# Bytewise Fellowship



## Task 8: Theoretical Exploration of Data Processing Concepts

| | |
|---|---|
| Name | Muhammad Ahsan Saleem |
| Date | 17th July, 2024 |
| Submitted to: | Muhammad Bilal |

# Introduction:

In this task, we will look into two key concepts in the realm of data processing: ELT (Extract, Load, Transform) vs. ETL (Extract, Transform, Load) and Batch vs. Streaming Pipelines. We will explore the main differences between these concepts, identify when to use each approach, and provide use-case demonstrations to illustrate their application. Finally, we will justify the chosen solutions to prove their effectiveness for the specific use-case.

## Task 8 - Theoretical Exploration of Data Processing Concepts

**Task Statement:** In this task, we will explore and compare key data processing concepts: ELT (Extract, Load, Transform) versus ETL (Extract, Transform, Load), and Batch versus Streaming Pipelines. The aim is to understand the fundamental differences, appropriate use cases, and demonstrate practical scenarios for each concept.

**Requirements:**

1. **ELT vs ETL:**
   - Explain the main differences between ELT and ETL.
   - Discuss when to use each approach and provide practical examples.

2. **Batch vs Streaming Pipeline:**
   - Explain the main differences between batch and streaming data pipelines.
   - Discuss when to use each approach and demonstrate a proper use case with practical examples.

**Submission:** Prepare a detailed demonstration with the help of a use case for each concept. Add supporting arguments to justify why the chosen solution is the best for the given scenario.

# Task Documentation:

In this section, I will provide documentation of Task 8, covering the theoretical concepts, practical demonstrations, and justifications for the chosen approaches. This documentation aims to offer a comprehensive understanding of ELT vs. ETL and Batch vs. Streaming Pipelines, complete with use-case demonstrations to illustrate the practical applications and benefits of each approach.

## ETL vs. ELT:

### Main Differences:

#### ETL (Extract, Transform, Load):
- **Process:**
  - In ETL, data is first extracted from the source
  - Then, it is transformed into the desired format.
    - The source data is extracted to a staging area. In the staging area, the data undergoes a transformation process that organizes and cleans all data types. This transformation process allows for the now structured data to be compatible with the target data storage systems.
  - Finally, it is loaded onto the destination where it is stored.
- **When to Use:**
  - ETL is suitable for smaller data sets
  - When data quality and transformation requirements are complex and need to be handled before loading into the data warehouse.
- **Benefits:**
  - Results in cleaner data
  - Works with cloud data warehouses by using cloud-based SaaS platforms and onsite data warehouses.

#### ELT (Extract, Load, Transform):
- **Process:**
  - In ELT, data is first extracted from the source.
  - Then, it is loaded into the warehouse.
  - Finally, it is transformed within the data warehouse.
- **When to Use:**
  - ELT is suitable for large datasets and modern cloud-based data warehouses.

- o When leveraging the processing power of the data warehouse for transformation tasks.
- • Benefits:
  - o Enables faster implementation
    - ▪ The transformation occurs after the load function, preventing the migration slowdown that can occur during this process.
  - o ELT avoids server scaling issues by using the processing power and size of the data warehouse to enable transformation (or scalable computing) on a large scale.

## Use Cases:

### ETL:

Companies that have multiple ventures may have multiple consumers, suppliers and partners in common. As this data can be stored at multiple locations in different formats, ETL enables the transformation of this data in a unified format before loading it.
Two companies, Company A and Company B, are merging their operations. Both companies have their own sets of consumers, suppliers, and partners stored in separate databases. The data in these repositories is formatted differently, posing a challenge for integration. For instance, date formats might differ between the two companies (MM-DD-YYYY vs. DD-MM-YYYY). An ETL process is necessary to transform and unify the data before loading it into a central data warehouse.

### ELT:

Meteorological systems, such as national weather services, collect vast amounts of data from various sources like satellites, weather stations, and ocean buoys. This data needs to be processed and analyzed quickly to provide accurate weather forecasts and warnings. An ELT process can be used to load the raw data into a data warehouse and then transform it as needed for different analyses.

## Batch vs. Streaming Pipeline:

### Main Differences:

#### *Batch Processing:*

- **Process:**
  - Data is collected over a specified period and processed in groups or batches.
  - Typically involves aggregating data at intervals, such as hourly or daily, and processing the batch during a scheduled time.
- **When to Use:**
  - Suitable for scenarios where real-time processing is not critical.
  - Ideal for processing large volumes of data that can be collected and processed together at specific intervals.
  - Commonly used in applications like financial report generation, periodic data updates, and data archiving.
- **Benefits:**
  - Cost-effective, as it can be scheduled during off-peak hours.
  - Simplifies data management by processing data in bulk.
  - Efficient for handling large datasets that do not require immediate processing.

#### *Streaming Processing:*

- **Process:**
  - Data is processed in real-time as it is generated or received.
  - Continuous input and processing, with minimal latency between data arrival and processing.
- **When to Use:**
  - Suitable for scenarios requiring immediate processing and response.
  - Ideal for applications where data needs to be processed as soon as it arrives, such as monitoring systems, fraud detection, and real-time analytics.
- **Benefits:**
  - Enables real-time insights and immediate action on data.
  - Reduces latency, providing up-to-date information.

### Use Cases:

Consider a scenario where Spotflix, a popular video streaming service, needs to manage how it processes user activity data. Spotflix has two main types of user interactions: offline and

online streaming. These interactions generate different types of data that need to be processed efficiently.

## Batch Processing:

1. **Scenario:** *Spotflix needs to update its library with new uploads.*
   **Process:**
   - Videos uploaded are collected throughout the day.
   - Every ten minutes, the system batches the new uploads and processes them together to update the database.
   - This approach allows efficient handling of multiple uploads without overwhelming the system with continuous processing.

2. **Scenario:** *Spotflix's Offline Streaming*
   **Process:**
   - When a user downloads the video for offline watching, the entire video is batched and saved to the user's device in one go.
   - This approach ensures the user can access the video without interruption when offline.

## Streaming Processing:

1. **Scenario:** *A user signs up for the streaming service and wants to start using the service immediately.*
   **Process:**
   - The user's sign-up data is sent through the pipeline as soon as they complete the registration.
   - The system processes the data in real-time, updating the database and enabling the user to access the service without delay.
   - This approach ensures a seamless and immediate experience for new users.

2. **Scenario:** *Spotflix's Online Streaming*
   **Process:**
   - When a user watches a video online, it is streamed in real-time, with data sent and processed continuously.
   - This approach provides an immediate watching experience without waiting for the entire video to download.

## Conclusion:

In this document, we have explored essential concepts in data engineering, focusing on the differences and appropriate use cases for ETL (Extract, Transform, Load) and ELT (Extract, Load, Transform) processes, as well as batch and streaming pipelines. By understanding these fundamental processes and their practical applications, we can make informed decisions when designing and implementing data processing systems.

## References:

1. IBM. ELT vs. ETL: What's the Difference? Retrieved from https://www.ibm.com/blog/elt-vs-etl-whats-the-difference/
2. DataCamp. Understanding Data Engineering. Retrieved from https://app.datacamp.com/learn/courses/understanding-data-engineering