

# ETL Fundamentals

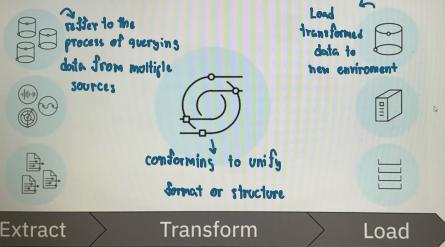
## Objectives

- ETL isn't automated data pipeline

After watching this video, you will be able to:

- Describe what an ETL process is
- Describe what data extraction means
- Describe what data transformation means
- Describe what data loading means
- List use cases for ETL processes

### What is an ETL process?



### What is an ETL process?

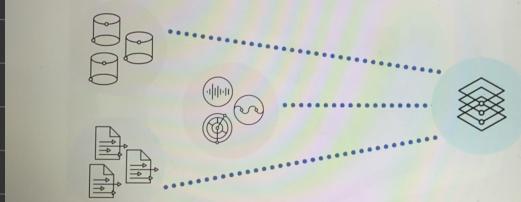
T=Transformation:



Transforming data into the format for the output

### What is an ETL process?

E=Extract: Extracting data from a source



### What is an ETL process?

L=Load: Loading data into a database, data warehouse or other storage



### What is Extraction?

- Configuring access to data and reading it into an application: normally, this isn't automated process
  - Web scraping using python or R
  - Connecting programmatically via APIs
- The data may be static or streaming online

### What is data transformation?

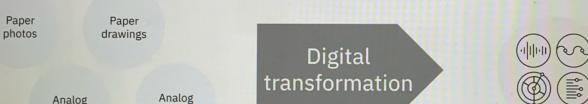
- Processing data
- Conforming to target systems and use cases
- Cleaning
- Filtering
- Joining
- Feature engineering
- Formatting and data typing

### What is data loading?

- Moving data into a new environment
- Examples: a database, data warehouse, or data mart
- Making the data readily available for analytics, dashboards, reports

### Use cases for ETL pipelines

- Digitizing analog media



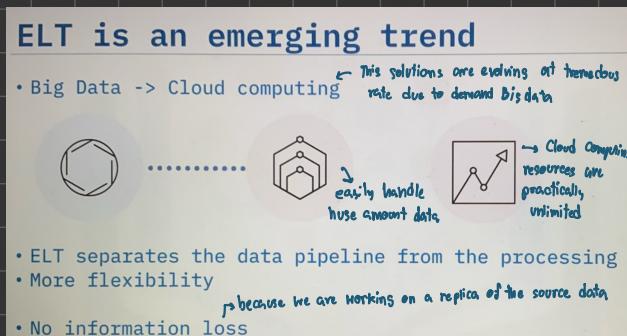
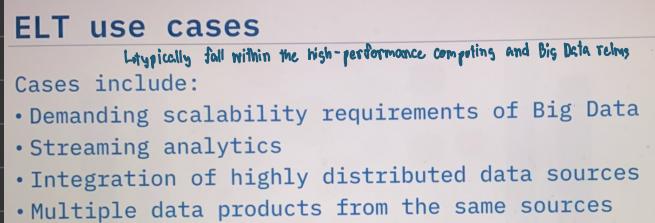
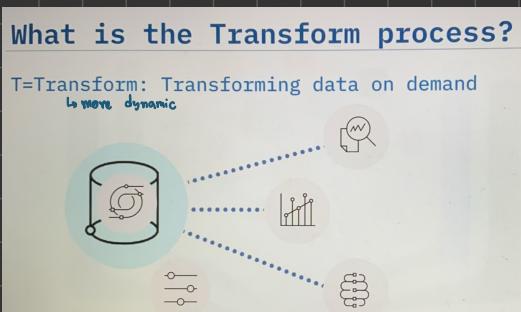
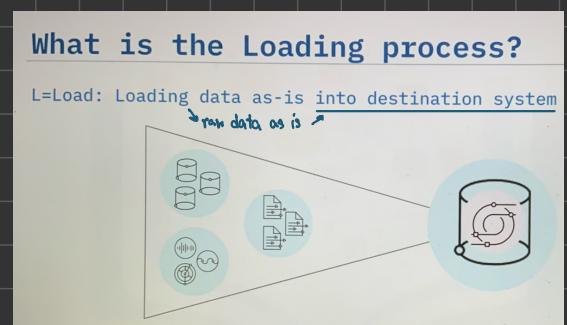
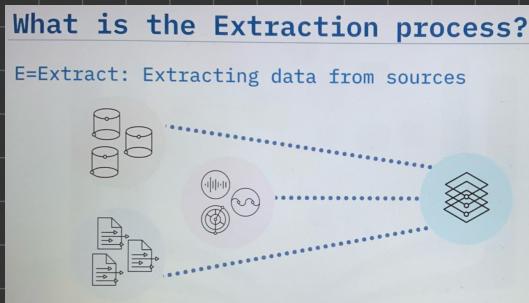
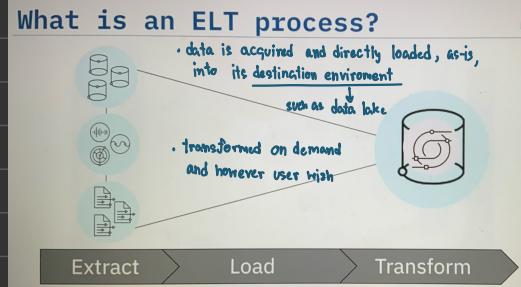
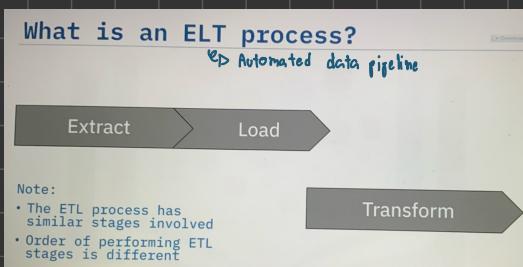
- Moving data from OLTP systems to OLAP systems
- Dashboards
- Machine learning

