# FACTORED DATATHON 2023

Cristina Gomez
Daverson Arenas

22-08-2023

Factored

**CRISTINA GOMEZ**

**DAVERSON ARENAS**

# PAISA GENIOUS TEAM

## Bioengineer and Data Scientist:

As the lead data scientist in our project Cristina:

- Managed feature engineering, data cleansing, and conducted data mining to create insightful visualizations
- Constructed the gold data set, essential for visualization and ML models
- Executed machine learning models for sentiment analysis, enhancing project insights

## Mechanical Engineer and Data Engineer

As the lead data engineer in our project Daverson:

- Orchestrated Databricks workspace deployment on AWS, integrating Github for efficient CI/CD.
- Designed lakehouse data architecture, unifying data warehouses and data lakes.
- Created ELT pipelines for batch and streaming incremental data ingestion

# Datathon Challenge Roadmap

**Factored**

### The Challenge

Create an innovative data solution (web apps, chatbots, dashboards, model interfaces...) to empower businesses with insights from product reviews.

## Stage 1

Understand the problem and define a solution approach.

Day 1    Day 2

## Stage 2

Choose technologies, deploy services, set up workspace and create GitHub repo. First batch ingestion

Day 3    Day 5

## Stage 3

Data architecture design, data engineering for batch and stream ingestion. Exploratory data analysis

Day 6    Day 8

## Stage 4

Data engineering for streaming, data cleansing, feature engineering and data visualization. Combine the data from both sources
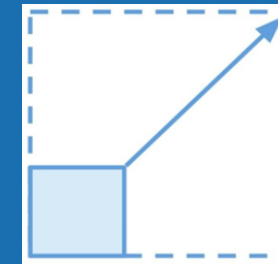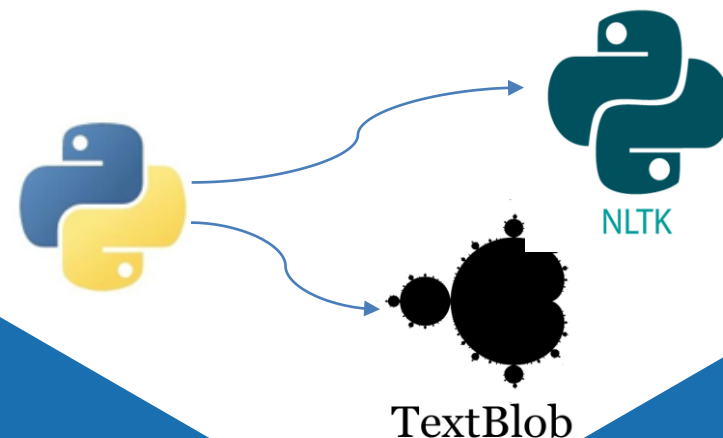
Day 9    Day 11

## Stage 5

Machine learning models, frontend design, final design dashboard, documentation and presentation
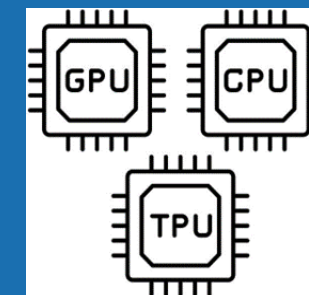
Day 12    Day 13

# MAIN TECHNOLOGIES



**DATABRICKS + APACHE SPARK**
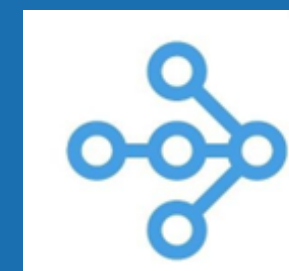
Scalability, Parallelism and Speed

- Databricks combines data warehouses & data lakes into a lakehouse architecture.

