

Before we start

→ Phone on **silent** mode

→ Feel free to ask questions

→ All doubts will be answered, so don't worry.

→ We will have 5 min break

→ Please complete all assignments

→ Join whatsapp group

Intro to Data Engineering

Agenda

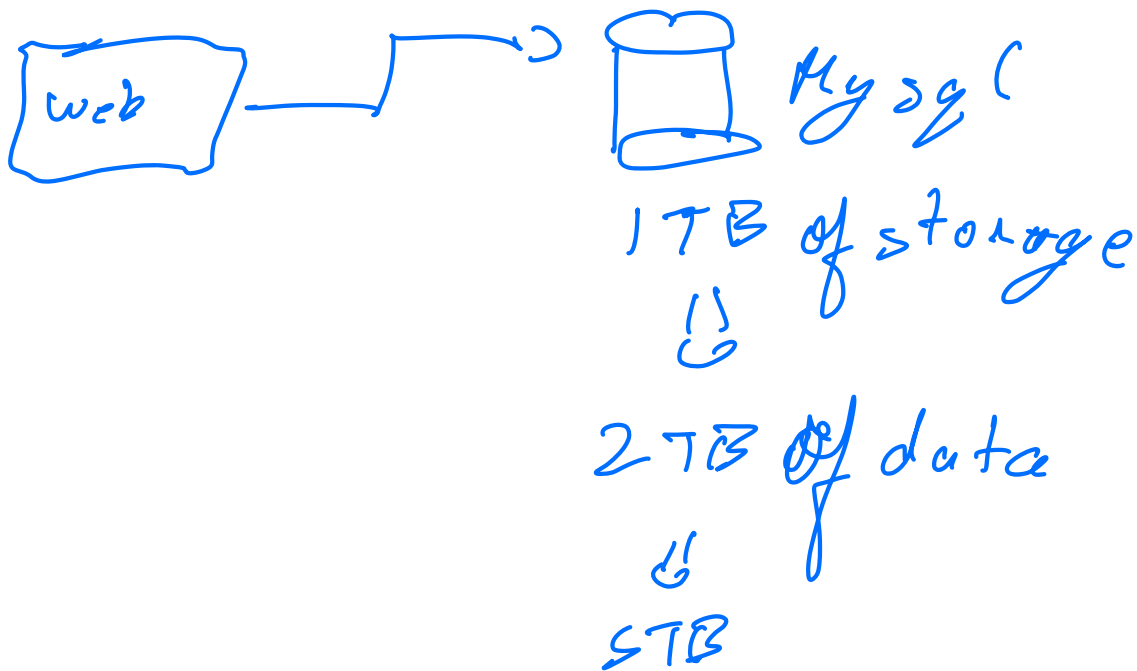
- Introduction to Course
- Big Data Story
- What defines Big Data
- Data Engineer's Role

