# End to End Monitoring of Data Pipeline
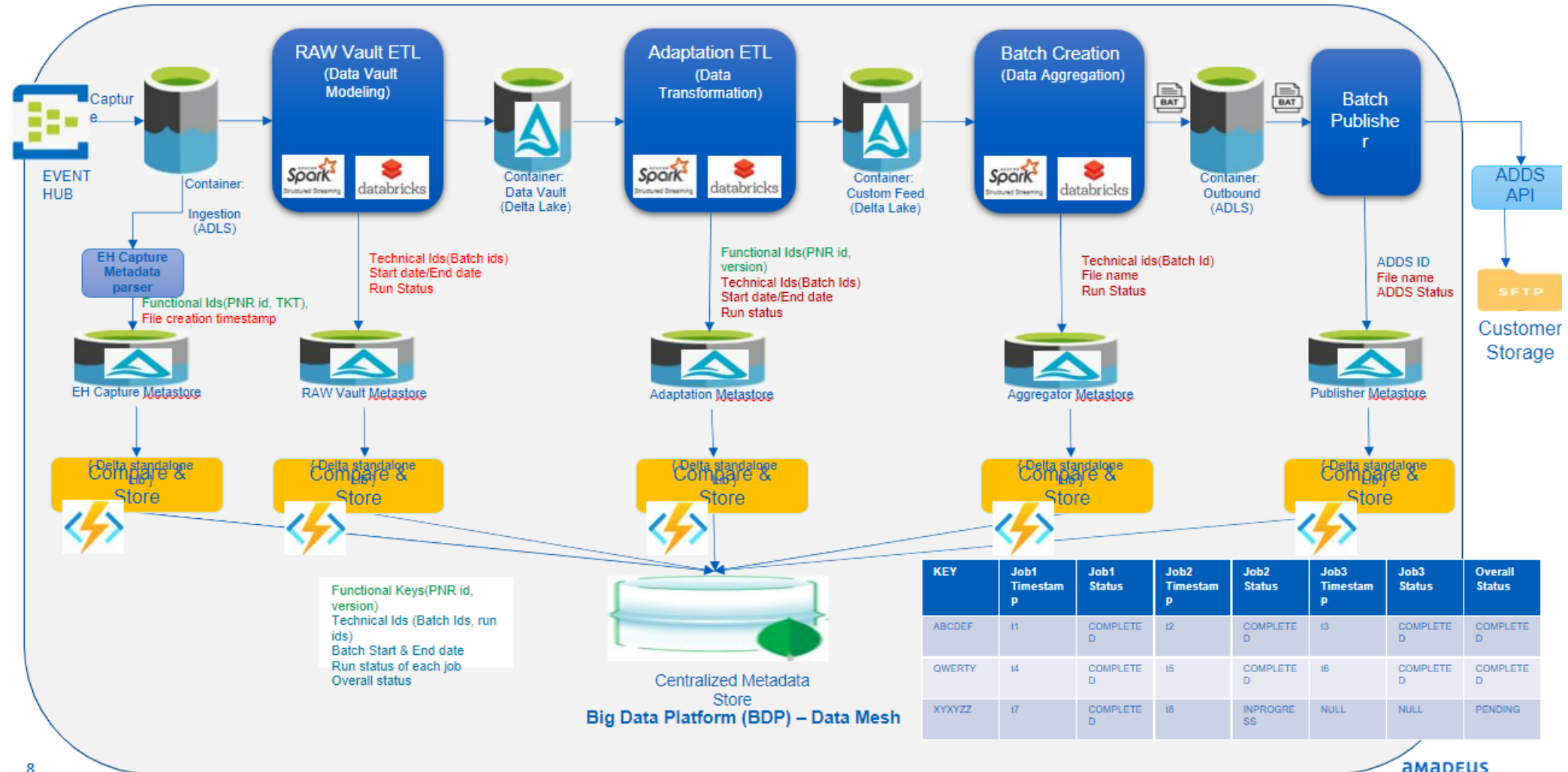
amadeus

# Agenda

**aMaDEUS**

# 1.
# Introduction

amadeus

# Context and Goal of Internship

➤ Create an end-to-end data