

Lambda Architecture Assignment

1. Explain the factors leading to Big Data, List and Explain major sources of Big data?
- The rise in technology has led to the production & storage of large amounts of Data. & Affordability of Smart phones.
- * Earlier MetaBytes of data were used but now a days Petabytes of data are used for analyzing, discovering new facts & knowledge.
- * Traditional systems were not built to handle large amount of data so there is a need for new systems to analyse.
- * The Big Data is a high-volume, high-velocity, high-variety & a high veracity in formation set.
- The Major Sources of Big Data are:-
- i) Social Networks & web Data:- Such as Facebook, Twitter, Emails, blogs & youtube on an average per day Billions of Bytes of data are produced by Social Networking sites making it a large contributor of Big Data.

- * Transaction data & Business Processes:- Such as credit card, transactions, flight bookings etc. & medical records.
- * Increased Data Generation:- As activities like online shopping, social media produce more data generation.
- * Devices and Sensors continuously collect and transmit data.

* E-commerce & Transaction data:-
generated online shopping transactions, payment processing data.

- * Health Care and Medical Records:-
EHR's and medical imaging
- * Telecommunications:- call records, network traffic data
- * Government and public records:- census data, public health records

② List and Explain the characteristic of Big Data

→ Big data is defined by several key characteristics that distinguish it from traditional data. These characteristics are often summarized by the term 3Vs. Given below are the characteristics list

and their descriptions based on its nature and

1. Volume:- It represents the big size of the data. It comes from large pieces of data or collection of small data over a period of time.

* The name Big Data itself is related to its size which is enormous.

Volume is a huge amount of data to determine the value of data. Size of data plays a crucial role.

* If the volume of data is very large, then it is actually considered as a 'Big Data'.

Ex:- Facebook can generate approximately a billion messages 4.5 billion times that the 'Like' button is recorded, & more than 350 millions new posts are uploaded each day.

2. Velocity:-

* It refers to the high speed of accumulation of data.

* In big data velocity data flow in form sources like machine, network, social media, mobile phone etc.

* There is a massive & continuous flow of data.

* It contains the linking of incoming data lets speed, rate of change & activity burst.

Veracity: It means how much the data is reliable.
It has many ways to filter & translate the data.
Veracity is the process of being able to handle & manage
data efficiently.

Ex: Facebook post with hashtag.

Variety: Big data can be structured, unstructured & semi-structured that come from different sources.

- * In earlier the data will only be collected from databases & sheets.
- * But these days the data will come in array forms, that are PDFs, Emails, Audios, Photos.
- ③ List and Explain the major challenges of Big Data systems.

⇒ ① Data Security & Privacy: Security is a big concern for organizations. Non-Encrypted information is at risk of theft or damage by cyber-criminals. therefore, data security professionals must balance access to data against maintaining strict security protocols.

- * Data Quality Issues:- Data is the foundation of data analytics, & its quality directly impacts the insights derived.
 - * Missing information can lead to skewed results & hinder effective analysis.
 - * Errors in data entry or inconsistencies can lead to misleading conclusions.
 - * Ethical and Bias challenges:- Data & algorithms can inherit biases from the real world, leading to unfair or discriminatory outcomes. It's crucial to be vigilant.
 - * Algorithms trained on biased data can perpetuate those biases in their outputs.
- ④ Discuss the desired/required properties for Big Data systems
- Big data systems must meet certain properties to efficiently handle the challenges posed by the large volume, velocity, variety and veracity of data. Here are few required properties.

1. Scalability:- Big data systems must be scalable to handle the ever-growing volumes of data. Scalability can be horizontal (adding more machines or nodes to a system) or vertical (adding more resources to existing model).
2. Fault Tolerance and Robustness:- Big data systems need to be fault tolerant, meaning they should continue to work even if some components fail. Reliability ensures that data processing is accurate and consistent, even after hardware failure.
3. Low latency and High throughput:- Latency is the time delay between data input and the corresponding output. Low-latency systems can provide fast responses or quicker data processing. High throughput refers to the system's ability to process a large amount of data at a given time.
4. Generalization:- A general system can support a wide variety of applications based on functions of all data, it generalizes to all applications like finance management systems, social networking or computer analysis etc.

