# Agenda →

1. what is data Engg ?

2. In AI world → why DE is important ? ✓ !

3. what does DE actually do ?

4. Skills ?

5. SQL → ?

# 1. what is D E ?

→

## 1. Structured ?

Table →  ✓

→ Purchase ?
?

## 2. Semi structured

← PDf ,

Email , Hashtag

| id | doc loc |
|----|---------|
| 1  | . ~ ~   |

C JSON
event log)

## 3. unstructured ← Review

Image

→ img → Text
→ Videos → voice

Text

→ Reviews

→ <u>Extract</u> ✓

    ↳ <u>website</u> ✓

    ↳ <u>API</u> → Paytm
                → GPAY
                →

    ↳ <u>files</u>
             ↘

→ <u>Store</u>
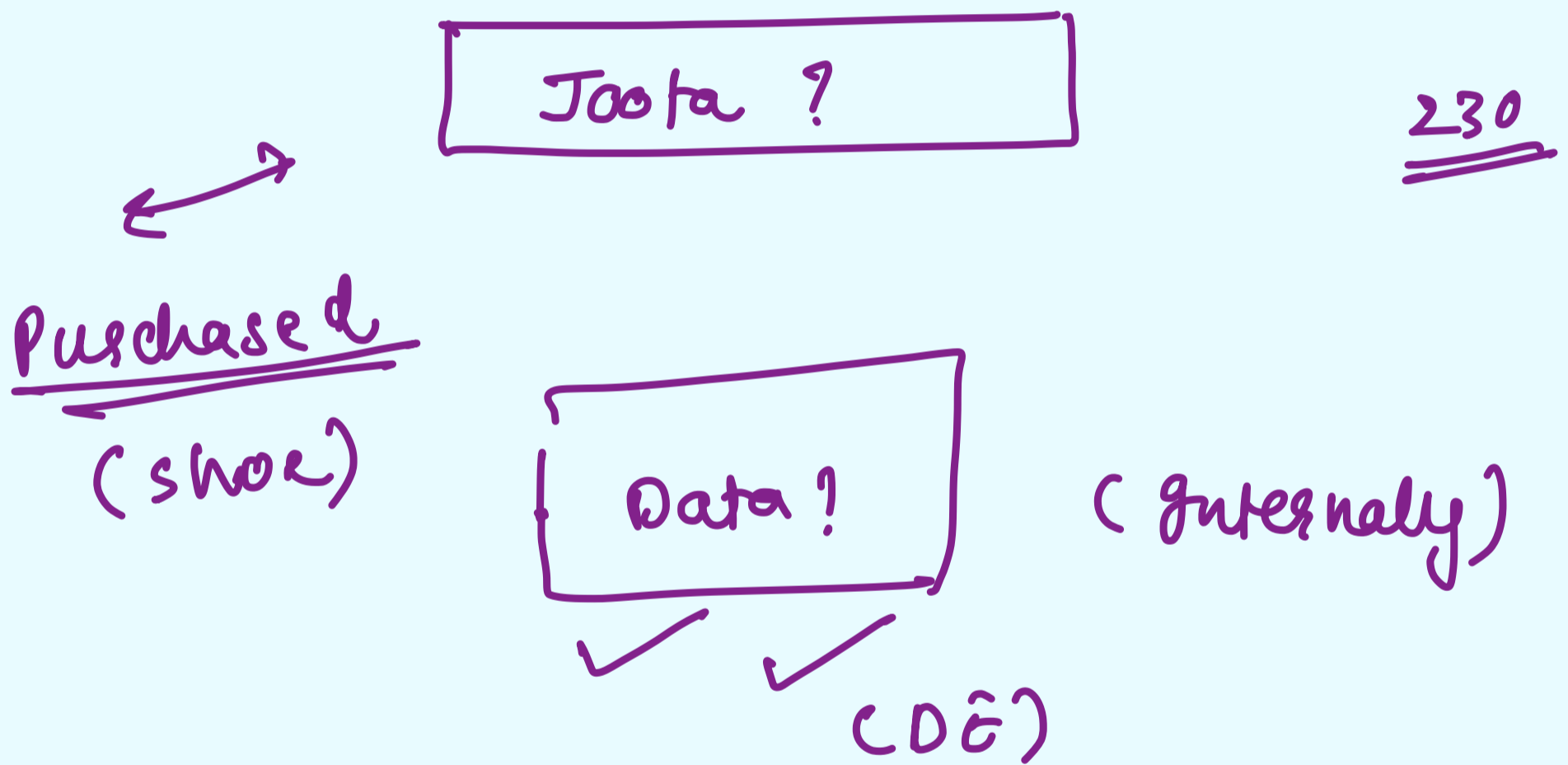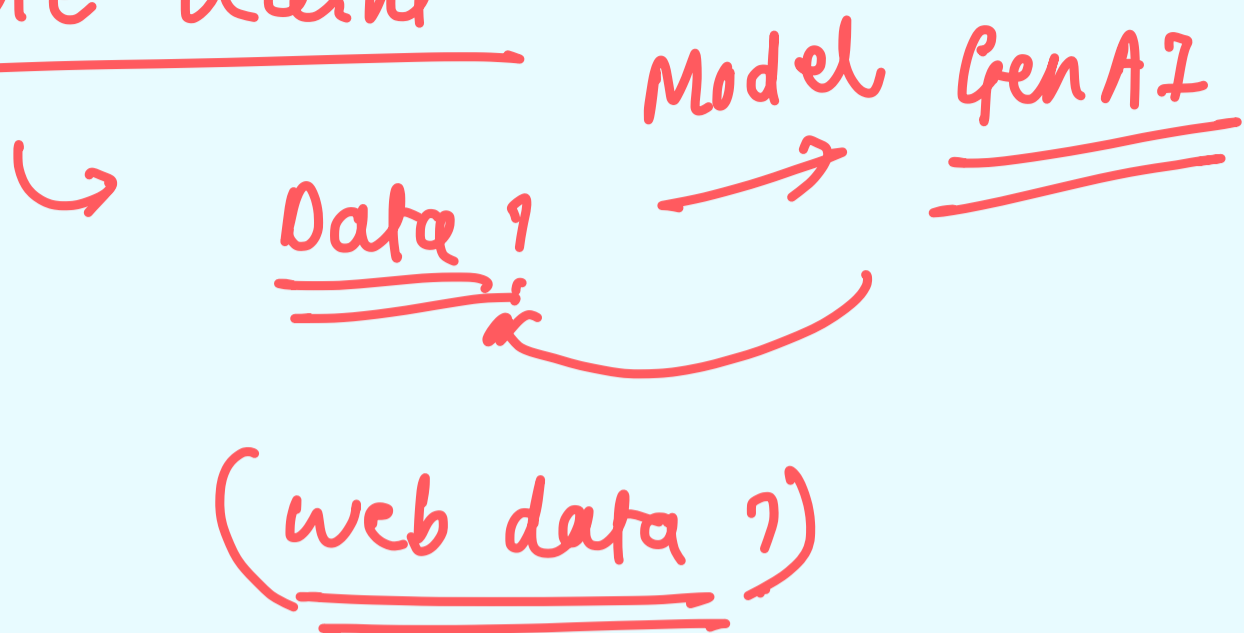
                  ?
_____

Data warehouse     Database

       ↗    Data lake ?

→ <u>Transforming</u> → clean
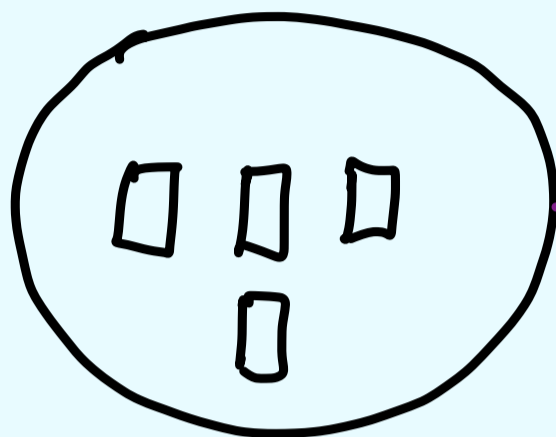
→ <u>Delivering</u> for analysis

Machine learn

Model GenAI

Data ?

