

# Data Processing in Databricks

Leveraging Pandas, PySpark, and SQL

Marcelino Mayorga Quesada



## Marcelino Mayorga Quesada

- 14 Years of Experience in Software in finance, marketing and video games sectors, and over the last 3 years in AI.
- Experienced in technical and delivery management roles.
- Clients: Blackstone, Cambridge Associates, EA Sports, Citibank and BAC Credomatic.
- Technical Instructor on GCP and .NET courses.
- Passionate about Artificial Intelligence 🤖, music 🎹, and video games 🎮.

# Agenda

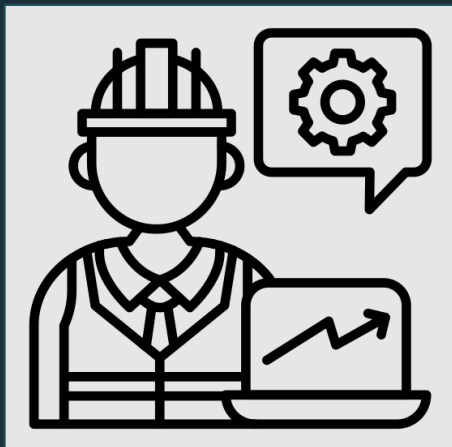
1. Introductions
  - Instructor, Topic, Audience
2. Data Processing
  - Concept, Operations
3. Databricks
  - Solution, Pyspark, Pandas, SQL
4. Demo
  - Community Version
  - Use case
5. Key Takeaways
6. Q & A

QR

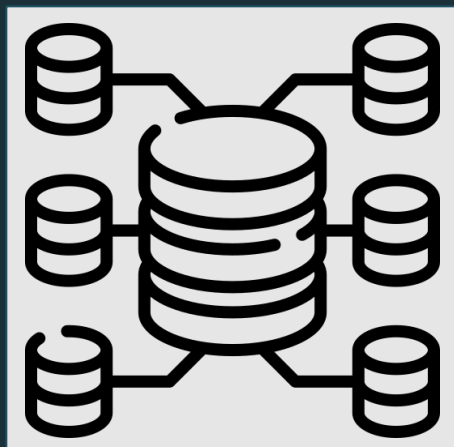
# Data Processing

## Concept

A series of operations to convert raw data into **meaningful information**.



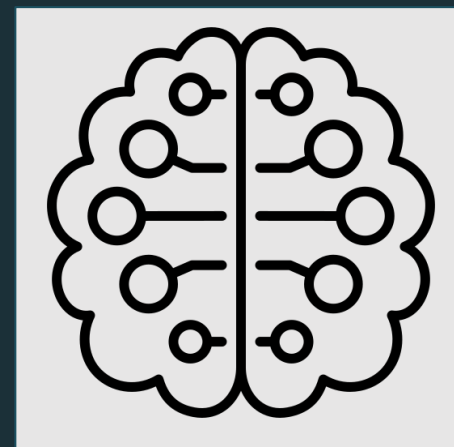
Data Engineers



Storage



Analytics



DL & ML Models

# Data Processing

## Operations

### 1. Cleaning

- Removing duplicates
- Impute or delete missing values
- Correct errors and inconsistencies

### 2. Integration

- ETL (Extract Transform Load)
- Merge and Join data warehousing
- Augmentation

### 3. Transformation

- Normalization and Standardization
- Aggregation (Summing, Averaging)
- Pivoting tables
- Encoding categorical values

### 4. Reduction

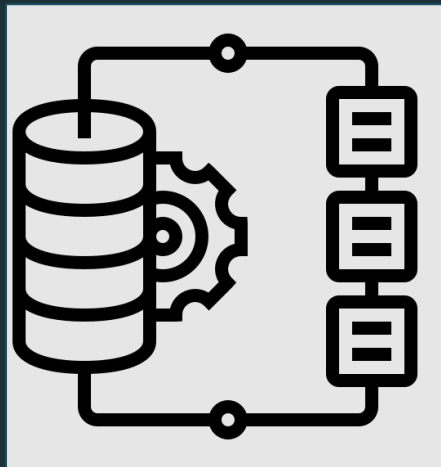
- Dimensionality Reduction: PCA, t-SNE
- Feature Selection & Extraction
- Sampling
- Compression