5. Extensibility:- Extensible systems allow functionality to be added with a minimal development cost.
Ex:- Change / add of a new feature required a migration of old data into a new format but as big data systems are extensible it is easy and quick to do big migration.

6. Adhoc Queries:- Being able to do adhoc queries on your datalets has increased opportunities for business optimization and new applications.

7. Minimal Maintenance:- Maintenance is the work required to keep a system running smoothly. This includes anticipating when to add machines to scale, keeping processes up and running, debugging anything that goes wrong in production.

8. Debuggability:- A big data system must provide the information necessary to debug the system when things go wrong. The key is to be able to trace for each time & value in the system.

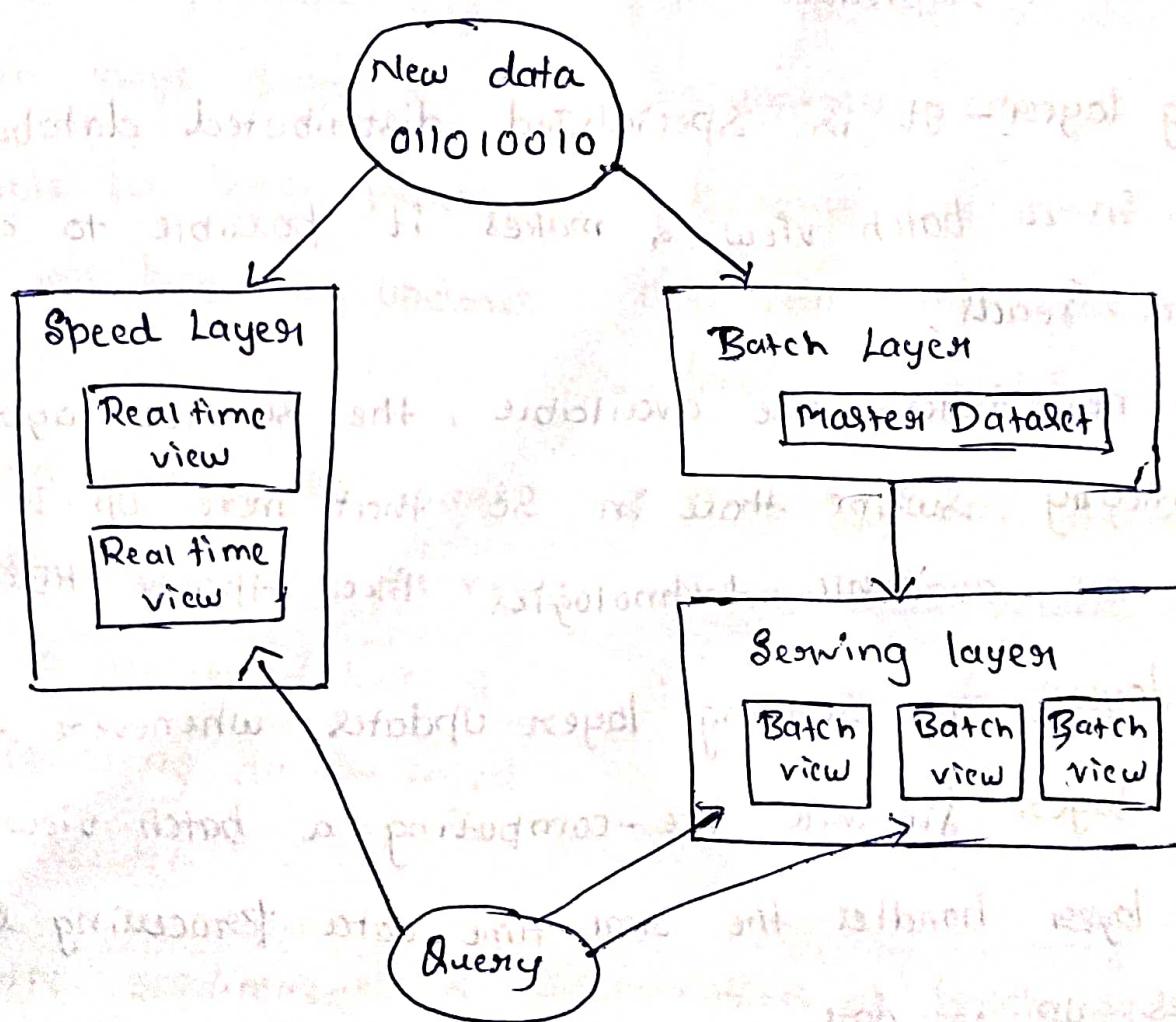
Q) Discuss the problem faced by traditional database system.

