# Agenda

1. Curriculum To be Covered
2. Introduction to Data Engineering
3. What do Data Engineers do, and Where does all the data come from
4. How do we manage different sources
5. Big Data Examples
6. What are the main challenges involved in handling Big Data
7. Various methods of storing Data, based on use cases
8. Data Platform Architecture
9. ETL pipeline
10. Distributed Systems
11. Leveraging the Big Data technology into building our own platform.

## Rules:-

1. class will start at 9:02 PM, no waiting for others.

2. Class will of 2 hrs ( 9 to 11 pm) and 30 minutes (11 to 11:30 PM)

3. Put your current topic doubts is chat window, also any other [off topic on topic or anything) put that is Questions tab (11 to 11:30) PM

4. Break at sharp 10:00 PM and that will of 5 minutes.

5. Most of the Concepts will be revised by default 2 times. and for more complex topics it will 3 times.

6. Feedback | Assignments

Surprise Gift

$\overline{\phantom{ab}}$ $\overline{\phantom{ab}}$

( Data ) Analyst ✓
   ↳ Scientists ✓
   ↳ ( Engineers )

| Analyst | Scientist | Engineers |
|---|---|---|
| → Python / Java<br>+<br>SQL<br>+<br>Any Visualization tool (tab,<br>Power BI, looker<br>quick Sight etc) | → Python<br>+ SQL<br>+ ML Algo<br>+ DL Algo<br>+ MLOPS<br>+ Gen AI | → Python / Java / Scala<br>+<br>SQL<br>+<br>cloud ( Aws /<br>GCP /<br>Azure )<br>+<br>( Data Pipelining Skills ) |
| Roles<br>→ who is responsible for deriving the insights from raw data. | Roles<br>→ who is responsible for cleaning, manipulation, Statistical analysis, predictive, Prescriptive, building models and delivering results that have a impact on business | Roles<br>→ who is responsibly for ETL on any data ( Batch / stream ), to ( Create Data Pipeline. ) |

··at DO DE do? + Curriculum

→ wha↴

↳ Skill Set

↳① Infrastructure Components
↳ VM
↳ N/W
↳ load balancing
↳ Application performance
Monitoring

② cloud - based Services
↳ AWS | GCP | Azure / IBM |
Oracle . . . .

③ Databases | Data warehouses
↳ RDBMS = Mysql / oracle / Postges
MS SQL ) . . . .

↳ NoSQL = Redis / MongoDB
DB               c* / Neptune . . .

↳ DwH = oracle exadata,
GCP BQ,
AWS Redshift...

④ Proficiency working with date pipelines :-
↳ Apache Beam | Airflow |
GCP Data flow

- - - - tools                    l - d  tre ]

✓ ⑤ ETL tools
                    ↳ Aws Glue / Informatica /
                       IBM infosphere

✗ ⑥ Languages :- → ⬭SQL⬭
                    → PL = Python / Java
                    → shell = Linux shell

⑦ Big Data processing tools
                    ↳ Hadoop
                    ↳ Spark
                    ↳ MR
                    ↳ Kafka

Curriculum :-
      ↳ SQL → 7 classes
   ③ ↳ Data modelling    | DWH → BO / Redshift
   ④ ↳ Hadoop Components → Hive / MR / HDFS / Yarn
      Spark → RDD / DF / Stream
   ⑥ ↳
   ① ↳ Airflow
   ① ↳ Kafka
   ① ↳ No SQL D B → ⬭AWS⬭
                       AWS Glue / S3

① ↳ ( Data lane ) →

---

# Big Data ?

↳ Volume
↳ Velocity      ⟩ (3V)
↳ Variety
↳ Veracity
↳ Value
↳ Variability

**6 V's**

① Volume :- Amount of Data that is generated
| Stored per Day.
✓

$$mB \rightarrow GB \rightarrow TB \rightarrow PB \rightarrow EB \rightarrow ZB$$
....

$$1 PB = 1024 TB \xleftarrow{} \boxed{2 PB}$$

② Variety :-   ⓐ Structured
✓           ⓑ Semi Structured
           ⓒ Unstructured

| Structured | Semi Structured | Unstructured |
|---|---|---|
| ↳ Data which follows a [rigid format] that can be organized into rows and columns very neatly. | ↳ is a mix of data that has consistant characteristics but data that does not conform to a [rigid structure]. | ↳ Data that does not have an easily identifiable structure, and cannot be organized into any format. |
| ↳ RDBMS (SQL)  ↳ schemas (structure) | ↳ XML  [Parquet]  JSON  CSV  Avro  TSV  Node  → logs | → Pdf → Videos  → text → Images  → Audio |

