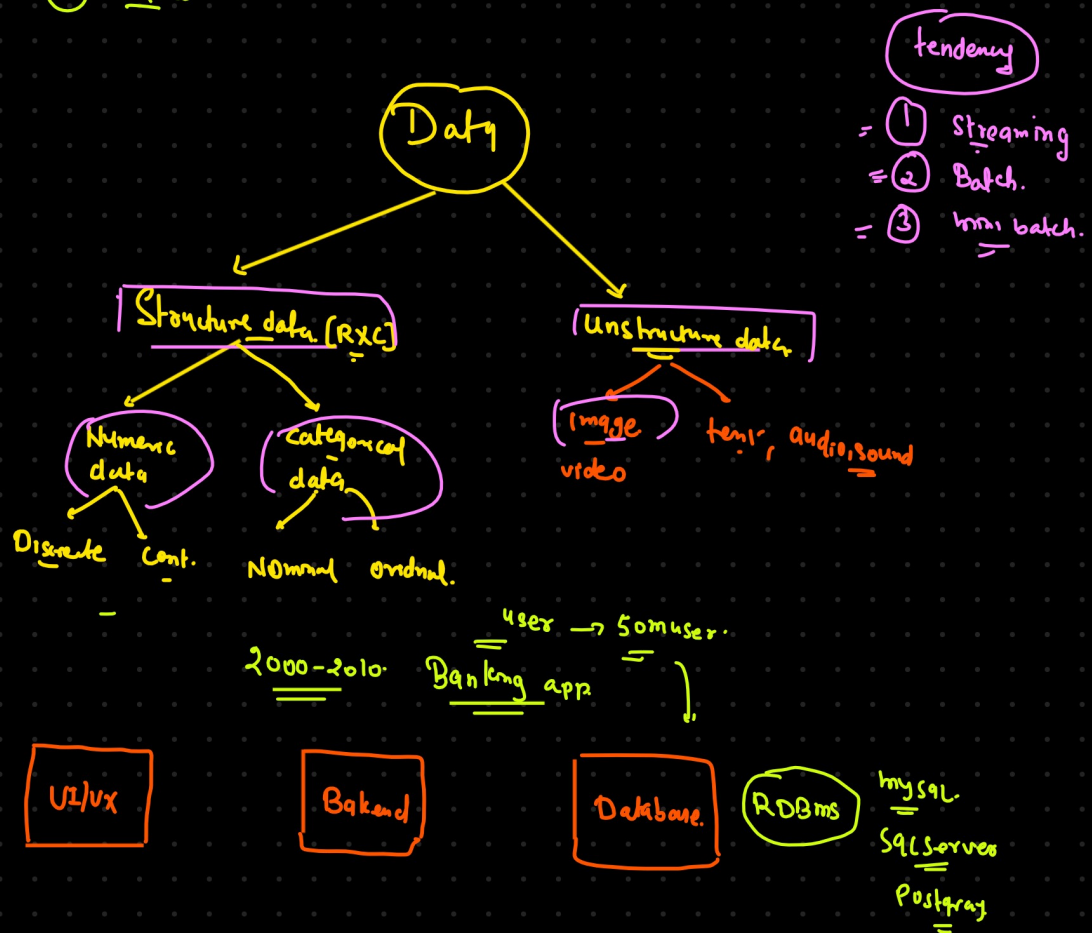
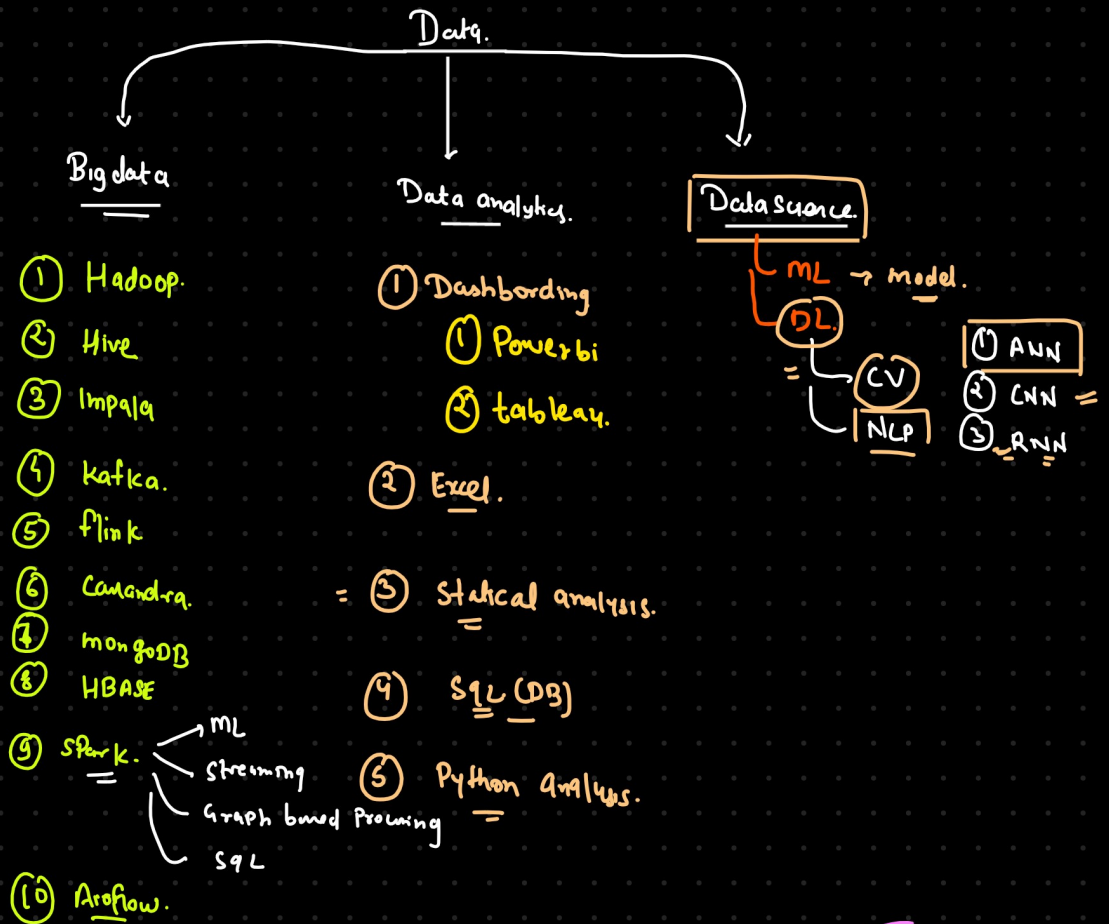
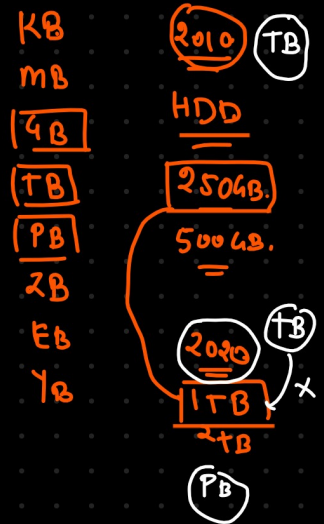


