

## Imp Questions of Big Data

- 2 Ques-1 Discuss RDBMS, NoSQL and difference ✓  
2 Ques-2 what are Big Data Processing tools? ✓  
1 Ques-3 Discuss the Applications and characteristics of Big Data. ✓  
2 Ques-4 Explain Data Marts, Data Lakes, Data pipelines etc types of Data. ✓  
1 Ques-5 Explain 6V's of Big Data. ✓  
2 Ques-6 what are different file formats of Big Data? ✓  
Ques-7 How will you improve Quality and storage of data. ✓  
3 Ques-8 what is the process of security of BD? ✓  
2 Ques-9 what is usability in BD. → Applications ✓  
1 Ques-10 Explain ETL (in detail). → Applications ✓

### (1) RDBMS and NoSQL.

- RDBMS :-
- Relational Database Management Systems.
  - Most Popular Database.
  - data is stored in form of rows ~~and~~ tuples
  - there are numbers of tables in it.
  - Data can be easily accessed bcoz it is stored in the tables.
  - This Model was prepared by E.F. Codd.

### → NO SQL :- non-SQL Database.

- doesn't use tables to store the data
- used for storing and fetching the data in database.
- generally used to store large amount of data.

GOOD WRITE

- It supports query language and provides better performance.

~~DBMS~~ (Expensive, need big servers)

### RDBMS

- (1) used to handle data coming in low velocity.

- (2) gives only read scalability.

- (3) manages structured data.

- (4) Data arrives from one or few locations.

- (5) supports complex transactions.

- (6) has single point of failure.

- (7) handles data in less volume.

- (8) Transactions written in one location.

- (9) Support ACID properties.

- (10) difficult to make changes in Database once it is defined.

GOOD WRITE

DATE:  
PAGE:

easy, uses cheap servers

No SQL

- (1) high Velocity.

- (2) gives both read and write scalability.

- (3) manages all type of data.

- (4) many locations.

- (5) simple transactions.

- (6) No single point of failure.

- (7) high volume.

- (8) many locations.

- (9) doesn't support.

- (10) Enables easy and frequent changes to database.

(11) Schema is mandatory to store the data.

(12) Deployed in Vertical fashion.

(11) Schema design is not required.

(12) Deployed in Horizontal fashion.

Q2

## Big Data Processing Tools :-

→ Big Data Tool is a software that extracts info from various complex data types and sets and then processes them to provide meaningful insights/insights. Traditional databases cannot process huge data hence best big data tools that manage big data easily are used by businesses.

→ Tools are :-

### Apache Hadoop:

- open source software
- processes data sets of big data with the help of MapReduce programming model.
- written in Java.
- provides cross-platform support.
- Most popular big data Tool.
- used by IBM, Amazon Web Services, Microsoft and Facebook etc.

Pros :-

- offers flexibility and faster data processing.
- Highly scalable, provides fast access to data.
- HDFS (Hadoop Distributed File System), which is its core strength, has ability to hold all data types like video, Images, XML, JSON and ~~text~~ plain text over the same file system.

- Cons +
- Gives disk space problems sometimes.

Pricing: free to use under Apache Licence.

### (a) Apache Cassandra

- open source (free of Cost).
- it employs CQL (Cassandra Structure language) to interact with database.
- used by Fortune 500 Companies including Facebook, Honeywell, Yahoo, Accenture and American Express etc.

Pros +

- Handles huge data volume very fast without any single point of failure.
- comes with simple ring architecture and log-structured storage.
- offers automated replication.

Cons +

- Clustering needs improvement.
- Troubleshooting and maintenance are not easy.

Pricing: free of Cost.

### (b) MongoDB

- free open-source.
- supports multiple operating systems, including windows Vista(), OS X, Linux, Solaris and freeBSD.
- it is a Document-oriented database.
- is written in C, C++, and Java Script.
- used by brands like Facebook, ebay, Google.

- Pros
- Reliable, low-cost, easy to learn
  - Smooth installation and maintenance.

Cons

- Has limited analytics and it's slow for certain use cases.

Pricing: Pricing for enterprise and SMB Versions is available on request.