# ELT Basics

## Objectives

After watching this video, you will be able to:

- Describe what an ELT process is
- List use cases for ELT processes
- Describe why ELT is an emergent trend



# ETL vs ELT

## Objectives

- After watching this video, you will be able to:
- List key differences between ETL and ELT
  - Describe ELT as an evolution of ETL
  - Describe the trending shift from ETL to ELT

## Differences between ETL and ELT

When and where the transformations happen:

- Transformations for ETL happen within the data pipeline before data reached to destination
- Transformations for ELT happen in the destination environment

Flexibility: fixed process

- ETL is rigid pipelines are engineered to user specifications
- ELT is flexible – end users build their own transformations

## The evolution of ETL to ELT

- Increasing demand for access to raw data
  - Intermediate storage facility → holding area for raw extracted data
  - then load to data marts/mart
- In ELT, the staging area fits the description of a data lake
- Staging areas – private ETL landing zones
- Self-serve data platforms are the new “staging area”

## Differences between ETL and ELT

Support for Big Data:

- Organizations use ETL for relational data, on-premise – scalability is difficult
- ELT solves scalability problems, handling both structured and unstructured Big Data in the cloud

Time-to-insight:

- ETL workflows take time to specify and develop
- ELT supports self-serve, interactive analytics in real time

## The shift from ETL to ELT

ETL still has its place for many applications

ETL ..... Shift ..... ELT

ELT addresses key pain points:

- Lengthy time-to-insight
- Challenges imposed by Big Data
- Demand for access to siloed information

# Data Extraction

## Objectives

After watching this video, you will be able to:

- List examples of raw data sources
- Describe data extraction techniques
- Relate use cases with data sources and extraction techniques

### Examples of raw data sources



### Examples of raw data sources



- Data is everywhere

### Techniques for extracting data

Data extraction techniques include:

- OCR Optical Character Recognition : used for text scanned from paper
- ADC sampling, CCD sampling → charge-coupled devices : capture and digitize images
- Mail, phone, or in-person surveys and polls
- Cookies, user logs for tracking human or system behavior

Analog → Digital

### Techniques for extracting data

More techniques include:

- Web scraping
- APIs
- Database querying SQL, NoSQL
- Edge computing
- Biomedical devices

