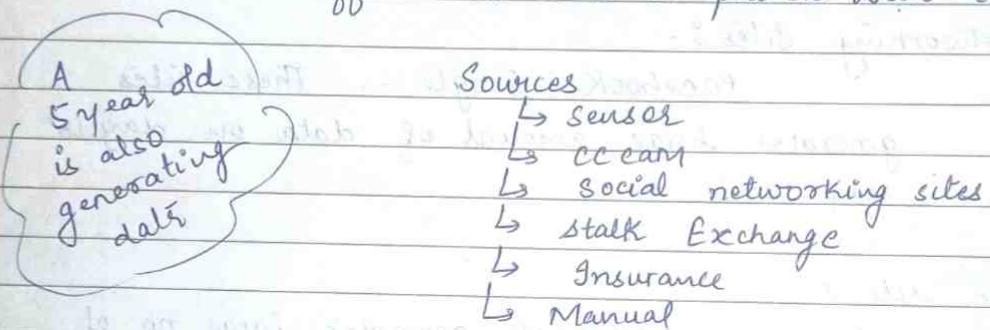


BIG DATA

UNIT - I

- B → Business
 - I → Intelligence
 - G → Gathering of data
- data having intelligence for growth of business.

Ex Email - Beyond 25MB data will be Big data which is difficult to store or process with existing system.



If data is generating rapidly, we need to take care of it.

What is Big Data?

↳ Data which is large in size.

↳ Data represent quantitative or qualitative values that are variable and usually come in tables, with rows and columns, in tree like format such as a set of nodes with parent children relationships or a graph with a set of interconnected nodes.

↳ It is a collection of large datasets that can't be processed using traditional and management tools efficiently.

Ex

- Social Media
- NYSE
- Google Maps

- Medical Industry
- Transportation

→ "Big data" is high volume, velocity and variety information assets that demand cost effective, innovative forms of information processing for enhanced insight and decision making.

Sources of Big Data

1. Social Networking Sites :-

Facebook, Google ... These sites generates huge amount of data on day to day basis.

2. E-commerce sites :-

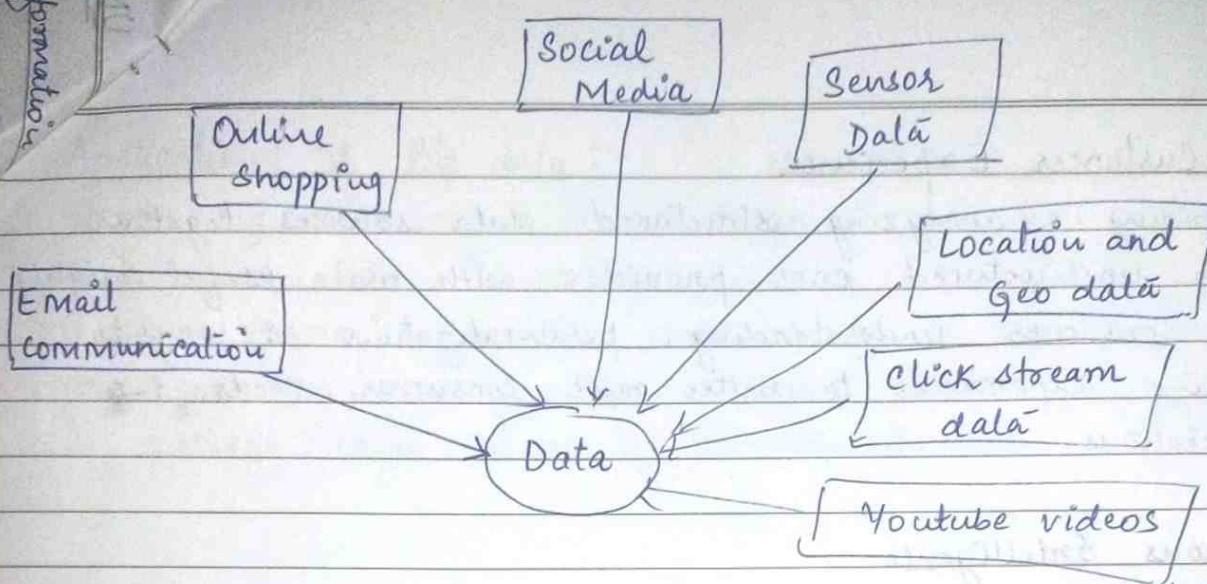
Amazon, Flipkart - generates large no. of logs from which user's buying trends can be traced.

3. Weather Station :- Gives huge data that is stored and manipulated for forecast weather.

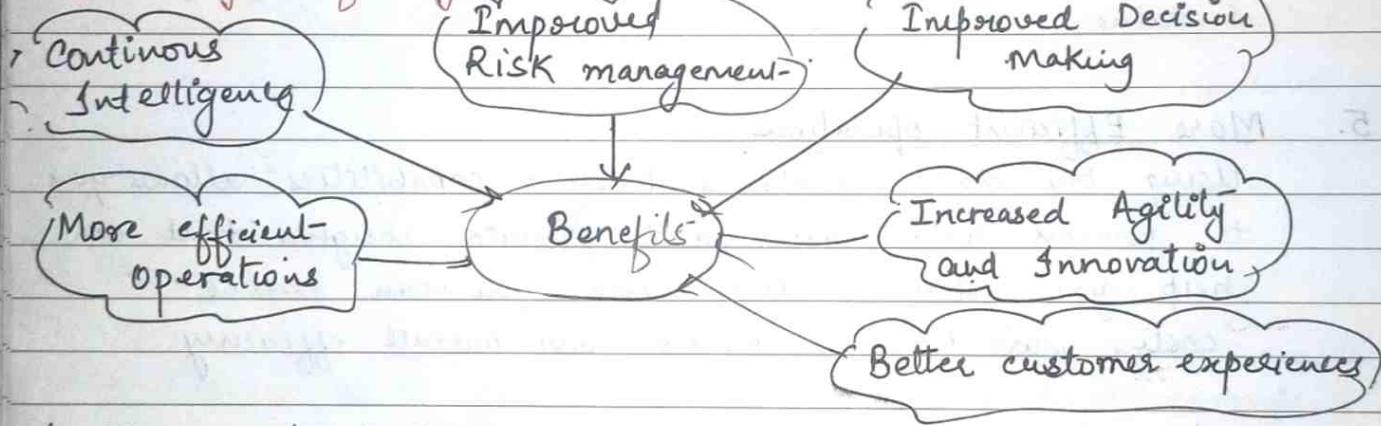
4. Telecom company :- Airtel, Vodafone study the user trends & publish their plan.

