

Assignment -1

Date:

Page:

Lambda Architecture

- Explain the factors leads to Big Data, List and explain Major source of Big Data.
-
- 1) Technological Advancements:- The exponential growth in computing power & storage capabilities led to the growth of Big Data.
 - 2) Digitalisation:- Business operations, Social interaction research etc. more has led to the increase in the amount of Big Data.
 - 3) Internet & Connectivity:- The wide spread availability of Internet and the proliferation of connected devices.
 - 4) Social Media & User Generated Content:- The rise of social media platforms, bloggers and other online platforms.
 - 5) Sensor Networks:- Advances in sensor technology have enabled the collection of data from various sources.
 - 6) Data Storage & Cloud Computing:- The advent of cloud computing platform has made it easier & more cost effective.
 - 7) Data Analytics tools:- The development of Data Analytics tools have made the organisations to store & process large volume of Data at given time.

huge amount of data stores them
without loss of quality is a challenge

→ Fault tolerance - Computers Fault
tolerance is extremely hard.

→ Scalability - Big Data projects can grow
and evolve rapidly. The scalability issue of
big Data has lead toward cloud
computing.

4) Discuss the problem faced by traditional
Database Systems.

→ * Scalability

Traditional Database Systems often struggle
to handle massive volumes of Data
generated by modern applications

* Performance

As Data size and query complexity
increase traditional Database may experience
performance is compromised. Complex
queries, joins and aggregation can
lead to slow response times, making
realtime analytics.

* Data Model Rigidity

Traditional Database are typically designed
with a fixed Schema. makes it
challenging to adapt to changing data Requirements.

* Data Consistency

Maintaining data Consistency in
distributed environment can be difficult

Traditional Database use Acid transaction which can be limits in distributed Systems

Lack of Support for Unstructured Data

The traditional Database are primarily Designed for structured data & may not handle unstructured or semi structured data types.

Security Concerns:

Data security is critical and traditional databases may face vulnerability if not configured and Managed properly.

5) Discuss The required properties of Big Data Systems.

→ * Robustness & fault tolerance

Systems need to behave correctly despite machine goes down randomly. The complex semantics of consistency in distributed database.

* ^{read} low latency & updates :- The vast majority of applications require reads to be satisfied with very low latency.

* Scalability

Scalability is the ability to maintain performance in the face of increasing data or load by adding resources to the System.

* Extensibility

Extensible Systems allow functions to be added with a minimal development cost.

* Adhoc queries:-

Nearly every large dataset has unwanted values within it. Been able to mine a dataset arbitrary gets opportunities optimize ad new applications

* Minimal Maintenance

Maintenance is the work required to keep the System Running smoothly in all environments.

c) Explain Different layers of Lambda Architecture:-

→ There are 3 layers.

1) Batch layer:-

This layer is responsible for processing large volume of Data in a fault tolerant & scalable manner. It stores all the data in its raw form and precomputes batch views on top of it.

2) Speed layer:-

This layer only responsible for processing real time data streams. It stores only the most relevant and recent

data & compute real time views on top of the existing one.

3) Service Layer:

This layer is responsible for serving queries on the data. It merges the batch view & real time views to provide a complete view of the data.

7) Differentiate between re-computation Algorithm and Incremental Algorithm.

→ Incremental Algorithm

Update a batch view by using the new data and the current state of batch view to perform an update. This approach is more efficient in terms of resource usage because it only needs to process the new data & update the relevant parts of the batch view.

Re-computation Algorithm

Update a batch view by looking at the entire master dataset, including the new data. This approach is more general and can handle any type of update, but it is less efficient in terms of resource.

8) List & explain Requirements of Responsibilities of batch layer.

→ Requirements

→ Process large volume of Data in a fault tolerant and scalable manner.

- store all data in its raw form
- Precompute batch views on top of the raw data
- Batch views should be immutable & recomputed periodically ^{to} reflect the latest Data.

Responsibilities:

- * Ingest raw data from various sources
- * Validate and clean the Data
- * Store the Data in a Distributed File System or Database
- * Precompute Batch views on top of the raw Data using Batch processing Algorithm.
- * Handle updates to the Data by recomputing the batch views.
- * Ensure consistency across the batch views by merging them with the real-time views in the serving layer.

9) Explain the responsibilities of serving layer

* low latency access

The serving layer must provide low latency access to the precomputed batch views.

* Scalability :- The serving layer must be scalable enough to handle the large volume of Data.

* Random Read The Server layer must be supports random read. with indexes provides direct access to small portion of the views.

* Batch writable

The server layer must be batch writable, meaning that it can be updated in bulk.

* Error tolerance

The server layer must be able to correct the errors, This is achieved by resampling the server layer views.

10) With Example slow latency and high throughput can be achieved in server layer
 → To achieve low latency the data must be precomputed and indexed in a way that allows for fast queries. This can be achieved by denormalizing the data & tailoring the views to the specific queries they serve.

Let's consider the example of large Dataset of customer transactions for an e-commerce website. We want to be able to quickly query the total revenue for a specific product categories over a certain time period. To achieve this we can denormalize the data by only precomputing batch row level aggregates the revenue for each product categories over

time. The core layer can then index the batch view and provide an interface for querying the data. By using a distributed database such as Cassandra we can achieve the given request.

11) List the reqs & responsibl of Speed Layer.

→

Requirements:-

- low latency updates
- Incremental Computation
- Random writes.
- Fault tolerance.
- Scalability.

Responsibilities:-

- Ingest Real-time Data.
- Process Real time Data
- Store Real time views.
- Merge batch & real time views.
- Expires Data.

12) Differentiate between Batch & Speed Layer.

Batch Layer.

* The Batch layer is responsible for precomputing batch views from the master Dataset.

Speed Layer

* The Speed Layer is responsible for providing updates to the real time views.

* Uses Batch Computation to process the entire Dataset and generate batch view
* Incremental computation process only the new data is update the relative view.

* Designed to handle large volume of data & is optimised for throughput rather than latency
* Designed to handle relative data is optimised for low latency update rather than throughput.

* Fault tolerant & can recover from failure by recomputing the batch view
* Recover from failure by recomputing the new Data only the relative view.