### Use cases

- Integrating disparate structured data sources via APIs
- Capturing events via APIs and recording them in history to capture periodic or asynchronous events to store
- Monitoring or surveillance with edge computing devices
- Data migration (direct to storage) for further processing
- Diagnosing health problems with medical devices

# Data Transformation

## Objectives

After watching this video, you will be able to:

- Name data transformation techniques
- Compare schema-on-write vs. schema-on-read
- List ways information can be “lost in transformation”

## Data transformation techniques

↳ mainly about formatting the data to suit the application

Data transformations can involve various operations, such as:

- Data typing converting one data format to another, such as JSON, CSV to table
- Data structuring
- Anonymizing, encrypting



## Data transformation techniques

Other types of transformations include:

- Cleaning: duplicate records, missing values
- Normalizing: converting data to common units
- Filtering, sorting, aggregating, binning
- Joining data sources

## Schema-on-write vs. schema-on-read

Schema-on-write is the conventional ETL approach:

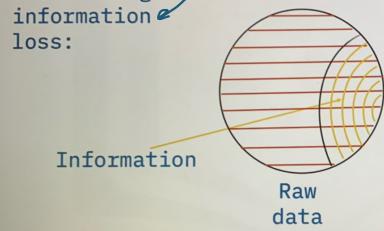
- Consistency and efficiency
- Limited versatility

Schema-on-read applies to the modern ELT approach:

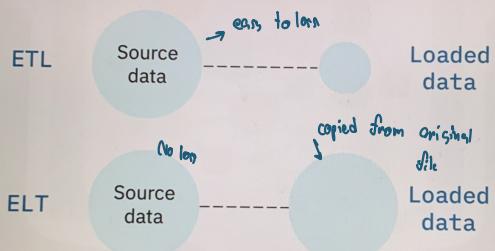
- Versatility
- Enhanced storage flexibility = more data

## Information loss in transformation

Visualizing information loss: there are many ways in which information can be ‘lost’ in transformation



## Information loss in transformation



## Information loss in transformation

Examples of ways information can be lost in transformation processes include:

- Lossy data compression
- Filtering
- Aggregation
- Edge computing devices



# Data Loading

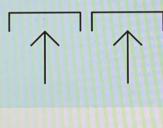
## Objectives

After watching this video, you will be able to:

- List data loading techniques
- Differentiate batch loading from stream loading
- Explain push vs. pull
- Describe parallel loading

## Data loading techniques

- Full
  - ↳ load an initial history into database
  - Full is applied to insert new data or to update
- Incremental already loaded data
  - ↳ schedule data loading
  - Scheduled
    - ↳ load it as required
  - On-demand
- Batch and stream
  - ↳ push to server or clients by server
- Push and pull
- Parallel and serial



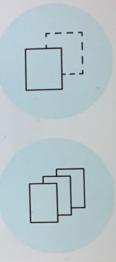
## Scheduled vs. on-demand loading

- ↳ Data is often loaded on schedule
- Scheduled loading
- Periodic loading, like daily transactions to database
  - Windows Task Scheduler, cron tools
- On-demand loading, triggered by ↳
- Measures such as data size
  - Event detection, like motion, sound, or temperature change
  - User requests, like video or music streaming, web pages



## Full vs. incremental loading

- Full loading → loading data in one large batch
- ↳ Start tracking transactions in a new data warehouse
- Used for porting over transaction history to the new system
- Incremental loading → to ensure the transaction history is tracked
- Data is appended to, not overwritten
  - Used for accumulating transaction history
  - Depending on the volume and velocity of data, can be batch loaded or stream loaded



## Batch vs. stream loading

- Batch loading
- Periodic updates using windows of data
- Stream loading
- Continuous updates as data arrives
- Micro-batch loading
- Short time windows used to access older data



## Push vs. pull technology

- basis on ↳
- Client-server** model
- then server responds to ↳ clients request
- Pull - requests for data originate from the client
  - For example: RSS feeds, email ↳ email box
- Push - server pushes data to clients
- For example: push notifications, instant messaging



## Parallel loading

- ↳ employed to
- Multiple data streams
- ↳ to boost loading efficiency
- ↳ particularly when data is big and long distances
- 
- Destination

## Parallel loading

- ↓
- File partitioning
- ↳ splits file into smaller chunks, the chunks can be loaded simultaneously

