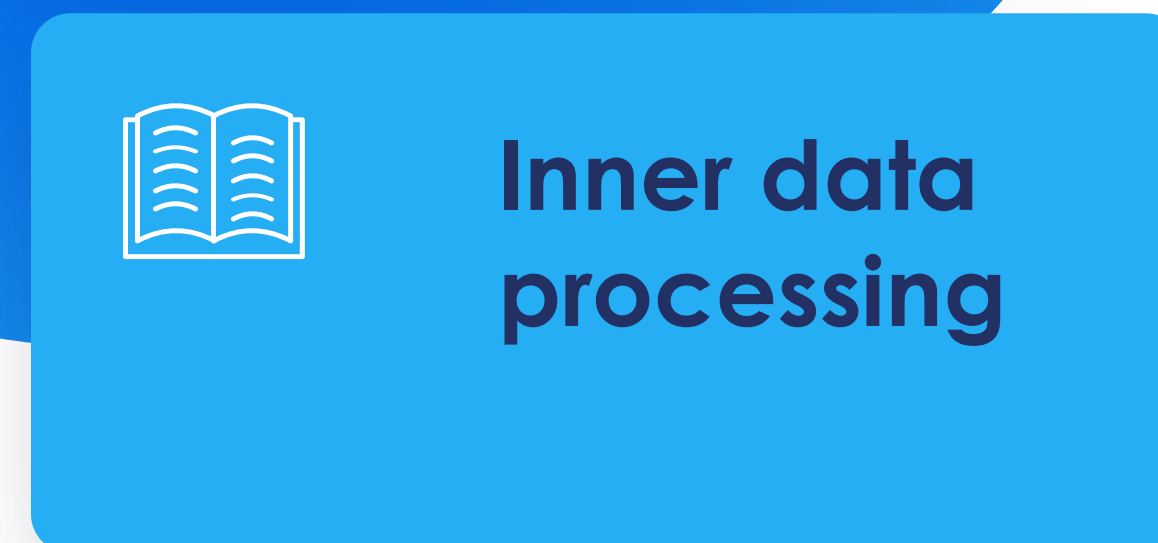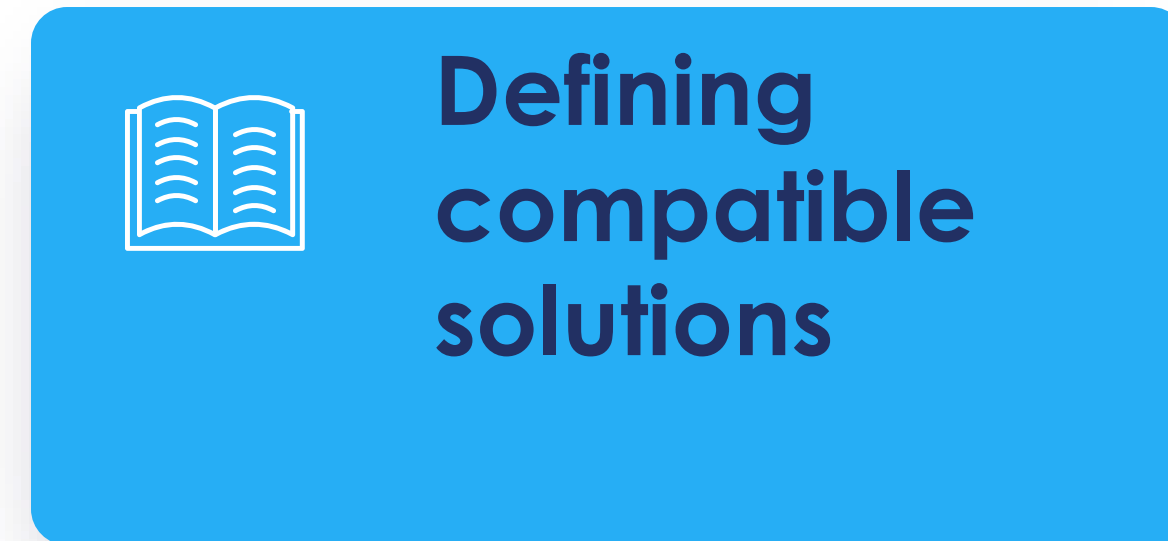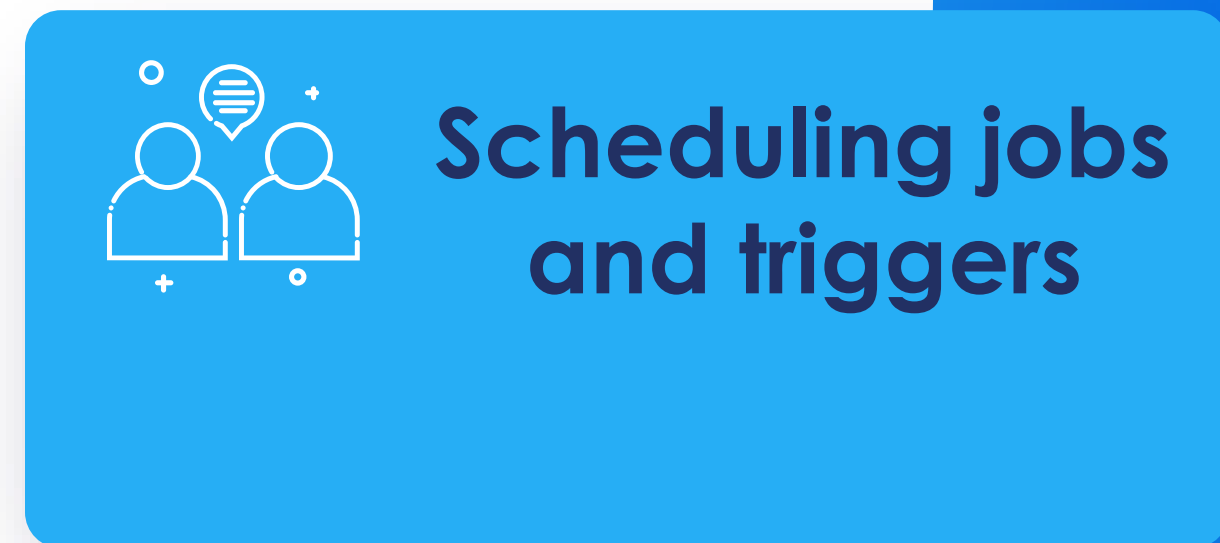| Icon | Label |
|------|-------|
| | Experiment Request |
| | CloudFormation Script |
| | S3 Data Bucket |
| | Experimental Environment |
| | Bitbucket Repository |
| | CodePipeline |
| | Model Packager |
| | Model Trainer |
| | Data Catalog |
| | Model Mart |
| | Model Hosting |
| | Model Monitor |
| | Model Retraining |