#### (4) Tableau

- Famous Software with three diff. options:
  - A) Tableau Desktop (for analyst).
  - B) Tableau Server (for enterprise use).
  - C) Tableau Online (cloud-based)Tableau Reader and Tableau Public are recent additions.
- used mainly for data visualization, exploration and understanding.
- can handle all sizes of data and quite easy to handle.

- Pros
- offers huge flexibility to create various types of visualizations as desired.
  - provides superb data blending.
  - Mobile friendly.

Cons

- Scope of Improvement areas include formatting controls and built-in tools for deployment and migration b/w other tableau environments servers.

Pricing: It is not free and Pricing starts from \$35/month.

#### (5) Datawrapper

open source

- helps in Data Visualization and prepares precise and simpler charts fast.
- used by many big brands like Twitter, Bloomberg, The times and Fortune to name a few.

### Pros

- Device independent and works well on all types of devices - mobile, tablet or desktop.
- Very fast, fully responsive, interactive.
- No Coding required.

### Cons

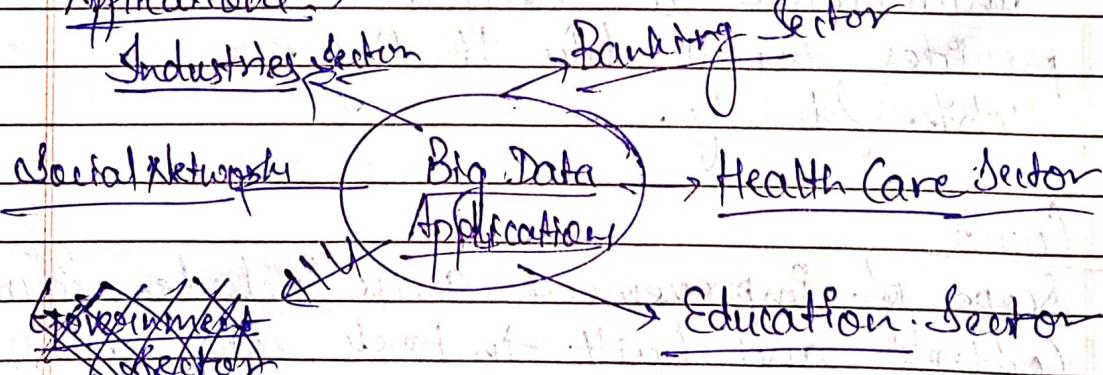
- Nothing versatile. Just offers a limited color palette.

Price Offer both free and pricing models.

③

### Application and characteristics of Big Data

#### Applications



④

### Banking Sector

- Banks are creating a Digital Database in which all data related to deposits or withdrawals, transactions and various customer service records is stored.

GOOD WRITE

By analyzing the whole data they can enhance cybersecurity and can provide innovative and personal offers to each individual.

(1) Big Data use Big data to Track Customer spending habit, and shopping behavior to provide recommendations.

(2) Smart traffic system + Data about condition of traffic of diff. road, collected through Camera kept beside road, GPS device placed in vehicle (ola, uber, cab etc). All these data are analyzed and jam-free or less jam way, less time-taking ways are recommended.

(3) Virtual Personal Assistant Tool + Big Data helps Virtual Personal Assistant Tools (Siri, Google Assistant etc) to provide answer to various questions asked by user.

(4) Education Sector + Online Educational Course Conducting organization utilize big data to research Candidate, interested for that course.

(5) Media and Entertainment Sector + Netflix, Spotify, Amazon prime do analysis on data collected from their user. to set the next business strategy.

(6) Healthcare Sector + Big data is used to improve services and better understand the overall health situation to solve pt quickly.

(7) Industries + Big Data helps Companies enhance the quality of and efficiency of their various products.

helps companies make more flexible business decisions and quickly solve problems.

### (9) Social Networks →

- By this data Companies can predict the customers demand before they request it.

### (10) Government sector →

→ Characteristics of Big Data and 6 V's of Big Data (One - S).

- Volume, Velocity, Variety, Veracity, Value, Variability.