Schema → Col name
       → Col Datatypes

$\{a : b\}$   $\{ a>$
$\{b : c\}$   $<b?$
              $c|>>$
              $c|a>$

---

③ Velocity :- the Speed at which data is getting generated. Here we apply

2 types of Processing.

| Batch | Stream |
|---|---|
| ↳ Predefined Data | ↳ Continous Data |
| 25h  Mgr → [CSV] → HR → them  ↓  delay ← Bank  [2000]  Size. | ↳ Size will be less but overall it will huge. (24 hrs)  [1 mik5] × 60 × 60 × 24 |

⤷ Data is of big ~~...~~ | $1 s = \overparen{(1\ldots)}$
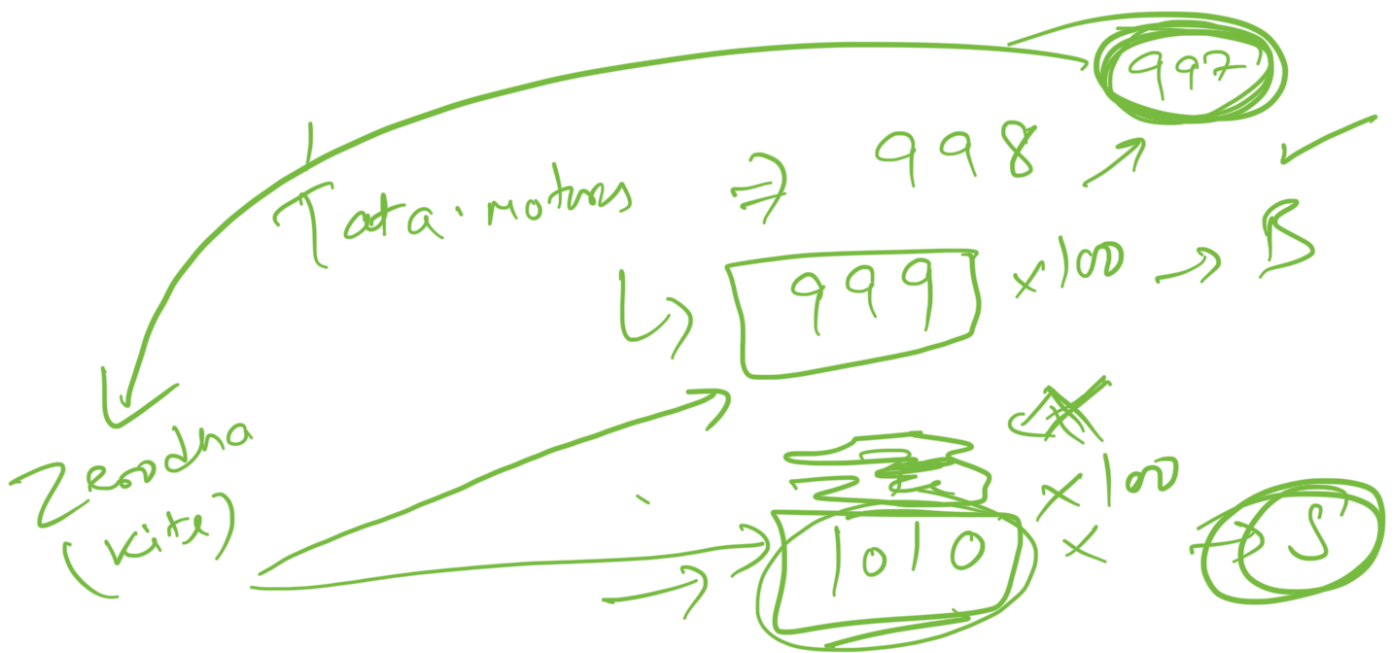
④ Veracity:- quality, accuracy and trustworthiness of Data.

✓
C1 → Internal Sales DB with regular audits
✓
C2 → Social media posts with slang, Sarcasm and fake accounts.

* Being able to identify the relevance, correctness or accuracy of data and apply it to appropriate purposes

Tata motors ⇒ 998 ↗ 997

⤷ 999 ×100 ⇒ B

Zerodha (kite)

1010 ×100 × S

⑤ Value :- Usefulness or benefit derived

for Data.

→ Analyze Customer purchase history
  to recommend few products

✗ → outdated Customer Surveys. (1mm)

⑥ Variability :— Degree to which data is
                 subject to change and
                 inconsistent

→ Stock market data fluctuating throughout
  the day

→ historical population Census data, data54
  over long periods

Methods of Store Data

↳ OLTP   ⇒ Database

↳ OLAP   ⇒ Datawarehouse



reporting
&
analysis.

5 7 10 15...

DWH

day h day

OLTP = online txn Processing

OLAP = online Analytical Processing