



# Class-1 - Big Data

Audio ✓  
Video ✓  
Screen ✓

# Zoom

raise your hand

## Other people chat

chat

$\leftarrow \rightarrow Q_n A$

← → Raise hand

( lot active )

(very less)

## clearing doubts

[ Doubt on  
concept ]

## Community Chat

@ MayonR

## Agenda

→ Who am I ?

## → Introduction & Getting into Big Data

Mayank Aggarwal

→ 2012 → Coding  
love maths → Coders

2013

NSIT, Delhi → 1<sup>st</sup> year coded

# 2<sup>nd</sup> year Data Science

→ IIT-1

3<sup>rd</sup> → ML, DL, AI, DSA

$\rightarrow 0\gamma 0$

$$4^m \rightarrow \text{PPD} +$$

Professional

DYO → 5 months

Goldman Sachs → Big Data or Data engineering

Mindhive → Spark + AI/ML

Ineuron

Scaler + Coding Ninjas

Course

1. Notes + Code + hints + Recording + files + pdf

↳ Critique  
↳ Resource section

2. Class 8 - 11 Doubt session  till 12

Few Instruction / Good practice

1. Be a leanR paper
2. Make your own notes. [use pens]
3. Don't jump to Future
4. Theory + Practical → Projects
5. Great Teacher → Everything + how to learn +  
vs  
Good Teacher → everything
6. Dumb with doubts cleared than smart with having doubts

Feedback → -ve skewed

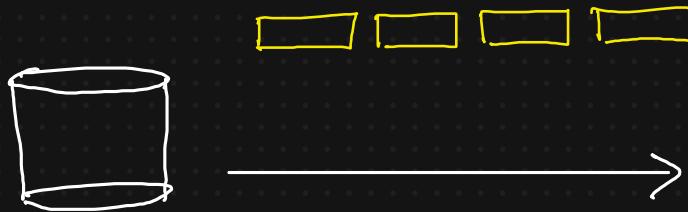
1 we have syllables which is shared

# Big Data

## Data Science

Kaggle, mail, pendrive, sqe

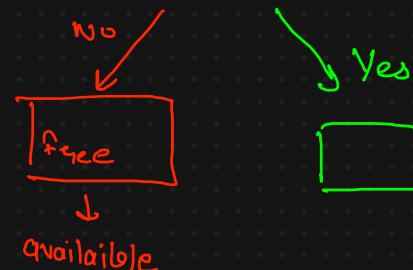
### WhatsApp Chat Automation



#### Problems

- transaction
- small volume

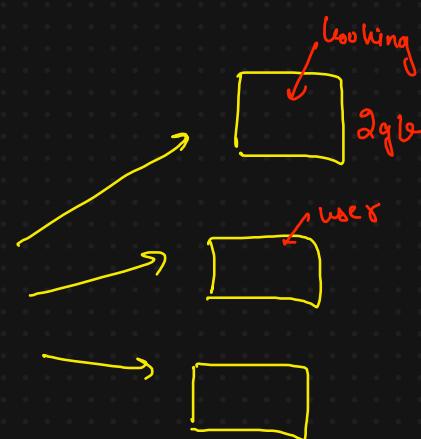
Hi Mayank, we are looking forward to host you for your booking ID SCTP8012 at OYO 12191 Hotel Airlift, Delhi. To confirm your arrival, please reply 'confirm' to this message. If your plans have changed and you need to cancel this booking, please reply 'cancel'.