# Data Processing

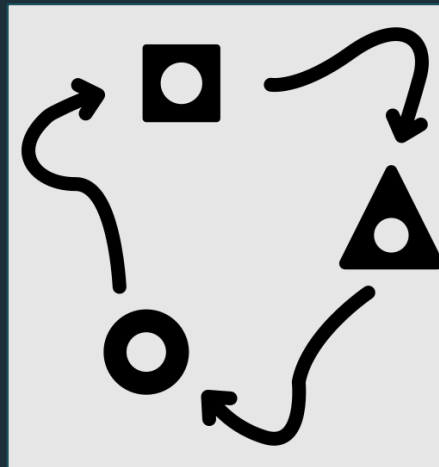
## Operations



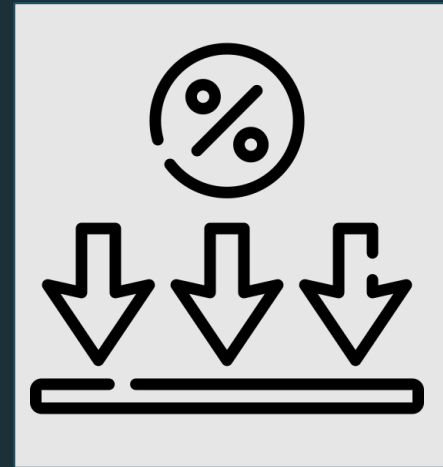
Cleaning



Integration



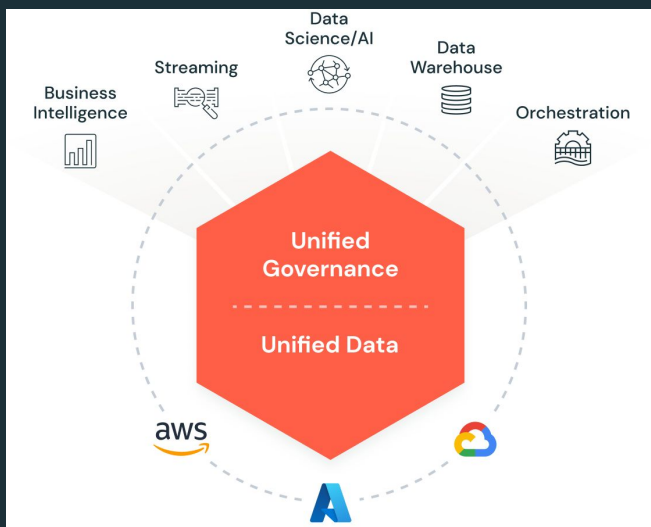
Transformation



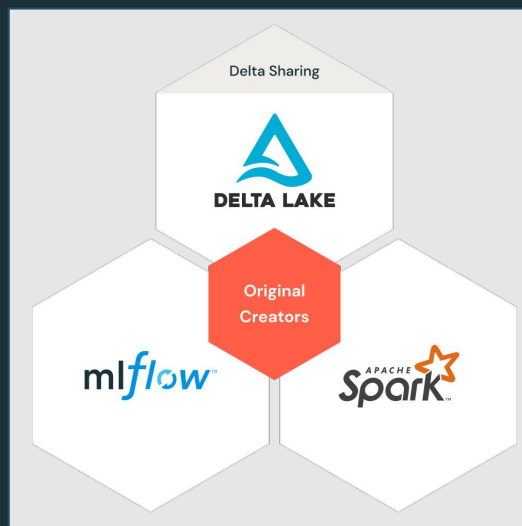
Reduction

# Databricks

## Data Lakehouse Architecture



Unified



Open



Scalable

# Pandas

## Data Analysis and Manipulation Library

- Data cleaning, transformation and visualization
- **Low volume** data for single computing
- Prototyping
- Available:
  - Stand alone library
  - Pandas API in Spark
  - **Pandas Dataframes**





# PySpark

## Big Data Processing Framework

- Python API for Apache Spark
- **High volume** data (TB, PB)
- Distributed computing (Clusters)
  - Parallel processing
  - Lazy Evaluation
  - Fault Tolerance
- Complex data transformations
- Used via **Spark Session & Context** and **Spark Dataframes**



# SQL

## Structure Query Language

- Managing and manipulating **relational databases**
- Queries and Transactions over multiple tables.
- Available:
  - Stand alone DBMS (mysql)
  - SQL in Spark
  - **SQLContext**