**AWS + DELTA LAKE**

heterogeneous hardware

- Amazon S3 serves as the data lake, coupled with Delta Lake, which functions as the storage layer
- AWS EC2 instances as the compute resources for Databricks clusters
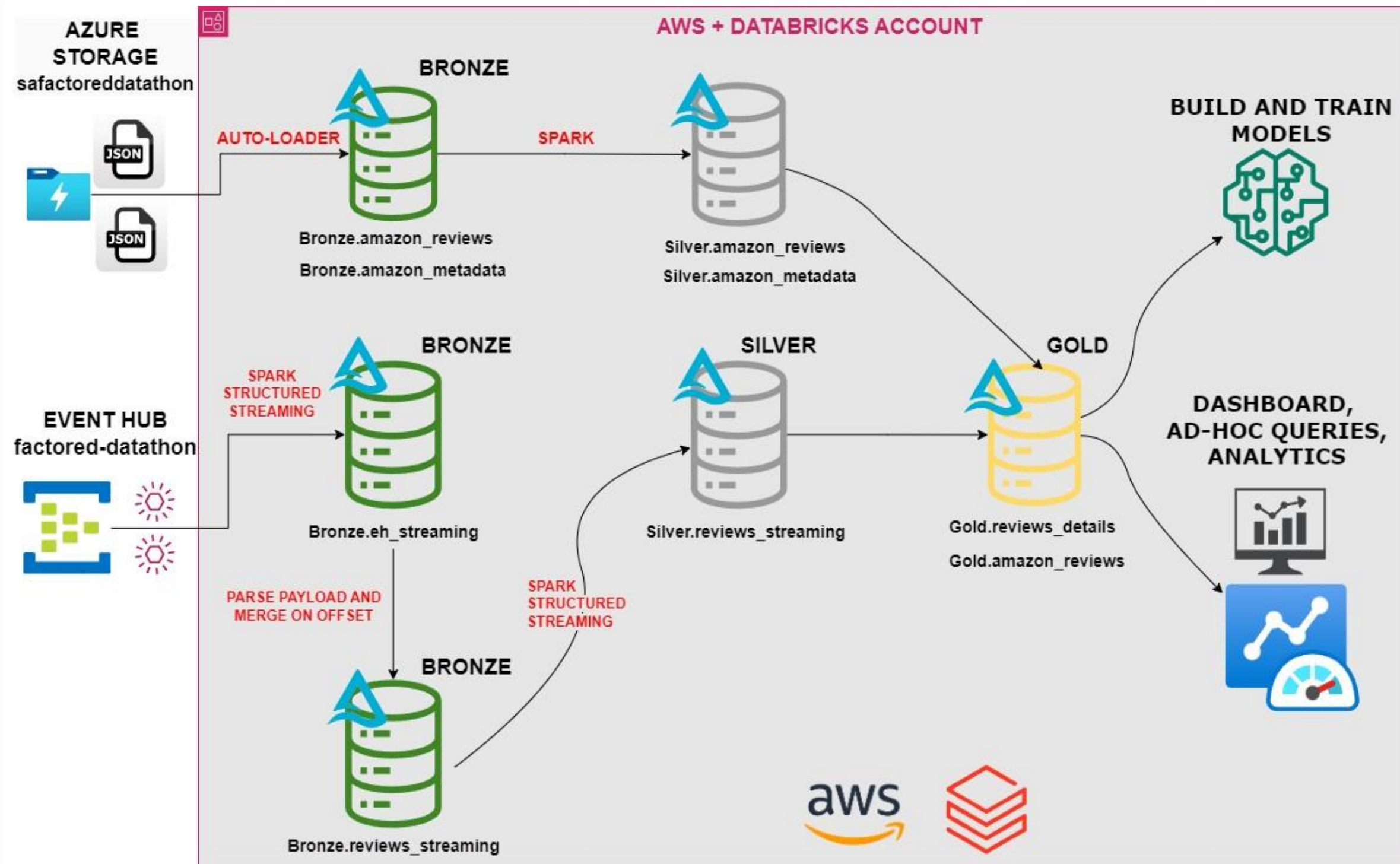
**MLFLOW + TEXTBLOB + NLTK**

