

# Big Data

Session - 1

## Big Data

Collection of large, complex and diverse data sets that are difficult to manage.

Data can be structured, unstructured and semi-structured data that continues to grow exponentially over time.

### Why to use Big Data

① Large No. of Users on internet ~~platforms~~, IOT devices, social media,

### Global data generation

Over ~~1000~~<sup>150</sup> zettabytes expected by 2025

$1 \text{ zettabyte} = 1 \text{ billion Terabytes}$

$$= 10^{21} \text{ bytes}$$

## Categories of Data

- ① Structured Data
- ② Semi-Structured Data
- ③ Unstructured Data

### ① Structured Data

Data which is organized in rows and columns,

~~④~~ Data containing a undefined data type, format, and structure.

ex → CSV files, simple spreadsheets, traditional RDBMS.

This type of data can be easily stored and queried in database.

~~⑤~~ SQL databases like MySQL, PostgreSQL, Oracle.

### Real-World App

- ① Banks having transactions & balances

Q)

## Semi- Structured Data

This type of data does not follow the rigid structure of traditional databases yet it contains some level of organization and can be processed.

flexible Schema

Schema is not strict, predefined

More human readable

JSON files → Key-Value pair  
XML

③)

## Unstructured Data

Does not follow any definite Schema or set of rules

→ Photos, videos, text documents, Pdf

→ YouTube Story Videos

## Data Storage Mechanisms

- ① Flat files
- ② Relational Databases
- ③ NoSQL

### Flat files

Data stored in plain text or binary files  
(ex. .txt, .csv)

name, sapid, batch

X	,	123	,	B1
Y	,	345	,	B2

Rows represents records, & column (if present) are separated by delimiter like commas or tabs.

#### Adv-

- ① Easy to create & use
- ② Light weight & portable
- ③ Good for simple app for small datasets or

Disadv.

- ① Not good for large dataset because no indexing
- ② Querying is limited
- ③ No relationships b/w data  
Because of no foreign keys

Ex -~~Book~~

CSV reports for reporting  
Configuration files

Adv.②

## Relational Databases (RDBMS)

Database that store data in a structured format using tables, rows and columns.

~~These~~ These tables are related to each other using keys (primary keys, foreign keys)

Employee Table

ID	Name	Dept
1	Vish	CSE
2	Lam	SOCs

Limit.

Adv.

① Support for complex queries using SQL.

② Data integrity & consistency  
(Data normalization, 1NF, 2NF, 3NF, BCNF)

③ Relationships b/w tables for efficient data organization.

### Limitations

- ① Predefined schema so it is less flexible
- ② Not well suited for unstructured data or semi-structured data.

Can be used in Banking Systems,  
E-commerce platforms

MySQL, PostgreSQL, Oracle, Microsoft SQL Server

### ③ NoSQL Databases

These are non-relational & designed for handling large volumes of unstructured or semi-structured data.

Does not rely on a fixed schema & are designed to scale horizontally.

#### Types of NoSQL

- ① Document-based → MongoDB
- ② Key-value pair → Redis, DynamoDB
- ③ Column-based → Cassandra, HBase
- ④ Graph-based → Neo4j

Adv

- ① flexible schema
- ② Horizontal scalability for large scale
- ③ handle semi-struct. or unstr. data

use

- ① Social Media app

Comparison

<u>Storage</u>	Flat files Simple	RDBMS Predefined schema	NoSQL Flexible schema
<u>Data relationships</u>	None	Strong (via keys)	Limited or schema-less
<u>Scalability</u>	Limited	Vertical scaling	Horizontal scaling
<u>Users</u>	Small dataset	Structured data	Unstructured data

Characteristics of Big DataThe 5 V's of Big Data

- ① Volume
- ② Velocity
- ③ Variety
- ④ Value
- ⑤ Veracity

Volume :

The enormous size of data generated & stored

Storage: cloud storage like Amazon S3

Processing: Tools like Hadoop & Spark help process large datasets

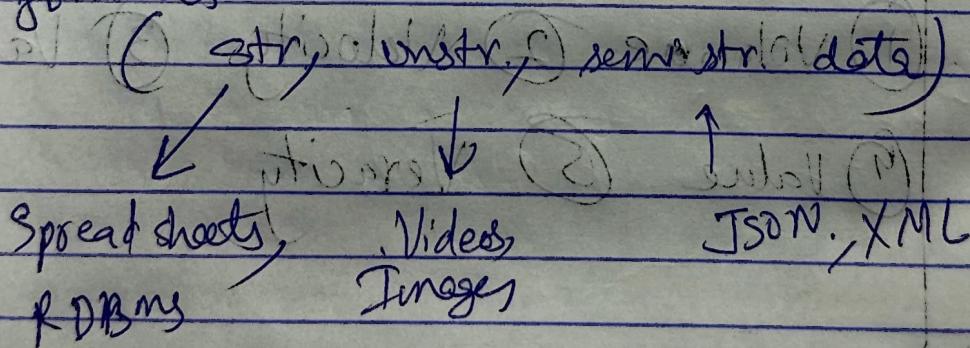
## ② Velocity : Speed of Data Generation & processing

The rate at which data is generated, collected & processed in real-time

- + Real-time processing requires tools like Apache Kafka or Flink
- + Network bandwidth for high speed transmission

## ③ Variety : Diversity of Data formats

Big data can store data in various formats



## ④ Value : Insights Derived from data

Importance or worth of data in driving decisions & generating business value.

