Before we start lets take a quick recap of what we learnt in the last lecture. =>

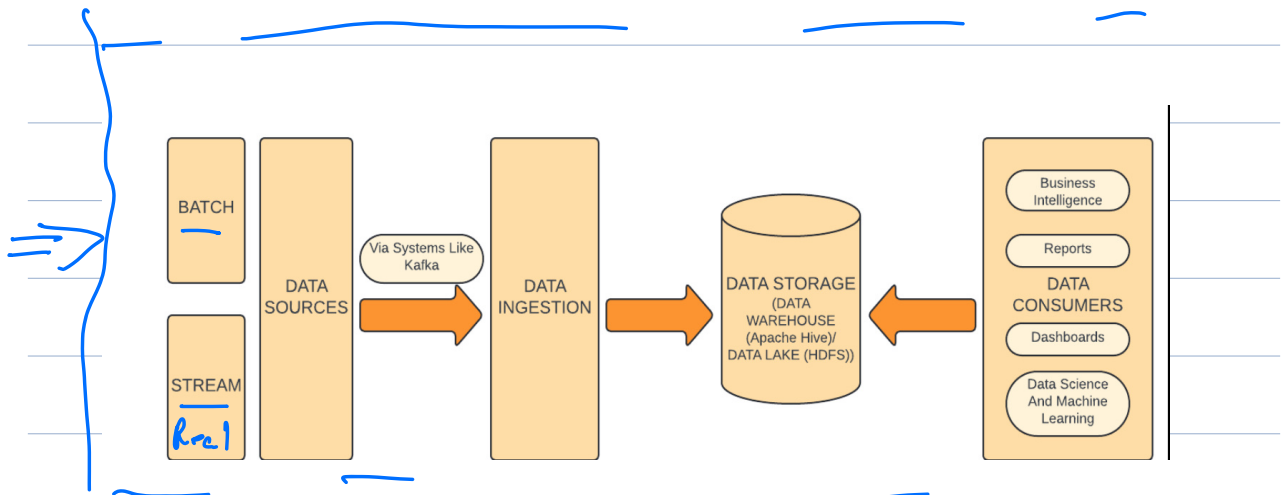① Why do we need big data technologies?

② OLTP vs OLAP

③ 6 v's of Big data -> [Volume] Variety, Velocity, Value, Veracity & Variability

④ DE role & tech stack required

Uber Tech stack

[2-3 PB/day]

| | | | | |
|---|---|---|---|---|
| BATCH | DATA SOURCES | Via Systems Like Kafka → DATA INGESTION → | DATA STORAGE (DATA WAREHOUSE (Apache Hive)/ DATA LAKE (HDFS)) ← | DATA CONSUMERS · Business Intelligence · Reports · Dashboards · Data Science And Machine Learning |
| STREAM Real | | | | |