- i) Scalability Issues:- Traditional databases struggle to scale horizontally, meaning they are designed to run on a single server or clusters of tightly-coupled servers.
- ii) challenges with Fault Tolerance:- Ensuring fault tolerance in traditional databases, especially in distributed systems can be complex, requiring mechanisms like replication, failover systems and backups.
- iii) Bigdata is too big for traditional storage.
Though in theory traditional database systems ~~can~~ can handle large amounts of data, it simply can't keep up with demands of modern data to deliver the efficiency & insight we need.
- iv) Big data is too complex for traditional storage.
[Though in theory traditional database systems can handle large amount of data, it simply] this is wrong.

Since big data consists of more than rows & columns some may be structured but large part of it is unstructured. which can be an issue as they can't properly categorize it.

- v) Big data is too fast for traditional storage.
→ An RDBMS designed for rapid fluctuations but

6. Explain different layers of Lambda Architecture.



Batch layers: The batch layer stores the master copy of the dataset & precomputes batch views on that master dataset which can be thought of as a very large list of records.

- Batch layer should be able to store an immutable, constantly growing master data set and,
- compare this is immutable append-only data set that stores the entire historical data.
- typically using distributed systems like Hadoop or Spark & it produces batch views handles complex computations for dataset analysis.

Ex: Hadoop, Mapreduce, Apache Spark.

Serving layer: - It is specialized distributed database that loads in a batch view & makes it possible to do random reads.

* When new batches are available, the service layer automatically swipe those in so that more up-to-date results are available. technologies like Apache HBase

Speed layer: - The serving layer updates whenever the batch layer finished pre-computing a batch view.

* Speed layer handles the real-time data processing & provides up-to-date.

Realtime view = function (real time view, new data)

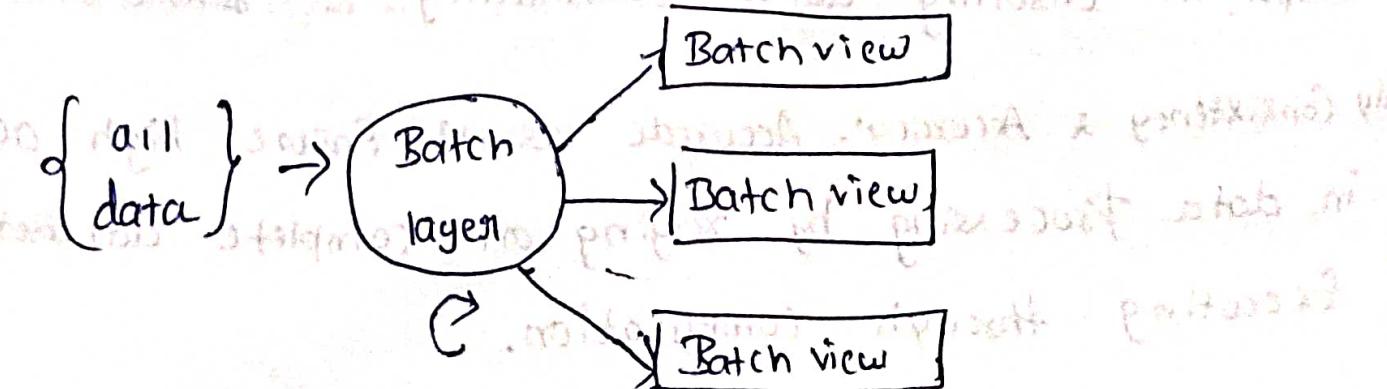
batch view = function (all data)

Realtime view = function (realtime view, new data)

query = function (batch view, realtime view)

- ⑦ Differentiate b/w re-computation algorithm and increment algorithm
- re-computation algorithm Increment algorithm
- * It re-computes the entire dataset from scratch
 - * It updates only the parts of the data that changed
 - * Lower efficiency compared to Increment
 - * High latency, especially with large datasets
 - * Suitable for batch processing and less frequency updates
 - * It will update the count by first appending the new data to master dataset & then counting all the records from scratch
 - * It will count the number of new data records & add it to the existing count.