# DQ Goals and opportunities

Consistency
Sed non felis non leo pulvinar dictum vitae eu elit.

Timeliness
Sed non felis non leo pulvinar dictum vitae eu elit.

**Uniqueness**

DESCRIPTION

**DQ vision**

**Completeness**

DESCRIPTION

Accuracy
Sed non felis non leo pulvinar dictum vitae eu elit. Sed euismod id risus quis scelerisque.

**Validity**

DESCRIPTION

# Amazon Deequ

# What could be better than Amazon (Deequ)?



You're exceptionally awesome but we at least exist in this reality

Here we go!

**Amazon Deequ**

Deequ Implementation Scheme

Functions and jar-file

Pros and Cons

POC

Scala code structure

Output example

# Deequ Implementation Scheme



Sequence execution

- Source data is stored in AWS Data Lake S3

- Load and install .jar-files for Deequ libs

- Save scala code in S3 via Jupyter Notebook cells executed on spark-shell

- Use EMR Terminal via Spark-shell to execute stored .scala
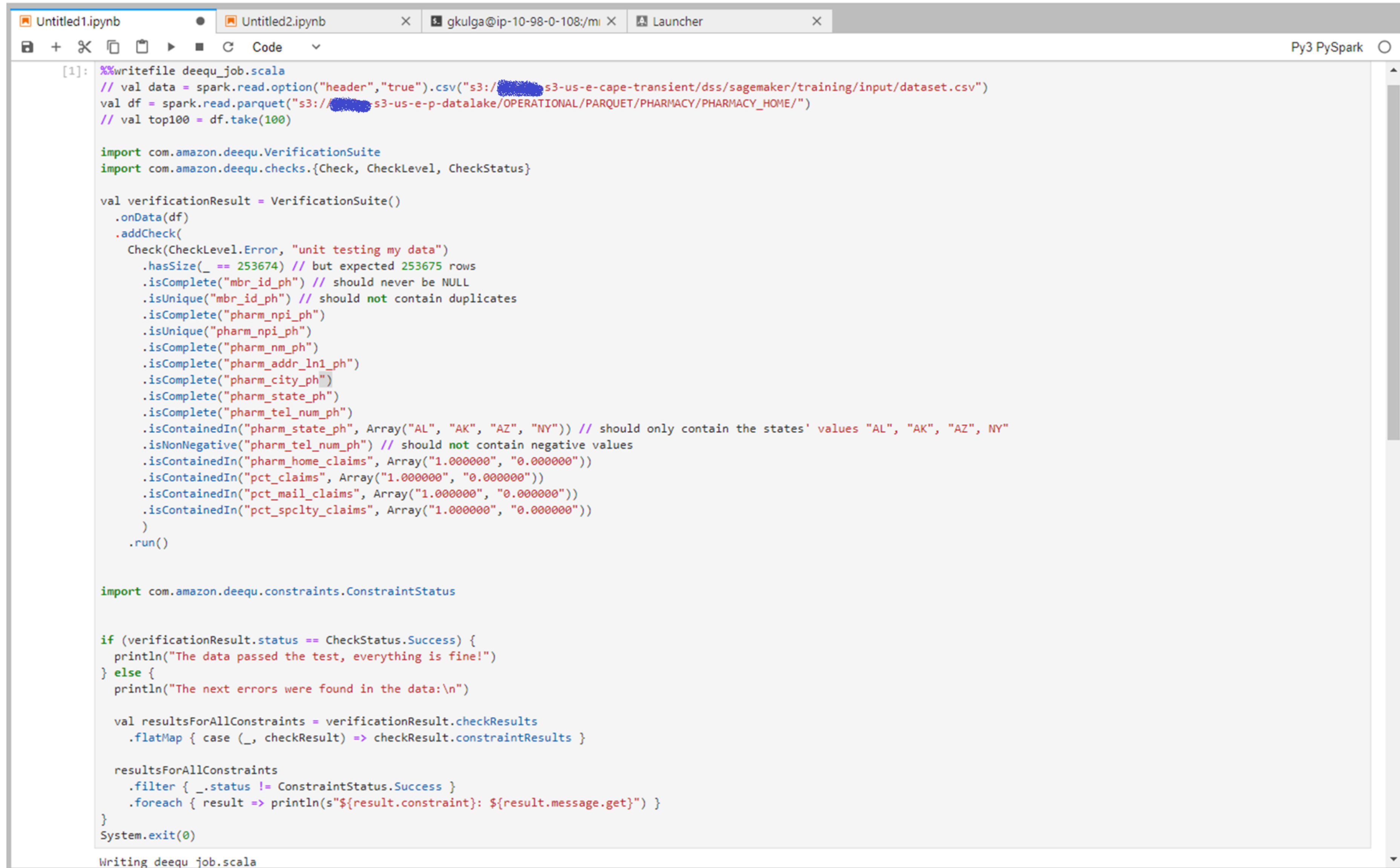
- Get the result into EMR Terminal

