# Fundamentals of Data Engineering :-

## What the fuck is Data Engineering ?
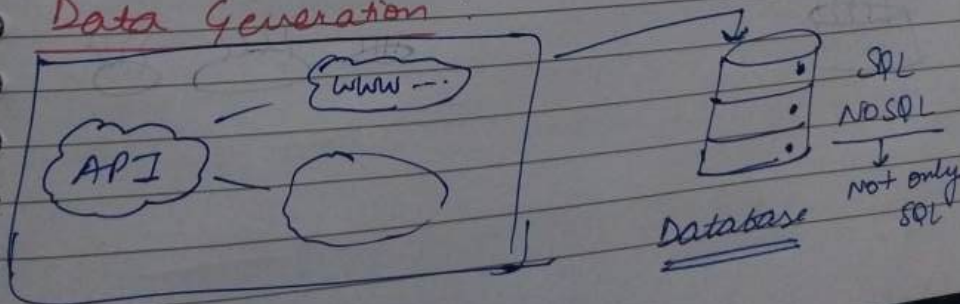
```
┌─────────┐      ┌──────────┐      ┌─────────────┐
│  SQL    │ ──→  │  Data    │ ──→  │ Stakeholders│
│         │      │ Engineer │      │             │
└─────────┘      └──────────┘      └─────────────┘
                      │ does the
                      ↓
Raw Data ──────→ Transformation ──→ Stakeholders.
```

## Data Engineering Workflow :-

Data Production / ──→ Data ──→ Data
Data Generation.     Transformation    Serving

[www. amazon - Com]

## Data Generation



```
SQL
NOSQL
  ↓
Not only
  SQL
```

Database

## Data Transformation :-



Transformations

Cleaned
Database

DBA

Raw Data ⟶ Curated Data

## Data Serving

## Upstream - Downstream :-

Software Dev / DBA

Data Engineer

Data Analyst

| 🛢 | API | □ | ⬠ | → Upstream |

↓ Data

| Data Engineer | Transformation |

↓ Transformed Data

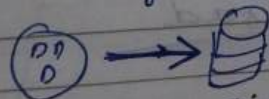| Manager. | Data Analyst |

Downstream.        Downstream.

---

## OLTP

- Online Transaction Processing

EX : MySQL, PostgresSQL

- Managed by DBAs

Transaction
- Writes and Updates efficient

## OLAP.

Online Analytical Processing.

- Also called Data warehouse

EX :

- Read Efficient.

| OLTP | OLAP | E |
|---|---|---|
| Modelling | EX :- Wheres cape Snowflake. | |

Modelling

↳ Normalization → 1NF
→ 2NF
→ 3NF

Dimensional Modelling for reading data.

OLTP
(Write and update Efficient
Difficult in read)

X Read is difficult

OLAP
(Read Efficient
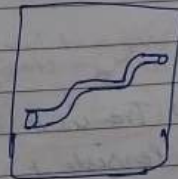
[Mmg]

Ex

[OLTP] → [ETL] → [OLAP.]

Extract        Transform        Load

## ETL :-



| Extract | Transform (Pipeline Creation) | Load |
|---------|-------------------------------|------|

Extract

↓

OLTP
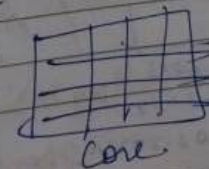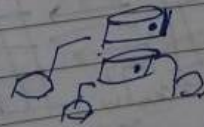
OLAP

Data Warehouse

for

## Data Warehouse and Layers :-



Core

OLTP

Drop and insert
date in staging
layer.

→ Transient

→ Persistent
↓
When we
want to
store history
of data in
staging

Staging

Core.

Incremental loading :-

API

Day 1

Day 2
Jan 02
2024 core

Day 1

Day 2
Jan 02
2024 Dimensional
Model.

Staging

Jan 2023
Jan 2024 ⟩ Day 1
Jan 02 2024 ⟩ Day 2

Day 2
Pull only
latest data

## Dimensional Modelling :-

Dimension 2 → Fact ← Dimension 1

Fact
STAR

Dimension 4 → ← Dimension 3

Stores only numeric values.