# Databricks

## Summary

Name	Type	Purpose	Usage via	Ideal
Pandas	Data Analysis and Manipulation <b>Library</b>	<ul style="list-style-type: none"><li>• Data Handling</li><li>• Transformation</li><li>• Eager Execution</li></ul>	<ul style="list-style-type: none"><li>• Standalone Library</li><li>• Pandas API on Spark</li><li>• Dataframes</li></ul>	<ul style="list-style-type: none"><li>• Low Volume Data</li><li>• Prototyping</li></ul>
PySpark	Python Interface for Big Data Processing <b>Framework</b>	<ul style="list-style-type: none"><li>• Distributed Computing</li><li>• Parallel Processing</li><li>• Lazy Evaluation</li><li>• Fault Tolerance</li></ul>	<ul style="list-style-type: none"><li>• Pyspark</li><li>• Context &amp; Session</li><li>• Dataframes</li></ul>	<ul style="list-style-type: none"><li>• Big Data</li><li>• Scalability</li><li>• Performance</li><li>• Integration with multiple sources</li></ul>
SQL	Structured Query <b>Language</b>	<ul style="list-style-type: none"><li>• Querying</li><li>• Transformations</li><li>• Transactions</li><li>• Storage</li></ul>	<ul style="list-style-type: none"><li>• SQLContext</li><li>• Datasets</li></ul>	<ul style="list-style-type: none"><li>• Querying</li><li>• Managing relational databases</li></ul>