# Functions Variety

| Metric | Description | Usage Example |
|---|---|---|
| ApproxCountDistinct | Approximate number of distinct value, computed with HyperLogLogPlusPlus sketches. | ApproxCountDistinct("review_id") |
| ApproxQuantile | Approximate quantile of a distribution. | ApproxQuantile("star_rating", quantile = 0.5) |
| ApproxQuantiles | Approximate quantiles of a distribution. | ApproxQuantiles("star_rating", quantiles = Seq(0.1, 0.5, 0.9)) |
| Completeness | Fraction of non-null values in a column. | Completeness("review_id") |
| Compliance | Fraction of rows that comply with the given column constraint. | Compliance("top star_rating", "star_rating >= 4.0") |
| Correlation | Pearson correlation coefficient, measures the linear correlation between two columns. The result is in the range [-1, 1], where 1 means positive linear correlation, -1 means negative linear correlation, and 0 means no correlation. | Correlation("total_votes", "star_rating") |
| CountDistinct | Number of distinct values. | CountDistinct("review_id") |
| DataType | Distribution of data types such as Boolean, Fractional, Integral, and String. The resulting histogram allows filtering by relative or absolute fractions. | DataType("year") |
| Distinctness | Fraction of distinct values of a column over the number of all values of a column. Distinct values occur at least once. Example: [a, a, b] contains two distinct values a and b, so distinctness is 2/3. | Distinctness("review_id") |
| Maximum | Maximum value. | Maximum("star_rating") |
| Mean | Mean value; null values are excluded. | Mean("star_rating") |
| Minimum | Minimum value. | Minimum("star_rating") |

# Deequ scala code



```scala
%%writefile deequ_job.scala
// val data = spark.read.option("header","true").csv("s3://xxxxxxx-s3-us-e-cape-transient/dss/sagemaker/training/input/dataset.csv")
val df = spark.read.parquet("s3://xxxxxx-s3-us-e-p-datalake/OPERATIONAL/PARQUET/PHARMACY/PHARMACY_HOME/")
// val top100 = df.take(100)

import com.amazon.deequ.VerificationSuite
import com.amazon.deequ.checks.{Check, CheckLevel, CheckStatus}

val verificationResult = VerificationSuite()
  .onData(df)
  .addCheck(
    Check(CheckLevel.Error, "unit testing my data")
      .hasSize(_ == 253674) // but expected 253675 rows
      .isComplete("mbr_id_ph") // should never be NULL
      .isUnique("mbr_id_ph") // should not contain duplicates
      .isComplete("pharm_npi_ph")
      .isUnique("pharm_npi_ph")
      .isComplete("pharm_nm_ph")
      .isComplete("pharm_addr_ln1_ph")
      .isComplete("pharm_city_ph")
      .isComplete("pharm_state_ph")
      .isComplete("pharm_tel_num_ph")
      .isContainedIn("pharm_state_ph", Array("AL", "AK", "AZ", "NY")) // should only contain the states' values "AL", "AK", "AZ", NY"
      .isNonNegative("pharm_tel_num_ph") // should not contain negative values
      .isContainedIn("pharm_home_claims", Array("1.000000", "0.000000"))
      .isContainedIn("pct_claims", Array("1.000000", "0.000000"))
      .isContainedIn("pct_mail_claims", Array("1.000000", "0.000000"))
      .isContainedIn("pct_spclty_claims", Array("1.000000", "0.000000"))
    )
    .run()


import com.amazon.deequ.constraints.ConstraintStatus


if (verificationResult.status == CheckStatus.Success) {
  println("The data passed the test, everything is fine!")
} else {
  println("The next errors were found in the data:\n")

  val resultsForAllConstraints = verificationResult.checkResults
    .flatMap { case (_, checkResult) => checkResult.constraintResults }

  resultsForAllConstraints
    .filter { _.status != ConstraintStatus.Success }
    .foreach { result => println(s"${result.constraint}: ${result.message.get}") }
}
System.exit(0)

Writing deequ_job.scala
```

