

Himalayan Peaks of Testing Data Pipelines

Talk about quality of pipelines

Pasha Finkelshteyn, JetBrains

Pasha Finkelshteyn



Developer  for Big Data @ JetBrains

 @asm0di0
 @asm0dey@fosstodon.org

Data processing



kafka

topic: actions

modify data

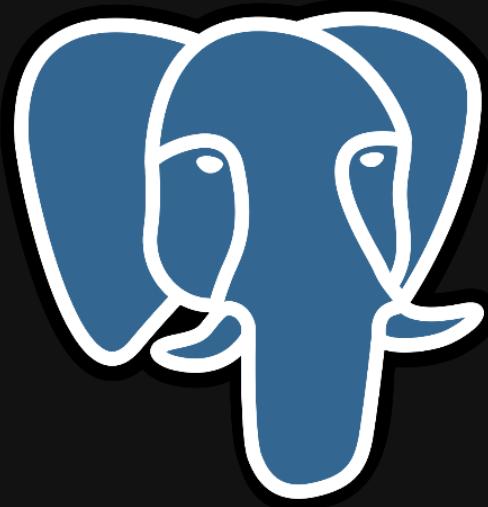
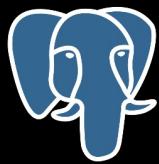
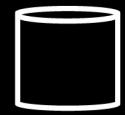


table: actions

Data lake?



PostgreSQL



MySQL™



```
{  
  "data": {  
    "pk": 2,  
    "aa": "text"  
  },  
  "op": "c",  
  "ts_ms":  
  "1580390884335"  
}
```

pk	aa	datetime
2	text	2020-01-01

data	aa_count
2020-01-01	231



Who needs pipelines

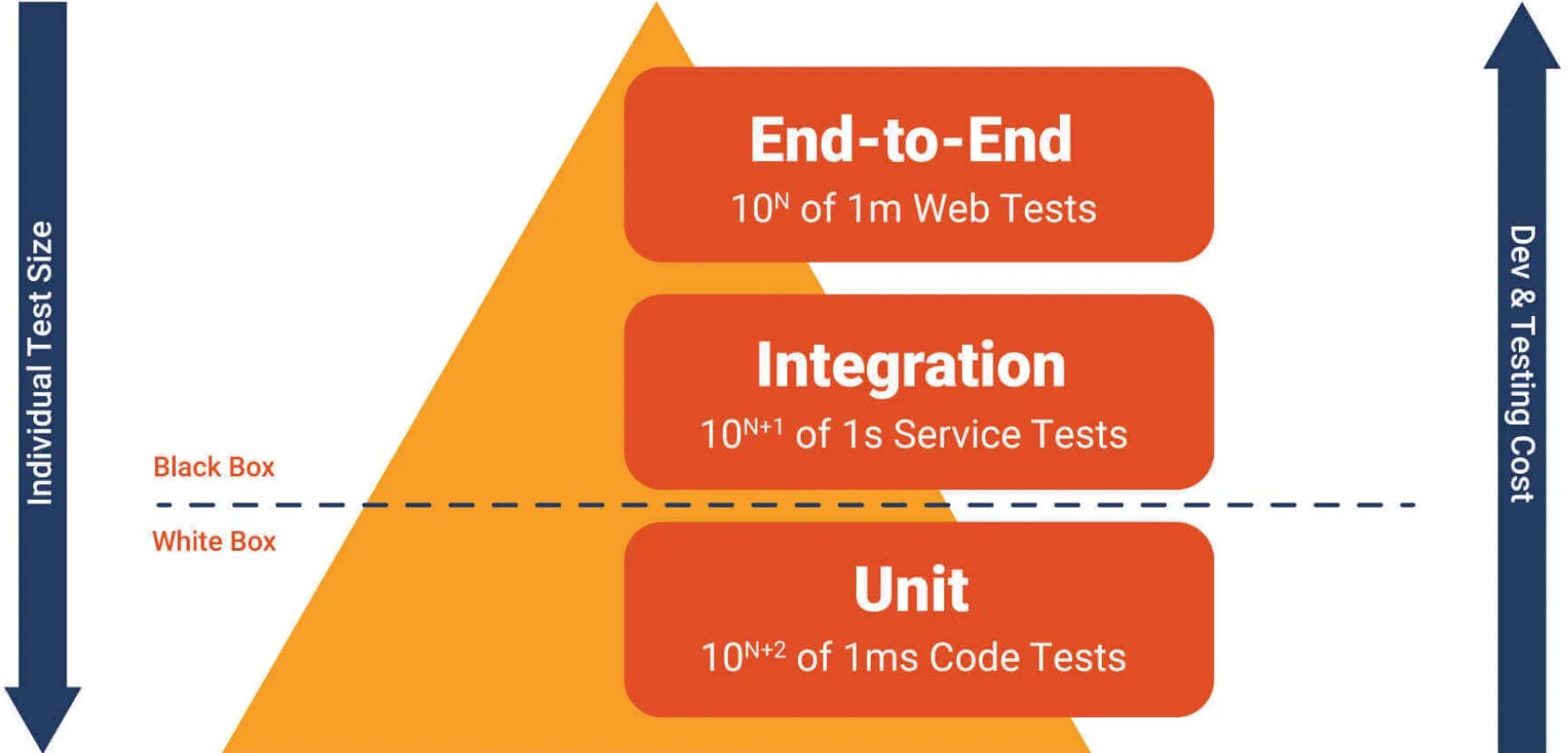
- Data Scientists
- Data Analytics
- Marketing
- PO