(web data ?)

Joota ?

230

Purchased
(shoe)

Data ?

(Internally)

(Dᴇ)

R - II

↳ Data ?

||  ✓



Data storage

Dᾱ

↓

~~~~~~~~~~~~~

Process data

clean data

Analysis

Data → structured → Data Analyst ✓

(Structured + unstructured) → DS ✓

Icecream

$$\underline{\text{DS equation}}$$
$$+$$
$$\underline{\text{New data}} \qquad = \qquad \underline{\underline{ML}}$$

$$(\underline{\underline{ML}} + (\text{Decision} \checkmark$$
$$\qquad \underline{\underline{Algo}} \,) \qquad = \quad \underline{\underline{AI}} \checkmark$$

→ <u>Skills</u>

→ <u>Python - scripting</u>

→ SQL → query (Tables)

→ Java or Scala ( Big data)
         —x— spark

2) <u>Data tools</u>

→ ETL → <u>Apache Airflow / ADF</u>
     [Extract | Transform | Load ]
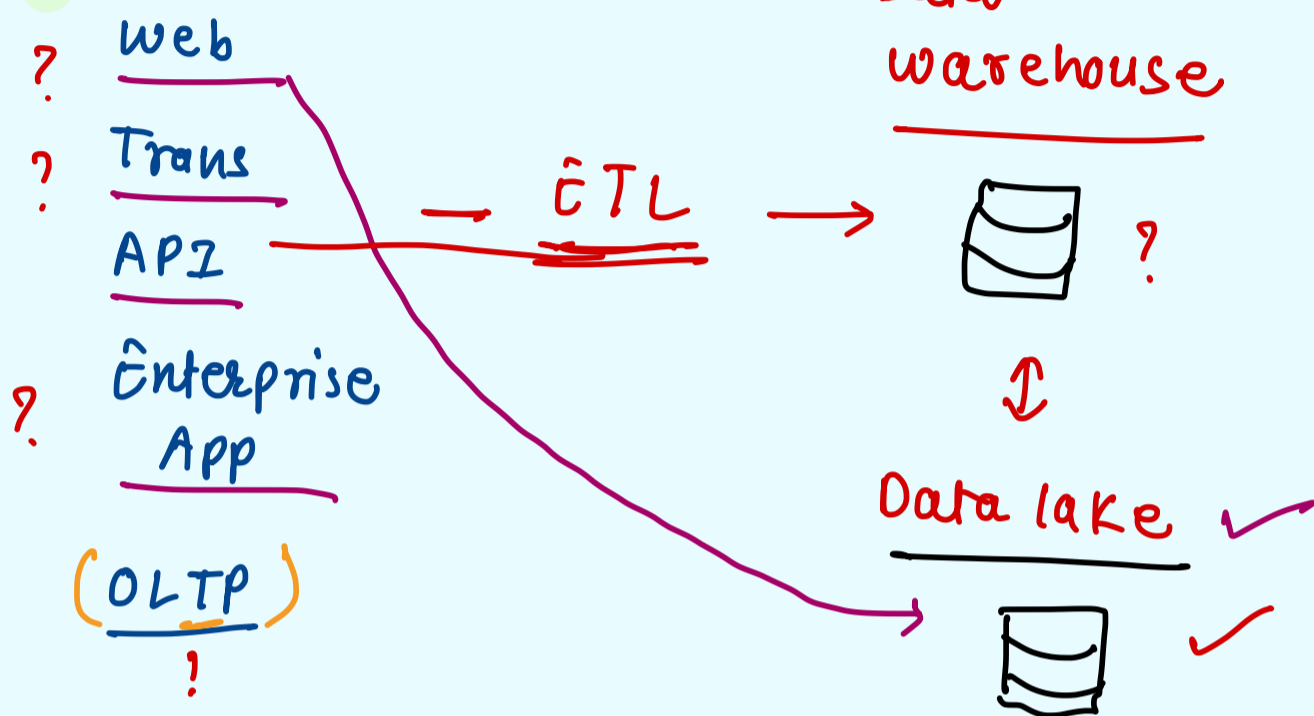
→ <u>Hadoop | spark</u>

→ AWS / GCP / Azure (Cloud)

## Database / storage

- <u>Structured</u>    —    SQL

- semi structure → <u>No SQL</u>
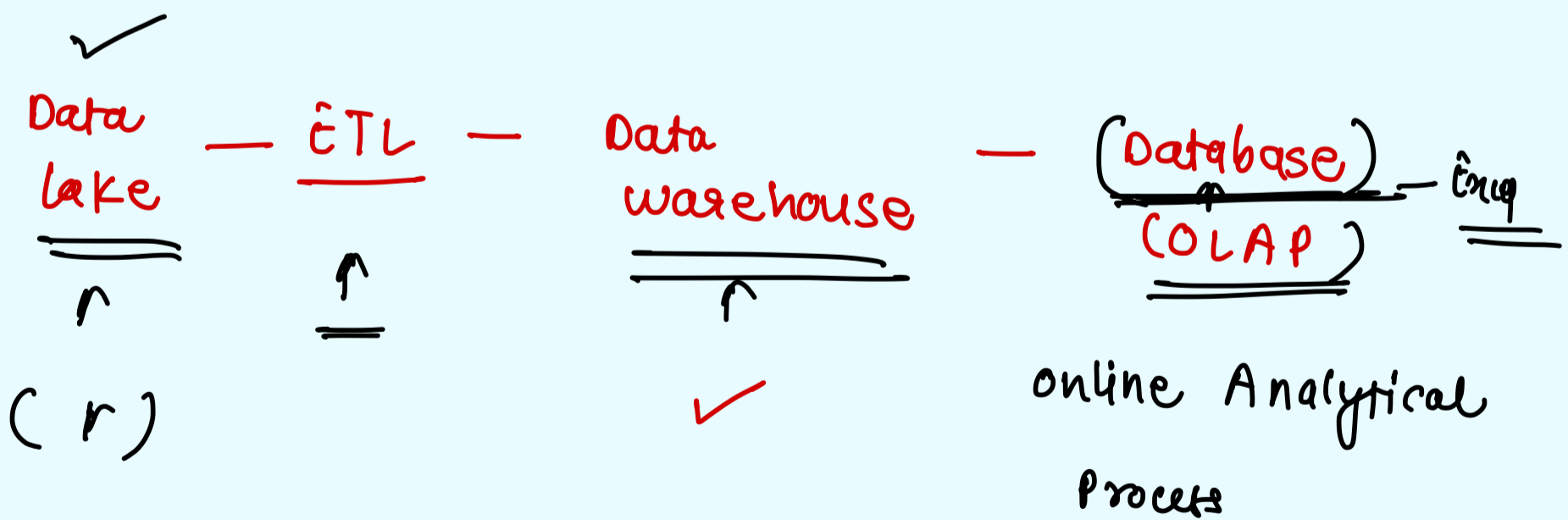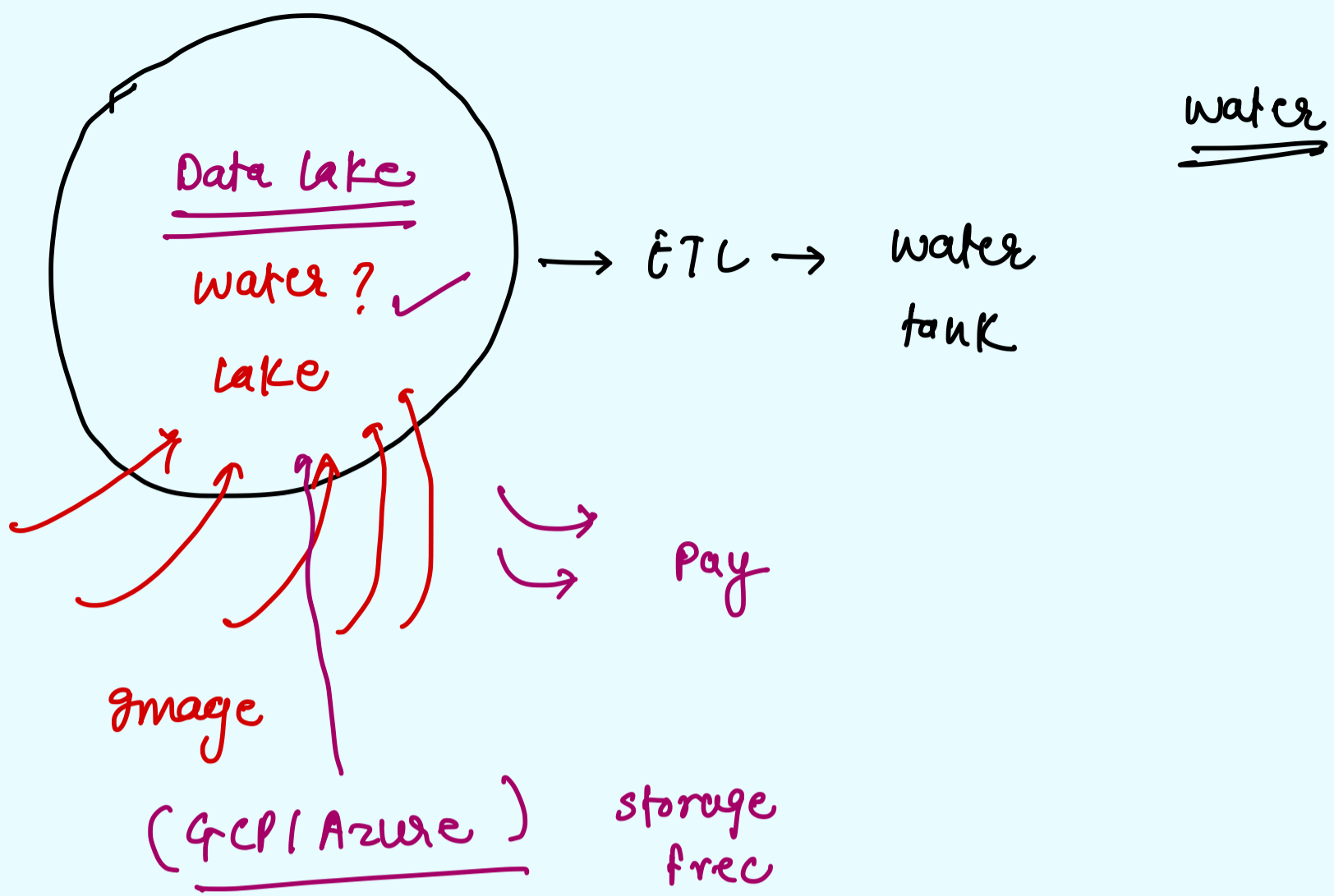
# Myntra

## Data Arch

Data Source

? web

? Trans

API

? Enterprise
App

(OLTP)
!

online Transcation
procesing ✓

Data
warehouse

→ ETL →

?

↕

Data lake ✓

✓

OLTP → Database ?

→ (OLTP)

Pos ? →

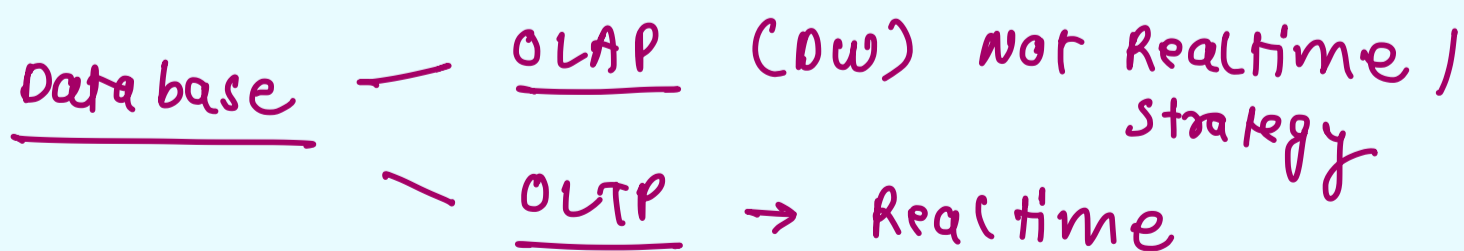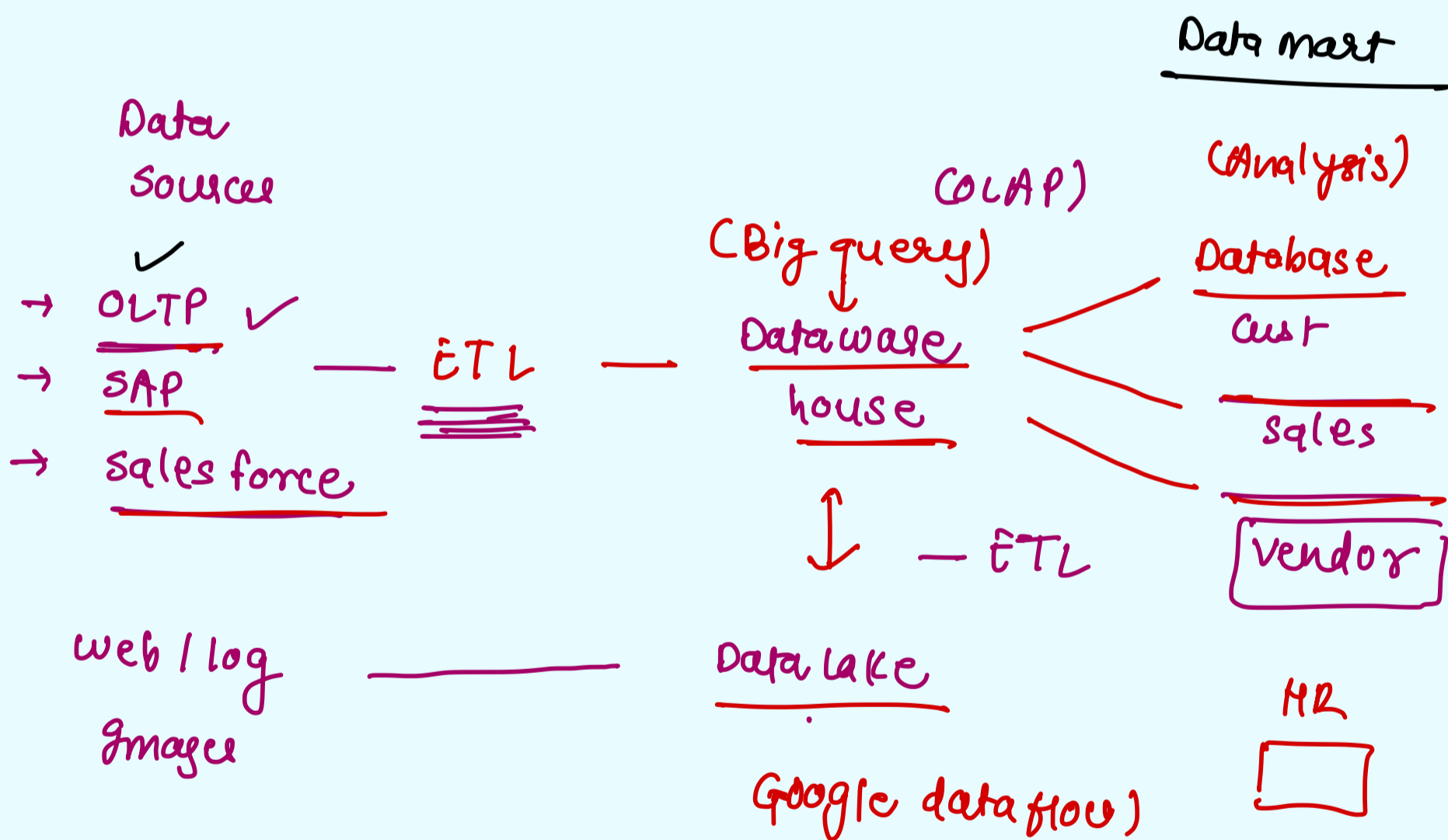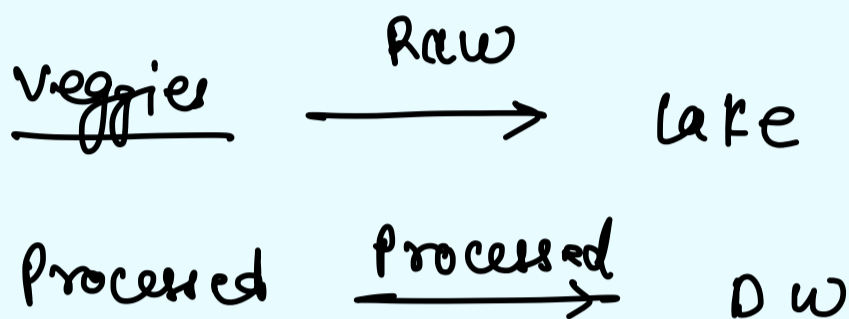| Row | cus id | order | order item | unit | Price |