Optimized resource utilization

- Distributed capabilities for large-scale experiments
- Efficient experiment tracking, reproducibility and scalability
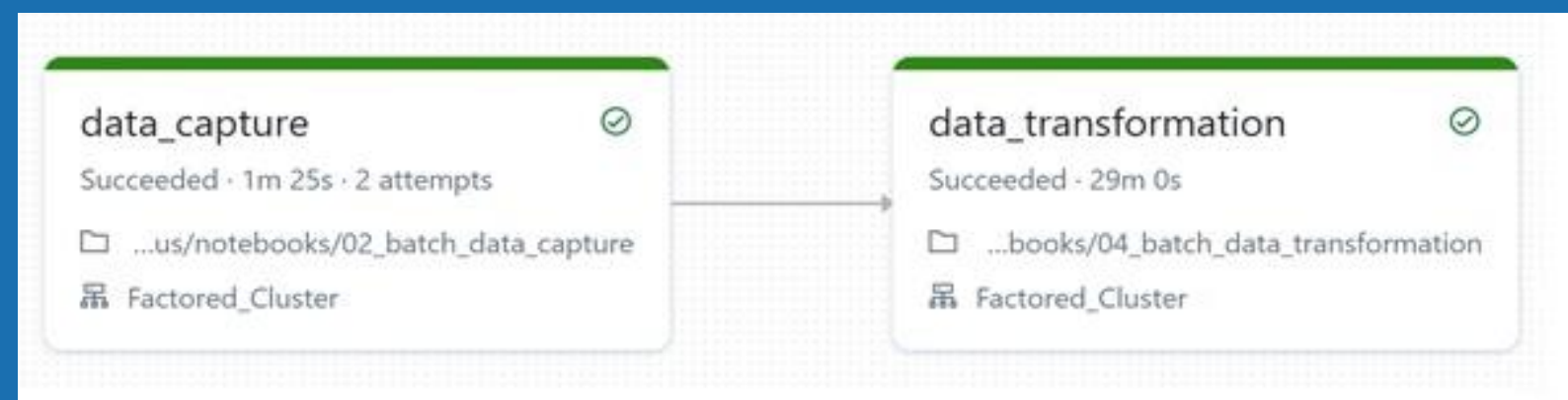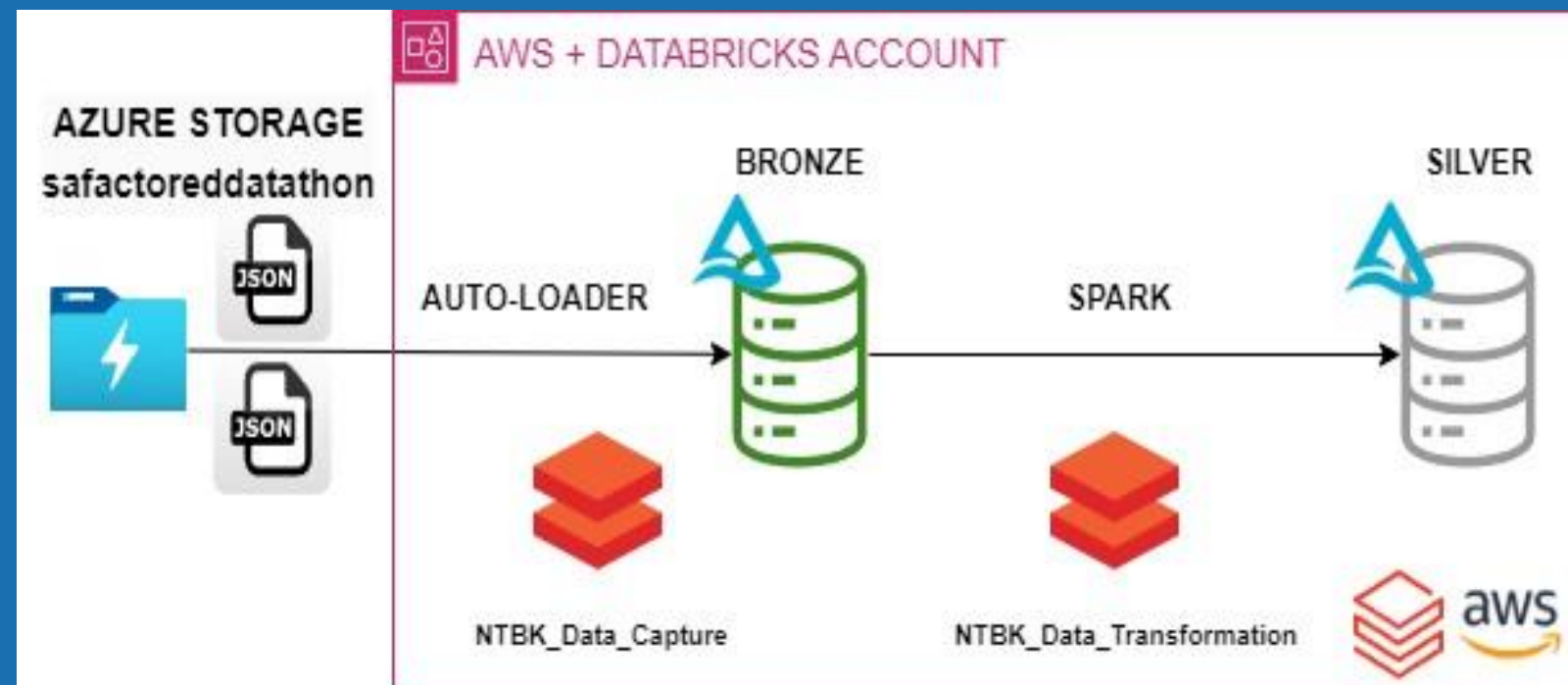
# DATA ENGINEERING WITH DATABRICKS

