# BigData — Intro

* Python
* SQL

Career

→ role of a D.E ?

→ where BigData fits in the data engineering ecosystem

```
                    ┌─────────┐
                    │ Big Data │
                    └─────────┘
     ┌───────────┬───────┴────┬──────────────┬──────────────┐
┌──────────┐ ┌───────────┐ ┌────────┐ ┌──────────────┐ ┌────────────┐
│Core Concept│ │Infrastructure│ │Processes│ │Storage Solutions│ │Career Content│
└──────────┘ └───────────┘ └────────┘ └──────────────┘ └────────────┘
  ├ Big Data?   ├ on-premise   ├ ETL      ├ datalake        ├ Data Engineer role
  ├ five v's    └ cloud        └ ELT      ├ data warehouse  └ Big Data Applications
                                          └ Data lake
```

## What is Big Data ?

Bigdata refers to datasets that are too large, fast, or complex for traditional databases to process efficiently

A problem —— handling massive data

A solution —— technologies developed to address these challenges

## Traditional RDMS limitations & (SQL)

Case 1    small online store : customers order
                              inventory
                              basic analytics

$\left.\begin{array}{l} \\ \\ \end{array}\right\}$ RDMS

1, well struct data {
   limited volume
2, queries efficiently retrieve
   data → order history
          product details ...

## Case 2  Company like walmart

↳ Billion of transactions across thousands of stores

↳ need to analyze customer behaviour, inventory trends & product performance

✱ hadoop / BigQuery instead of RDMS to process

# Why RDBMS fails for BigData !  (RTD) ✓ 100GB

64
QTD     (100GB) → 2 RTB

## 1, Scaling Limitation :

→ Vertical scale has physical/practical limits

→ Cost is expensive if we infinitely expand storage

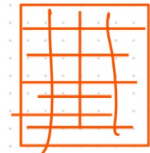## 2, Performance Issues :

* Large datasets require complex partitioning

* joins ——→ time-consuming

→ prone to error

* operations becomes extremely slow with growth of data

[Structured]

## 3, Media Storage Challenges

poor to handle/process → images, videos &c-y (unstructured)

## 4, Cost Inefficiency

hardware costs increases for enterprise-grade solutions

# Why Big Data Emerged :

Internet Growth → Social Media ———→ IOT Devices ——→ Big Data Need
(Fb, wh, Instagram, Youtube, LinkedIn)

Data Penetration :
                        Storage

Facebook : 4+ petabytes of data generated daily

Twitter : 500+ million tweets daily

Youtube : 500+ hours of video uploaded every minute

requires specialized
Storage & processing
solutions

IOT Devices : estimated to reach 75 billion by 2025

# Processing

**Netflix:** 450 + billion events daily for content recommendation

**Uber:** handles 100+ petabytes of data for real-time data

**Weather Forecasting:** Analyzes 15+ terabytes of data daily from satellites & sensors

**Genomics:** human genome → generates 200+ gigabytes per person

**Big Data Scale:** → nearly half of all the world's population in just last 3-5 years

→ grow to 175+ zettabytes by 2025

## Five V's of Bigdata

* framework to understand the challenges & characteristics of BigData
* helps us to determine whether a problem requires Big Data technologies