IT Home

# Results via Amazon Deequ

IT Home

```
20/04/06 17:02:39 WARN Utils: Truncated the string representation of a plan since it was too large. This behavior can be adjusted by setting 'spark.debug.maxToStringFields' in SparkEnv.conf.
The next errors were found in the data:

SizeConstraint(Size(None)): Value: 255021 does not meet the constraint requirement!
UniquenessConstraint(Uniqueness(List(pharm_npi_ph))): Value: 0.00861889805153301 does not meet the constraint requirement!
ComplianceConstraint(Compliance(pharm_state_ph contained in AL,AK,AZ,NY,`pharm_state_ph` IS NULL OR `pharm_state_ph` IN ('AL','AK','AZ','NY'),None)): Value: 0.9786252896820262 does not meet the constraint requirement!
ComplianceConstraint(Compliance(pharm_tel_num_ph is non-negative,COALESCE(pharm_tel_num_ph, 0.0) >= 0,None)): Value: 0.002403723614917987 does not meet the constraint requirement!
ComplianceConstraint(Compliance(pharm_home_claims contained in 1.000000,0.000000,`pharm_home_claims` IS NULL OR `pharm_home_claims` IN ('1.000000','0.000000'),None)): Value: 0.0 does not meet the constraint requirement!
ComplianceConstraint(Compliance(pct_claims contained in 1.000000,0.000000,`pct_claims` IS NULL OR `pct_claims` IN ('1.000000','0.000000'),None)): Value: 0.715592049282216 does not meet the constraint requirement!
ComplianceConstraint(Compliance(pct_mail_claims contained in 1.000000,0.000000,`pct_mail_claims` IS NULL OR `pct_mail_claims` IN ('1.000000','0.000000'),None)): Value: 0.9480434944573192 does not meet the constraint requirement!
ComplianceConstraint(Compliance(pct_spclty_claims contained in 1.000000,0.000000,`pct_spclty_claims` IS NULL OR `pct_spclty_claims` IN ('1.000000','0.000000'),None)): Value: 0.9854051234996334 does not meet the constraint requirement!
[gkulga@ip-10-98-0-108 gkulga]$
```

