

Lambda Architecture Assignment

1. Explain the factors leading to Big Data. List & explain major sources of Big Data.

- The rise in technology has led to the production and storage of Large amounts of Data. & Affordability of smart phones.
- Earlier MetaBytes of data were used but nowadays 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 analyze the fast growing complex and unstructured Data.
- The Big Data is a high-volume, high-velocity, high variety & a high veracity information set.

The Major Sources of Big Data are:

- 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 contribution of Big Data.
- Transaction data & Business Processes: Such as credit card, transactions, flight bookings etc. & Medical records.

- Machine generated Data: From IoT sensors, Home automation, traffic sensors etc. Data from computer systems like Logs, web logs etc.
- Human-generated Data: biometrics data, human Machine interaction data.

Q. List and explain the characteristics of Big Data.

- Volume: It represents the big size of the data. It comes from large piece of data or collection of small data over a period of time.
- Velocity: The term Velocity refers to the speed of generation of data. Velocity is so fast to meet the demands & the challenges of processing Big Data, the velocity of generation of data plays a crucial role.
- Variety: Big Data comprises of variety of data. Data is generated from multiple sources in a system. This introduces variety in data & therefore introduces 'complexity'. It is caused to the availability of large number of heterogeneous platforms in the industry.

→ Veracity: It defines the quality of data captured, which can vary greatly, affecting its accurate analysis.

3. List & explain the major challenges of Big Data Systems.

→ Scaling up & down Big Data according to current Demand.

④ Due to data verity it is difficult to create common storage & difficult to integrate.

→ Collecting and Integrating massive & Diverse Datasets.

④ Collection of data might be difficult as a lot of institutions would refuse to share & available data consists of diverse datasets. Finding a common ground to integrate would be difficult.

→ Picking the Right NoSQL Tools:

→ Depending on the type of data choosing the right NoSQL tools can be difficult as it might take trial & error approach.

→ Maintaining Data Integrity, Security and Privacy:

→ Securing the Data for the analysis can be tough as some data can be sensitive like patient information.

→ Overcoming Big Data talent & Resource constraint.

- ④ Due to increasing tools & technologies in big data it is hard to find the right talent as people are still working on traditional methods.
- ⑤ To maintain & process data of large volume it is near impossible due to hardware constraints.

4. Discuss the problems faced by traditional database Systems.

→ Big data is too big for traditional storage.

- ⑥ Though in theory traditional database systems can handle large amount of data. It simply can't keep up with demands of modern data to deliver the efficiency & insights we need.

→ Scaling up vs scaling out.

- traditional data storage systems are centralized, so we are forced to scale up which is less resource efficient than scaling out.

→ Big data is too complex for traditional storage:

- 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.

- Big data is too fast for traditional storage.
- An RDBMS isn't designed for rapid fluctuations. But Big data grows almost instantaneously & analysis needs to occur in real time. for example sensor data, which produce at concurrent rate, which may be difficult to analyze.

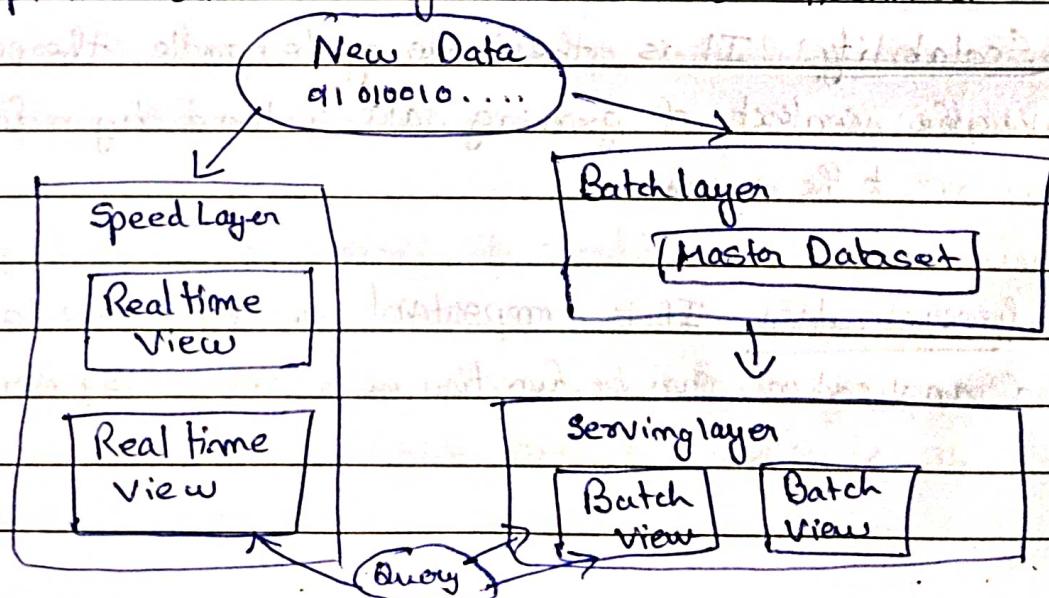
5. Discuss the required Properties for Big Data Systems.