## BATCH

**1. Extract and Load data to the Lake house**
**Auto-loader:** Incrementally load new data files as they arrive

**2. Transform, Clean and Filter**
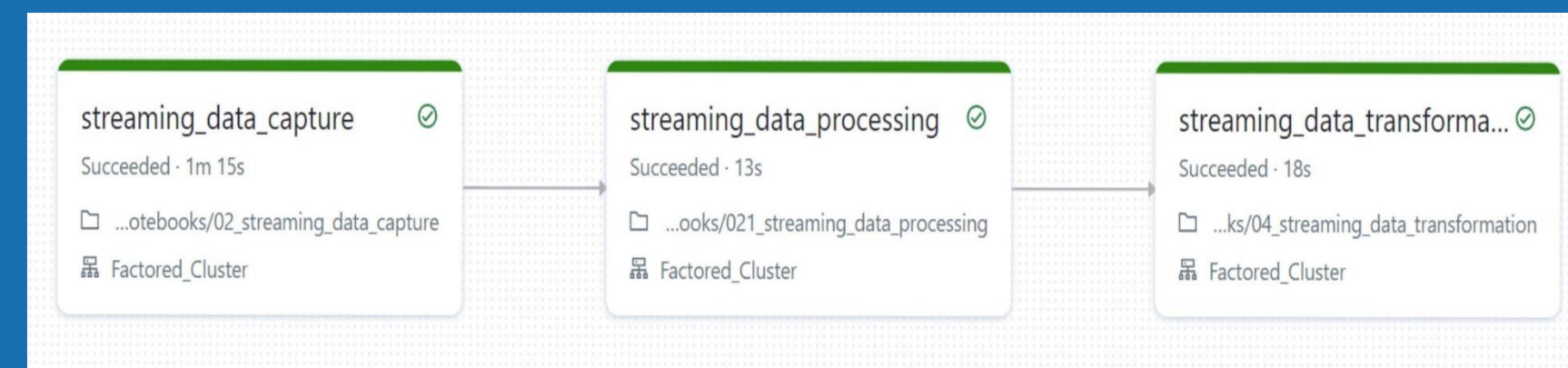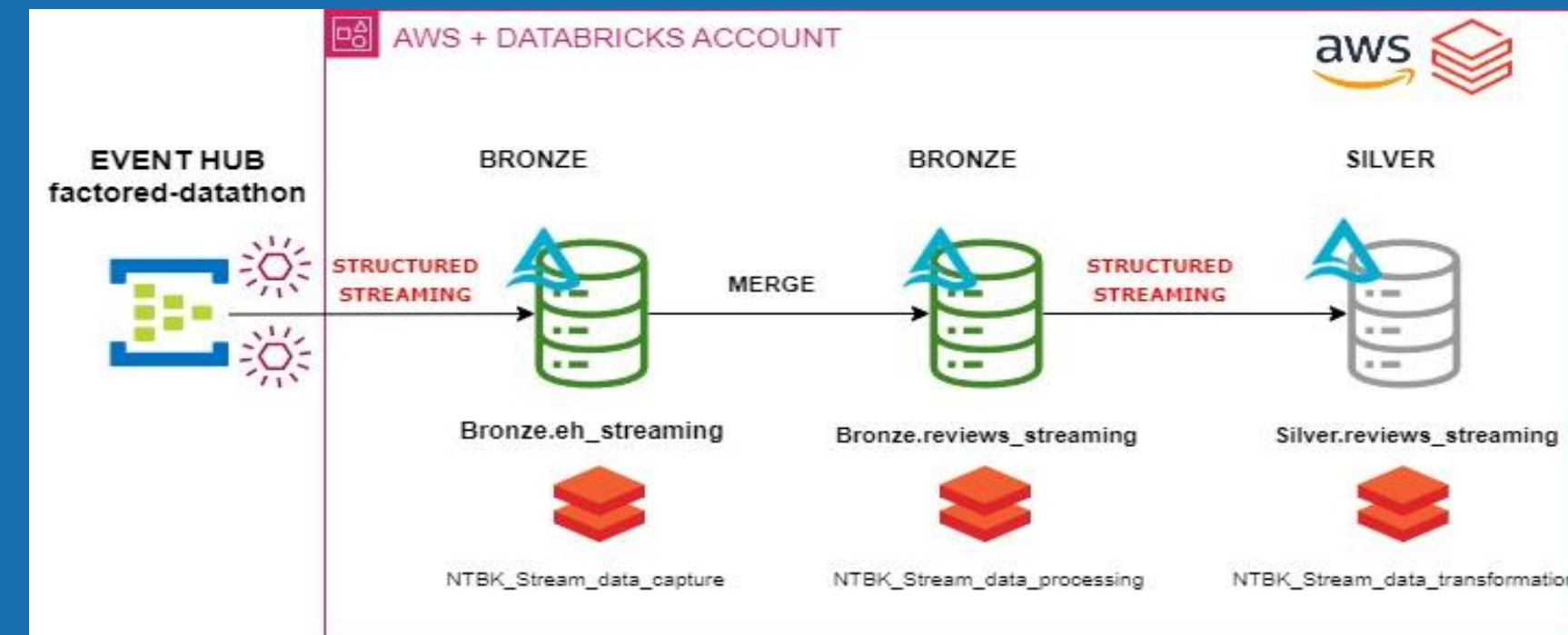**Spark structured streaming:** Processing and transformation tasks with a batch-like behavior