*Results could be stored in S3*

C
b
c

# Amazon Deequ Pros and Cons

## Pros

- **Small and easy-supportable library**

- **SQL-like functions supported via scala syntax**

- **Fast deliverable code: one script – one file – one run**

## Cons

- **Console output**

- **No documentation with alerts, time-postponed functions**

- **Poorly configurable options**

IT Home

IT HOME
CONFERENCE

# Random squirrel

- If you tired - look right

- If you distracted - look right
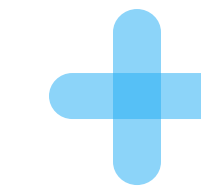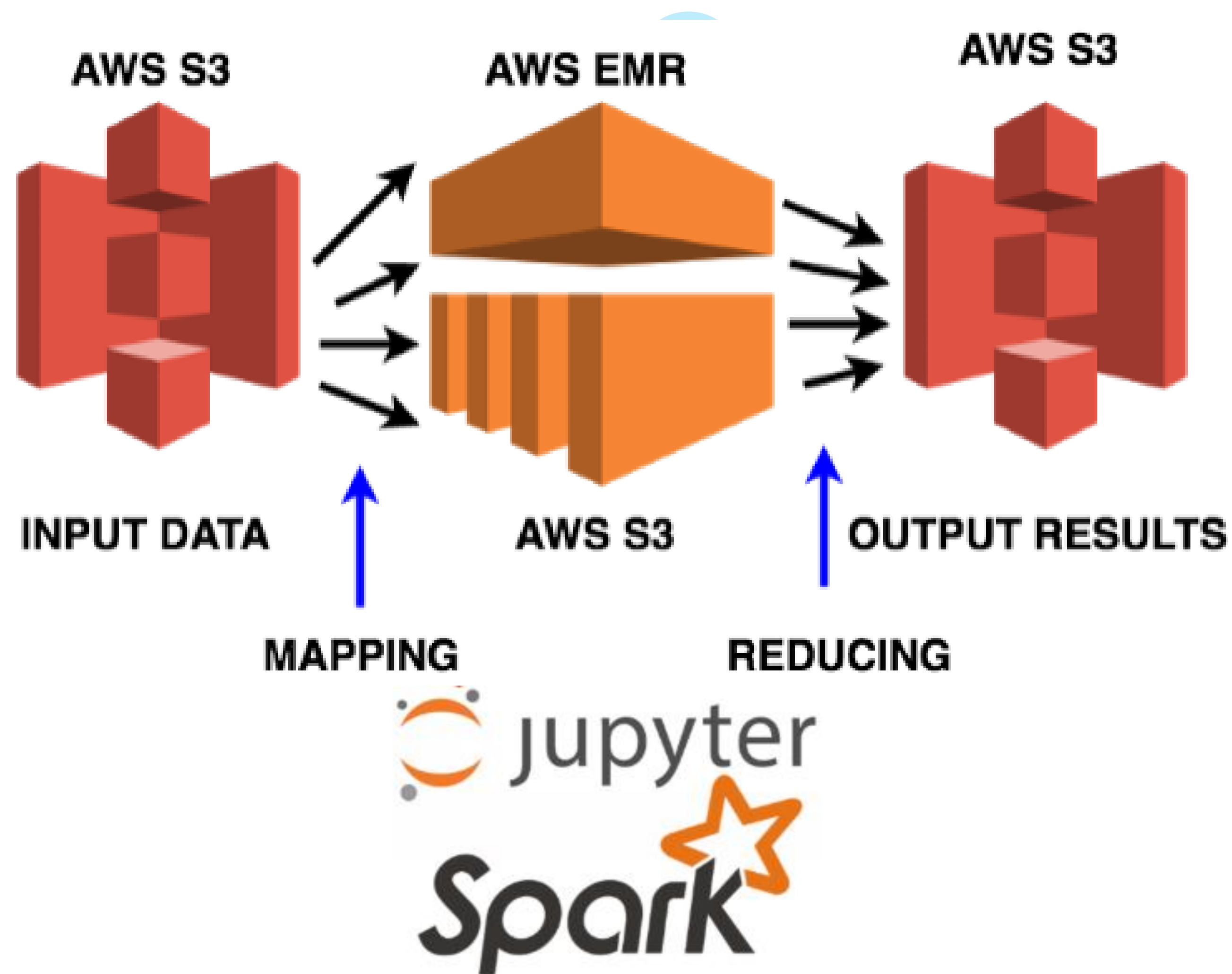
Don't be rude – wave back!

# Apache Griffin

# Apache Griffin

- Production implementation schema

- Configuration files structure

- Pros and Cons

# Griffin Implementation Scheme



Sequence execution

- Source data is stored in AWS Data Lake S3

- Load and install .jar-files for Griffin libs

- Save json-files in S3 via Jupyter Notebook cells executed on spark-shell

- Use EMR Terminal via Spark-shell to execute stored jsons

- Get the result into S3 storage or console output

# Possible Features

IT Home

- Batch and Streaming are included both

- HDFS, Cloud, Console, Elasticsearch ways as options for output

- Combining functions rule created

- Scheme must be described: in mismatching case run will fail with relevant error