5. Share Market :-

Stock Exchange across the world generates huge amount of data through its daily transaction.



Advantages of Big Data :-



1. Improved decision making :-

↳ When we manage and analyze big data, we discover patterns and unlock insights that improve & drive better operational and strategic decisions.

2. Increased Agility & Innovation :-

Big data allows to collect & process real-time data points & analyze them to adapt quickly & gain a competitive advantage.

These insights can guide & accelerate the planning, production & launch of new product, features & updates.

3. Better Customer Experiences

↳ combining & analyzing structured data sources together with unstructured ones provides more useful insights for consumer understanding, personalization & ways to optimize experiences to better meet consumer needs & expectations.

4. Continuous Intelligence

↳ Allows to integrate automated, real-time data streaming with advanced data analytics to continuously collect data, find new insights, discover new opportunities for growth & value.

5. More Efficient operations

Using big data analytic tools & capabilities allows you to process data faster and generate insights that can help you determine areas where you can reduce costs, save time, & increase your overall efficiency.

6. Improved Risk Management

Analyzing vast amount of data helps companies evaluate risk better - making it easier to identify & monitor all potential threats & report insights that lead to more robust control & mitigation strategies.

Applications of Big data :-

Term Big data refers to as large amount of complex and unprocessed data. Now a days company's use Big Data to make business more informative and allows to take business decisions by enabling data scientists and other professionals to analyze large volume of transactional data.

Applications
of
Big Data

Travel & Tourism

1. Travel & Tourism :-

- ↳ Enables us to forecast travel facilities requirements at multiple locations, improve business through dynamic pricing & many more.

2. Financial & Banking sector :-

- ↳ These people use big data technology extensively.
- ↳ Help banks and customer behaviour on the basis of investment patterns, shopping trends, motivation to invest and inputs that are obtained from personal or financial backgrounds.

3. Health care :-

- ↳ Big data has started making a massive difference in healthcare sector, with the help of predictive analytics, medical professionals & health care personnel.
- ↳ It can also produce personalized healthcare & solo patients also.

4. Telecommunication Media :-

- ↳ Telecommunications and the multimedia sector are the main users of Big Data.
- ↳ There are zettabytes to be generated every day and handling large scale data that require big data technologies.

5. Government & Military :-

- ↳ Government & military also used technology at high rates.
- ↳ In military, a fighter plane requires to process petabytes of data.
- ↳ Government agencies use Big Data & run many agencies managing utilities, dealing with traffic jams and the effect of crime like hacking and online fraud.
- ↳ Aadhar card - Government has a record of 1.21 billion citizens. This vast data is analyzed & stored to find things like no. of youth in the country.
- ↳ Big data cannot store in a traditional database, so it stores & analyze data using Big data analytics.

6 E-commerce :-

- ↳ Application of big data.
- ↳ Maintains relationships with customers that is essential for ecommerce industry.
- ↳ E-commerce websites have many marketing ideas to retail merchandise customers, manage transactions, & implement better strategies of innovative ideas to improve businesses with Big data.

7 Social Media :-

- ↳ Largest data generator.
- ↳ statistics have shown that around 500+ terabytes of fresh data generated from social media daily, particularly on Facebook.
- ↳ Single activity on social media site generates many stored data & gets processed when required.
- ↳ Data stored is in terabytes it takes a lot of time for processing. Big data is a solution to the problem.

8 Academia :-

- ↳ Numerous online educational courses to learn from.
- ↳ Academic institutions are investing in digital courses powered by Big Data technologies to aid all-round development of budding learners.

Big Data Analysis Tools and Software :-

↳ Tools that are used to store and analyze a large no. of data sets and processing these complex data are known as big data tools.

1. Xplenty

- ↳ Data analytic tool for building a data pipeline by using minimal codes in it.
- ↳ Offers wide range of solutions for sales, marketing, and support.
- ↳ with the help of its interactive graphical interface, it provides solutions for ETL, ELT.
- ↳ Best part is its low investment in hardware & software & it offers support via email, chat, telephonic ---.
- ↳ Platform to process data for analysis over cloud and segregates all data together.

Features

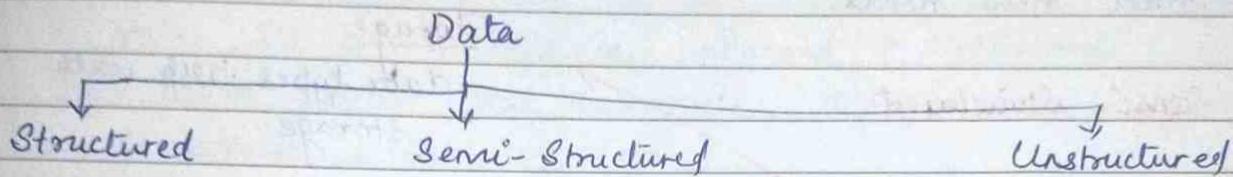
• Rest API : User can possibly do anything by implementing API.