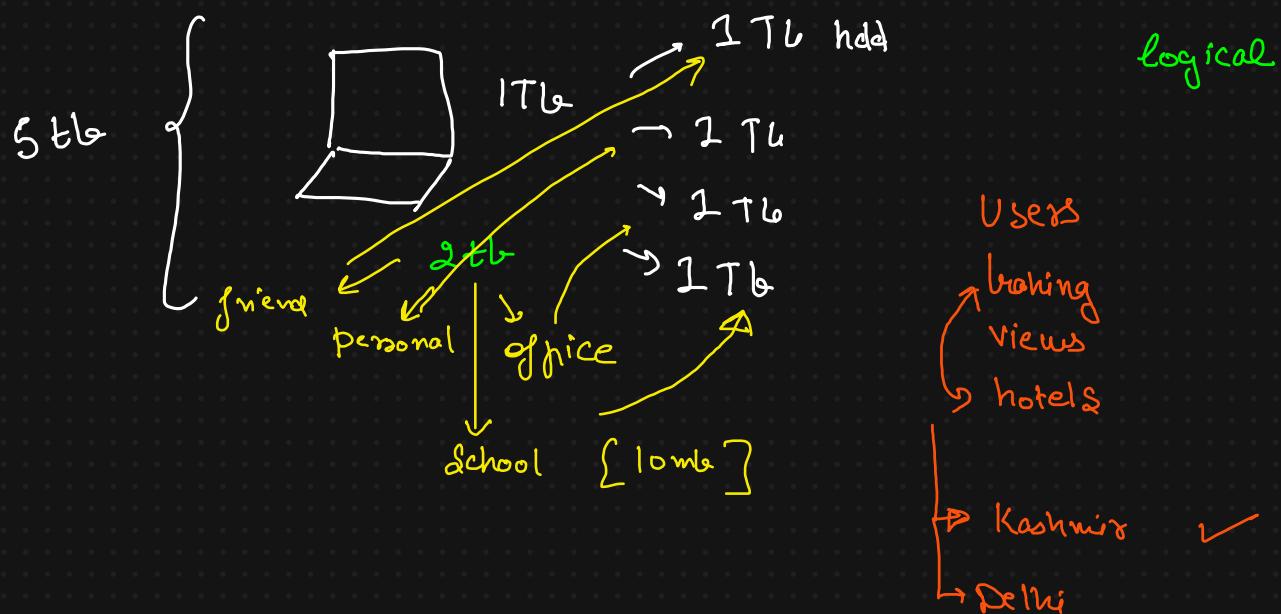
2nd project → Delhi - 2x



gb/tb can be stored



- ① 1. There is still a limit which we can expand
- 2. as data increases we have to position our data
- 3. Unstructured data ✓
- 4. Costly.



**Big data :** huge amount of data that cannot be handled by traditional systems.

These systems deal with data that is too large & complex for traditional system to handle

Big Data is problem  
& solution

Why Big data emerge?

1. Internet
2. Social Media
3. IoT

# 5 V's of Big Data

5 Vs of Big Data are crucial / provide a framework for understanding challenges & characteristic of Big Data.

1. Volume :- Size of data getting generated

OYO vs GS

there is no set volume  
when your problem  
becomes a big data

2. Velocity: Speed at which data is getting generated & needs to be processed.

Real Time

Bank Transaction

Fraud

live feed

Near Real Time

data is continuously  
generated but we are  
taking 1min, 2min  
rather

Amazon

→ Confirm in 2min

live tracking

Batch

Credit card

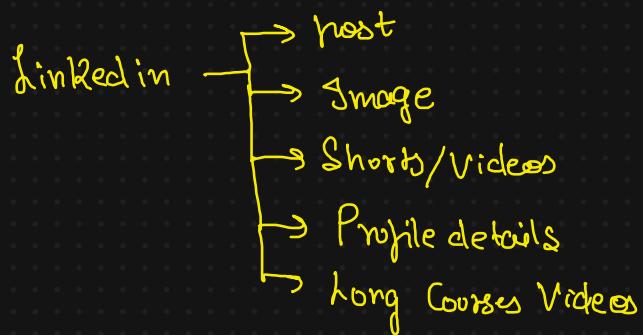
[data engineer]

Once a month



specific time period

3. Variety data can be in different formats & we have to deal with all of them.



Structured

Row & columns

Schema is enforced

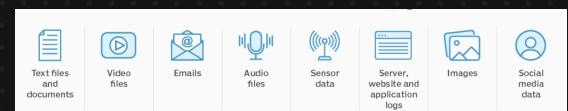
Semistructured

JSON

XML

CSV

Unstructured



#### Structured data

ID	Name	Age	Degree
1	John	18	B.Sc.
2	David	31	Ph.D.
3	Robert	51	Ph.D.
4	Rick	26	M.Sc.
5	Michael	19	B.Sc.

```
<University>
<Student ID="1">
  <Name>John</Name>
  <Age>18</Age>
  <Degree>B.Sc.</Degree>
</Student>
<Student ID="2">
  <Name>David</Name>
  <Age>31</Age>
  <Degree>Ph.D. </Degree>
</Student>
...
</University>
```

Parquet

```
{
  "customer": "John Doe",
  "gender": "Male",
  "items": [
    {
      "SSN": "123-45-6789",
      "exp": "10/01/2025",
      "moresubnesting": {
        "SSN": "123-45-6789",
        "newfield2": "123 Main Street"
      }
    }
  ]
}
```

4. Veracity: Trustworthiness or Quality of data

→ messy data needs to be corrected.