#### (1) Volume →

- Big Data = enormous data.
- Volume is a huge amount of data.
- To determine volume of data, size of data plays a very crucial role.  
→ if volume of data is very large, then it is actually considered as a 'Big Data'.
- Hence, while dealing with 'Big Data' it is necessary to consider a characteristic 'Volume'.

#### (2) Velocity →

- refers to high speed of accumulation of data.
- In Big Data, Data flows in from sources like machines, networks, social media, mobile phones etc.
- There is massive and continuous flow of data. This determines the potential of data that how fast the data is.

generated and procured to meet the demands.

- Sampling data can help in dealing with the issue like 'Velocity'.

### (3) Variety

- refers to nature of data that is structured, semi-structured and unstructured data.
- also refers to heterogeneous sources.
- it is basically the arrival of data from new sources that are both inside and outside of an enterprise.
  - Structured - organized data.
  - Semi-structured - semi-organized data.  
Ex - log files.
  - Unstructured - unorganized data.  
Ex - Text, Images, video etc.

### (4) Veracity

- refers to inconsistencies and uncertainty in data, that is data which is available can sometimes get messy and quality and accuracy are difficult to control.
- Big Data is also vulnerable becoz of the multitude of data dimensions resulting from multiple disparate data types and sources.
- Ex - Data in bulk can create confusion whereas less amount of data could convey half or incomplete info.

### (5) Value

- The bulk of Data having no value is of no use to company, unless you turn it into something useful.
- Data in itself is of no use or importance but it needs to be converted into something valuable to extract information.
- Hence it is the most important V of all 6 V's.

## 6. Variability (tendency to change or shift)

- How fast or available data
- How fast your data is changing or shifting.
- How often the meaning or shape of your data changes.
- Ex: If you are eating James Pre-cream daily and the taste just keep changing.

## Ques-4: Types of Data

### Structured Data

### Unstructured data

### Semi-structured data

(1) Well-organized data.	(2) Not organized at all.	(3) Partially organized.
(1) less flexible and difficult to scale. if it is schema dependent.	(2) is flexible and scalable. it is schema independent.	(3) is more flexible than unstructured data, but lesser than unstructured data.
(3) based on RDBMS.	(3) based on characters and binary data.	(3) based on XML and RDF.
(4) Easy analysis.	(4) Difficult analysis.	(4) Difficult analysis compared to structured data but easier when compared to unstructured data.
(5) Ex: financial data, bank codes.	(5) Ex: Media logs, videos, audios, images.	(5) Ex: Tweets organized by hashtags, folders organized by topics.

## → Reasons for Creating Data Mart

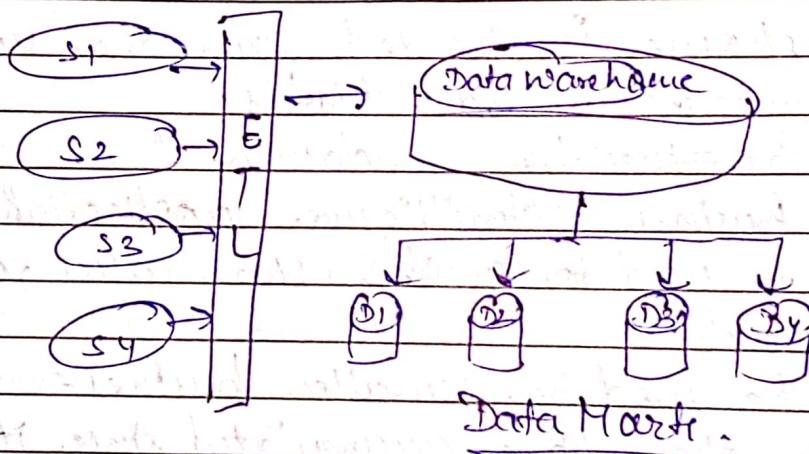
- Easy to Create
- Creates collective data by group of users
- Easy access to frequently needed data
- Improves end-user response time
- Lower cost than data warehouse
- Contains only essential business data