Environment Parameters

```
{
  "spark": {
    "log.level": "WARN",
    "checkpoint.dir": "hdfs:///griffin/streaming/cp",
    "batch.interval": "1m",
    "process.interval": "5m",
    "config": {
      "spark.default.parallelism": 5,
      "spark.task.maxFailures": 5,
      "spark.streaming.kafkaMaxRatePerPartition": 1000,
      "spark.streaming.concurrentJobs": 4,
      "spark.yarn.maxAppAttempts": 5,
      "spark.yarn.am.attemptFailuresValidityInterval": "1h",
      "spark.yarn.max.executor.failures": 120,
      "spark.yarn.executor.failuresValidityInterval": "1h",
      "spark.hadoop.fs.hdfs.impl.disable.cache": true
    }
  },
  "sinks": [
    {
      "type": "console",
      "config": {
        "max.log.lines": 100
      }
    }, {
      "type": "hdfs",
      "config": {
        "path": "hdfs:///griffin/streaming/persist",
        "max.lines.per.file": 10000
      }
    }
  ],
  "griffin.checkpoint": [
    {
      "type": "zk",
      "config": {
        "hosts": "<zookeeper host ip>:2181",
        "namespace": "griffin/infocache",
        "lock.path": "lock",
        "mode": "persist",
        "init.clear": true,
        "close.clear": false
      }
    }
  ]
}
```

# :Data.json

```
1  %%writefile dq_hdfs_mod.json
2  {
3    "name": "uniq_batch",
4    "process.type": "BATCH",
5    "data.sources": [
6      {
7        "name": "src",
8        "connector": {
9            "type": "parquet",
10           "config": {
11             "file.name": "s3://        s3-us-e-p-datalake/OPERATIONAL/PARQUET/PHARMACY/PHARMACY_HOME/part-00000-f2fe82ae-258c-4da6-a55f-81a95905231d-c000.gz.parquet",
12             "skipOnError": "false",
13             "schema": [
14               {
15                 "name": "mbr_id_ph",
16                 "type": "string",
17                 "nullable": "true"
18               },
19               {
20                 "name": "pharm_npi_ph",
21                 "type": "string",
22                 "nullable": "false"
23               },
24               {
25                 "name": "pharm_nm_ph",
26                 "type": "string",
27                 "nullable": "false"
28               },
29               {
30                 "name": "pharm_addr_ln1_ph",
31                 "type": "string",
32                 "nullable": "false"
33               },
34               {
35                 "name": "pharm_city_ph",
36                 "type": "string",
37                 "nullable": "false"
38               },
```

"process.type":"BATCH"

Could be Streaming

"file_name":"s3://…"

Source file path

{
"name":"pharm_addr_lnl_ph",
"type":"string",
"nullable":"false"
}

29

# Environment.json

- Structure "env" defined in main file separated to another file and described in it;

- Technology and path according to it should be written;

- Output file size in lines must be defined in structure.

```
1   %%writefile env_hdfs.json
2   {
3     "spark": {
4       "log.level": "WARN"
5     },
6     "sinks": [
7       {
8         "type": "hdfs",
9         "config": {
10          "path": "s3://        s3-us-e-cape-transient/dss/sagemaker/training/output/",
11          "max.lines.per.file": 10000
12        }
13      }
14    ]
15  }
```

Output path

# **Apache Griffin Pros and Cons**

## Pros

- Wide functionality

- Data-driven tool

- Clear mapping

- Ability to create alerts and notifications

- Two deeply configurable jsons

## Cons

- Open-source tool

- Documentation isn't fully updated

- Difficulties with manage and support

IT Home

# Overall Estimation

# Amazon Deequ vs Apache Griffin

## Amazon Deequ

- Small and easy-supportable library

- SQL-like functions syntax via scala

- Fast deliverable code:

  one script – one file – one run

- No alerts and notifications at start

- No mapping

## Apache Griffin

- Complex library

- Structured functions

- Two deeply configurable .json-files:

  first - for data, second - for environment

- Ability to create alerts and notifications

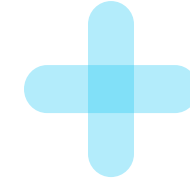- Clear mapping at start

IT Home

# Amazon Deequ vs Apache Griffin

- Amazon Deequ

- Apache Griffin

IT Home



34

Thank you
for attention

IT Home