

## Agenda

- a. Introduction to Data Engineering ✓
- b. What do Data Engineers do? ✓
- c. How do we manage different sources? ✓
- d. Big Data Examples ✓
- e. What are the main challenges involved in handling Big Data ✓
- f. Various methods of storing Data, based on use cases
- g. What is an ETL pipeline?
- h. Why Can't Standalone Systems Suffice?
- i. Curriculum To be Covered



## Rules :-

- ① class will start at 9:02 PM ↴
- ② class will be of 2 hrs + 30 min DCS
- ③ put your Doubts in chat box,  
Question tab
- ④ Most of the Concepts will be revised  
2 times default and upto 3 times, in  
case of a complex topic
- ⑤ Break of 5 mins around 10:00 PM
- ⑥ Feedback + Assignments
- ⑦ We will a Project / Case study at the  
end of module.

Data Engineering :- It is a part of Engineering  
where a user build the  
pipelines that collect and  
deliver data for Data Scientists.

DS



- ① are the people who analyze data, create algorithms and make prediction.
- ② DJ develop models that is used for Prediction & prescriptive Modeling
- ① build the pipelines to collect and deliver data for DS
- ② Extract, organize & integrate raw data from different sources & then extract value from flat Dbs.

### Skill Set

- ↳ Infrastructure Components
  - ↳ Virtual M/C
  - ↳ Networking
  - ↳ Load balancing
  - ↳ Monitor a application
- ↳ cloud based services
  - ↳ AWS / GCP / Azure / IBM
- ↳ Database & Datawarehouse
- ↳ Data Pipelines
- ↳ ETL tools → AWS Glue / Datacatalog
- ↳ Language      QnQ

↳ ↳ ↳ ↳ ↳ ↳ ↳ ↳

- ↳ Python | Terra
- ↳ Python
- ↳ Linux

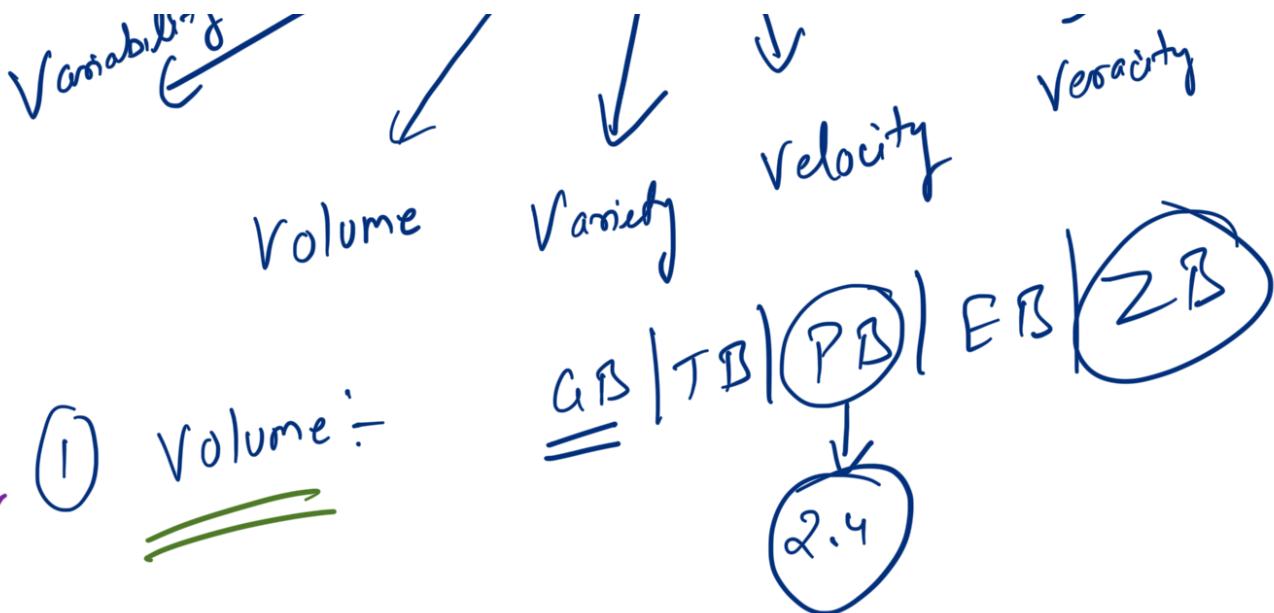
↳ Big Data Processing tools : Hadoop,

↳ Orchestration : Airflow ~~Spark, Kafka~~

### Curriculum

- 7-8 days
- ① SQL Database (MySQL)
  - ② Data Modelling & Data Warehousing  
AWS, Hadoop
  - ③ Batch Pipeline
  - ④ Real time Pipeline
  - ⑤ Orchestration with Airflow
  - ⑥ Git and GitHub
  - ⑦ CI | CD  $\Rightarrow$  Docker





✓ ① Volume :

✓ ② Variety :

Structure	Semi Structure	Unstructured
<p>① Data which is present in rows &amp; columns   ACID Properties</p> <p>→ RDBMS mysql   oracle   Postgres SQLite   DB2 ...</p>	<p>① which can be stored in row   column   key value pairs or any specific file format</p> <p>→ JSON, CSV → XLSX, Parquet → Avro, XML ....</p>	<p>① Don't have any schema</p> <p>→ Video → Audio → PDF → Text → Images → Streams ...</p>
<p>② Schema = Column name + Column datatype</p>	<p>② Schema → Column name Column datatype</p>	

③ Velocity = amount of data generated by sec

④ Veracity = the degree to which big data can be trusted.

⑤ Value = the business value of the data collected.

⑥ Variability = the ways in which the big data can be used and formatted.

### Method of storing Data

↳ Structure  $\Rightarrow$

OLTP

↳ online transactional Processing

$\Rightarrow$  CRUD



RDBMS / OLTP

= 100GB | 400GB | 10TB / 100TB

200TB | 150TB

↳ 1 hr (processing)

