

Video 10 Summary: Understanding Data Pipelines Through Analogy

Data as Raw Material

In cloud data analysis, data is like raw materials in a factory—it must be processed **efficiently** and **consistently** to be useful.

What Is a Data Pipeline?

A **data pipeline** is a series of processes that move data from **source** to **destination** for **storage** and **analysis**.

It consists of **three main stages**:

1. **Extract**
 2. **Transform**
 3. **Load**
-

Factory Analogy: Toy Car Assembly Line

- **Extract:** A worker gathers all parts and places them on a tray.
- **Transform:** Another worker assembles the toy car.
- **Load:** A third worker packages the car.

This mirrors how data is handled in a pipeline.

Detailed Breakdown of ETL Stages

1. Extract

- **Definition:** Retrieving data from one or more sources.
- **Action:** Move raw data to a **staging area**.
- **Example:** An animal rescue organization extracts data from:
 - A CSV file with microchip registrations.
 - A database with pet owner contact info.

2. Transform

- **Definition:** Cleaning and formatting data.
- **Action:** Remove duplicates, fix errors, standardize formats.
- **Goal:** Make data usable for analysis.

3. Load

- **Definition:** Inserting data into a **target system** (e.g., database, warehouse, lake).
 - **Action:** Store transformed data for future use.
 - **Example:** The rescue organization loads data into a **data warehouse**.
-

Pro Tips

- Data pipelines are **customizable**—they may not follow the same order or steps every time.
- Like assembly lines, they are tailored to specific tasks and data types.
- Pipelines help:
 - **Automate** data handling
 - **Save time and resources**
 - **Improve accuracy**
 - **Increase data value**