## STREAMING

**1. Extract and Load data to the Lake house**
**Spark structured streaming:** Stream data from event hub
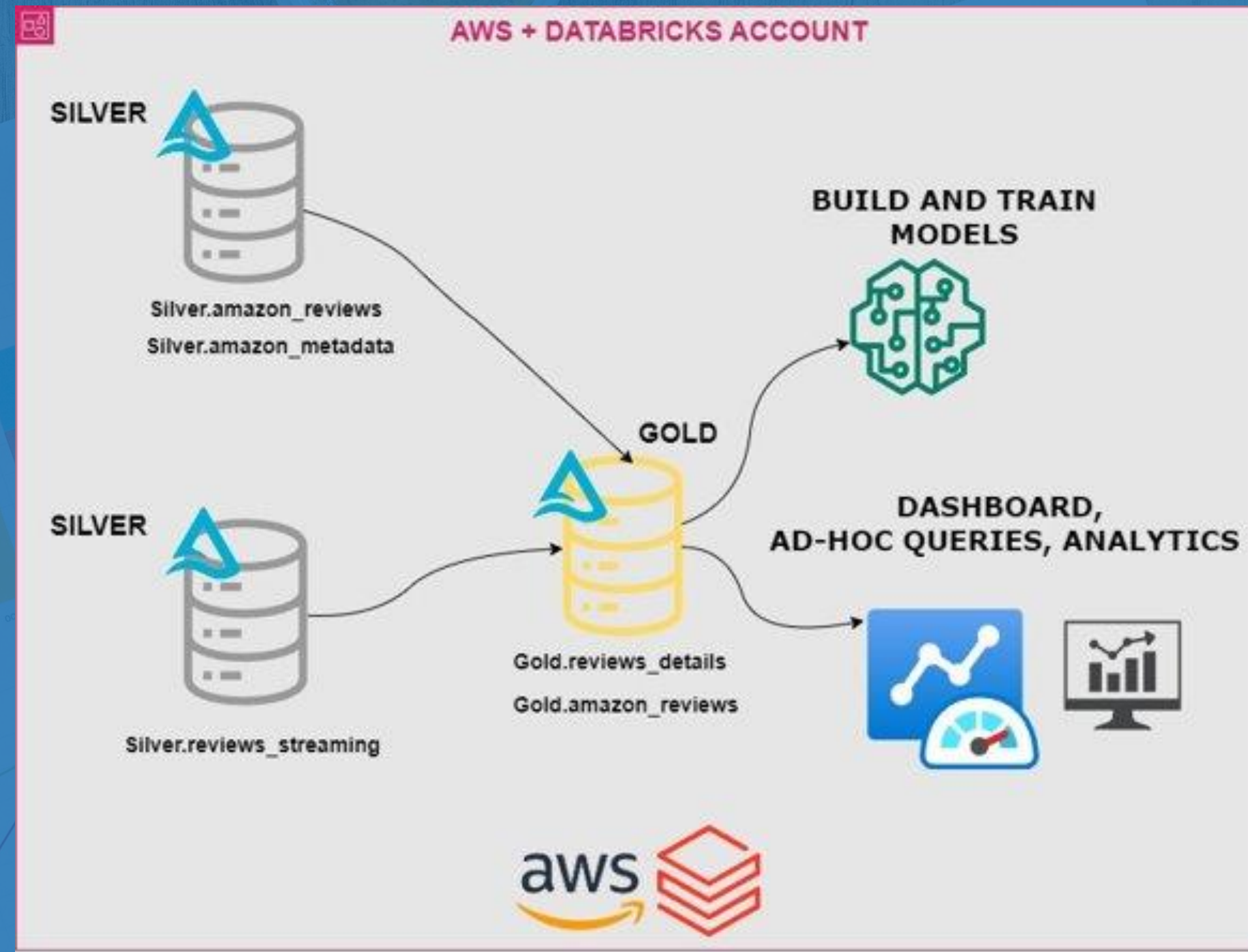
**2. Transform Data, Clean and Filter**
**Spark structured streaming:** Processing and transforming tasks as data becomes available