|-----|--------|-------|------------|------|-------|
| 1 | C1 | 01 | ② | | 100 |

✓

Myntra → _____

| Inv | item-id | In stock | datetime |
|-----|---------|----------|----------|
| | 1 | 10 | d1 t1 |
| | 1 | 8 | dt t2 |

**Data lake**

water ? ✓
lake

→ ETL → water
tank

water

⇉ Pay

gmage

( GCP / Azure )   storage
free

---

Lake → Treatment   water   overhead    Tap
        Plant   ‾ᴡ  tank  ‾ᴡ   tank

---

✓

Data — ETL — Data   — (Database) — Enq
lake           warehouse        (COLAP)

( r )              ✓         online Analytical

                             Procts

lake
 ↑
veg

        ⤴ Data warehouse
           ↑
        (Processed)

---

(Database)

OLTP        ETL _____    [OLAP]  → Analysis
(Trans)                        ↑
                        Data warehouse

event ⟶

veggies    Raw    ⟶    lake

Processed   Processed    DW

Data mart

Data
Sources
      ✓
→ OLTP ✓                        (OLAP)        (Analysis)
                    (Big query)        Database
→ SAP    — ETL —    Data ware              Cust
                      house
→ Sales force                            Sales

                        ↕  — ETL        [vendor]

web / log    _____    Data lake        HR
Images                                  ▭
              (Google data flow)

Database  —  OLAP  (DW) Not Realtime /
                                Strategy
          ↘  OLTP  → Realtime

OLTP (DB)

Order system

POS

OM

OLAP (DW) ✓

DB     vs     D

C drive

System

D-drive

Movie

E

study

gcp

ETL —

SQL

DB ✓

→ Real time

→ detailed
  Trans ‿

→ Limited ‗

→ Change frq ‗

DW

Mom order

Trans per day

Many

DC

Raw

Orginal

DW

Processed

Agg

✓

Costlier

DW

?

Processed    id   Image   location
                        Type

E
unstructured
        ↑
Not processed → Data lake ✓

Data source

x1
x2 } — [ETL] — Data
x3                              warehouse
                                                          Data mart

x2
x5 } — no process —

Data
Lake

Extraction → Hadoop Ecosystem