## → Types of Data Mart

- Dependent Data Mart
- Independent Data Mart

### (1) Dependent Data Mart

External Source

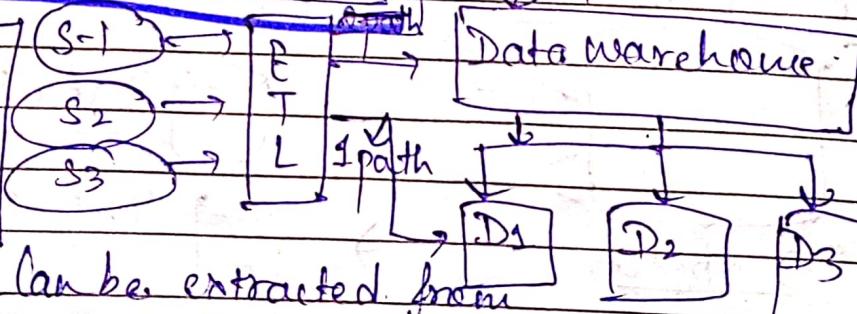


- This model is used by big organizations
- it is created in Top-down approach

### (2)

### Hybrid Data Mart

E path reflects  
dependent data  
model,

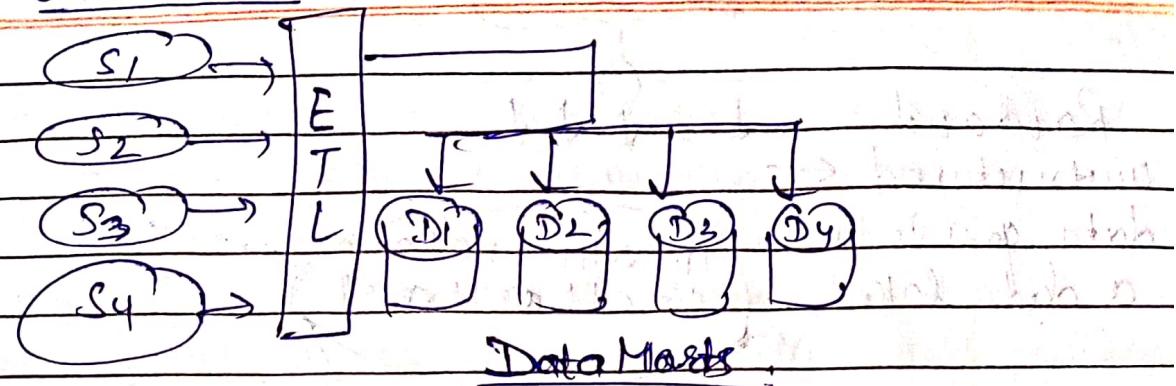


- data can be extracted from operational source or from data warehouse, from external source.
- I path reflects accessing data directly

## Independent data Mart

DATE: / /  
PAGE

External source



- Used by small organizations.

- It is Cost effective.

## Data Lakes

- Stores large volumes of unstructured, semi-structured and semi-structured data in its native format.
- Design approach is bottom-up.
- Data lake captures anything the organization deems valuable for future use.
- Data can be images, videos, PDFs, anything.
- Can be used for data analytics and report creation.
- Technology used in Data lake is much more complex than in a data warehouse.
- Because of level of complexity and skill required to leverage, a data lake requires users who are experienced in programming languages and data science techniques.

## Characteristics of Data Lake

- Collects all data from many disparate data sources over an extended period.
- Meets the needs of various users in organization.
- If it is uploaded without an established methodology.
- Processes and cleans data and stores it in the data lake.

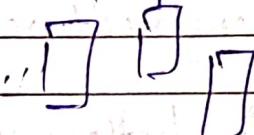
GOOD WRITING

Raw and unstructured data goes into a data lake

Positional  
Binary data

0 0 0 0 0 0 0 0  
0 0 1 1 1 1 0 0 1  
0 0 0 0 0 0 0 0 0  
0 0 0 0 0 0 0 0 0  
0 0 0 0 0 0 0 0 0  
0 0 0 0 0 0 0 0 0  
0 0 0 0 0 0 0 0 0  
0 0 0 0 0 0 0 0 0

Data is selected and organized as needed.