f v's are

Volume
Velocity
Variety
Veracity
Value

a, **Volume:** refers to the size or amount of data being generated

* Volume determines if BigData technologies are required
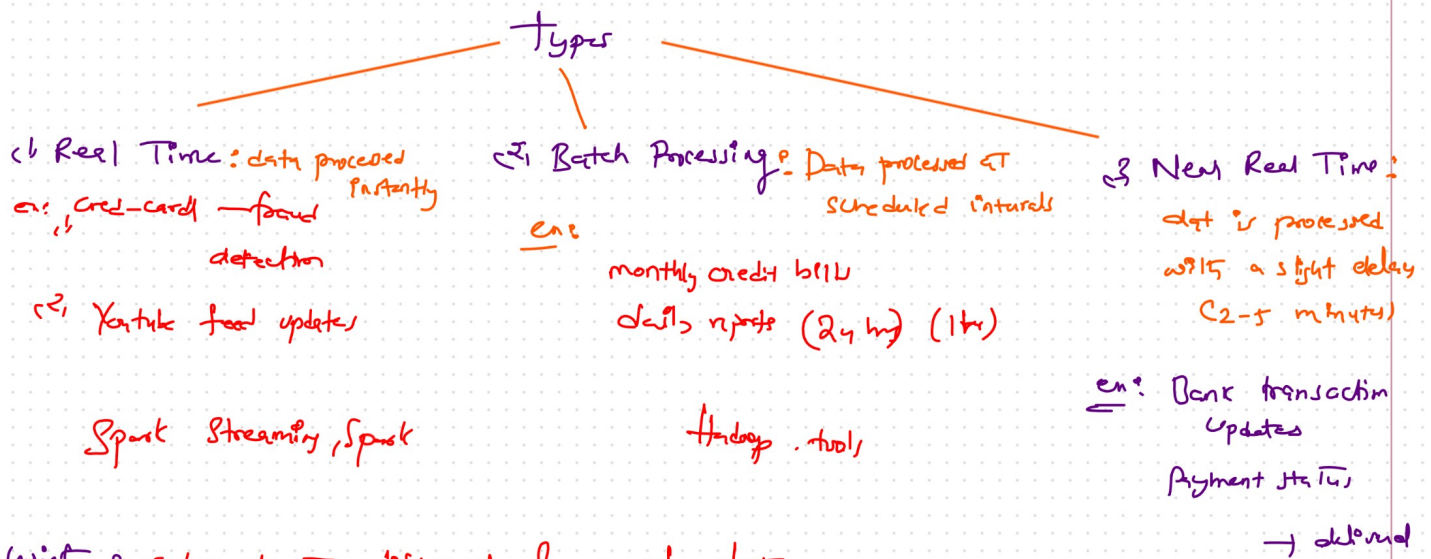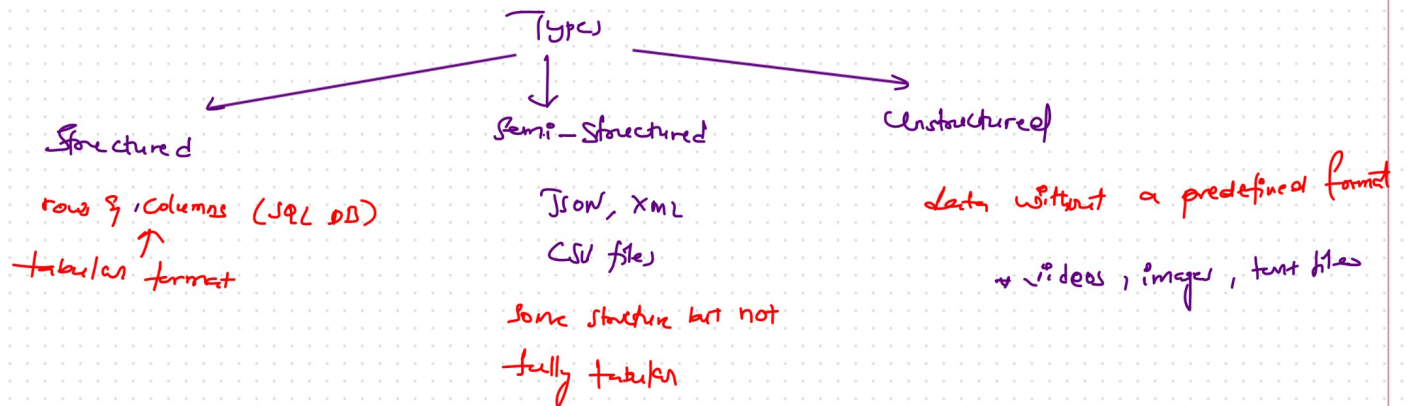  (Hadoop) → Store & process volume of data

example

Grocey store ──→ RDMS          Company → 30GB (Startup)

Amazon ──→ Big Data

* volume alone does not define BigData problem

(2) **Velocity:** The speed at which data is generated & processed

**Types**

**c) Real Time:** data processed instantly
ex: Cred-card fraud detection
c²) Youtube feed updates

Spark Streaming, Spark

**c²) Batch Processing:** Data processed at scheduled intervals
ex:
monthly credit bill
daily reports (24 hr) (1hr)

Hadoop tools,

**c³) Near Real Time:**
data is processed with a slight delay
(2-5 minutes)

ex: Bank transaction Updates
Payment status
→ delivered

3. **Variety:** refers to the different forms of data

**Types**

**Structured**
rows & columns (SQL DB)
↑
tabular format

**Semi-Structured**
Json, XML
CSV files

Some structure but not fully tabular

**Unstructured**
data without a predefined format

* videos, images, text files

(4) **Veracity:** trustworthiness or quality of data

* data can have issues like inaccuracies / inconsistencies

eg: -ve age value

* poor data quality leads to bad decisions

**Solution:** Data Cleaning & validation process

ex: Student / employee records

(5) **Value:** The insightful or business value derived from data.
✓

* It should provide insights & help in decision-making

ex: Sales data → analysis to identify trends [ iphones ____

* Big Data technologies turn raw data into meaningful insights
* Data without value is meaningless, even if it has high volume / velocity