→ -ve age, Quality issues

5. Value: data should be useful & provide some value

→ help in business decision

→ provides insights

Volume  $\neq$  Big data

For a problem to be a big data problem, not all V's need to be satisfied. Generally a combination of 2~3 is enough

# Big Data & Distributed Systems.

gta 3 /  $\longrightarrow$  gta 5  
Sam Andreas

We need more resources

1. Storage
2. Memory (RAM)
3. Processors

## Single System

Monolithic Systems  
 $\downarrow$   
1



Single, self contained unit

$X \rightarrow 2X \rightarrow 10X$

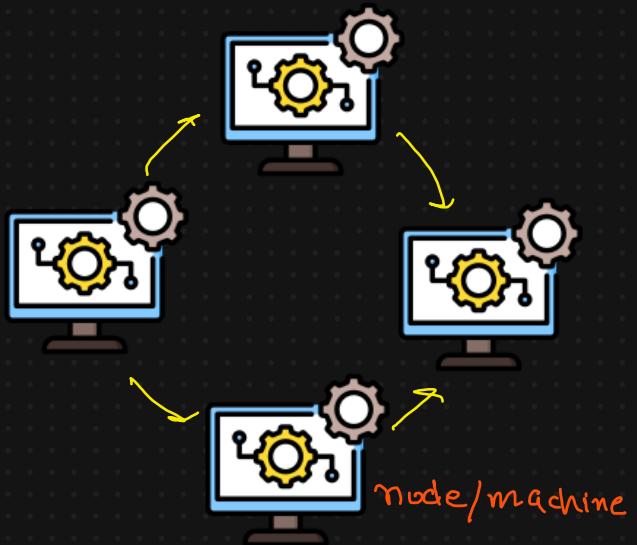
$P \rightarrow 2P \rightarrow 6P$

Scalability problem

cost & performance

## Multiple Systems

Distributed architectures



True Scaling

$X \rightarrow 2X \rightarrow 10X$

$P \rightarrow 2P \approx 10P$

they don't have problem of scalability

## Vertical Scaling

## Horizontal Scaling

All good big data systems are based on distributed architecture

Designing a Good Big Data System

Cloud vs On-Premise

Database vs Data Warehouse vs

Data Lake

ETL vs ELT

Data Engineer vs Big Data

Engineer

Hadoop Introduction

## Class 2 : Terms in Big Data

### Designing a good Big data System

1. Scalability: Flipkart Big billion

(a) Storage scalability: store increased data

(b) Processing scalability: process large / complex data faster by adding more m/c.

2. Reliability and Fault tolerance: it should continue to work / operate even if some component fails.



3. Cost effectiveness: system should be able to balance performance & scalability with cost

4. Security and Data Privacy: data is saved from unauthorized access.

## On Premise vs Cloud

When choosing to deploy Big Data solutions, companies face an issue between On-premise vs Cloud.

### On-Premise

Buying an office

- space
- initial investment
- procure all machines
- install software, OS, security updates
- setting up → operating expenditure

$$\begin{array}{c} \square \\ \triangle \end{array} * 20 = 20 \text{ nodes}$$

More machines → scaling up

+ 10 more machines

### 1. Security

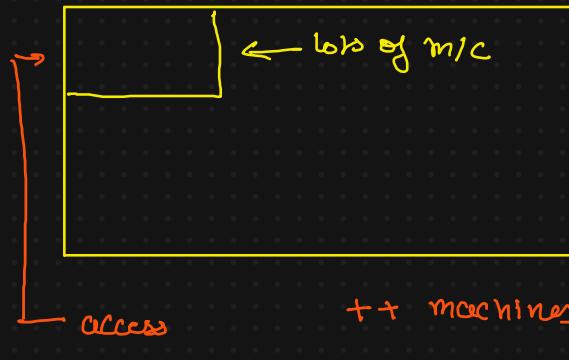
### Cloud

aws  
gcp  
azure  
IBM



Renting

pay-as-you-go  
model



++ machines  
-- machines

Info on rent

### Advantages.

1. Scalability
2. Operating expenses: More operating expenses  
Capex vs Opex
3. Fault tolerant

Aspect	On-Premise	Cloud
<b>Deployment</b>	Hardware and software are hosted within the organization's facilities.	Resources and services are hosted on the provider's servers and accessed via the internet.
<b>Cost Model</b>	High upfront costs for hardware, maintenance, and IT staff. <i>capex</i>	Pay-as-you-go pricing, with minimal upfront costs. <i>opex</i>
<b>Scalability</b>	Limited by the organization's hardware capacity; scaling requires purchasing and installing new equipment.	Highly scalable—add resources on demand instantly.
<b>Maintenance</b>	The organization is responsible for managing and maintaining hardware, software, and updates.	Managed by the cloud provider (e.g., AWS, Azure, GCP).
<b>Flexibility</b>	Fixed capacity with little room for dynamic needs.	Highly flexible, allowing resources to scale up or down.
<b>Security</b>	Data remains within the organization's premises, offering greater control.	Security is managed by the provider; often meets global compliance standards but might raise concerns for sensitive data.
<b>Disaster Recovery</b>	Requires internal backups and disaster recovery systems.	Cloud providers offer built-in disaster recovery and redundancy.