- Robustness & error tolerance: It's important to have a system that can handle errors and continue to function even if parts of the system fail.
- Systems are required to behave in a right manner despite machines going down randomly. Robustness of big data System is the solution to overcome the obstacles associated with it.
- De buggability: A system must be debug when something happens by the required information, delivered by the big data system.
- Scalability: It is the tendency to handle the performance in the context of growing data & load by adding resource to the system.
- Generalization: It is important to generalize a wide range of application can be function in a general system.

- Ad hoc queries: Every large dataset contains an anticipated value in it. The ability to perform ad hoc queries on data is significant.
- Extensivity: Extensible system enables to function to be added cost effectively.
- Low latency reads & updates: In big data system, there is a need of applications low latency on updates propagated shortly.
- Minimal maintenance: Selecting components with probably little complexity plays a significant role in minimal maintenance.

Most of these problems are solved by Lambda Architecture.

6. Explain different layers of Lambda architecture.



Speed Layer: Batch Layer:

- 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 dataset and,
- Compute arbitrary functions on the dataset.

Serving Layer:

- The batch layer emits batch views as the result of its functions
- The serving layer is a specialized distributed database that loads in a batch view & makes it possible to do random reads.

Speed Layer:

- The serving layer updates whenever the batch layer finishes precomputing a batch view
- The goal is to ensure new data is represented in query functions as quickly as needed for the application requirement

7. Differentiate between re-computation algorithm & increment algorithm.

re-computation algorithm

→ When new data arrives it throws away the old batch views and recomputing functions over the entire master dataset.

→ It will update the count by first appending the new data to master dataset & then counting all the records from scratch.

→ Lower efficiency compared to Increment Algorithm

→ extremely tolerant of human errors

→ Requires computational effort to process the entire master dataset

Increment Algorithm

→ When new data arrives it will update the batch views directly when new data arrives.

→ It will count the number of new data records & add it to the existing count.

→ Higher efficiency compared to re-computation Algorithm.

→ Doesn't facilitate repairing errors in the batch views.

→ Requires less computational resources, but may generate much larger batch views.

8. List the requirements & responsibilities of batch layer.

Responsibilities of the Batch Layer:

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

Requirements of Batch layer:

1. Scalability: It must be able to scale horizontally to handle large volumes of data.
2. Fault tolerance: It should be designed to recover gracefully from failures.
3. Performance: The system should efficiently perform large scale data computations.
4. Data Integrity: Ensuring that the results produced are consistent & accurate, even in the face of failures.

9. 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 direct access.
- Fault tolerant: Because a serving layer database is distributed, it must be tolerant of machine failures.

10. With example Show how low latency & high throughput can be achieved in serving layer of Lambda architecture.



11. List the requirements & responsibilities of Speed layer.

- Random reads - A realtime view should support fast random reads to answer queries quickly. This means the data it contains must be indexed.
- Random writes - To support incremental algorithms, it must also be possible to modify a realtime view with low latency.
- Scalability: As with the serving layer views, the realtime views should scale with amount of data they store and the read/write rates required by the application.
- Fault tolerance: If a disk or a machine crashes, a realtime view should continue to function normally.

12. Differentiate between batch and speed layers.