① Real time → Immediately

② Batch → Bulk

① Real time

→ Latency↓ → Throughput↓
      Cab pricing

② Batch

→ Latency ↑      Throughput ↑

      Forecasting of cabs required
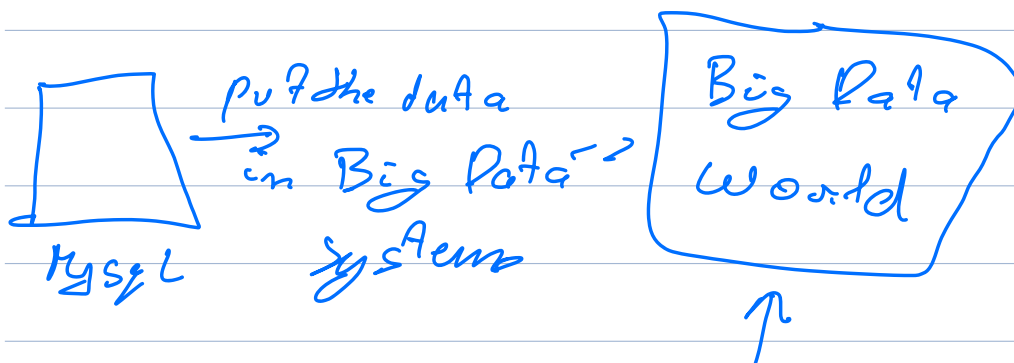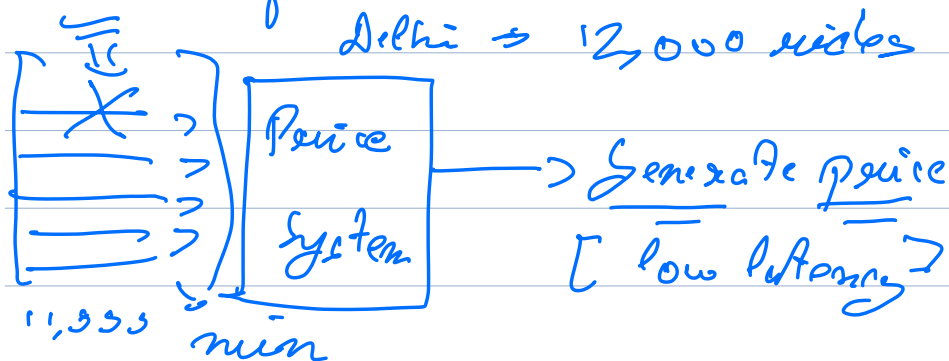
# ① Data Sources

=> From where data arise

=>

## OLTP

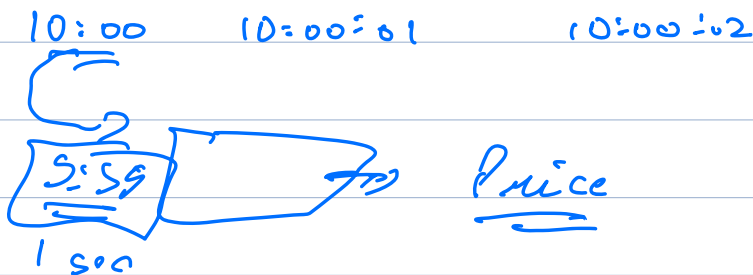Book a cab -> give me my cab details
Amount I paid, etc...  -    -

## OLAP

=> Price of a customer ride

Delhi => 12,000 rides

| Price System | -> Generate price |

[ low latency ]

1,999  min

mysql   put the data in Big Data System   Big Data World

Analyst

1 month

how many cabs we need ?

10:00        10:00:01              10:00:02

9:59     →     Price

1 sec

---

OLTP → Transactional Data
    ↳ All interactions that happen
        on App

OLAP ⎰→ Realtime → Data for cab
                        pricing
       ⎱→ Batch
              ↳ data process latter
         ex → cab count forecasting

Realtime -> Flink, Elastic Search,
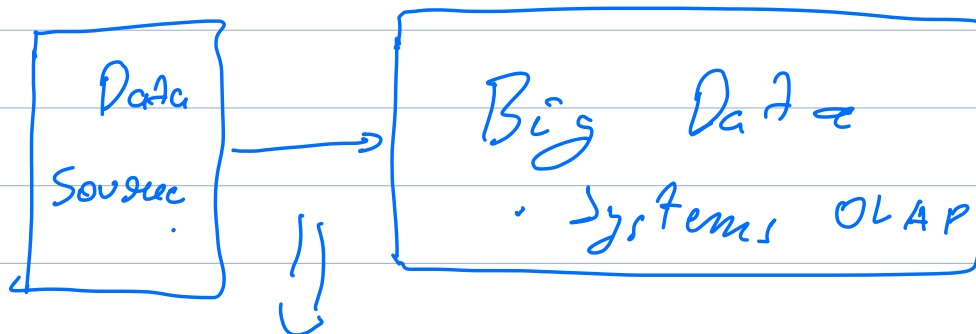Pinot, Presto