## Types of cloud

1. Public : AWS, Azure, GCP

2. Private: infra set up by company. *Compliance & Security*

3. Hybrid : Private + Public

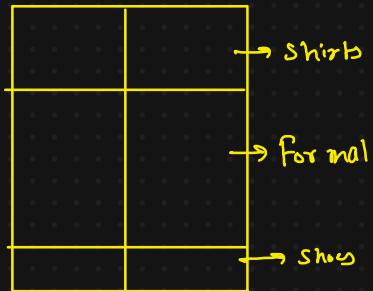
4. Community : Shared by organization with common concern.

*eg: University  
Hospitals*

# Data Base vs Data warehouse vs Data Lake

## Database

- system which stores the data
- organized i.e. has a structure
- Structured data with defined schema



Transaction data = Day to day operations  
e.g. (updated balance  
order,  
inventory)

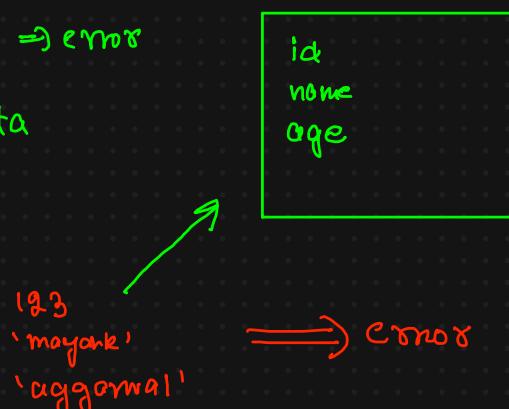
## OLTP (Online transaction processing)

We store just the recent data due to performance & cost

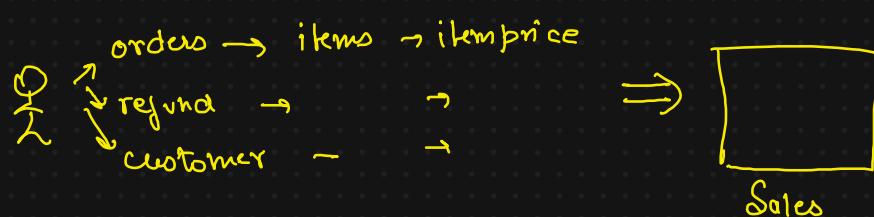
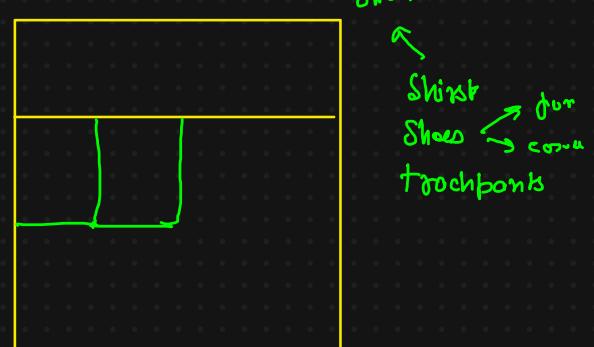
don't hold years of data, recent data for perform.

Schema on write      only mismatch  $\Rightarrow$  error

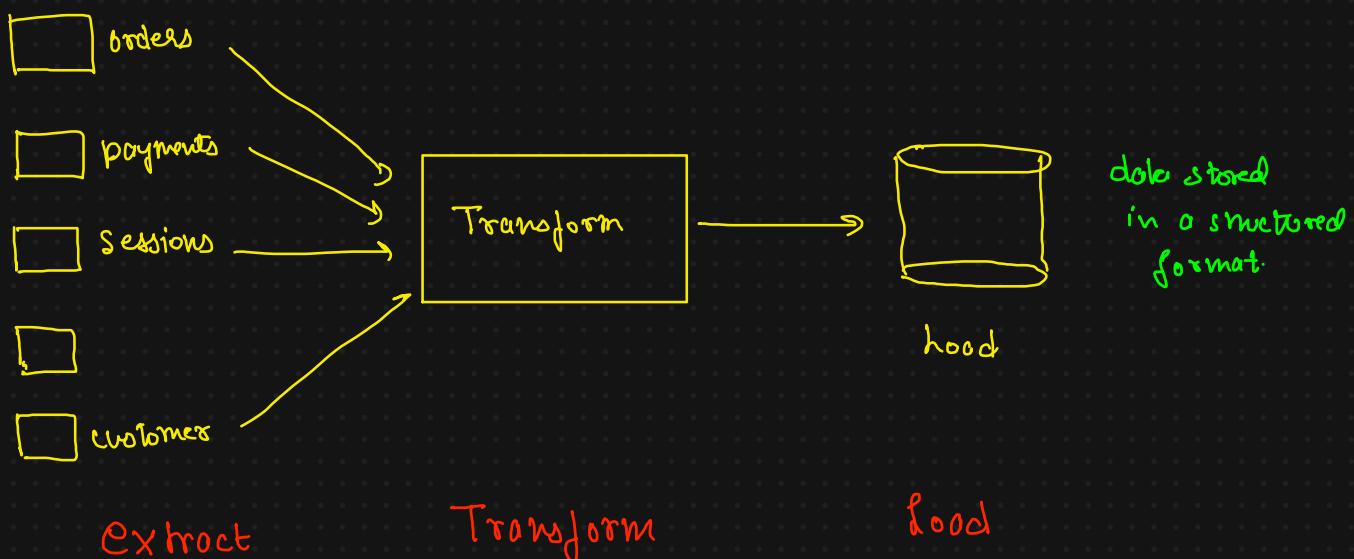
Overall cost high  $\Rightarrow$  just for recent data