### → Data pipelines

- Can work with any type of data.
- It is a means of moving data from one place (the source) to a destination (such as data warehouse).

Along the way data is transformed and optimized, partitioned in a state that can be analyzed and used to develop business insights.

- It is a step involved in aggregating, organizing and moving data.

### → Elements

Data pipeline consists of three essential elements

- (1) Sources - Sources are where data comes from.

GOOD WRITE

- Common sources include RDBMS systems like MySQL, CRMs such as Salesforce and HubSpot.

(2) processing steps → transformation, augmentation, filtering, grouping, and aggregation.

(3) Destination → where data arrives at the end of the processing, typically a data lake or data warehouse for analysis.

→ characteristics of (to look when considering a data pipeline include)

- continuous and extensible data processing.
- The elasticity and agility of cloud.
- Isolated and independent resources for data processing.
- High availability and disaster recovery.

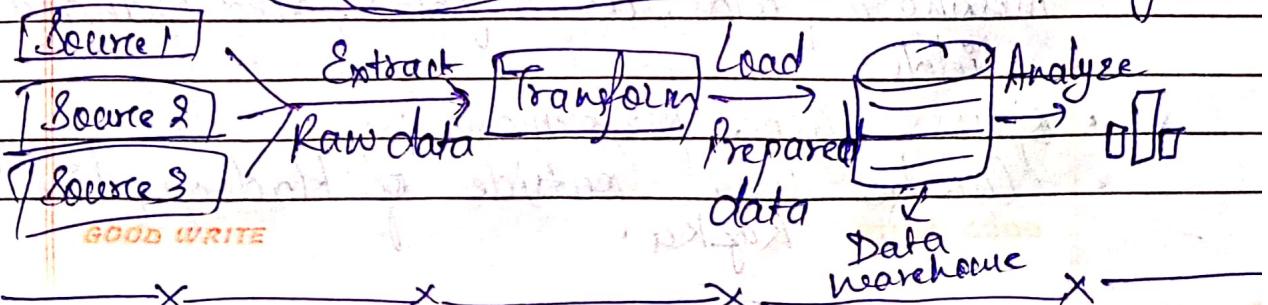
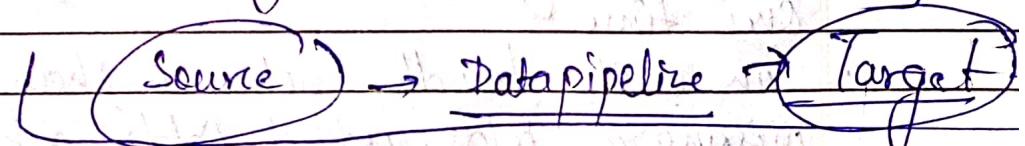
→ challenges to building Data Pipelines

- Netflix has built its own data pipeline. However, Building your own pipeline is very difficult and time consuming.

challenges are

- Connection.
- Flexibility.
- Centralization → Ex: McDonald's (exact no. of pickles is put on each burger no matter where you are in the world).
- Latency.

(delay)



## Ques. Diff. file formats of Big Data

### (1) CSV →

- Comma-separated values.
- here data is stored in row-based file format
- mostly used for exchanging tabular data.
- easier to use
- human-readable.
- widely used format for tabular data representation, but lacks many other capabilities which other formats provides.
- very slow.

### (2) PARQUET →

- open-source file format for Hadoop.
- helps to achieve efficient storage and performance.
- is column-oriented format.
- it is default format for spark, Spark.
- widely used in Hadoop ecosystem for querying data.

### (3) AVRO →

- row-based storage format.
- widely used for serialization.
- supports dynamic data schemas that change over time.
- can easily handle schema changes like missing fields, added fields, and changed fields.
- offers good performance.
- can be used outside of Hadoop, like in Kafka.

#### (4) ORC

- optimized Row Columnar.
- provides highly efficient way to store data.
- were designed to overcome the limitations of other file formats.
- improves overall performance when Hive (SQL kind of interface, built on hadoop) reads, writes, and processes the data.
- stores Collection of Rows.
- Mostly used in Hive.

