# Data Engineering & Pipelines

## What is Data Engineering?

Data engineering is the practice of designing and building systems for collecting, storing, and analyzing data at scale. The primary goal is to make data available and usable for various purposes, such as business intelligence, machine learning, and reporting.

## The Stages of a Data Pipeline

A data pipeline is a set of processes that moves data from one or more sources to a destination. The stages are commonly described by two acronyms: **ETL** and **ELT**.

### ETL (Extract, Transform, Load)

- **Extract:** Data is pulled from a source (e.g., a database, API, or CSV file).
- **Transform:** Data is cleaned, filtered, and aggregated *before* being loaded.
- **Load:** The transformed data is loaded into the final destination.

### ELT (Extract, Load, Transform)

- **Extract:** Data is pulled from the source.
- **Load:** The raw data is loaded directly into the destination (e.g., a data warehouse or data lake).
- **Transform:** Data is cleaned and prepared for use *after* it has been loaded.

ELT is often favored in cloud environments because modern data warehouses can handle the heavy lifting of the transformation step more efficiently.

## Key Pipeline Components

When viewed from a high level, a data pipeline can be broken down into four essential components:

1. **Ingestion:** This is the starting point, where raw data is collected from its source. This step is equivalent to the "Extract" stage.
2. **Processing:** The core of the pipeline, where raw data is cleaned and structured. This is the "Transform" stage. Common tasks include:
   - **Data Cleansing:** Handling missing values, correcting errors.
   - **Data Aggregation:** Summarizing data (e.g., calculating daily sales totals).
   - **Data Formatting:** Converting data types or standardizing text.

3. **Storage:** The final destination for the processed data. This is the "Load" stage, and the storage system is chosen based on how the data will be used. Examples include data warehouses, data lakes, or databases for applications.
4. **Output:** How the final data is consumed. This could be a business intelligence dashboard, a machine learning model, or a report for business users.

# The Role of Workflow Automation

Workflow automation is the practice of scheduling, running, and monitoring the tasks within a data pipeline automatically. This is handled by a dedicated tool called a workflow orchestrator.

### Why is Automation Important?

Automation is crucial for:

- **Reliability:** It ensures tasks run in the correct sequence, with automatic retries for failed jobs.
- **Efficiency:** It removes the need for manual, time-consuming tasks.
- **Scalability:** It allows for the management of hundreds or thousands of pipelines at once.
- **Visibility:** It provides a central hub to monitor pipeline health and troubleshoot issues.

# Real-Life Examples of Automation

Here are some examples of how workflow automation is used in real-world data pipelines:

1. **E-commerce Sales Reporting:** A workflow orchestrator (like **Apache Airflow**) schedules a pipeline to run every morning at 3:00 AM. It pulls sales data from a database, transforms it into a daily summary, and loads the results into a data warehouse. If the database connection fails, the orchestrator automatically retries the job.
2. **Marketing Analytics Dashboard:** A tool like **Prefect** manages the dependencies for a pipeline that pulls data from various ad platforms (Google Ads, Facebook Ads). The tool ensures the data processing only begins *after* all source data has been successfully ingested, preventing the dashboard from displaying incomplete or inaccurate information.
3. **Real-Time IoT Sensor Monitoring:** An orchestrator monitors a continuous stream of data from factory sensors. If the data flow from a specific sensor stops, the automation tool immediately triggers an alert, allowing the operations team to quickly identify and fix the issue.