

# Lambda Architecture

Q1) Explain the factors leading to Big Data. List and explain major sources of Big Data?

A) Factors leading to Big Data: The emergence of Big Data is driven by several key factors that have collectively transformed the way data is generated, collected, and processed:

- i) Increase in Data Volume: The exponential growth in data generated from various sources, such as social media, sensors, and transaction records, has led to the creation of massive datasets.
- ii) Variety of Data: Data comes in various forms, including structured, semi-structured, and unstructured data.
- iii) Velocity of Data Generation: The speed at which data is generated and needs to be processed is crucial factor.
- iv) Advances in Storage Technology: The reduction in the cost of storage and the development of scalable storage solutions.
- v) Development of Advanced Analytics: The need for more complex analytics & machine learning algorithms to derive insights from large datasets.
- vi) Cloud Computing: The advent of cloud computing has provided the necessary infrastructure to store and process large datasets.
- vii) Increased Connectivity: The proliferation of internet connectivity & mobile devices has led to a surge in data generation.  
\* Major sources of Big Data:
  - i) Social media platforms
  - ii) IoT Devices and Sensors
  - iii) Transactional data
  - iv) Healthcare
  - v) Telecommunications

## E-commerce Data : Scientific Research

### Q1) List and explain the characteristics of Big Data.

Big data is typically defined by its unique characteristics, often referred to as the '3Vs'.

i) Volume: refers to the vast amount of data generated every second from various sources.

The sheer size of data being produced is enormous, ranging from terabytes to petabytes and even zettabytes.

ii) Variety: refers to the different types of data that are generated from multiple sources.

Data comes in various forms, including structured data, semi-structured data, and unstructured data.

iii) Velocity: refers to the speed at which data is generated, collected, and processed.

Real-time data streams from social media feeds, financial markets, and sensor networks require immediate processing and analysis to extract actionable insights.

"True, all the 3's of Big data but over time, additional 'V' has been added to encapsulate its complexity."

iv) Veracity: refers to uncertainty or truthworthiness of data. Not all data is accurate or reliable. Veracity deals with the quality & accuracy of data as well as the challenges posed by inconsistencies, biases & noise in data.

v) Value: refers to the potential economic or informational

benefits derived from data. The ultimate goal of analyzing Big Data is to extract valuable insights that can lead to better decision making.

Variability: refers to the inconsistencies & fluctuation in data over time. Data can be highly inconsistent with peaks and troughs in data flows.

vi) Visualization: refers to the ability to represent complex data in a visual format that is easy to understand & interpret. With the vast amounts of data involved, traditional data visualization methods often fall short.

### Q2) Major challenges of Big Data Systems.

i) Data volume Management: one of the primary challenges is handling the sheer volume of data.

ii) Data variety: Managing this variety is challenging because different data types require different processing techniques.

iii) Data velocity: Realtime data sources require immediate analysis to provide actionable insights.

iv) Data quality & veracity: Big data often includes incomplete, noisy and inconsistent data, which led to incorrect analysis & poor decision making.

v) Scalability: This includes scaling up storage capacity, computation resources & data processing capabilities.

vi) Data security & privacy: Ensuring privacy is challenging due to distributed nature of Big Data Systems.

- vii) Data Integration: Data may come from different databases, even with its own found schema & semantics.
- viii) Data Governance: It is challenging due to the distributed & complex nature of the data.
- ix) Data Processing & Analysis: Traditional analytical tools often struggle with the scale & diversity of Big data.
- x) Cost Management: Building & maintaining BDs can be expensive.

**Q14)** Discuss the problems faced by traditional database systems?

- i) Scalability limitations: RDBMS often involves expensive & complex hardware upgrades (vertical scaling).
- ii) Handling unstructured data: RDBMS can't handle unstructured data because of rigid schema structure of traditional databases.
- iii) Performance bottlenecks: RDBMS can be overwhelmed for real-time data processing & may slowdown or delayed.
- iv) Complex Data Integration: RDBMS struggle to integrate and manage data from different diverse sources.
- v) Lack of flexibility and agility: Any change to schema, such as adding new fields or modifying existing ones, can be time-consuming & disruptive.
- vi) High cost of maintenance: The need for powerful hardware, specialized database administrators, and regular upgrades add to the overall cost.
- vii) Limited support for distributed computing: RDBMS lacks support for distributed computing.