Batch -> Hadoop, Spark, Hive, MR

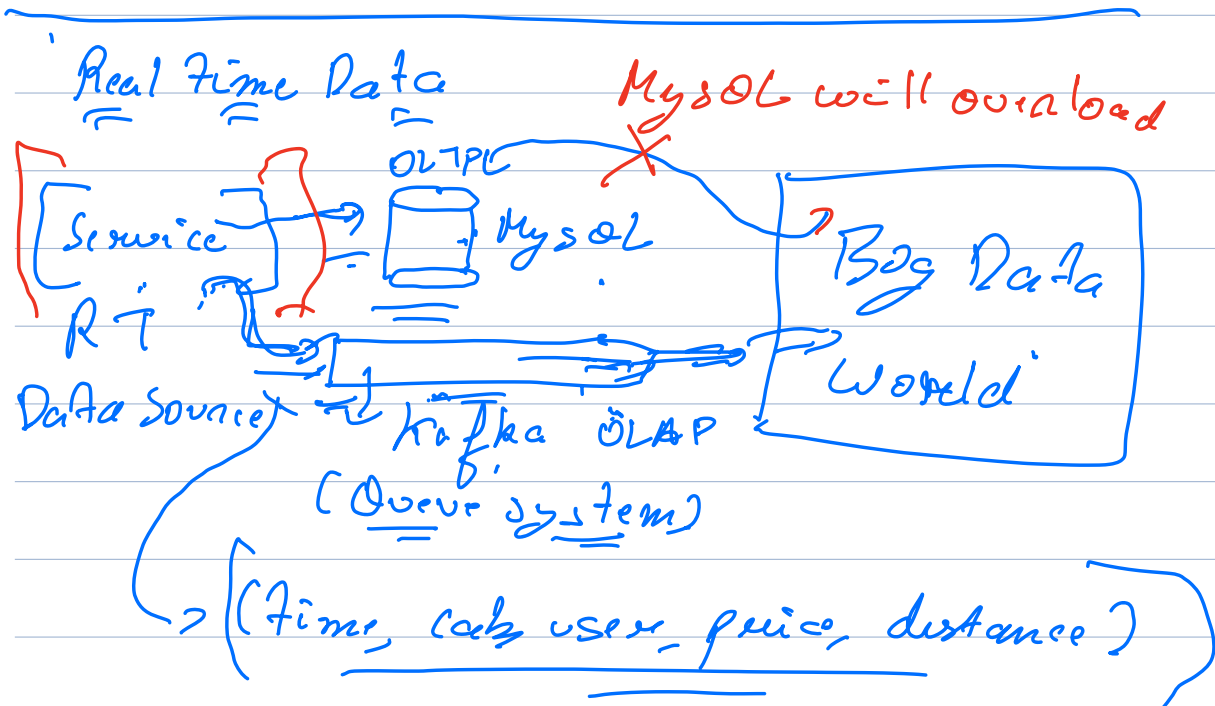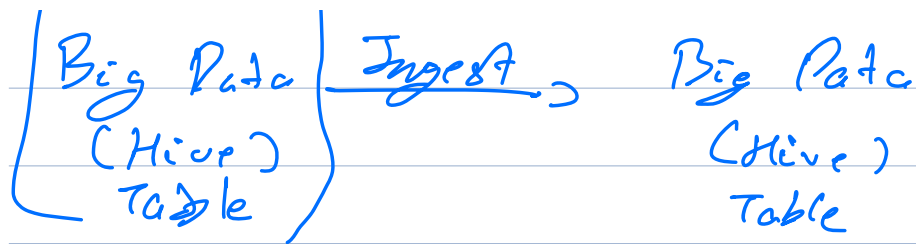OLTP -> MySQL, NoSQL [Mongo ....]