pipeline monitoring framework, which can help to monitor batch execution status based on functional keys across different jobs in a big data pipeline.
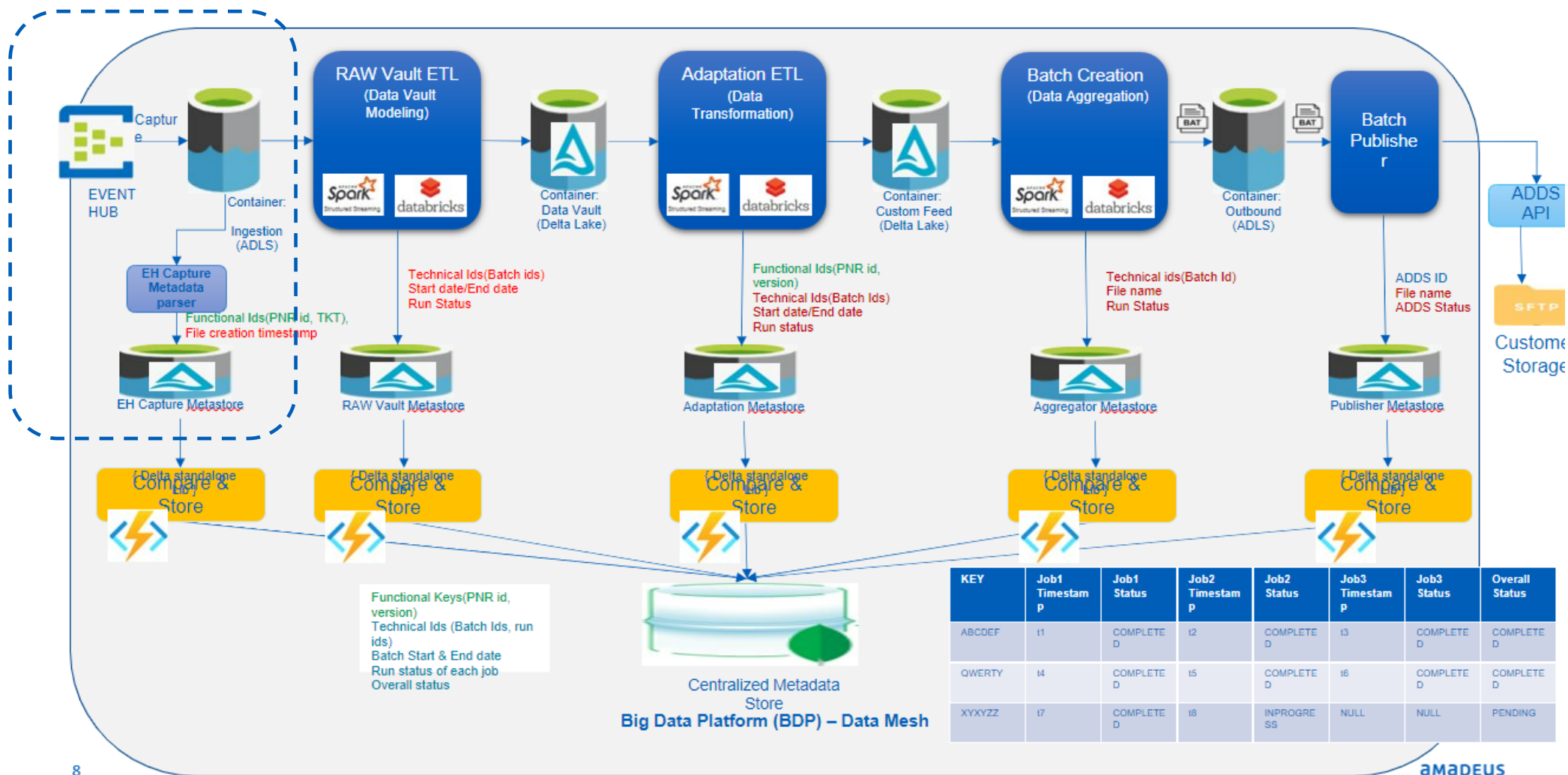
# Introduction