Rest of things
will go to dimension

EX: Dimension 1 :- Details of customers.
(customer Dimension)

Dimension 2 :- Details of Product
(Product Dimension)

→ Data is Denormalised / less normalised
in Fact / Dimension tables.

## Dimensional Modelling

### Star Schema

| Dimension1 |          | Dimension |

| Fact |

| Dimension |    | Dimension 4 |

### Snowflake Schema

| Dimension |    | Dimension |

| Fact |

| Dimension 7 |    | Dimension 7 |

| Dimension |

• / Heirarchy of Dimensions

### Slowly changing Dimensions :—

① Type - 0
② Type - 1
③ Type - 2
④ Type - 3

## Type-0 Slowly changing Dimensions :-

| Product_ID | Name | Prod.cat |
|---|---|---|
| 1 | Honey | Food |
| 2 | Shirt | clothing |
| 3 | Comb | clothing |

// No change in Dimensional Table.

## Type-1 Slowly changing Dimensions :-

UPSERT (update + insert)

| Product_ID | Name | Prod_cat | |
|---|---|---|---|
| 1 | Honey | Food | Update |
| 2 | Shirt | Clothing | Update |
| 3 | Comb | Clothing | Update |
| 4 | Shoes | Footware | Insert |

| Product_ID | Name | Product_cat |
|---|---|---|
| → 1 | Honey | Food |
| → 2 | Shirt | clothing |
| → 3 | comb | Hair |
| → 4 | Shoes | Footware |

Before                      After

CHANGE HISTORY IS NOT KEPT IN THIS
DIMENSION

Preserving history.

# Type-2 Slowly Changing Dimension :-

**Before**

| Product_ID | Name | Prod_Cat | Eff StartDate | Eff EndDate | Inuse |
|---|---|---|---|---|---|
| 1 | Honey | Food | 1/1/2024 | 1/1/3000 | Yes |
| 2 | Shirt | Clothing | 1/1/2024 | 1/1/3000 | Yes |
| 3 | Comb. | Clothing | 1/1/2024 | 1/1/3000 | Yes |

Add 3 more columns to handle change history.

① Effective start date
② Effective End date
③ Inuse.

Comb Prod_Cat is changed to Hair.

After.

Product_ID

1
2
3
4

Type-

Produc
1
2
3

Befo

After

## After

| Product_ID | Name | Prod_cat | Eff start date | Eff end date | Issue |
|---|---|---|---|---|---|
| 1 | Honey | Food | 1/1/2024 | 1/1/3000 | Yes |
| 2 | Shirt | Clothing | 1/1/2024 | 1/1/3000 | Yes |
| 3 | Comb | Clothing | 1/1/2024 | 1/2/2024 | No |
| 4 | comb | Hair | 1/2/2024 | 1/1/3000 | Yes |

## Type-3 Slowly changing Dimensions :-
### (Previous values)

| Product_ID | Name | Prod_cat | Prev Prod_cat |
|---|---|---|---|
| 1 | Honey | food | food |
| 2 | shirt | clothing | clothing |
| 3 | comb | clothing | clothing |

(Before)

(After)

| Product_ID | Name | Prod_cat | Prev Prod_cat |
|---|---|---|---|
| 1 | Honey | Food | Food |
| 2 | Shirt | Clothing | Clothing |
| 3 | comb | Hair | Clothing |

Issue
date

Yes
Yes
Yes

# Data Lake



- Datalake can store semi-structured data formats (CSV, json) etc.

- We can deal with various file formats.