Big - data.



Big data.

- ① massive collection
- ② Increasing exponentially.
- ③ Complicated.
- ④ Store, manage and processing is very difficult.



SV →

- ① Variety. → Structure, Semi-structure, Unstructure. ⁷⁰⁻⁸⁰
- ② Velocity. → Speed. | rate at which data being generated.
- ③ Volume. → massive volume of data
- ④ Veracity. → trustworthy.
- ⑤ Value. → valuable and reliable.

Batch data. → Periodic data. [1 week, 1 month, 6 month, 1 year]

Mini batch data. → less frequent. [1 week, 10 day, 5 day, 2 day].

Streaming data. → very frequent data [less than second
(live) within second]

- ① Big data def
- ② S'v of big data
- ③ type of data
- ④ tendency of data.

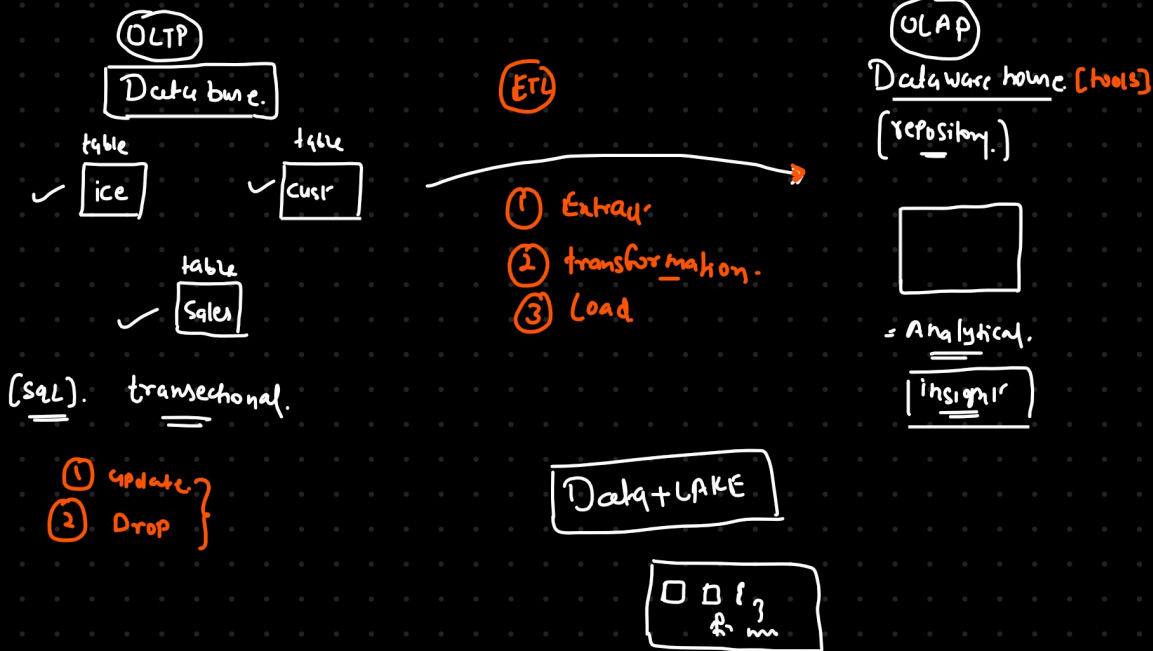
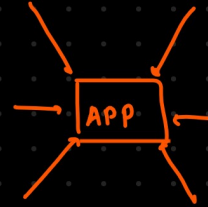
Evaluations of the tools.

- ① file system.
- ② DBMS (RDBMS)
- ③ Data ware house
- ④ Data lake

DFS (Hadoop)

<u>File system</u>	<u>Database.</u>
entire data	specific data
(No Role) non secure.	secure (role)
not able to handle.	handle.
not able to handle	handle

- ① Size.
- ② Security.
- ③ Concurrency.
- ④ Redundancy



ETL → Extract, transform, load.

RDBMS → HDFS → Spark → HBASE

= (E) (T) (L)

- ① Data industry.
- ② introduction of Big data.
- ③ file system.
- ④ Database
- ⑤ DW and DL

Hadoop ⇒ DFS
HDFS

- ① History of Hadoop
- ② HDFS Architecture.
- ③ HDFS → Storing the data.
- ④ mapreduce → Processing of data.
- ⑤ component of Hadoop, HDFS
- ⑥ HDFS/LFS
- ⑦ YARN
- ⑧ Hadoop 1.x v/s 2.x
- ⑨ all the component of YARN
- ⑩ Read and write
- ⑪ Rack awareness.
- ⑫ HDFS commands.

- | | |
|-----------|-------------------------|
| ① HDFS. | ⑥ - ⑦ <u>clustering</u> |
| ② HIVE | |
| ③ Kafka | |
| ④ Spark | |
| ⑤ Airflow | |

Project