# DATA ENGINEERING WITH DATABRICKS
## Gold Layer: Business level Aggregates

The Gold layer aims to deliver continuously updated, clean data to downstream users and applications, including machine learning models, ad-hoc queries, and analytics tools.

# EXPLORATORY DATA ANALYSIS

**Data cleaning**

**Feature exploration**

**New variables**

**Pattern recognition**

**Report generation**

FINAL DASHBOARD

-Missing values and duplicate records.
-Variables imputation: price, main_cat, brand, title.
-Removing unwanted characters, converting to lowercase, and handling special cases: reviewText, title, main_cat, brand

-Number of unique customers,.
-Number of unique products.
-Number of unique reviews

-"Month"
-"Year"
-"Sentiment": positive, negative and neutral
-Number of words per review

Time Series analysis:
-Number of review per year.
Number of review per month.
-average of overall ratings per year
-Setiment analysis per year

# REPORT FROM HISTORICAL DATA

**Review Verification Impact:**
Does sentiment differ between verified and non-verified reviews, and do verified reviews tend to be more credible?
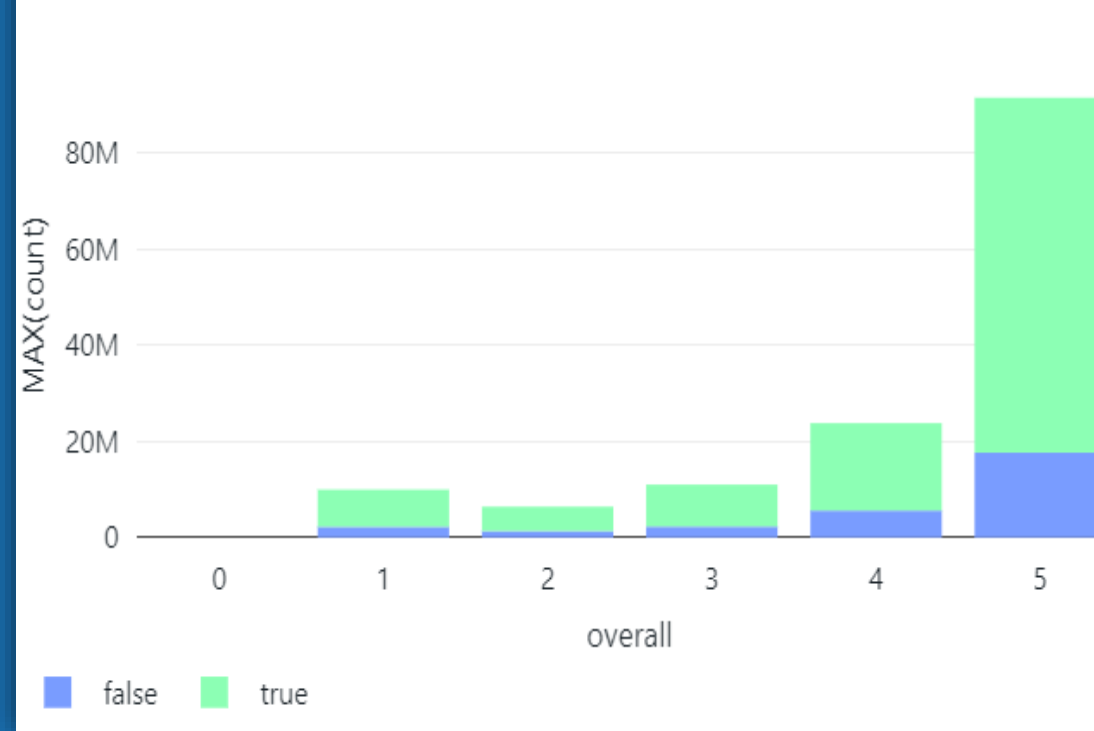
**Review Length and Sentiment:**
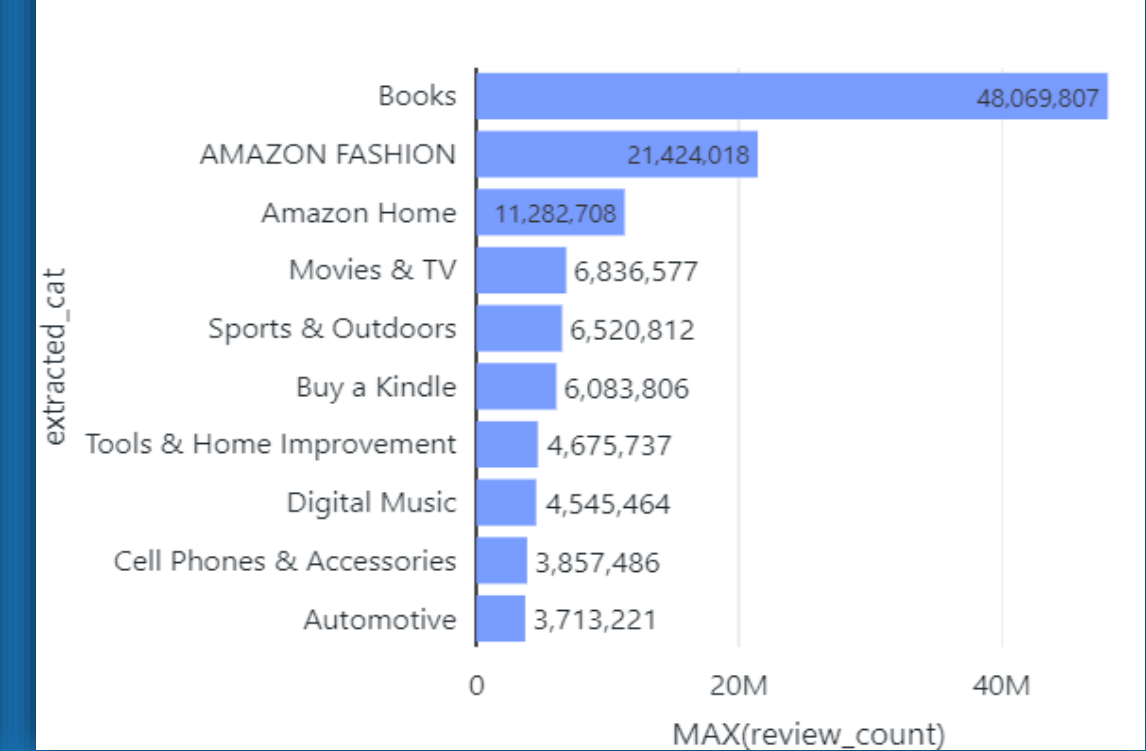Are longer or shorter reviews more likely to have a positive or negative sentiment?