• Flexibility : data can be sent, and pulled to databases, warehouses & salesforce.

• Data Security : offers SSL/TSL encryption .

• Deployment : offers integration apps for both cloud & in-house

Types of Big Data :-



1. Structured Data :-

- ↳ data that resides in a fixed field within a record.
- ↳ type of data most familiar to our everyday lives.
- ↳ Structured data is also called relational data.
- ↳ data is split into multiple tables to enhance integrity of data by creating a single record to depict an entity.

Structured data

→ Databases such as Oracle, DB2, MySQL etc.

→ Spreadsheets

→ Online transaction processing systems.

- ↳ Structured data is easy to enter, query & analyze.
- ↳ All of the data follows same format.
- ↳ Business data of an ecommerce website can be considered to be structured data.

Cons of Structured data :-

- ↳ Structured data can only be leveraged in cases of predefined functionalities.
- ↳ data is stored in data warehouse with rigid constraints & a defined schema. Any change in requirement

would mean updating all of that structured data to meet new needs.

Advantages of Structured Data

Storage

Data types help with storage

Ease with structured data

Scalability

Not an issue with increase in data

Update and delete

Updation & deletion is easy due to structured form.

Security

Data is secured

Retrieve information :- well defined structure helps in easy retrieval of data

Ease with structured data retrieval

Indexing & searching :- data can be indexed not only on a text string but other attributes as well

Mining data

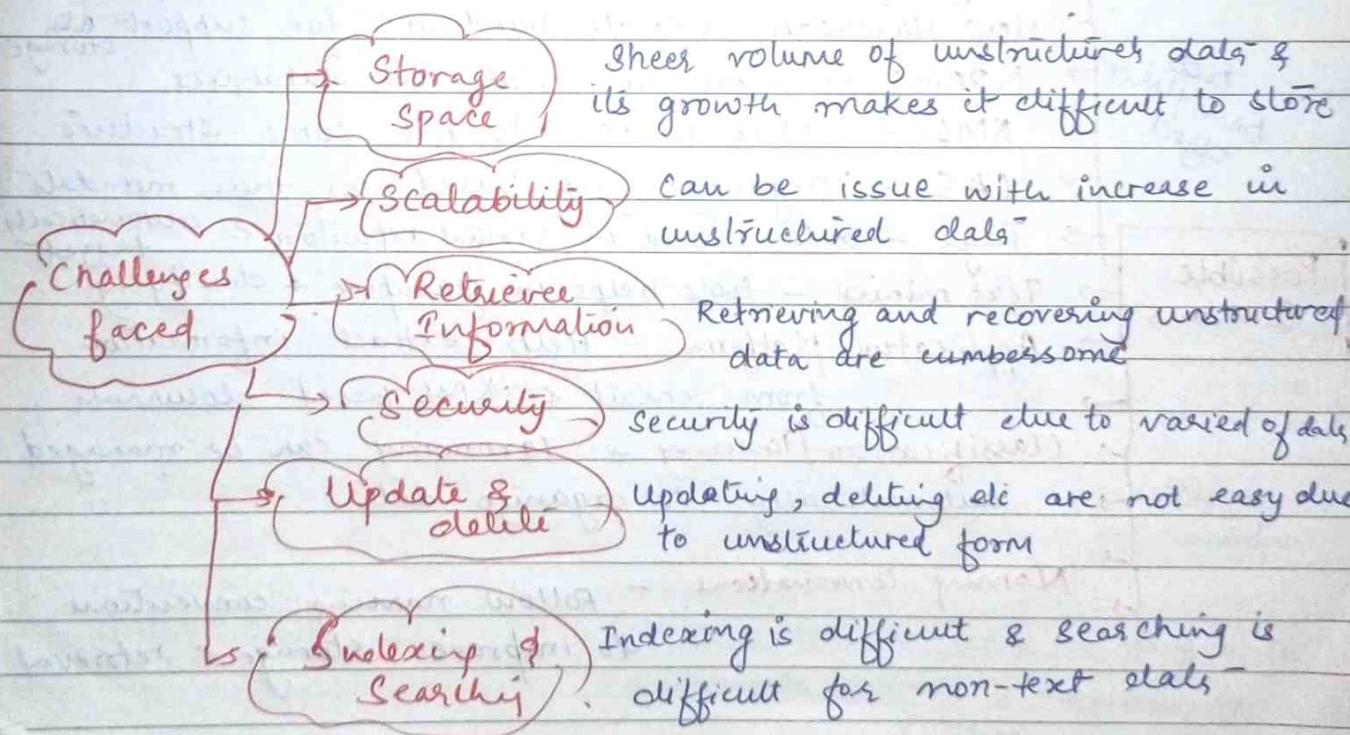
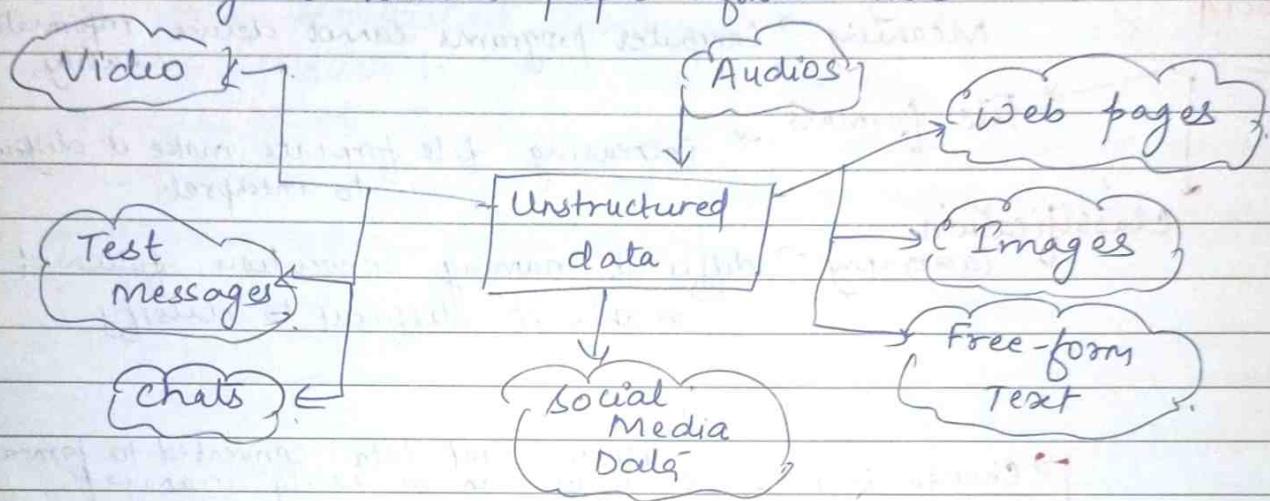
:- data can be easily mined & knowledge can be extracted