It have to be tested!

Pyramid of testing?

Pyramid of testing?

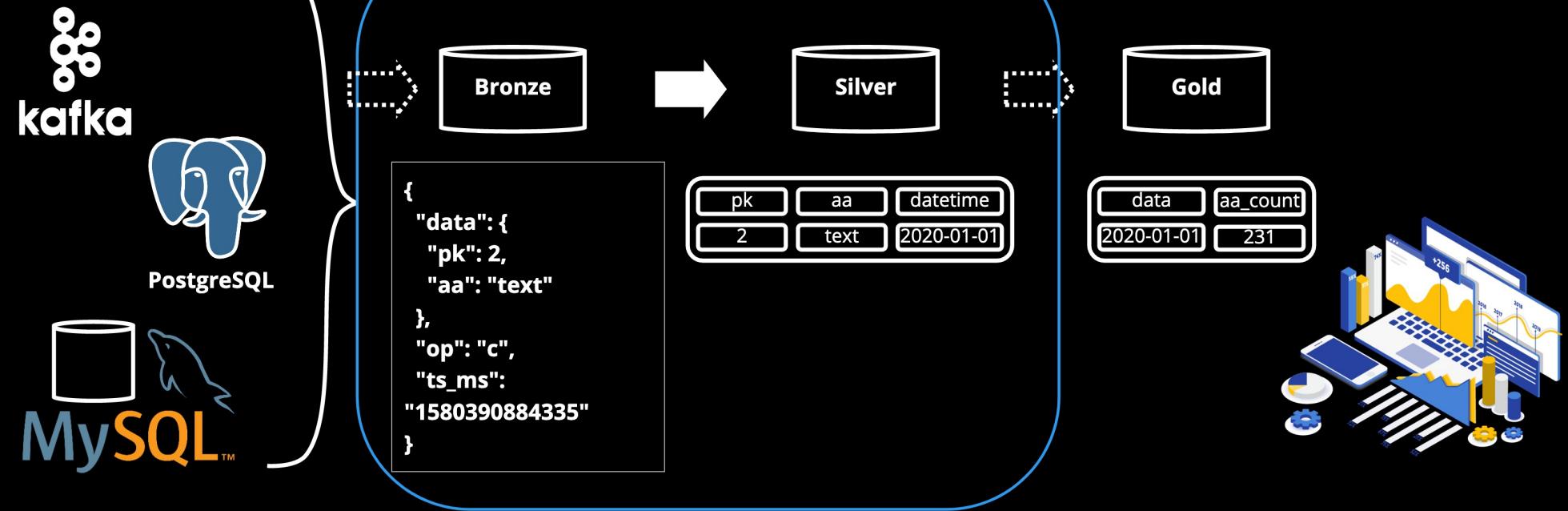




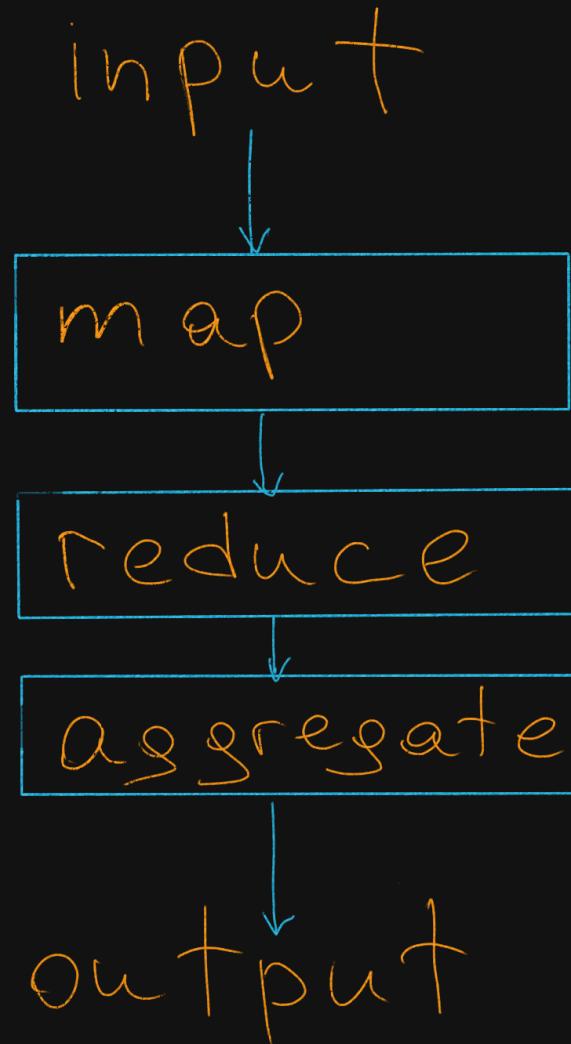
Pyramid of testing. Unit



UNIT



Typical pipeline



Typical pipeline

```
1 StructType schema = new StructType(new StructField[]{  
2     new StructField("pk", DataTuples.LongType, false, MetaData.empty()),  