## Data Warehouse



# Treppodato, Informatica, Tokend



Q: Why can't I use a database?

1] There are multiple sources of data

2] lots of data  $\Rightarrow$  cost  $\uparrow\uparrow$

3] We write lots of queries quite complex on DB

performance

↓↓↓

Complex  
Queries

transaction  
work

Read Heavy

OLAP  $\Rightarrow$  Online Analytical Processing

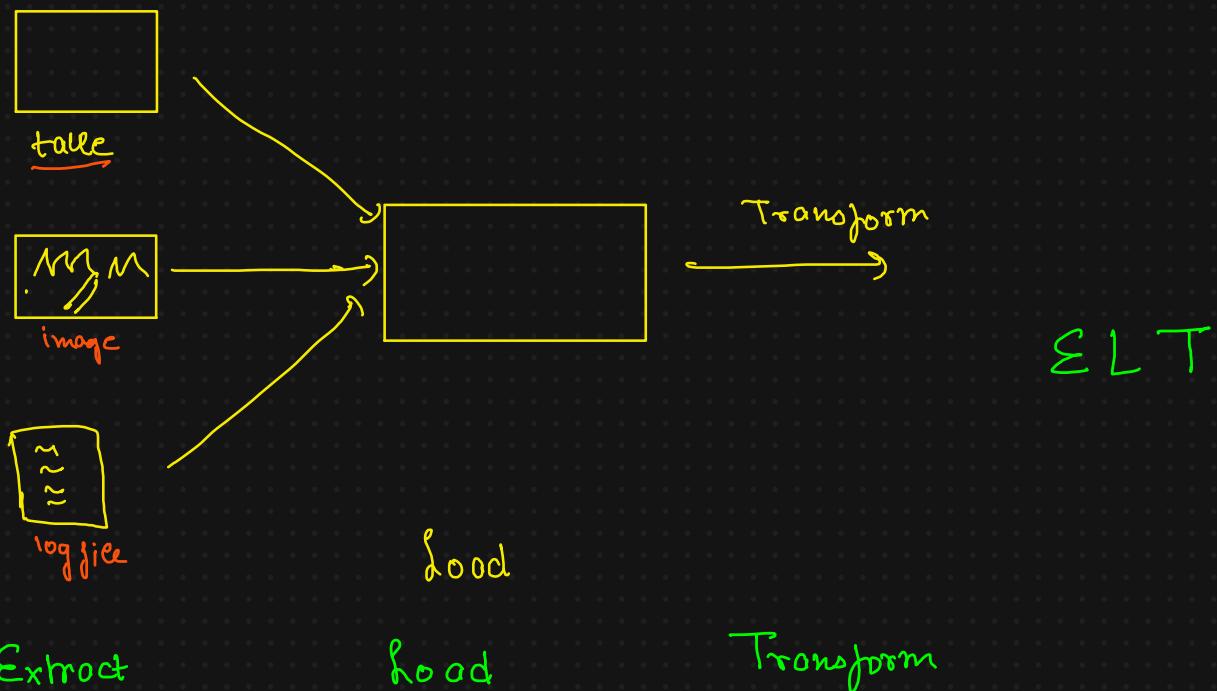
3. Data Lake 23, halfs

We are now mainly dealing with different kinds of data

↓      ↓      ↓  
Structur. unstru. semi-structured

A data lake is a centralized repository that stores all kind of raw data at scale, without any pre-structuring.