BI operations

:- BI works well with structured data

Unstructured data :-

- ↳ data doesn't adhere to any definite schema or set of rules. Its arrangement is unplanned & haphazard.
- ↳ Photos, videos, text documents, log files can be generally considered unstructured data.
- ↳ also known as "dark data" because it cannot be analyzed without proper software tools.



Other challenges

- Interpretation - data is not easily interpreted
- Tags - data grows, it is difficult to tag.
- Indexing - Algorithm construction is difficult for indexing.
- Deriving meaning → Computer programs cannot derive information easily.
- File formats → increasing file formats make it difficult to interpret.
- classification
Taxonomy → different naming convention followed makes it difficult to classify.

Big
Binary
large
object

Possible
Solutions

- Change formats → Unstructured data converted to formats which can be easily managed
- New hardware → Create hardware for support ~~size~~ storage
- RDBMS → store in relational databases.
- XML - store in XML to give some structure
- CAS → organize files based on their metadata
- Tags - data stored in virtual repositories & automatically tagged
- Text mining - tools helps in grouping & classifying.
- Application platforms - help extract information from email & XML based documents.
- Classification/Taxonomy → Taxonomy can be managed automatically to organize data.
- Naming conventions - Follow naming conventions to improve storage & retrieval.

Semistructured -

- ↳ data is not bound by any rigid schema for data storage & handling.
- ↳ data is not in relational format & is not neatly organized into rows & columns
- ↳ Doesn't need SQL rather called as NOSQL data.
- ↳ data serialization language is used to exchange semi-structured data across systems.
- ↳ often used to store metadata about business process.
- ↳ information typically comes from external sources such as social media platforms.

Audios → E-mails

Zip files ← Semistructured data → XML

JSON

other Markup languages

Challenges with Semi Structured Data

→ Storage Cost : Storing data with schemas ↑ cost.

→ RDBMS : Semistructured data cannot be stored in tables.

→ Irregular & Partial Structure - Some elements have extra info. while others not

→ Implicit Structure - In many cases, structure is implicit. Interpreting relationships is difficult.

→ Evolving Schemas - Schemas keep changing with requirements making it difficult to capture in database.

Challenges

- Distinction between schema & data - Vague distinction because schema & data exists at times making it difficult to capture data.
- Flat files - data stored in flat files difficult to search.
- Incomplete / irregular - Extracting structure when there is none & interpreting relations existing in structure is difficult task.

Properties

- XML - Allows to define tags & attributes to store data & store data in hierarchically nested structure.
- RDBMS - data stored in relational database by mapping data to relational schema which is then mapped to table.
- Special purpose DBMS - databases specifically designed to store semi-structured data
- OEM - data can be stored & exchanged in form of graph where entities are represented as objects which are vertices in graph.
- Indexing - Indexing data in graph enables quick search