3     new StructField("aa", DataTuples.StringType, false, MetaData.empty()),  
4 })  
5  
6 spark.read  
7     .schema(schema)  
8     .csv(/* path */)  
9     .map(/* mapper */)  
10    .write  
11    .parquet("path")
```

Unit testing of pipeline

What may we test here?

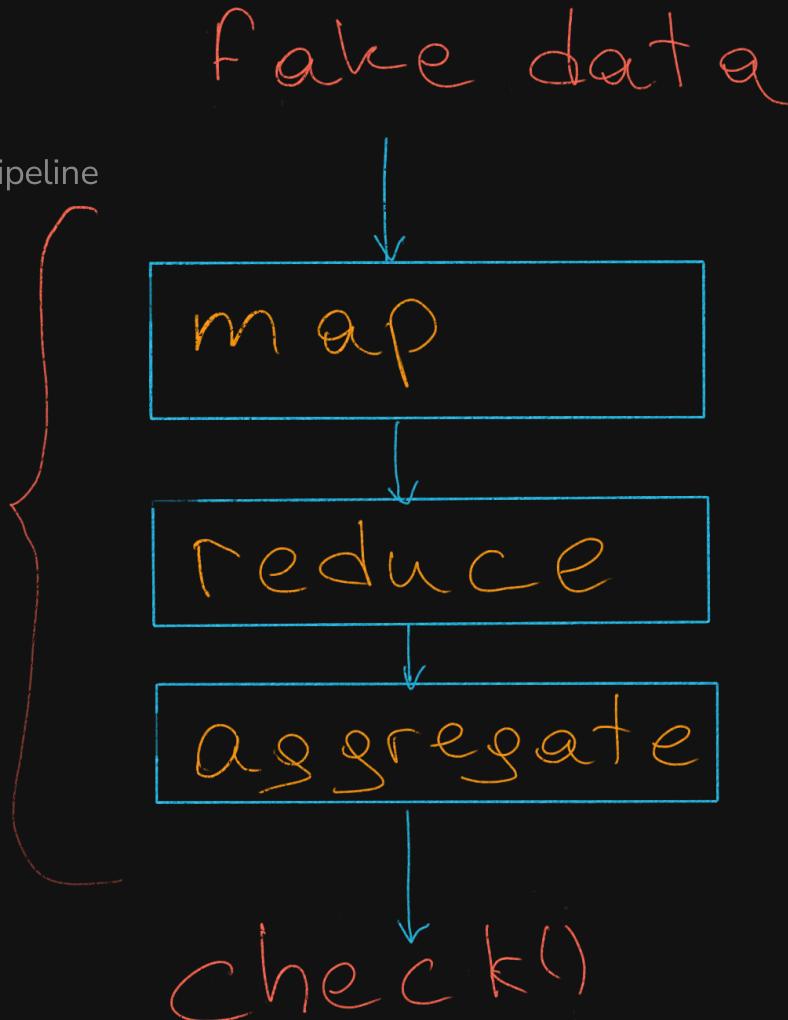
A pipeline should transform data correctly!

Correctness is a business term

Let's paste fakes!