Intro to DE

- Big Data Architecture
- HDFS
- Data Modelling
- HiveQL
- Hive

- Map Reduce
- Cloud offering
- Spark
- Flink, Ray, Pinot
- Kafka
- Airflow

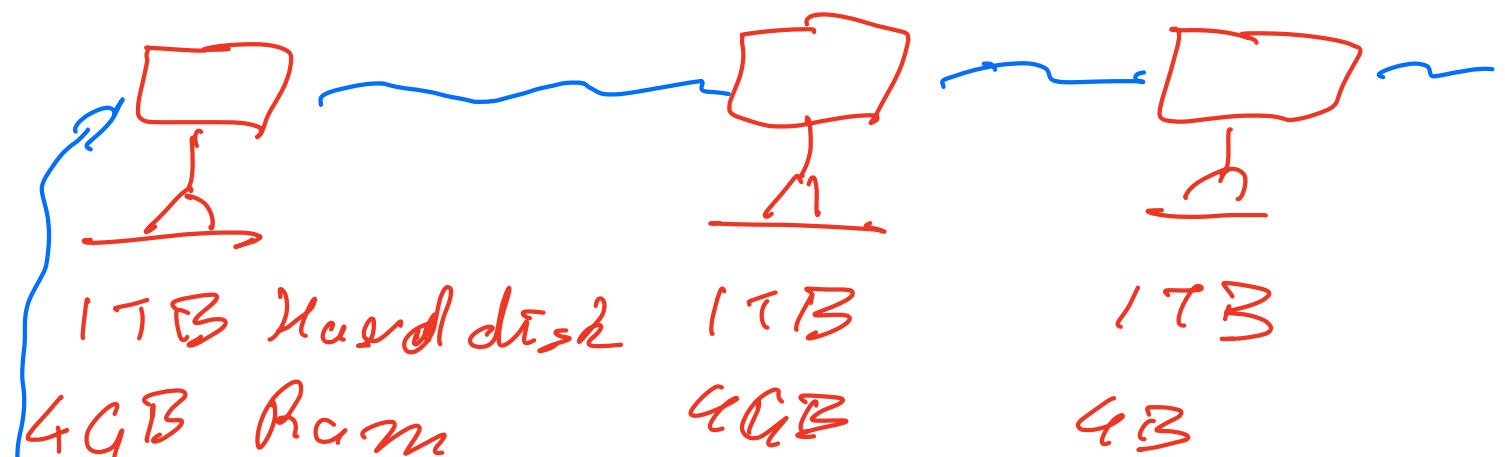
Flipkart - 2008



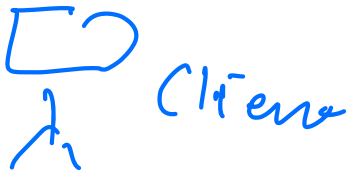
Vertical Scaling

Increase configuration of current system

Horizontal Scaling



100's of devices like this

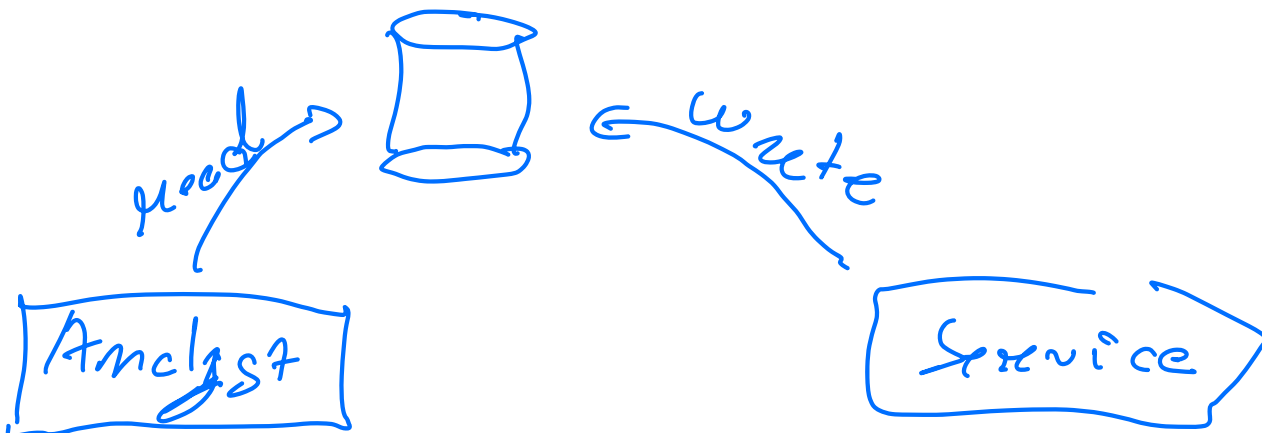


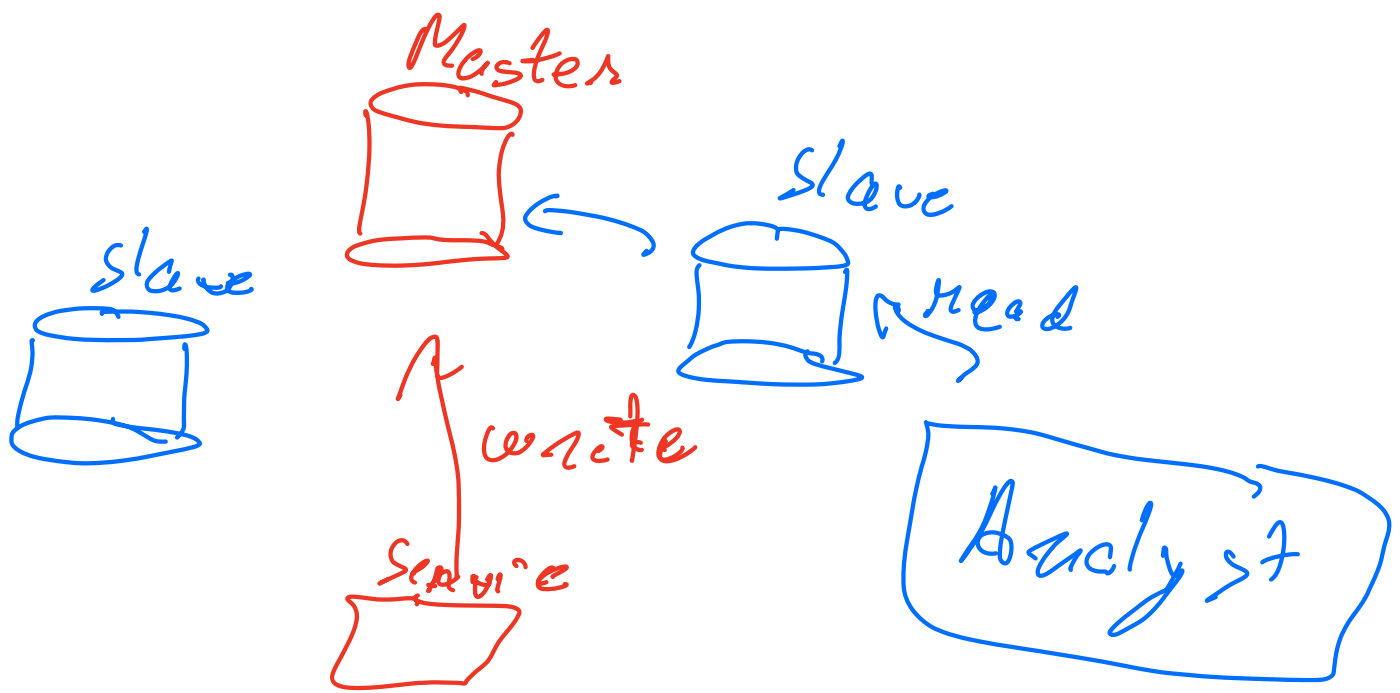
NoSQL →
↓

①

NoSQL ⇒ ACID

- ① Requirement of every the data
- Monthly sales
 - Best selling data





Node = Machine = Sensor

Service = Transactional Data

Analytics = Analytical Data

Service/Transactional \rightarrow OLTP

① MySQL ② NoSQL

Analytical \rightarrow OLAP

\Rightarrow BigQuery \Rightarrow Hive \Rightarrow Spark

OLTP

- ⇒ Online placement
- ⇒ Banking systems
- ⇒ Payment processing

OLAP

- ⇒ Analytical dashboards
- ⇒ Recommendation engine
- ⇒ Historical data

6 V's of Big Data

① Volume ⇒ 1 PB data/day

② Variety ⇒

Un-structured ordered product-id, date

Structured data

③ Σ Σ

JSON

3

"to": " —> "

"from": " , "

"body": " . "

3 "file": " . "

XML, HTML

Semi-Structured data

③ photo, file, pdf, video

Unstructured data

③ Velocity

⇒ Real Time

⇒ Batch

④ Value

Whether data is valuable for companies

⑤ Veracity

Data is trustworthy

⑥ Variability

Inconsistencies & evolving data

① Software Engineer \rightarrow Payments

② Data Engineer \rightarrow Store the data,
Manage the data, processed, data quality

③ Data Scientist \rightarrow models on top of data
recommender

④ Data Analyst \rightarrow query that data.

Software Engineer

\downarrow Data Produce

Data Engineer

↓ Play with data

Data Scientist / Data Analyst

Use this data

Skills required to become DE

⇒ Programming skills → Python / Java / Scala

⇒ SQL →

⇒ Big Data Tech → Hadoop, Spark, Hive, Kafka

⇒ ETL Tools → Extract Transform & Load

⇒ Cloud Platform → BigQuery

⇒ DBMS → SQL, NoSQL

⇒ Data Modelling → Design Schema

DE Interview

- DE Technologies/Skills
 - Basics of DSA
 - SQL
 - Soft skills / project
-