Possible Solutions

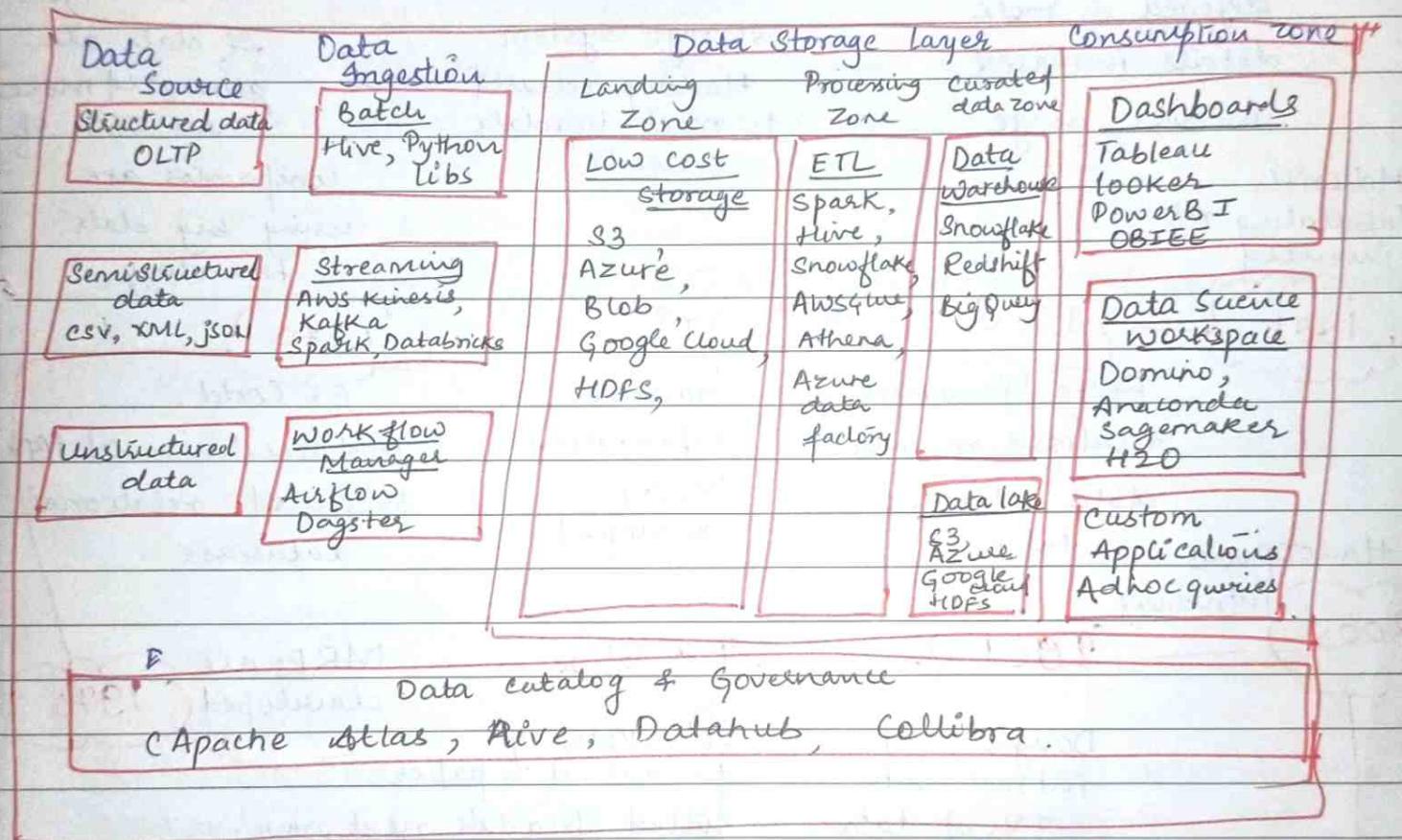
sufficient
at
between

difference between structured, semi-structured & unstructured

Properties	Structured Data	Semi Structured Data	Unstructured Data
Technology	Based on relational database table	Based on XML	Based on characters and binary data
Transaction management	Matured transaction management and various concurrency technique	Transaction is adapted from DBMS not matured	No transaction management and no concurrency
Version management	versioning over tuples, row, tables	versioning over tuples or graph is possible	versioned as a whole.
flexibility	Schema dependent and less flexible	More flexible than structured data but less flexible than unstructured	More flexible and there is absence of schema.
Scalability	Difficult to scale DB schema	Scaling is simpler than unstructured data	More scalable
Robust	very robust	New technology, not spread	
Query Performance	Structured query allow complex joining	Queries over anonymous nodes are possible	only textual queries are possible

Big Data Platform :-

- ↳ Integrated computing solution that combines numerous software systems, tools & hardware for big data management.
- ↳ One stop architecture that solves all data needs of a business regardless of volume & size of data at hand.
- ↳ Enterprises are increasingly adopting big data platforms to gather tons of data & convert it into structured.



Platforms

Hadoop Delta Lake Migration platform

↳ Open source platform by Apache Foundation

↳ used to manage & store large data set at a low cost & great efficiency

History of Big data & Innovation

Term big data was coined and used by Oxford dictionary

Introduction of Big data vs

MongoDB, Atlas
MongoDB, database service released

1980

1990

2001

2006

2016

2022

Big data defined in more details increased internet usage

Distributed

storage system

Hadoop released to meet big data needs

79 zetta bytes of data are generated more than 53% of

companies are using big data technology.

Hollerith Tabulating m/c invented

1881

1928

Fritz Pfleumer developed magnetic data storage on tape.

1948

Shanon's Information Theory developed

1970

EF Codd, mathematician at IBM presented relational database.

Hadoop was introduced

2005

2001

Dough Terry introduced 3Vs of data

1998

John Mashey present ed a paper titled "Big data...next wave".

1976

(2014)
Companies started moving ERP to cloud

② Data Catalog Platform :-

- ↳ Single self service environment to users to help them with find, understand & trust data source.
- ↳ Helps to discover new data sources
- ↳ use catalog discovery to find which data fits their needs
- ↳ Provides comprehensive, security, data governance & solution to protect data

③ Data Ingestion Platform :-

- ↳ First step for data coming from various sources
- ↳ Data is prioritized & categorized
- ↳ Importing & loading data into system

④ IOT Analytics Platform :-

- ↳ Provides wide range of tools to work upon big data

⑤ Big Data Integration & Management Platform :-

- ↳ Elixir Data provides highly customizable solution for Enterprise.
- ↳ Provides flexibility, security & stability for an Enterprise application to deploy public cloud.

⑥ ETL Data Transformation Platform :-

- ↳ Can be used to build pipelines & even schedule the running of the same for data transformation.

Essential components of Big Data Platform :-

There are so many essential components which are given as:-

① Data Ingestion, Management, ETL, & Warehouse

- ↳ Provides resources for effective data management and data warehousing.

2. Stream Computing :-
↳ Helps compute the streaming data used for real time analytics.
3. Analytics / Machine Learning :-
↳ features for advanced analytics & machine learning.
4. Integration :- integrate big data from any source with ease.
5. Data Governance - provides comprehensive security, data governance, and solutions to protect data.
6. Provides accurate data - Helps business to make right decision by providing accurate info.
7. Scalability - Helps scale application to analyze all climbing data ; offers scalable storage capacity.
8. Price optimization :- provides insight for B2C & B2B enterprises which helps business to optimize the prices.
9. Reduced Latency :- with warehouse, analytics tools & efficient data transformation, reduce data latency & increase throughput.

