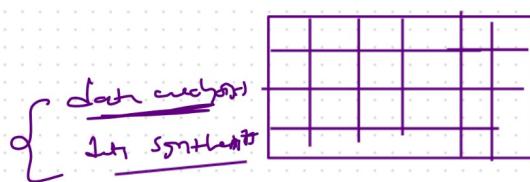


ETL vs ELT [key differences in Data Processing]

ETL [Extract, Transform, Load] ?

- ✓ Extract: Collect data from multiple sources [MySQL, JSON, XML, log file etc.]
- ✓ Transform: Clean, filter, and format the data
- ✓ Load: Store cleaned data into a data warehouse [data engineer]

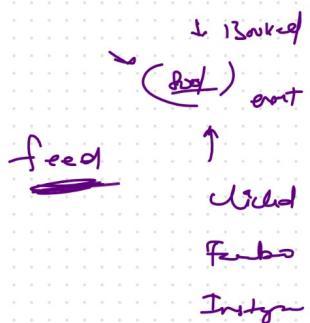
Structured data



ETL Workflow Examples:

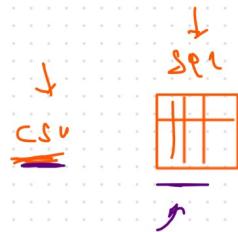
- c1) Collect data from multiple sources

- ex: At OYO Rooms data came from
 - * customer data in MySQL
 - * booking data from PostgreSQL
 - * user activity data
 - * marketing data from social platforms
 - * some data from XML format



- c2) Transform the data:

- * clean, filter & standardized data format
- * make different data sources compatible



- c3) Load into data warehouse

- * store in optimized structure for analysis
- * serves as the foundation for business intelligence

Q2 Food Delivery App

①

Collect data: Restaurants menus, Customer orders, tracking info

Transform: remove duplicate orders, invalid addrs

Load: Store clean data in data warehouse

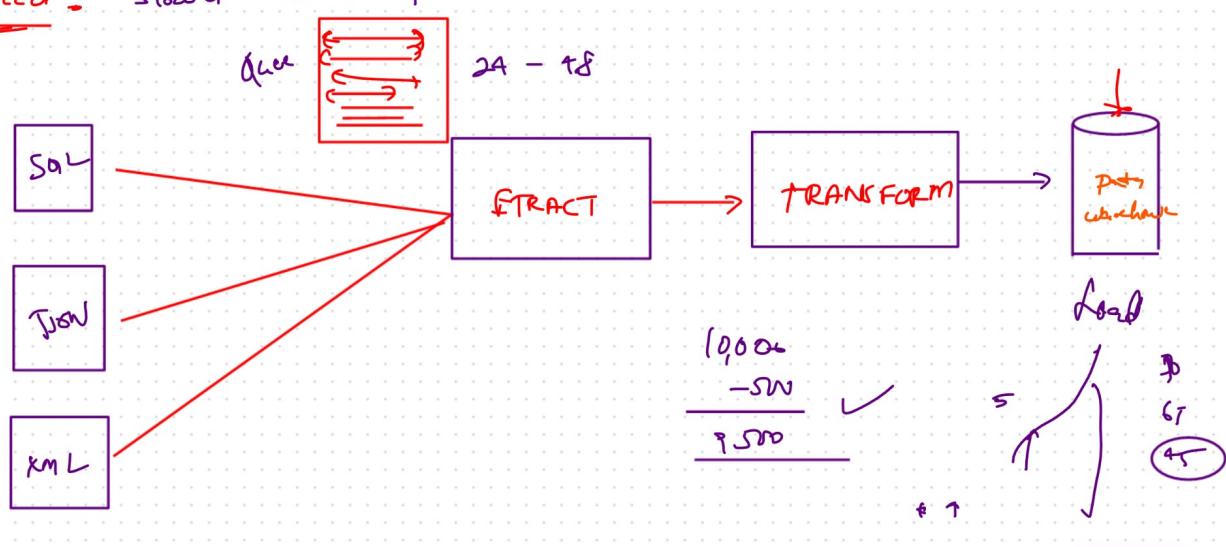
ETL Characteristics

Processing location: Data transformed before loading

Infrastructure: Traditional data warehouses with compute power

Data types: Best for structured data

Speed: Slower as transformation occurs beforehand



Use case: Bank transactions, order history, conversion cycles

Popular Tools

• Informatica

• Talend

• Pentaho

ELT (Extract, Load, Transform)