## Difference between Datalake vs Data warehouse

| Area | Datalake | Data warehouse |
|---|---|---|
| Data store. | It can capture and retain unstructured, semi-structured data in its raw | It can capture and retain only structured data. A data warehouse stores data in quantitative metrics |

| Area | Datalake | Data warehouse |
|---|---|---|
| | format A datalake can store all types of data irrespective of the sources and structure. | with their attributes Data is transformed and cleansed. |
| Schema Definition. | Typically the schema is defined after data is stored. This often offers high agility and data capture quite easily, but it requires lot of work at the end of the process (schema-on-read) | Typically, a schema is defined prior to when data is stored. It requires work at the start of the process, but it offers performance, security and integration (schema-on-write) |

| Area | Datalake | Data warehouse |
|---|---|---|
| Price and Performance | The storage cost is relatively low, compared to data warehouse and querying result is better. | The storage cost is high, and querying result is time consuming. |

### Lakehouse

```
                Lakehouse
            _____|_____
           |                 |
           ↓                 ↓
       Datalake          Warehouse
   Best for storage cost   Best of reporting
```

┌─────────────────────────────────┐
│  BI, Report, Data , ML           │
│            science               │
└─────────────────────────────────┘
                  ↑

┌──────────────────────────┐    ┌ Converts
│ Metadata and Governance layer │ │ data into
└──────────────────────────┘    │ Structured
                  ↑              └ schema
          
┌────────────────────────────┐
│ Structured, Semi-Structured, │
│ Unstructured Data.           │
└────────────────────────────┘

File-
(1) Row-
(2) Col-

Row-Ba

CSV, AVR
   ↓
used fo
faster
[ write

Row Based

Column Base

Select
Column

Row - ba

## File - Formats :

① Row - Based
② Column - Based.

**Row - Based**

CSV, AVRO
↓
used for
faster writes
[ write Efficient ]

| Product ID | Name | Product - Cat |
|---|---|---|
| 1 | Honey | Food |
| 2 | Shirt | Clothing |
| 3 | Comb | Clothing |

**Column - Based**

Parquet, ORC.

Used for faster reads.
[ Read Efficient ]

## Disk

**Row Based**  1 Honey Food 2 Shirt Clothing 3 Comb Clothing

**Column Based**  1 2 3 Honey Shirt Comb Food Clothing Clothing

**Select Product - ID from Dim Product**

**Column - Based :-** 1,2,3 Boom!!! we have all the data
(Read Efficient).

**Row - Based** 1 Honey Food 2 Shirt Clothing 3 comb Clothing,
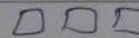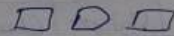↓ iterate over all the disk space.

## Delta Format                from files

- Open Table Format
- Built on top of Parquet.
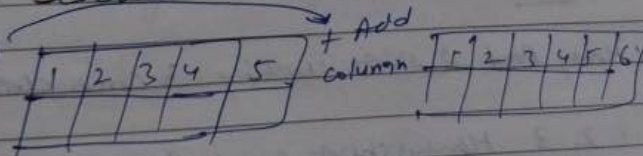
Parquet file  ⟶  Delta file

▢ ▢ ▢                 ▢ ▢ ▢        + Transaction
                                       log
                              ▢ ▢ ▢

⟶

① Helps in Versioning

② Schema Evolution

⑥ [1|2|3|4|5]  ⟶ + Add
                  column   [1|2|3|4|5|6]

⑤ Allows the ACID Capability.