Structured -> Hive => Cab ride data
Unstructured -> HDFS
Semi-Structured -> json, xml -> HDFS
=> specialized engine

=> [Ray = Apache Ray]

Data Source ----> Big Data
. Systems OLAP

Data Ingestion

Big Data | Ingest → Big Data
(Hive) (Hive)
Table Table

---

Real Time Data                    MySQL will overload

[Service]         OLTP
RT        →   [MySQL] ✗        Big Data
                                 World
Data Source →  Kafka  OLAP
              (Queue System)

→ (Time, cab, user, price, distance)

---

Batch Systems

[Service] | Hive Table | Big Data
          |  ⊞⊞⊞       |  ← →
          |            |  System

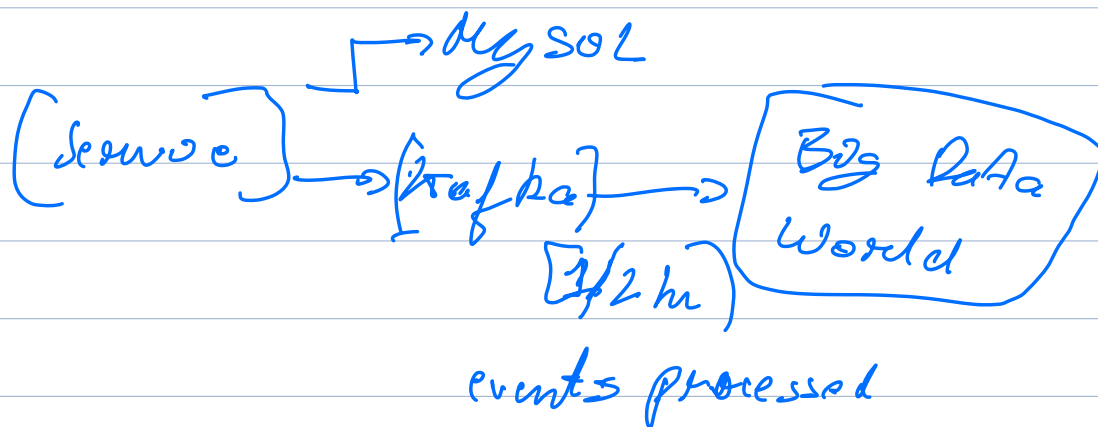Data Source

(Hive / Spark)

① Hive Table 1 => [message, cab_id, user_id, price
                                  50 more fields]

② Hive Table 2 => Hive Table 1 => Extracting
   the data
   =
                    ① [Transform]
              sum(price) date_str, . . .
   && user_id is not null

③ Insert into HiveTable 2 select * from
   HiveTable1 where <        ? .

                    => MySQL
   [service] —> [kafka] —> [ Big Data
                                World ]
                 [1/2 hr]
              events processed

ETL => Extract, Transform, Load.

kafka is a queue → capacity ⇒ 1 billion
                                    records
(1 million records) per day ⇒ our scale
                    ↑ 1/2 hr. ⇒ few thousands
                    ┌─────────────┐  data
                    │ Big Data System │
                    └─────────────┘

Data Lake                    Data Warehouse
   ↥                            ↥
Store the                    Organised
data                         data
≈ Organised data,            [Hive]
  picture, video
Hdfs, S3
Ray ⇒ Store in datalake & process it
      on Ray cluster

Data Processing

↳ Ingestion phase
↳ After ingestion

Hive Table → Hive Table

| Realtime | Batch |
|---|---|
| Latency ↓ | Latency ↑ |
| Throughput ↓ | Throughput ↑ |
| Ingested via | Ingested via |
| Kafka [Stream] | Kafka or Other Hive Table |

Data Lake → Any data
Data Warehouse → Structured.

Data Integration

Job → Clearly which runs at fixed time of a day → Cron job

⇒ [ Process my data at 12:00 am ⇒ ]
            ⇓
      updated every day

Apache Airflow

---

Data Consumers

ML engineers
Data Analyst
Reports / Dashboards To visualize data
                ⇓
        Tableau, Power BI

---

[ Table1 ] ⇒ userid, cabid, price, date
                    ⇓
    Table2 ⇒ date, sum(price)