This whole process of end-to-end monitoring is diving into four main modules

- EH Capture Metadata storage

- Metadata capture for ETL job (Generic Library)

- Batch Publisher metadata storage

- Compare and Store

aMaDEUS

# 3.
# Demo

aMaDEUS

# EH Capture Metadata storage

7

# EH Capture Metadata storage
## Introduction

### EVENT HUB

➢ Event Hub is an Azure service that enables in processing large amounts of event data from connected devices and applications.

### AUTO LOADER

➢ Databricks Autoloader is an **Optimized File Source** that can automatically perform incremental data loads from your Cloud storage as it arrives into the **Delta Lake Tables**.

➢ Databricks Autoloader presents a new Structured Streaming Source called **cloudFiles**.

### Uses of AUTO LOADER

➢ No file state management

➢ Scalable & Easy to use

➢ Schema inference & evolution support

**amaDEUS**

# EH Capture Metadata storage Demo Flow

_ Ingestion of files (AVRO) incrementally and efficiently from azure blob storage container

_ Consumption of the Avro files in the FrameWork(Spark Structured Streaming) and Performing deserialization and Uncompression into the Json Format.

_ Traversing the json document to derive columns for the meta datatable.

_ And pushing the Avro files after Consolidating them into the single record to those Related FeedType Database.

_ Creating the different Database-tables for related Functional feedTypes in the container in the form of delta table.

_ This Single application jar can work with the other feedtype Modules by Changing the config file while running the job.