**Product Prioritization:**
Which product categories receive the most feedback, and are they also the ones with higher satisfaction?
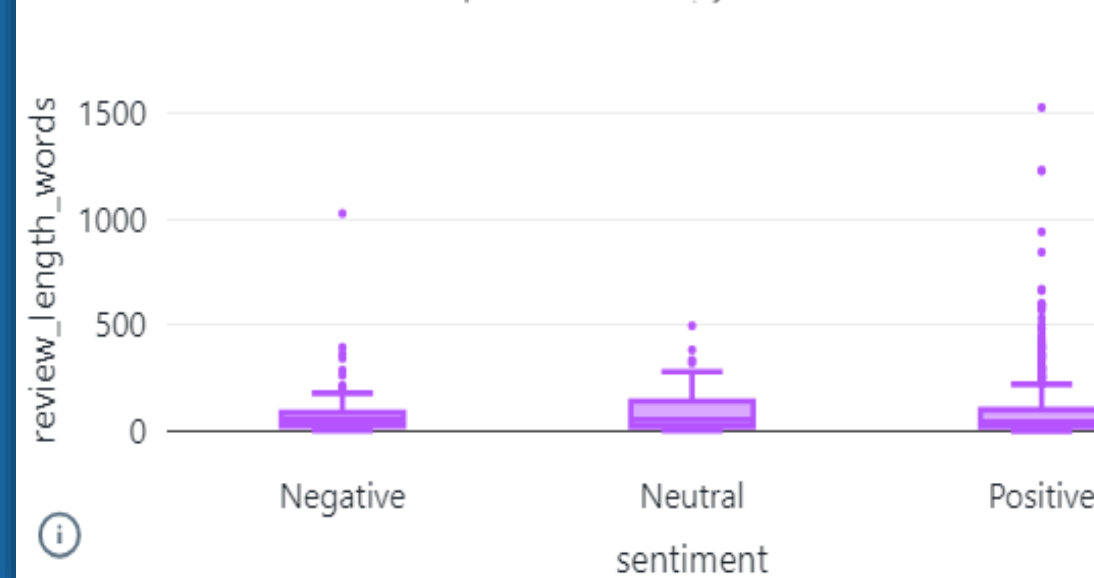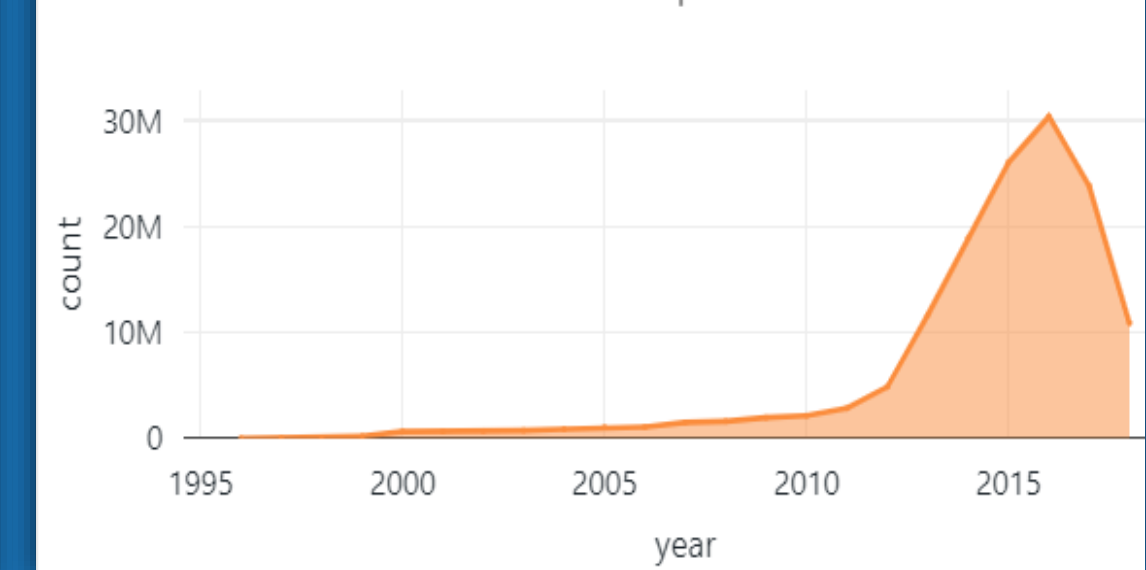