Big Da

① Apac
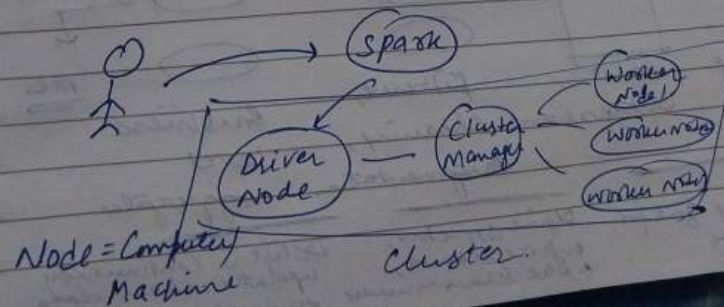
② Apa

③ Apa

④ Da

Distri

Ask
we d

Node = 
    M

## Big Data

Big Data frameworks :-

1. Apache kafka → Used for streaming Data, Real Time Data

2. Apache Airflow → Used for Data orchestration

3. Apache spark → "Discuss in future"

4. Databricks → Processing by using cluster

### Distributed computing with spark :-
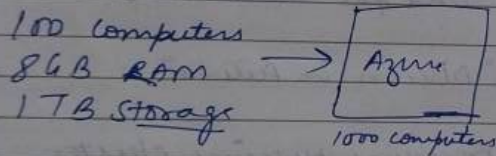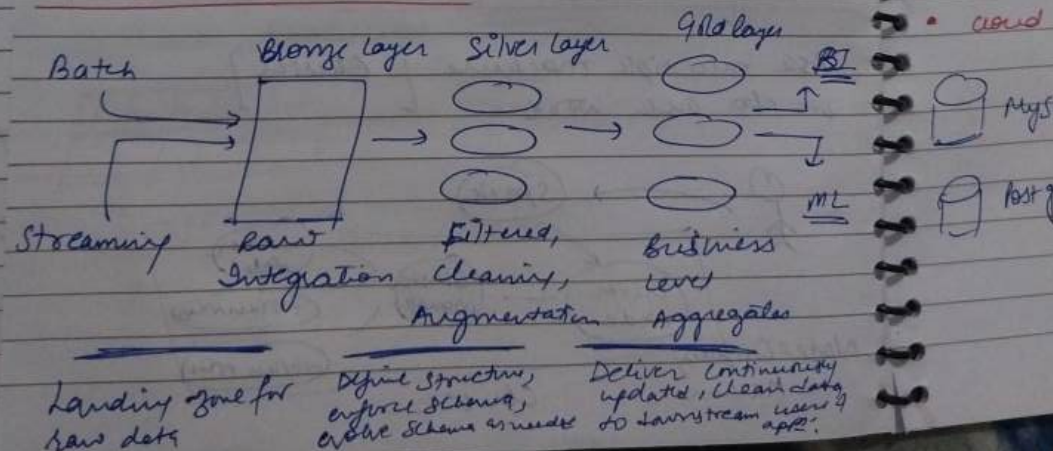
Ask multiple machine [ Cluster ]
we do our work



Node = Computer Machine        Cluster

# Cloud Data Engineering

| Azure | AWS | GCP |
|-------|-----|-----|
| ↓ | ↓ | ↓ |
| Microsoft | Amazon | Google cloud Platform |

```
100 computers
8GB RAM      →   [ Azure ]
1TB Storage
                1000 computers
```

## Medallion Architecture :-

```
Batch                Bronze Layer    Silver Layer      Gold Layer        BI
                     ┌──────┐          ◯                ◯          →
  →                  │      │          ◯          →     ◯          →
  →                  │      │          ◯                ◯
Streaming            Raw            Filtered,        Business        ML
                     Integration    Cleaning,        Level
                                    Augmentation     Aggregates
```

Landing zone for    Define structure,    Deliver continuously
raw data            enforce schema,      updated, clean data
                    evolve Schema as needed   to downstream users &
                                              apps.

# Cloud Data Engineering with Azure :-

- Data streaming source



Azure Event Hub

- Cloud sql database

SQL Statements
- Select
- Update
- Delete
- DDL
- TCL
- DCL

MySQL

Azure SQL DB

PostgresSQL

$\longrightarrow$ Storage Account $\longrightarrow$ 1. Datalake
2. Tables
3. File shares
4. ____

Datalake :-

Azure Datalake Storage Gen 2
(ADLS Gen 2)

Heirarchical Namespace
$\downarrow$

Containers (folder)
$\hookrightarrow$ Container (folders)
$\hookrightarrow$ ____

Cloud ETL Tool :-

• Azure Data Factory.

Big Data Transformation

Raw $\longrightarrow$ Databricks $\longrightarrow$ Silver.
$\downarrow$
Spark clusters

## Cloud Data Warehouse :-

- Azure Synapse Analytics
  Ex : Snowflake, Redshift.

## Data reporting

- Power BI

## Imp Terms

- Azure Purview :- Data Governance Tool.
- Azure devops :- CI/CD.
- Azure key vault :- Store confidential information
- Microsoft Entra ID :- Records all the users.
- Azure Monitor :- Monitor Production Runs.
- Cost Management :-