aMaDEUS

# EH Capture Metadata storage

# Generic library creation

8

**Big Data Platform (BDP) – Data Mesh**

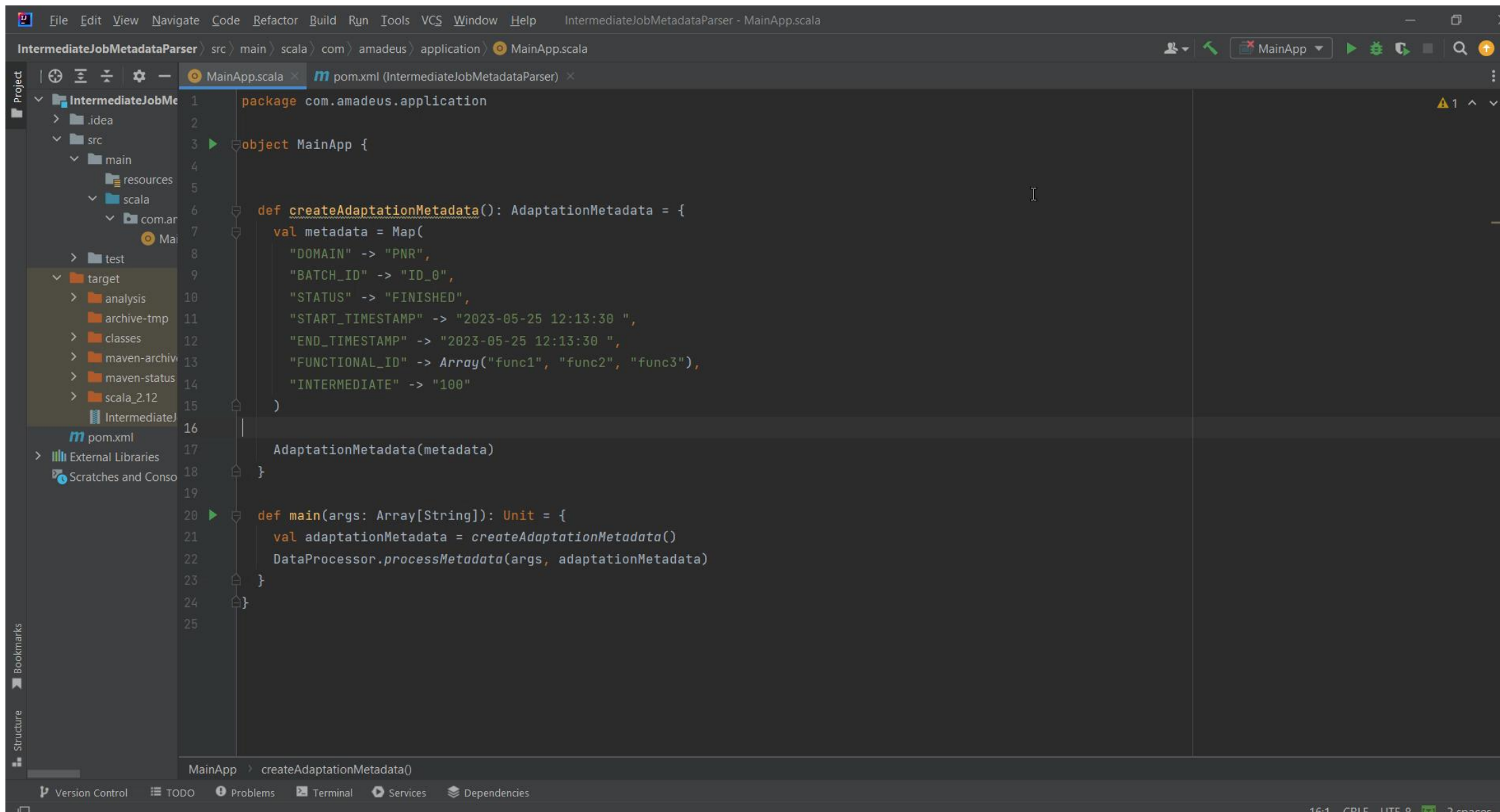| KEY | Job1 Timestamp | Job1 Status | Job2 Timestamp | Job2 Status | Job3 Timestamp | Job3 Status | Overall Status |
|------|------|------|------|------|------|------|------|
| ABCDEF | t1 | COMPLETED | t2 | COMPLETED | t3 | COMPLETED | COMPLETED |
| QWERTY | t4 | COMPLETED | t5 | COMPLETED | t6 | COMPLETED | COMPLETED |
| XYXYZZ | t7 | COMPLETED | t8 | INPROGRESS | NULL | NULL | PENDING |

11

# Generic library creation

## Introduction

_ Fetches data from modules

_ Takes schema for the fields from the configuration file

_ Writes the output inside a table after creating a data frame

```
<dependencies>
    <dependency>
        <groupId>org.dataprocessor</groupId>
        <artifactId>fetch_data</artifactId>
        <version>1.0-SNAPSHOT</version>
    </dependency>
```

```
File  Edit  Format  View  Help
schemaConfig {
 DOMAIN = "StringType"
 BATCH_ID = "StringType"
 STATUS = "StringType"
 START_TIMESTAMP = "TimestampType"
 END_TIMESTAMP = "TimestampType"
 FUNCTIONAL_ID = "ArrayType(StringType)"
 INTERMEDIATE = "IntegerType"
}
outputPath = "dbfs:/user/hive/warehouse/intermediate_e2e.db/source_table"
```

