

Week 4 :-
quiz

Data formats and streaming.

1. What is true bth data modeling and formatting of the data?
→ The data doesn't necessarily need to be formatted in way that represents the data model. Just so long as it can be extrapolated.

2. What is streaming?

→ Utilizing real time data to compute and change the state of an app. continuously.

Video 1 :-

Data model vs Data Format

Serialized representation of data and data model.

CSV doesn't mean Relational.

Video 2 :-

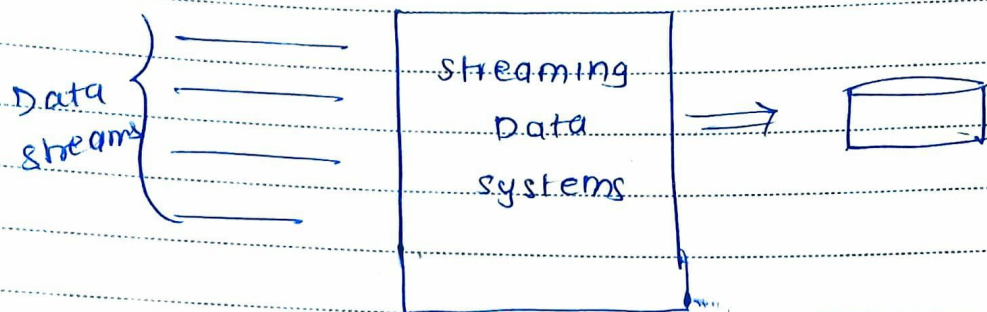
What is data stream.

Big data challenges
↓
velocity

- Key characteristics of a data stream
- requirements of streaming data systems.
- Recognize the data streams you use in your life.

⇒ social media :- Sales Trends + Sales Distribution
↓
Data-driven Marketing
Monitoring and Fault Detection.

A possibly unbounded sequence of data record.
↓
→ Time-stamped
→ Geo-tagged
may or may not be related



- # Streaming data systems → To overcome conventional data system
- Manage one record or small time window
 - near-real-time
 - independent computations.
 - Non-interactive.
- } many challenges.

Some streaming data systems :-

- amazon kinesis
- Apache storm
- Flink
- spark streaming

video 3:-

Why is streaming Data different

→ compare and contrast

"data-in-motion" and "data-at-rest"

→ differentiate btⁿ streaming and batch data processing

→ List management and processing challenges for streaming data

Data-at-Rest :- Mostly static data from one or more sources

- collected prior to analysis.

Data-in-Motion :- Analyzed as it is generated

Ex. sensor data from self-driving vehicles.

- Stream processing.

i) Static/Batch processing :- size determines time and space.

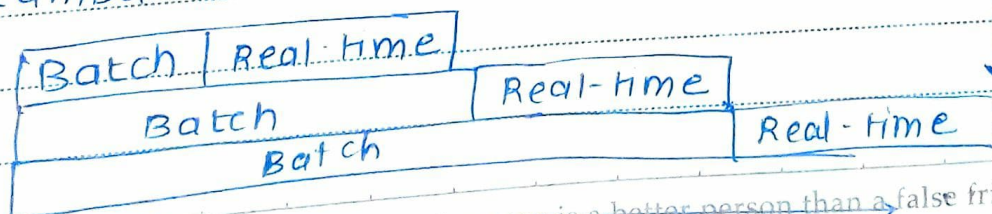
ii) streaming processing :- Unbounded sizes but finite time and space

Streaming data management and Processing

→ compute one data element or a small window of data elements at a time.

→ Relatively fast and simple computations.

→ No interaction with the data source

Lambda Architecture λ 

Priority

• A good enemy is a better person than a false friend.

Time

streaming data changes over
size + Frequency

M/T/W/T/F/S/SU/

changes can be periodic or
sporadic

i) periodic : evenings, weekends

ii) sporadic : major events (Breaking News)

Example of extreme changes :-

World Records for Tweets

- size → unbounded
- size and frequency → unpredictable
- processing → fast and simple

video 4: Understanding Data Lakes

→ Describe how data lakes enable batch processing
of streaming data

→ schema-on-write vs schema-on-read

→ Organizes data stream, data lakes and data warehouse
on a spectrum of big data management and storage

~~What~~ What is Data Lake ?

- click-streams
- social-media
- sensor data
- sales transaction
- Geo location

How Data Lake Works :-

Load data from source

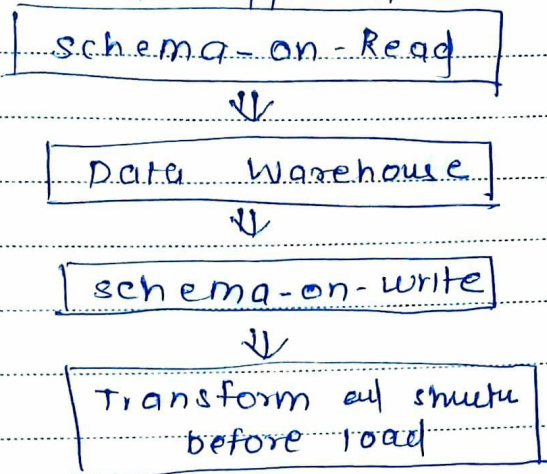
↓

store raw data

↓

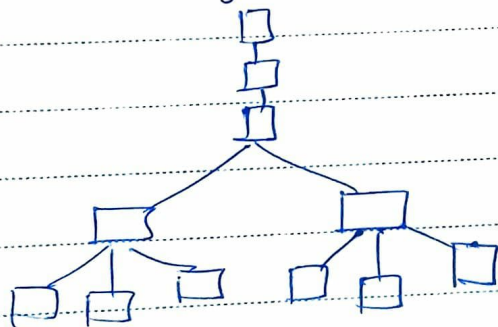
Add data model on read

Schema - On - Read Approach

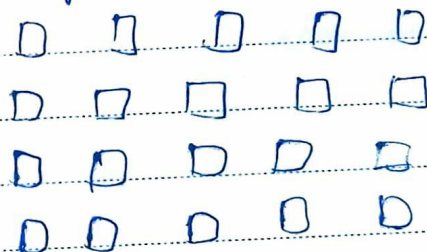


Data lake vs. Data Warehouse

Data Warehouse
Hierarchical File system



Data lake:
object storage



Data Lake object storage Application

API