# POP QUIZ!

QR

# Demo

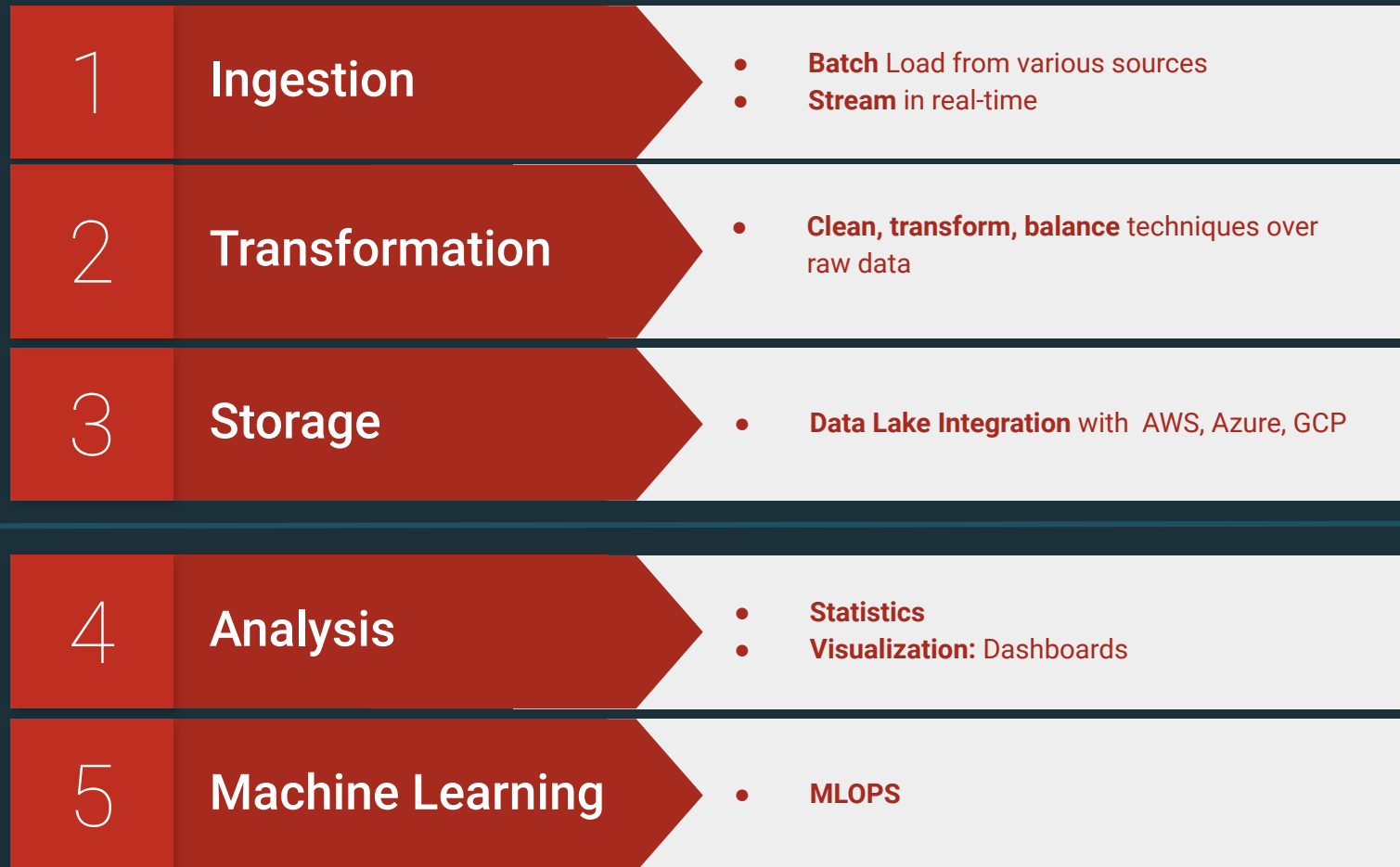
# Databricks

## Community Edition

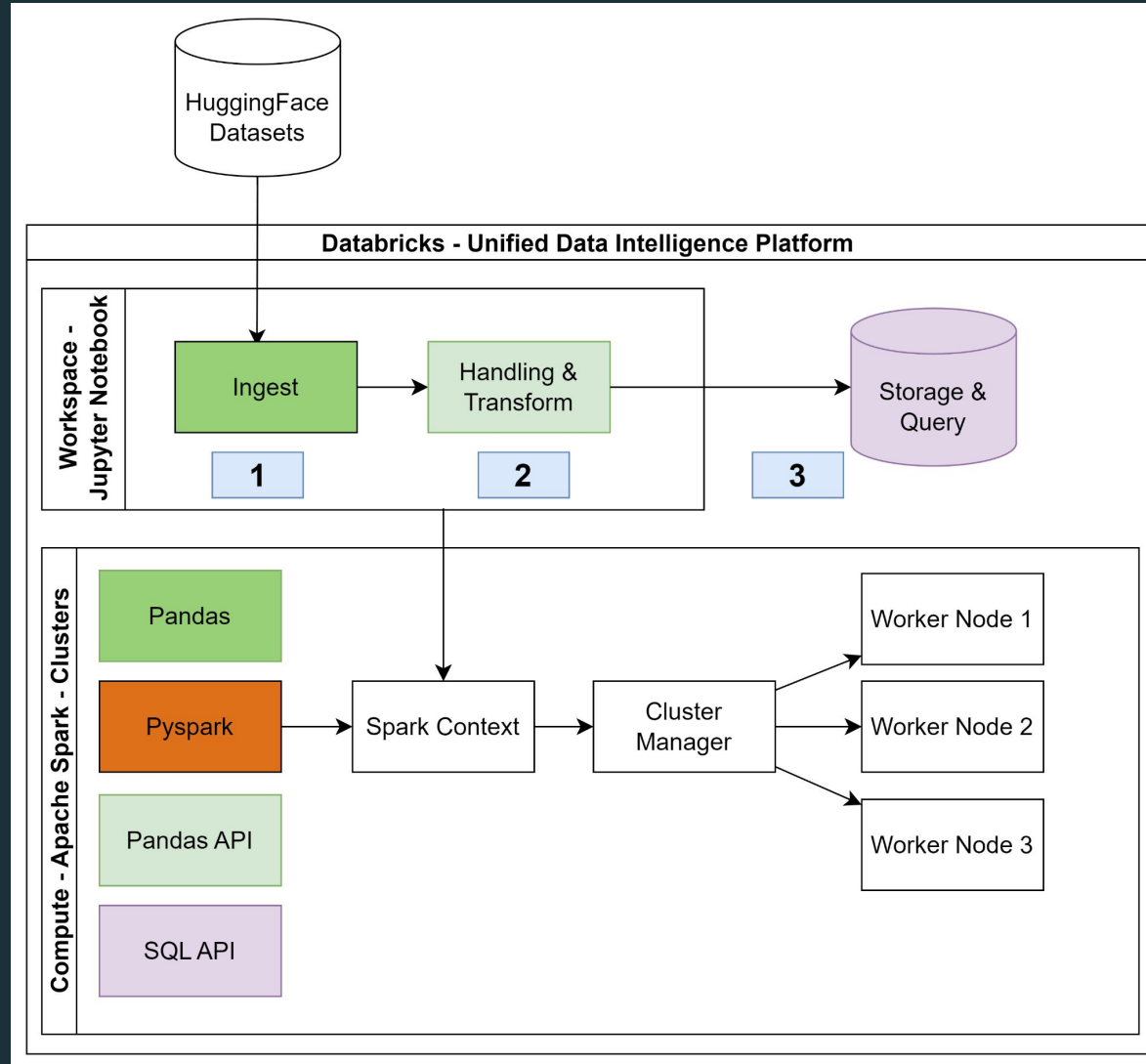
- Free access: <https://community.cloud.databricks.com/>
- Collaborative Notebook Environment
- Cluster Management
  - 1 Driver
  - 15.3 GB Memory, 2 Cores, 1 DBU
- Experiment Tracking

# Data Workflow

## Stages



# Use Case





# Key Takeaways

- Data Processing crucial step in Data Engineering
- Databricks Data Intelligence Platform
- Pandas is a powerful for data handling library aimed for single computing.
- Pyspark supports and enhances Pandas and SQL

# Q&A



databricks