- ⑧ List the requirements & responsibilities of batch layer



Requirements of Batch layers

- i) Scalability:- It should be horizontally scalable to handle massive volumes of data as more data is added over time.
Ex:- Apache hadoop & Apache spark.
- ii) Fault tolerance:- It should be designed to recover gradually from failures.
- iii) Performance:- The system should efficiently perform large scale data computations.
- iv) Data Integrity:- Ensuring that the results produced are consistent & accurate even in the face of failures.

Responsibility of the Batch layer

- i) Data Processing:- Batch layer processes large volume of data in batches, it is designed to handle comprehensive computations & generate results.
- ii) Data Storage:- It is immutable storage systems which helps in ensuring data consistency & reliability.
- iii) Consistency & Accuracy:- Accurate results ensure high accuracy in data processing by relying on complete datasets & executing through computation.

① Explain the requirements of serving layer in lambda architecture.

→ Batch writable:- The batch views for a serving layer are produced from scratch. When a new version of a view becomes available, it must be possible to completely swap out the older version.

Scalable

- A serving layer database must be capable of handling views of arbitrary size.

Random reads:- A serving layer database must support random reads, with indexes providing directness.

Fault Tolerance:- The serving layer must be fault tolerant, ensuring that even in the case of node or system failure, data remains available for

Querying:- It should replicate data across multiple nodes & regions to ensure high availability &

resilience against failure.

(10) with example show how low latency & high throughput can be achieved in serving layer of lambda

Architectures.

→ i) Low latency in serving layer.

Ex:- Suppose we're building a recommendation system

for a e-commerce site, where customers expect product recommendations as they browse the site.

Batch layer: Batch layer processes historical data such

as computed e-commerce site, customer purchase history, browsing history and product details. Generating precomputed recommendation models (batch view)

Speed layer: The speed layer process real-time data,

such as the customer's most recent click and page

views to adjust the recommendations on the fly.

Query workflow for above executing all three layers

* A customer lands on the site and starts browsing

* A query is made to the serving layer for their recommendations based on customer IDs, key techniques for low latency are.

Precomputed views,

Differentiate b/w batch layer and speed layer.

Batch Layer Speed Layer

- * Process large volume of historical data to produce accurate comprehensive results.
 - * It has Batch processing model.
 - * It has High latency rate.
 - * Provides accurate & complete results.
 - * Scaler are large dataset but not suitable for real time needs.
 - * typically uses distributed file systems like Hadoop.
 - * Suitable for complex, large scale queries on entire datasets.
- * It handles real time to provide low-latency, approximate results.
 - * It has Stream processing model.
 - * It has low latency rate.
 - * Produces approximate results due to real-time nature.
 - * Scaler to handle high throughput for real time data.
 - * used in memory systems like redis or kafka for fast access.
 - * Suitable for quick real time queries on recent data.

2. High throughput in the processing layer.

Ex:- Consider the financial analytics system that receives millions of queries per second for stock prices and market trends.

Batch layer:- Batch layer processes historical financial data, generating pre-computed views, such as moving averages, volume trends & stock price aggregation.

Speed layer:- The speed layer processes real-time stock price updates to provide the latest information.

(12) Requirements & Responsibilities of Speed layer:-

Low latency processing:- The primary requirement of the speed layer is to process data in near real-time, ensuring minimal delay from when data is ingested to when it can be queried.

Efficient Real-time Processing Framework:- It should use a stream processing framework capable of performing computations (eg- aggregations joining) in real time as data flows through the system.