- viii) Challenges with Real-Time Processing: RDBMS's inability to handle real-time data streams effectively limits the use.
- ix) Complexity in Handling High Volume Transactions: Ensuring ACID in high volume transactional systems may lead to bottlenecks.
- x) Data Silos: RDBMS can lead to the creation of data silos where data is isolated within different departments or systems, making it difficult to integrate.

**Q15)** Discuss the required properties for Big Data Systems?

- \* Scalability: It is the fundamental property of BDS, enabling them to handle increasing amounts of data & computational load.
- \* Flexibility: Ability to adapt to changing data schema, source of formats without requiring significant reconfiguration or downtime.
- \* Fault Tolerance: Designed to detect & recover from failures automatically minimizing downtime & data loss.
- \* High throughput: Ensure that the system can handle the ingestion, processing & analysis of data at a rate of real time.
- \* Low latency: Refers to minimal delay between data input & output enabling rapid processing & immediate response.
- \* Data consistency: That all users & applications see the same version of the data at any given time.
- \* Data security & privacy: This includes encryption, access control, authentication, & compliance with data privacy regulations.
- \* Efficient data storage and management: Use of distributed file systems, data compression, that balance speed & capacity.

- \* Advanced Data processing capabilities: enables organizations to extract valuable insights from large complex datasets.
- \* Interoperability & Integration: enabling data sharing, collaboration & comprehensive analysis across platforms.
- \* Cost Efficiency: systems should be designed to optimize resource utilization minimize operational costs.
- \* Ease of use & Management: includes intuitive interfaces, automated management tools & support for easy deployment.

### Q(6) Explain the different layers of lambda architecture.

Lambda Architecture is a data processing architecture designed to handle massive quantities of data by utilizing both batch & realtime processing methods.

1) Batch layer: is responsible for managing the historical data & processing it in large batches. This layer ensures that the data is processed in its entirety to produce a comprehensive accurate view of the data.

#### \* Components.

Master datasets: This is the immutable, append only dataset that stores the entire historical data. It acts as the source of truth of the system.

Batch processing systems: This component processes the master datasets periodically, typically using distributed computing frameworks like MapReduce or Spark. The output is a set of precomputed views or models, which

are stored and used to answer queries that require complete and accurate data analysis.

Output: The results of the Batch layer is usually a batch view or a precomputed model that can be queried for in-depth insights.

2) Speed layer: is designed to handle realtime data processing and provide immediate results. It complements the Batch layer by processing data as it arrives, allowing for quick responses to data events.

#### \* Components:

\* Real-time processing system: This system processes incoming data streams in realtime using technologies like Apache Storm, Apache Kafka Streams, or Apache Flink. It generates low latency realtime views of data, which may not be as accurate or comprehensive as those produced by the Batch layer but are available immediately.

Output: The speed layer outputs realtime views that are quickly accessible and provide immediate insights into recent data.

3) Serving layer: The serving layer is responsible for merging and presenting the output from both the batch layer & the speed layer. It enables querying of the data in a way that combines the accuracy of batch processing with the immediacy of realtime processing.

Query systems: This component allows users to query the merged views from both the batch layer & speed layer. It may involve databases or systems that can handle large volumes of data & provide fast query results.

View Merging: The serving layer must intelligently merge the results from the batch and speed layers, often giving precedence to the real-time data but ensuring that the batch data is used to provide accurate long-term insights.

Output: The serving layer provides users with a unified view of the data, combining the completeness of batch views with the timeliness of real-time views.

(79) Differentiate b/w Re-computation algorithm & incremental algorithm

Re-computation

Definition: Recalculates the entire dataset whenever an update occurs.

Process: Process all data from scratch upon every update.

Advantages: - Higher accuracy  
- Simpler to implement

Disadvantages: - Computationally expensive  
- High latency

Efficiency: low, especially for large datasets

Latency: High, due to reprocessing the entire dataset

Complexity: simpler to implement

Use cases: - Batch processing  
- Infrequent updates

Incremental

Process: updates only the affected part of the datasets.

Process: pushes only the new or modified data.

- more efficient
- low latency
- complex to implement
- risk of inconsistency

Efficiency: high due to processing only changes  
Latency: low as only a small portion of data is processed  
Complexity: more complex, requires tracking of changes

- Real-time Analytics
- Frequent data updates.