Fake input data Reference data at the end of the pipeline

Separate
Function

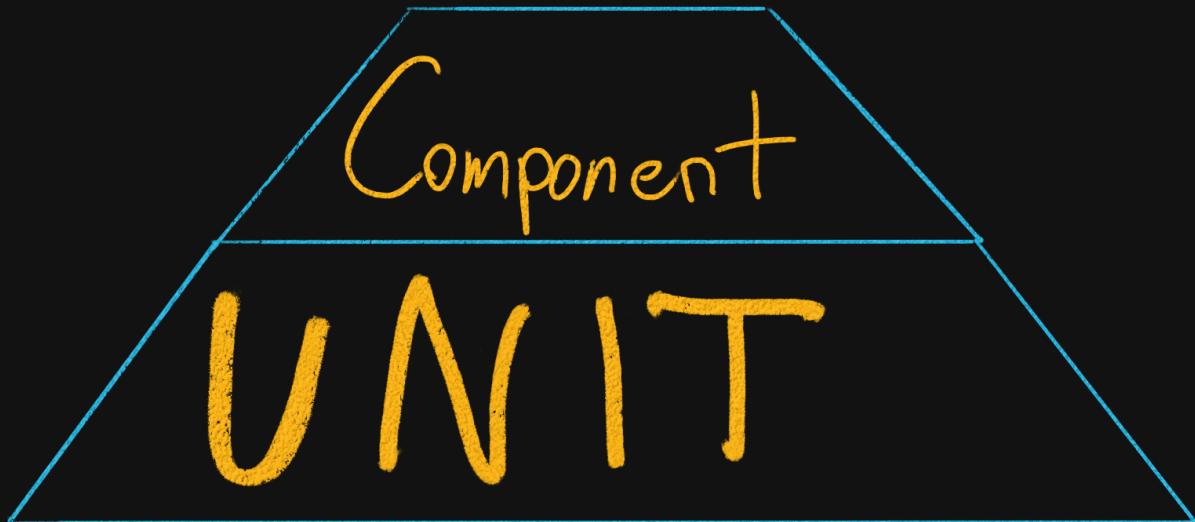


Tools

holdenk/spark-testing-base ← Tools to run tests

MrPowers/spark-daria ← tools to easily create test data

Component testing



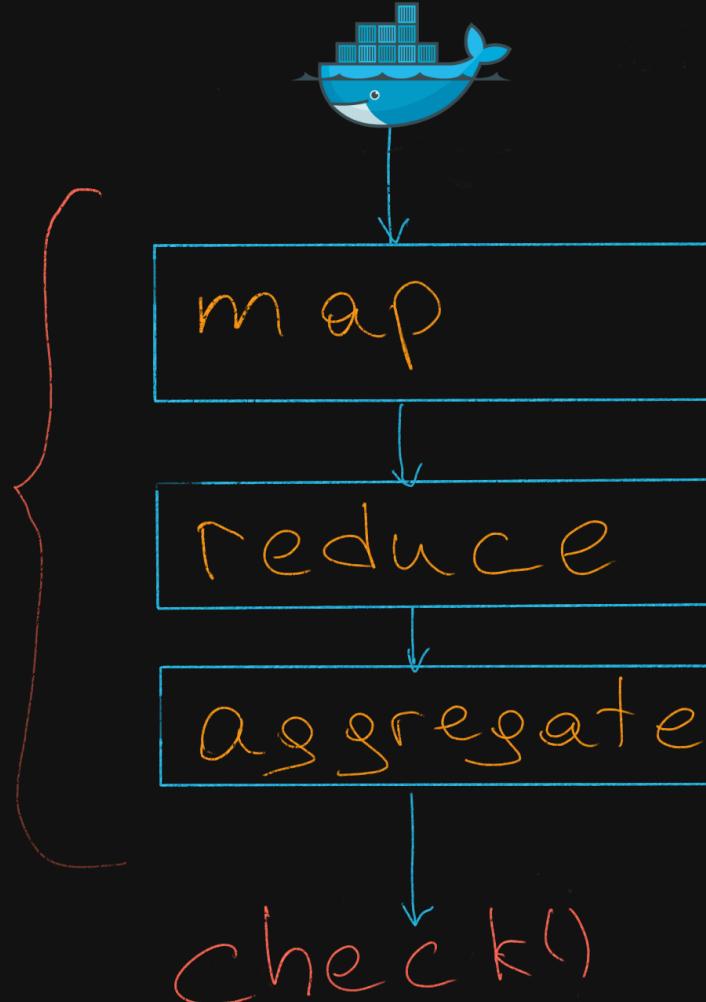
Testcontainers



TESTCONTAINERS

Testcontainers

Separate
Function



TestContainers

Supported languages:

- Java (and compatibles: Scala, Kotlin, etc.)
- Python
- Go
- Node.js
- Rust
- .NET

TestContainers

```
1  from testcontainers postgres import PostgresContainer
2  import sqlalchemy
3
4  postgres_container = PostgresContainer("postgres:13")
5  with postgres_container as postgres:
6      e = sqlalchemy.create_engine(postgres.get_connection_url())
7      result = e.execute("SELECT version()")
8      version, = result.fetchone()
9  print(version)
```