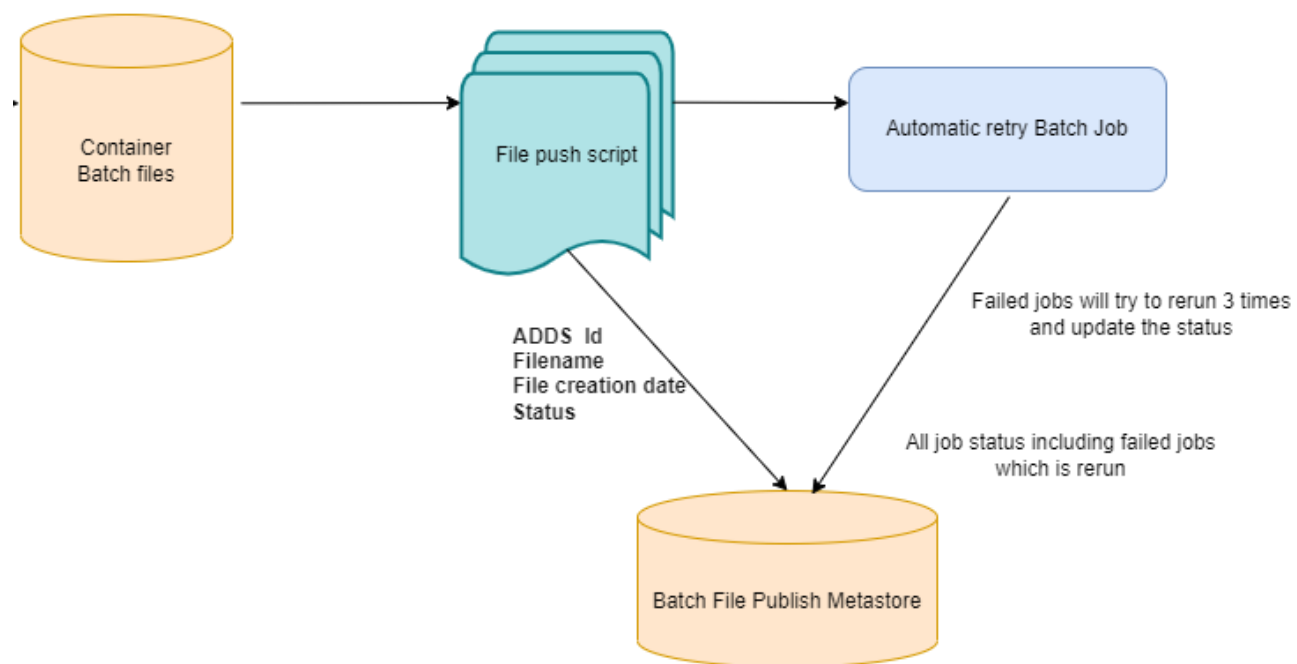# Generic library creation

# aDDS Batch Publisher Monitoring

## Introduction

_ aDDS for transferring files

_ aDDS provides different endpoints:

- aDDS API Stream
- aDDS API Notify

_ aDDS is separate team

_ Files status need to be monitored

**aMaDEUS**

# Batch Publisher Library

_ Already a batch publisher lib to transfer files from storage account to external sftp.

_ Now it will also monitor "Batch Publisher Metastore" metadata of files.



Batch File Monitoring

aMaDEUS

# Batch Publisher Library

# Compare & Store

## Delta Standalone

- The Delta Standalone library is a Java library that can be used to read from and write to Delta tables

- Standalone doesn't depend on Apache Spark

- ACID guarantees

## MongoDB

- MongoDB is document-oriented database

- Collection: Collection is a group of MongoDB documents, which is similar to table

- Document: Document is a set of key-value pairs, which is similar to row

## Azure Functions

- Azure functions allows you to schedule the execution of your functions

- Integrates with most of the development tools

- Automatically scales the execution environment

aMaDEUS

# Demo Flow

_ Reading Delta table data & schema by provided configuration

_ Creating MongoDB connection to access the maximum timestamp value

_ Comparing MongoDB timestamp with timestamps from delta table

_ Perform Transformations on the table record which has early timestamp

_ Check if the exploded ID is present in MongoDB

_ If ID exists, update the document with values

_ If ID doesn't exist, insert the values into new document

_ If the timestamp is later, end the workflow

**aMaDEUS**

# Compare & Store

# 4.
# Key Takeaways

aMaDEUS

# Key Takeaways

➢ Developed a deep understanding of Scala and Spark programming.

➢ Gained experience using Databricks and IntelliJ to run spark jobs.

➢ Developed proficiency in querying and manipulating data using SQL.

➢ Understood how data works with respect to high level

➢ Understood about real-time streaming using Autoloader

# Thank you!

aMaDEUS