The common platform :-

1. Apache Hadoop :-

- ↳ Open source programming architecture & server software
- ↳ store and analyze large dataset very fast.
- ↳ uses commodity hardware in a clustered computing environment.

2. Cloudera :-

- ↳ big data platform based on Apache's Hadoop.
- ↳ Handle huge volumes of data
- ↳ Enterprises store over 50 petabytes in this platform data warehouse.
- ↳ its dataflow enables real-time processing.

3. Amazon web Services :-

- ↳ popular as AWS
- ↳ Another Hadoop based big data platform
- ↳ Hosted as cloud environment.
- ↳ Using Amazon EMR, enterprises easily set up & scale other big data platforms like Spark, Aestio.

4. Oracle :-

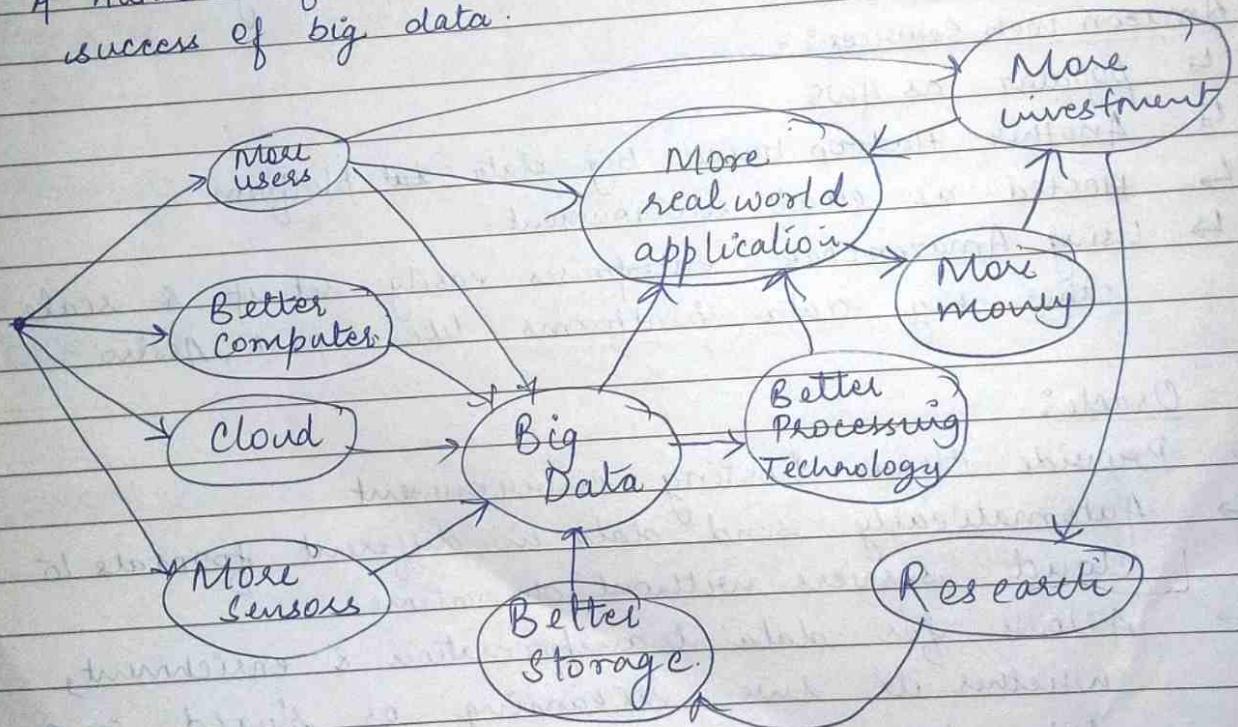
- ↳ Provide Cloud hosting environment
- ↳ Automatically send data in different formats to cloud servers without downtime.
- ↳ Allows for data transformation & enrichment, whether it's live streaming or stored in a data lake.

, performs and opinions of public on a

5. Snowflake :-
↳ data warehouse for storing, processing & analyzing data
↳ top cloud hosting frameworks & integrates with new SQL query engine.

6. Apache Storm :-
↳ Brainchild of Apache software Foundation.
↳ used in real-time analytics & distributed processing.
↳ supports all programming languages.

Drivers of Big Data :-
A number of business drivers are the core of success of big data.



1) Digitization of society :-

↳ Most people spend 4-6 hours per day consuming and generating data through variety of devices & social applications.

2) Plummeting (Reducing) Technology Costs :-

↳ Costs of data storage & processes need reducing, makes possible for small businesses & individuals to become involved with Big data.

3) Connectivity through Cloud computing :-

↳ organizations want to process massive quantities of data do not have to invest in large quantities of IT infrastructure. Instead they can license the storage & processing capacity they need.

4) Increased Knowledge about data Science :-

↳ Knowledge & education about data science has greatly professionalized & more information becomes available every day.

↳ data analysis & statistics is becoming a popular subject among students.

5) Social Media Applications :-

↳ Social media data can be used to identify customer preferences for product development, target new customers, behaviours, preferences and opinions of public on a

scale that has never been known before.

6. Upcoming internet of Things (IoT) :-

↳ increasingly gaining popularity as consumer goods providers start including smart sensors in household appliances.

Characteristics of Big data :-

How fast or available data that extent of data is the structure of data is changing?

Variability
often change the meaning & shape of data

Volume
Huge Amount of data

data is generated on daily basis from various sources. In 2016, data created was 8 zB by 2020, it rises to 40 zB

Variety

different formats of data from various sources

combination of data types that are being dumped into system - emails, PDFs, photos, audios, videos etc.