Overall ratings/Verified purchases

false    true



Total reviews per Main cat

| extracted_cat | MAX(review_count) |
| --- | --- |
| Books | 48,069,807 |
| AMAZON FASHION | 21,424,018 |
| Amazon Home | 11,282,708 |
| Movies & TV | 6,836,577 |
| Sports & Outdoors | 6,520,812 |
| Buy a Kindle | 6,083,806 |
| Tools & Home Improvement | 4,675,737 |
| Digital Music | 4,545,464 |
| Cell Phones & Accessories | 3,857,486 |
| Automotive | 3,713,221 |



Sentiment per Lenght of words



Total reviews per Year

# SENTIMENT ANALISYS DASHBOARD

- ❑ **Identifying Key Themes**
- ❑ **Prioritizing Focus Areas**
- ❑ **Sentiment Overview**
- ❑ **Monitoring Trends**
- ❑ **Customer Engagement**
- ❑ **Historical Insights**
- ❑ **Seasonal Patterns**

| SENTIMENT LEVEL | TOTAL REVIEWS | % Max |
|---|---|---|
| Total | 142,659,837 | 100% |
| Positive | 115,294,253 | 80.82% |
| Negative | 16,371,011 | 11.48% |
| Neutral | 10,994,573 | 7.71% |

**34,114,894** TOTAL USERS

**10,470,042** TOTAL PRODUCTS

# FRONTEND AND CONTINOUS DEPLOYMENT

# CHALLENGES AND CONCLUSIONS

SMALL TEAM

LIMITED TIME FRAME

LIMITED RESOURCES

## AWS DAILY COST RESOURCES

| | |
|---|---|
| Total cost | Average daily cost |
| $217.96 | $18.16 |

Costs ($)

60

45

30

15

0

Jul-25 Jul-26 Jul-27 Jul-28 Jul-29 Jul-30 Jul-31 Aug-01* Aug-02* Aug-03* Aug-04* Aug-05*

# THANKS FOR YOUR ATTENTION