nlp



⊕ cost effective  
enough flexibility

Many teams can use this

Feature	Database	Data Warehouse	Data Lake
Purpose	Real-time transactions	Historical data analysis	Raw data storage for diverse use cases
Data Structure	Structured	Structured	Structured, Semi-structured, Unstructured
Speed	High-speed for small queries	Optimized for analytical queries	Variable (depends on data processing)
Use Case	Operational systems (e.g., POS)	Business intelligence	Advanced analytics, machine learning
Scalability	Limited	Moderate	Highly scalable

Break  $\Rightarrow$  10min

10:00

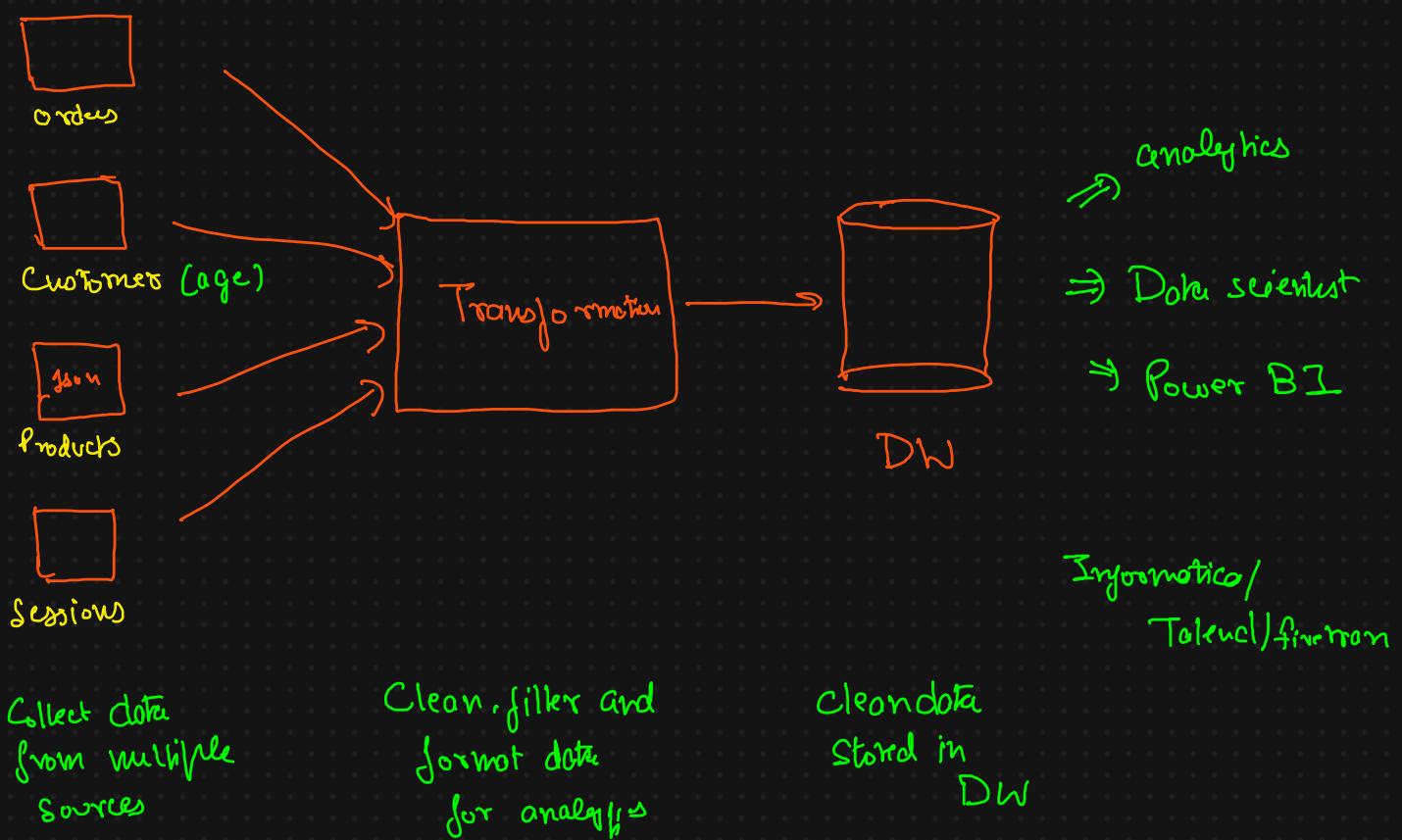
## Transformation :-

Aspect	OLTP	OLAP
Purpose	Handles day-to-day transactional data.	Handles historical data for analytical queries.
Data	Real-time, current operational data.	Historical, aggregated, or summarized data.
Operations	Frequent INSERT, UPDATE, DELETE operations.	Primarily SELECT queries for analysis.