Veracity
Inconsistencies & uncertainty in data

Big Data

value that can be derived from accessing & analysing big data.

Velocity
High speed of accumulation of data

Value
extract useful data

↳ Mechanism to derive correct meaning of data.

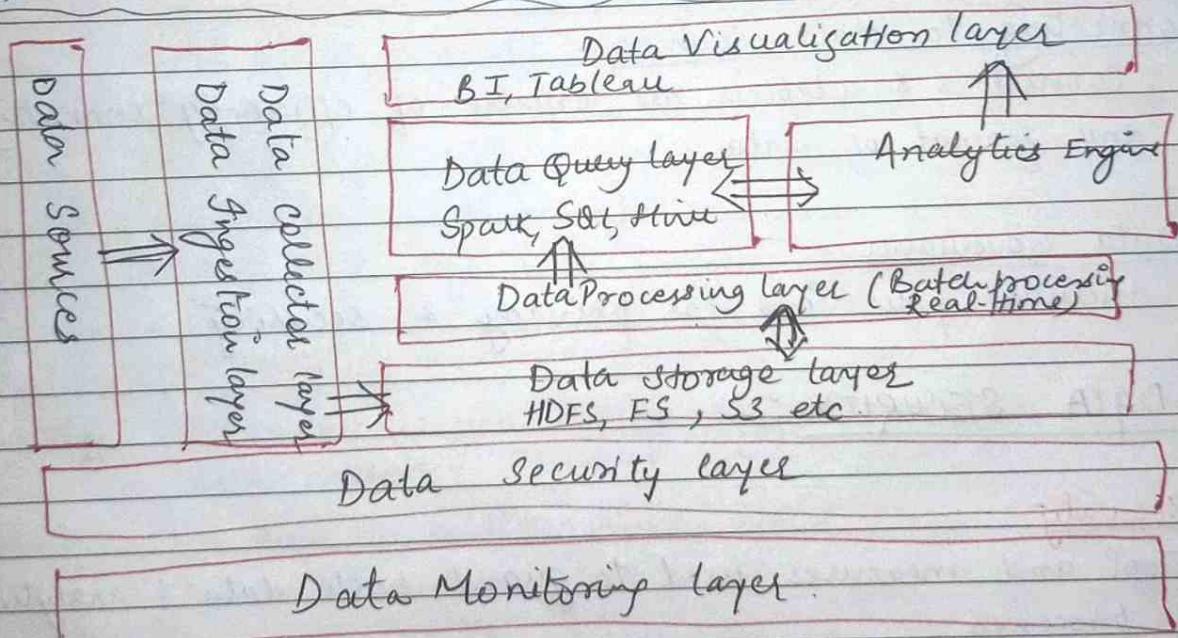
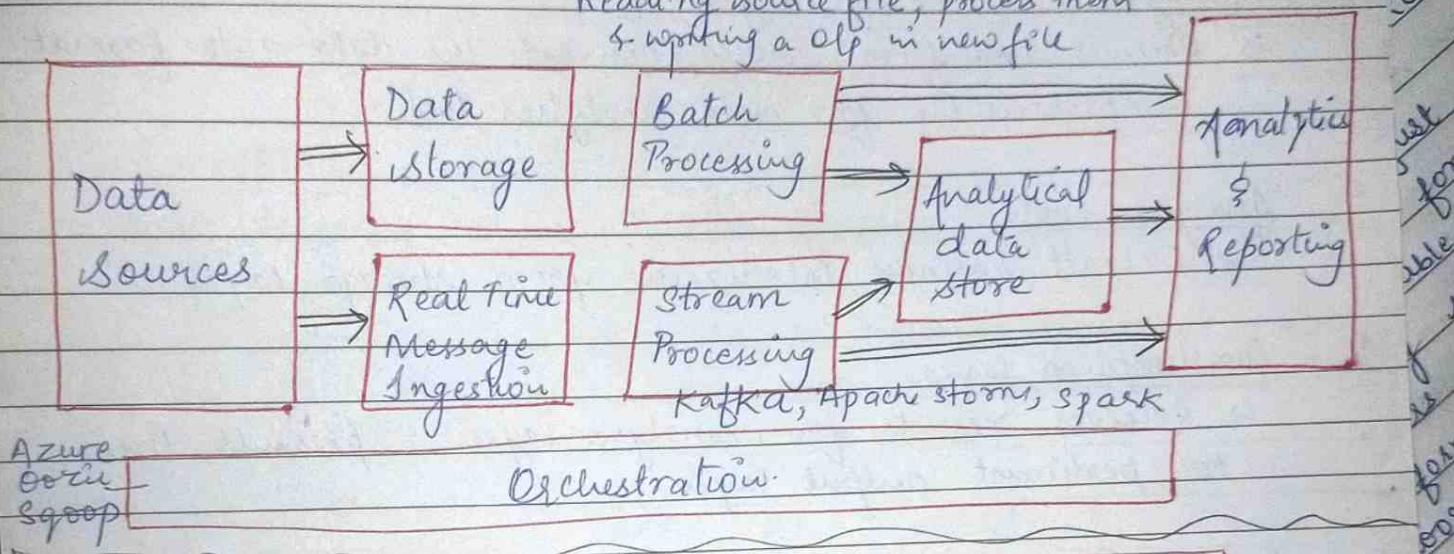
discrepancies found in data

Refers to speed of generation of data. How fast data is generated & processed to meet the demands.

Data Architecture :-

Big data architecture is designed to handle the ingestion, processing, and analysis of data that is too large or complex for traditional database system.