(80)

list the requirements and responsibilities of batch layer Requirements

- \* Data Storage: must store an immutable, append-only master dataset
- \* Scalability: able to scale horizontally, handle ever-increasing data
- \* Fault Tolerance: ensure data integrity & process continuity
- \* Comprehensive Data Processing: extensive data processing on the entire dataset to generate accurate & complete views
- \* Processing Time: should complete processing within an acceptable time

Responsibilities

- \* Storing the master dataset: Contain all historical data in an immutable format
- \* Batch processing: execute batch processing jobs on master dataset
- \* Generating batch views: accurate & complete batch views that reflect full state
- \* Handling large-scale data: ensuring that the system can handle vast amounts of information without performance degradation
- \* Data Integrity & Consistency: ensure that the data processing events are consistent and accurate, reflecting true state of dataset
- \* Providing Historical Data Analysis: offer deep insights & trends
- \* Resource Management: to optimize the processing of large datasets

(81)

Explain the requirements of serving layer in lambda architecture

- \* Low Latency: provide quick look-up and access
- \* Scalability: able to scale horizontally to handle a large number of queries
- \* Data Consistency & Merging: must intelligently merge results from both layers
- \* Fault Tolerance & Reliability: must be sufficient to failures
- \* Efficient Data Storage: support fast reads while optimizing storage space

- Versioning & Time Travel Queries: providing access to different versions of data.
- Handling Large Data Volumes: efficiently handle large data volumes.
- Multi-user Query Support: should support concurrent querying of data.
- Query flexibility: allow ad-hoc pre-defined complex queries.
- Support for Streaming Data: provide immediate access to real-time data.

10a) With example shows how low latency and high throughput can be achieved in Serving layer of Lambda architecture  
 • low latency ensure that user receive responses to their queries almost instantly :

- Precomputed Batch Views: suppose you are running a retail analytics system where you need to show the total sales for each product category every hour. The batch layer can precompute the hourly sales totals for each product category and store these precomputed results in a key-value store.

Low Latency Solutions: When a user queries the total sales for a specific product category, instead of recalculating from raw data, the serving layer can instantly retrieve the precomputed result from the key-value store, ensuring quick response times.

- Real time updates from the speed layer: In a social media platform the speed layer provides real time updates for user activity (likes, comments, shares). Whenever a user makes a query for recent activity, the serving

layer fetches these real time updates directly from the speed layer's in-memory store (like Redis or Memcached)  
Low Latency Solutions = By storing real time updates in fast in-memory cache, queries can retrieve the latest data with minimal delay providing low latency response.

Indexing for faster Querying: In a log monitoring system if users query for error logs in a specific time range, indexing the logs based on timestamps can significantly speed up retrieval. Instead of scanning the entire dataset, the index allows the system to directly access the relevant logs.

Low Latency Solutions: Indexing the data in the serving layer allows for quick lookups, reducing the time it takes to retrieve display results to the user.

Achieving high throughput in serving layer.

- Distributed Data Storage:

In a video streaming platform millions of users may be querying video recommendations at the same time. By storing precomputed recommendations in a distributed database like Cassandra or MongoDB, queries can be distributed across multiple nodes, preventing bottleneck.

High Throughput Solutions: A distributed key-value store allows the system to balance the load across different nodes, handling many concurrent queries while maintaining fast response times.

Sharding and partitioning: In an e-commerce platform the sharding layer stores product details across multiple partitions based on product categories. When a user searches for products the system only queries the relevant partition.

High throughput solution: sharding the dataset ensures that queries will distributed evenly across different partitions increasing the system's capacity to handle many queries at once.

- Q3) List the requirements and responsibilities of speed layer
- \* low latency data processing & must process incoming data quickly
  - \* efficient handling of real time data: handle continuous streams of data
  - \* Approximate and incomplete results
  - \* Scalability
  - \* fault tolerance
  - \* State Management
  - \* Responsibilities

real time data ingestion

real time processing

providing real time views

complementing the batch layer

ensuring eventual consistency

stateful computation

fault tolerant real time data processing

(20) Differentiate between batch & speed layer

~~Batch Layer~~

function: process large volumes of historical data

latency: Batch processing

ingestion: (large chunks of data)

latency: High latency (resulting in)

accuracy: provide accurate complete results

database works on large datasets

scalability: scales for large objects

but not suitable for real time

fault tolerance: fault tolerant for data can

be recomputed from dataset

responsibility: maintains master copy

storage: distributed file system DFS

operating: suitable for complex queries

use cases: historical trends,

ML models

Speed Layer

handles real time data to provide low latency

stream processing (continuous)

real time data

low latency (near instant results)

provide approximate results

process real time data streams

scales to handle high throughput

for real time

requires careful state management

real time views to fill the gap

use in memory

suitable for quick real time queries

real time monitoring

alarming systems