Query Complexity	Simple, fast queries (e.g., fetching single records).	Complex queries with joins, aggregations, and filters.
Schema Design	Normalized schema (e.g., 3NF) to minimize redundancy.	De-normalized schema (e.g., star or snowflake schema) for faster querying.
Transaction Type	Short, atomic transactions (e.g., order placement).	Long, batch-oriented queries (e.g., sales trends analysis).
Users	Operational staff, front-end applications.	Data analysts, decision-makers.
Data Volume	Smaller data sets, as it handles current operations.	Very large data sets (terabytes to petabytes).
Performance	Optimized for fast write operations.	Optimized for read-intensive operations.
Examples	Banking systems, e-commerce sites, CRMs.	Data warehouses, BI tools, analytical dashboards.
Tools/Databases	MySQL, PostgreSQL, MongoDB, Oracle (OLTP mode).	Snowflake, Amazon Redshift, Google BigQuery.

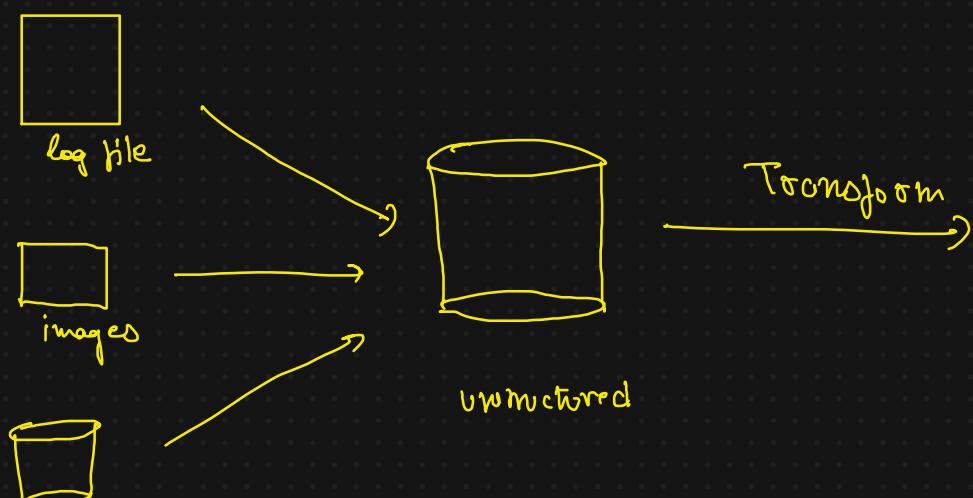
# ELT vs ETL

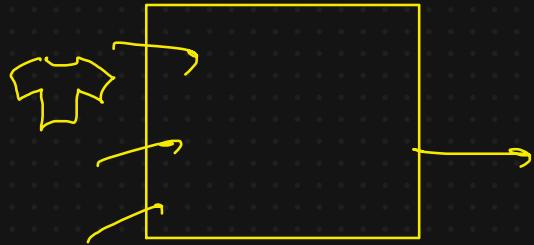
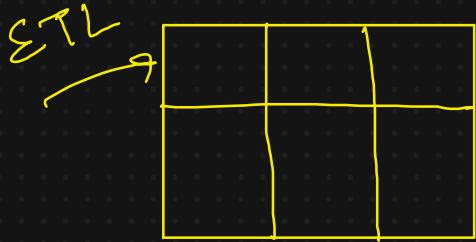
Extract Transform Load