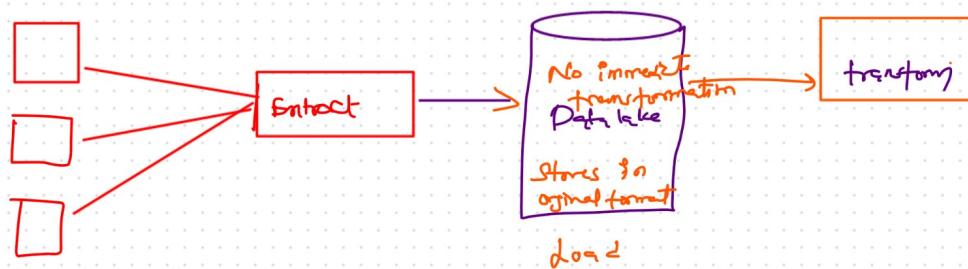
What is ELT?

Extract : Get data from multiple sources

Load : Store raw data in a data lake

Transform : Process data after it's stored, as needed

Workflow

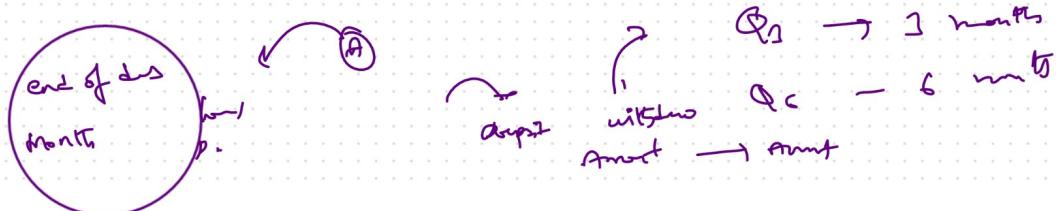


User Case Example : Food Delivery App

- * Raw data collections : Customer reviews (text), food photos(images), GPS data
- * Load into data lake : Store all formats without transformation
- * Transform data : Process as needed for specific analysis

ELT Characteristics

- * Processing location : Data transformed after loading
- * Infrastructure : Data lakes / cloud platforms with high compute power
- * Data types : Handles structured, semi-structured & unstructured data
- * Speed : Faster data ingestion as transformation happens after loading
- * Use cases : Social media analysis, storing raw data for future use



Common ETL Tools

- + Azure Data Factory
- Azure Data bricks
- + AWS Glue
- + AWS S3

Core Option



Data Engineer vs Big Data Engineering

Data Engineers — Core Concepts

Def: A data engineer designs, builds & maintains systems that can store pixels & make data accessible for analysis

Analogy: Water Pipeline Builders

Sources: Raw data from different origins (like water from rivers, lakes, oceans)

Process: Design & build pipelines [to collect → clean → store → distribute data]

End Users: Data analysts, data scientists, management, others

Primary Responsibilities

- Data pipeline creation
- ↓ Data transformation
- Ensuring data quality
- Automation

Traditional Tools & Approaches

- SQL, Python, Airflow, ETL tools
- ETL + ware warehouse
- focus on structured data, manageable data
- Storage: data lakes & data warehouses
 - ↳ monthly sales reports, standard business analytics

Big Data Engines

Definition: Specialized in tools & technologies for processing massive datasets
that traditional systems cannot handle effectively

Analogies DAM Engines

- * Handles "massive floods" of data
- * Uses specialized systems for extreme volume & complexity
- * required distributed computing & processing system

Characteristics of Big Data: (Five V's)

Volume: extremely large datasets

Variety: Structured, semi-structured, & unstructured data

Velocity: High-speed data processing

Veracity: Data quality & reliability

Value: Extracting meaningful insights

Tools & Technologies:

- * Hadoop, Spark, Hive, Kafka
- * NoSQL databases
- * Distributed Storage systems
- * Data lakes

Modern Approaches

- ELT + Big Data technologies
- Parallel processing for large-scale data
- e.g.: Analyzing twitter trends, real-time data processing