Reading source file, process them & writing a copy in new file



1. Big Data Sources layer - data warehouses, RDBMS, IoT devices
↳ Big data environment can manage both batch processing & real-time processing of big data sources
2. Management & storage layer -
↳ Receives data from source, converts the data into format comprehensible for data analytics tool.
3. Analysis tool
↳ Extract business intelligence from storage layer.
4. Consumption layer
↳ Receives results from analysis layer & presents them to pertinent output layer.
5. Connecting to data sources
↳ Connectors & adapters are capable of efficiently connecting any format of data.
6. Data Governance
↳ Includes provisions for privacy & security

BIG DATA SECURITY :-

Data Security

- ↳ Tools and measures used to guard both data & analytics processes.

Main purpose of big data security is to provide protection against attacks, thefts, & other malicious activities.

Data Security Technologies :-

- ↳ Encryption
- User Access Control
- ↳ Physical Security
- ↳ Centralized Key Management

1. Encryption :- - DES, RSA

- ↳ Secure a massive volume of data, different types of data.
- ↳ Can be user-generated or machine-generated code.
- ↳ Encryption tool along with different analytics tools to format or code the data.

2. User Access Control :- MAC, RBAC (oracle)

- ↳ Most basic security tool.
- ↳ Automated strong user access control is a must for organizations.
- ↳ Automation control manages complex user control levels that protect big data platform against the inside attack.

3. Physical Security :-

- ↳ Built in when you deploy big data platform in your own center.
- ↳ Can also be built around cloud provider's data center security.
- ↳ Video surveillance, security logs.

4. Centralized Key Management

↳ Applied in Big data environments, especially on those having wide geographical distribution.

↳ On-demand key delivery - logging, OTP, PIN

5. Intrusion detection & prevention - IPS behind the firewall.

Some of the companies practicing Big data securities are:-

1. Cloudwick

↳ KDAP (Cloudwick Data Analytics Platform) is a managed security hub that integrates security features from multiple analytics tool sets & machine learning projects.

2. IBM

↳ IBM Security Guardium to monitor Big data and NoSQL environments.

3. Gemalto

↳ Gemalto safenet protects Big data platforms

Big data Security Use Cases :-

Host sensitive data & also monitors Cloud & cloud hosted security infrastructure Monitoring

Network Traffic Analysis

Analyze in & infrastructure out of cloud environment Also identifies attack spaces that are hidden in infrastructure

Insider Threat Detection

Identify threats are dangerous as external threats

Incident - Provide contextualization to alerts Investigation

Use Cases

Threat Hunting

Helps to automate threat hunting

Monitor unusual User Behaviors of employees Behaviour Analysis DEBA

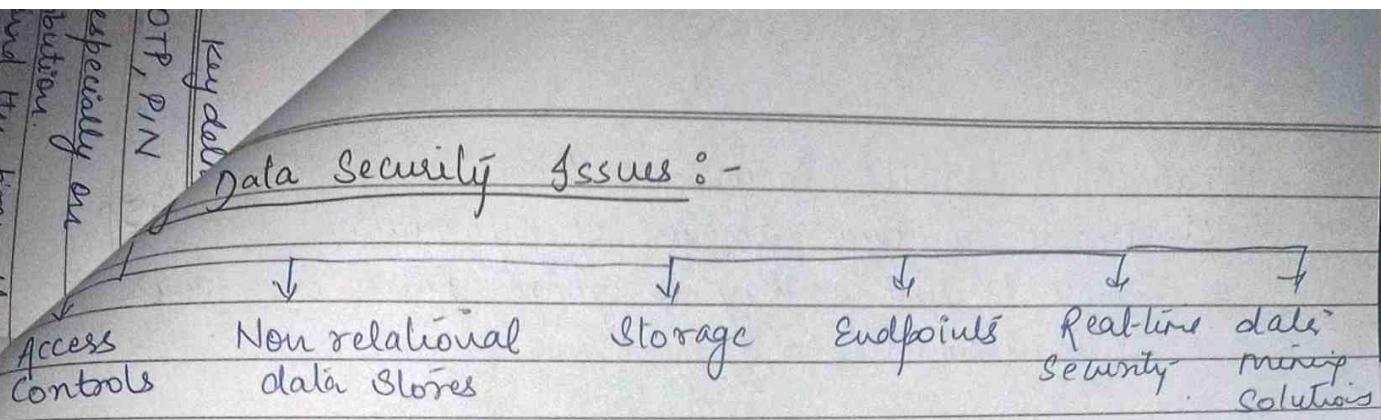
Data Exfiltration

Unauthorized movement of data

Detection

so that detecting malware is easy.

Detect leakage of data



1. Access Controls

↳ Permission to exchange data should be permitted to authenticated users only.

↳ Needs to be such that it would not get hacked by attackers, hackers or by any malicious activities.

2. Non-relational data stores

↳ Like NoSQL usually lack security by themselves

3. Storage

↳ we store data on multiple tiers

↳ Storage depends on business needs in terms of performance & cost.

↳ High-priority data is generally stored on flash media.

4. Endpoints

↳ Security solutions that usually draw logs from endpoints will need to validate the authenticity of those endpoints

5. Real-time Security / Compliance tools :-

- ↳ generate a large amount of information.
- ↳ Key to find a way to ignore false or rough information.

6. Data Mining solutions :-

- ↳ Find a pattern that suggests business strategies.
- ↳ Ensure security from both internal & external threats.

Big Data Privacy & Ethics :-

- ↳ Private customer data & identity should remain private.
- ↳ Share private information should be treated confidentially.
- ↳ Customers should have transparent view.
- ↳ Big data should not interfere with human will.
- ↳ Big data should not institutionalize unfair biases.

Reporting vs Analysis :-