homework



ELT process

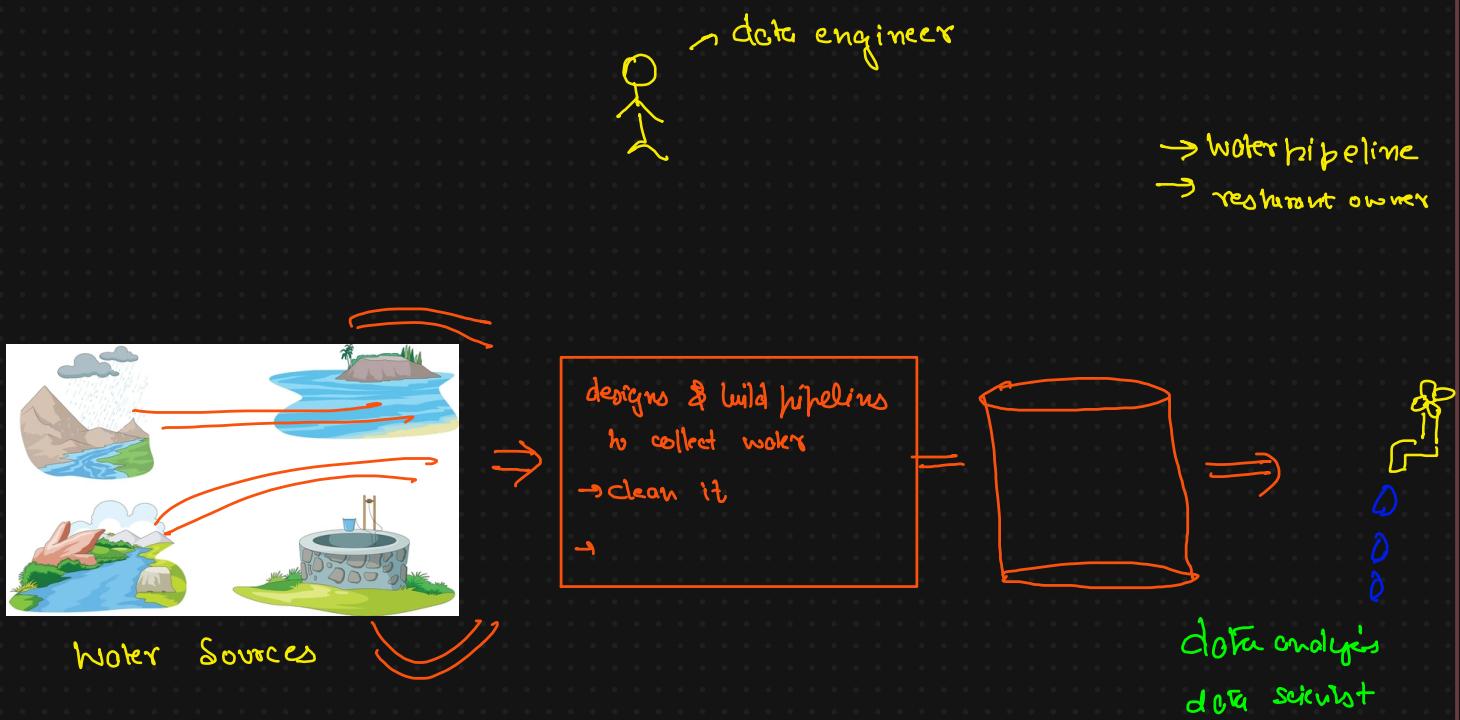




Aspect	ETL (Extract, Transform, Load)	ELT (Extract, Load, Transform)
<b>Processing Location</b>	Data is transformed before loading.	Data is transformed after loading.
<b>Target System</b>	Traditional data warehouses with limited compute power.	Modern data lakes or cloud platforms with high compute power.
<b>Data Type</b>	Structured data.	Structured, semi-structured, unstructured data.
<b>Speed</b>	Slower, as transformations occur beforehand.	Faster, as transformations happen after loading.
<b>Use Case</b>	Bank transactions (cleaned before storage).	Social media analysis (raw data stored for future use).



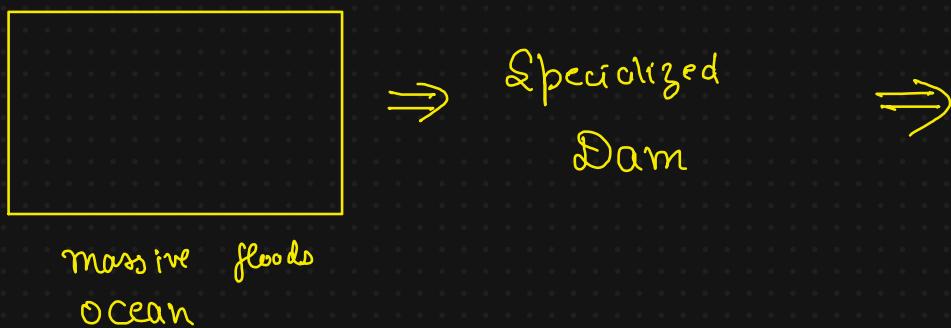
# Data engineer vs Big Data engineer



- ⇒ Data pipeline creation
- ⇒ Data transformation
- ⇒ Ensuring data Quality
- ⇒ Automation

ETL + DW

Q: Where does now Big data fits in?



Big data engineer specializes in technology that can process and manage massive data in distributed systems

ELT + D1

Aspect	Data Engineer	Big Data Engineer
Focus	Handles structured and manageable data.	Handles massive datasets (Big Data).
Tools	SQL, Python, Airflow, ETL tools.	Hadoop, Spark, Hive, Kafka, NoSQL databases.
Scalability	Works with traditional systems.	Works with distributed systems for scalability.
Storage	Data warehouses and databases.	Data lakes and distributed storage systems.
Processing	ETL pipelines for structured data.	Parallel processing for large-scale data.
Use Case	Preparing data for a monthly sales report.	Analyzing social media trends in real-time.