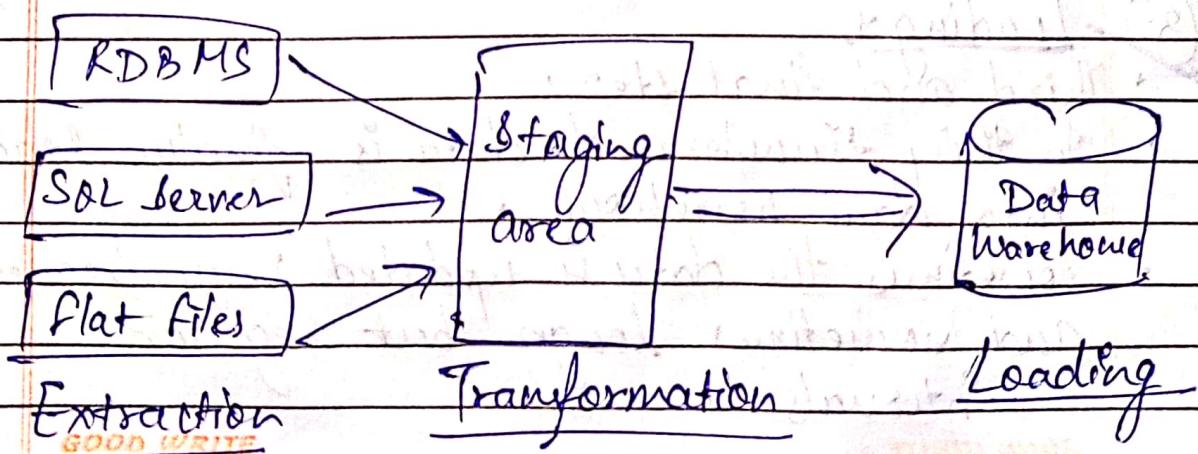
#### (5) Feather

- is a portable file format.
- used for storing Arrow tables or data frames.
- ~~Fast~~ is fast.
- Lightweight and easy to use binary file format for
- ~~easy to use~~ storing data frames.
- Files are generally faster in read and write performance when used with solid state drives, due to its simpler compression scheme.

### Ques. Explain ETL in detail

→ ETL stands for Extract, Transform, and Load.

- It is used to extract data from various sources, transform it into a format suitable for loading into a data warehouse, and then load it into the warehouse.



## (1) Extraction

- First step of ETL process.
  - In this data is extracted from various source systems which can be in various formats like Relational databases, NoSQL, XML and flat files in Staging Area.
  - It is important to extract data from various systems in Staging area because data can be corrupted so, directly loading it in data warehouse may damage it.
- Therefore, it is most important step of ETL.

## (2) Transformation

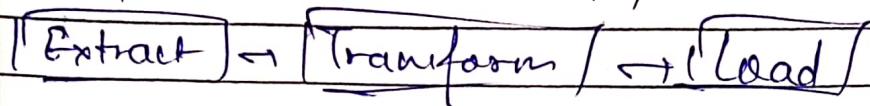
- Second step of ETL.
- In this, a set of rules or functions are applied on the extracted data to convert it into a single standard format.
- May involve following processes/steps:
  - Filtering
  - Cleaning
  - Joining (joining multiple attributes)
  - Splitting (splitting a single attribute into multiple attributes)
  - Sorting (sorting tuples on the basis of some attribute (generally key attribute))

## (3) Loading

- Third and final step.
- In this, transformed data is finally loaded into data warehouse.
- Sometimes the data is updated very frequently and sometimes longer but regular intervals.

- The state and period of loading solely depends on the requirements and varies from system to system.

→ ETL process can also use pipeline concept &



→ ETL Tools

- Hive, Sybase, Oracle Warehouse builder, Clover ETL, and MarkLogic.

→ Advantages

- Improved data Quality.
- Better data Integration.
- Increased data Security.
- Improved Scalability.
- Increased Automation.

Geeky.

→ Disadvantages

- High Cost.
- Complexity.
- Limited Flexibility.
- Limited Scalability.
- Data privacy concerns.

Geeky.