## Benefits of Big Data

Data mining

Data mining

Risks of Big data

- ① Privacy concerns
- ② Security risks
- ③ Data quality issues
- ④ High costs of infrastructure
- ⑤ Ethical & legal challenges
- ⑥ Complexity & skills gap

# Big Data Analytics

The ability to analyze and make sense of the data is Big Data Analytics.

It is a key driver of innovation, efficiency and competitive advantage.

Organizations that harness Big Data effectively can make more informed decisions, predict market trends, personalize customer experiences and uncover new opportunities.

## Importance of Big Data Analytics

### ① Transforming Data into Insights

Big data analytics is the process of examining large and varied datasets to uncover hidden patterns, unknown correlations, market trends & other valuable insights.

### Techniques Used:-

Data Mining: → Exploring datasets to find patterns

Statistical Analysis: Applying statistical methods to understand data.

Predictive Analytics: Using historical data to make predictions about future events

Machine Learning: Algorithms that improve through experience, identifying complex patterns

## ② Competitive Advantage

Organizations that adeptly leverage Big Data Analytics can significantly outperform their competitors

## ③ Enhancing Customer Experience → Personalization

Customer Engagement

## ④ Operational Efficiency

→ Process Optimization

→ Predictive Maintenance

## ⑤ Innovation & New Opportunities

→ Product Development

→ New Business Model

Page No. \_\_\_\_\_  
Date: 11

# Impact of Big Data Analytics Across Industries

## Healthcare

Personalized Medicine: → Using genetic data to tailor treatments

Predictive Analytics + Identifying patients at risk of disease like diabetes or heart conditions.

Resource Allocation:- Optimizing the use of medical staff & equipment.

Electronic Health Records (EHRs) (Electronic Health Records)  
Centralizing patient information for better care coordination.

Case

IBM Watson

Assists doctors in diagnosing diseases by analyzing vast amounts of medical literature.

P

## Finance

### Fraud Detection & Risk Management

- Real-time monitoring → Analyzing transaction to spot anomalies
- Credit Scoring → Incorporating diverse data sources for more accurate assessments.

### Personalized Banking Services

- Customer Segmentation: offering tailored financial products
- Investment Strategies: Algorithmic trading based on market data analysis

Case

JPMorgan Chase: → Uses machine learning for contract review, significantly reducing time and errors.

## Retail and E-commerce

### Customer Behavior Analysis

- Sales forecasting
- Marketing Strategies: → Targeting promotions to specific customer segments.

### Recommendation Systems

- Upselling and Cross-Selling

- Personalized shopping Exp.

Netflix → Recommends shows based on viewing history → 80%

## Manufacturing

Supply Chain Optimization

→ Demand Planning

→ Inventory Management

## Quality Control

→ Defect Detection

→ Process Improvement

## Public Sector

→ Smart Cities

→ Traffic Management

→ Energy Efficiency

→ Public Safety

→ Crime Analysis

→ Disaster Response

## Challenges in Big Data Analytics

### ① Data Quality & Management

Data Silos: → issue  
→ form

→ Data Cleanliness

→ Data Governance

### ② Privacy & Security Concerns

→ Data Protection

→ Compliance

→ Ethical Use of data

### ③ Skill gaps & Talent Shortage

### ④ Infrastructure

## Future Landscape of Big Data Analytics

### ① AI and ML

→ Enhanced Capabilities

→ AI Integration  
→ Automation

→ Applications

→ NLP  
→ Computer Vision

→ Impact

→ Decision Support

→ Efficiency

(2)

## Edge Computing & IoT

- Real-Time Processing
- IoT Expansion

(3)

## Ethical Considerations

(4)

## Data as a Service (DaaS)

## Emerging Trends in Big Data

What drives trends in BD

- (1) Growth of IoT, AI and machine learning
- (2) Increasing data generation from social media, sensors and digital platforms
- (3) Demand for real-time analytics & insights

## Why Emerging Trends Matter

- (1) Keep pace with evolving data needs
- (2) Leverage advancements to solve complex problems

Trends① Edge Computing

→ Processing data closer to the source instead of sending it to central server.

adv → ① Reduces latency, improves efficiency

② ex:- Autonomous vehicles process data locally for faster decision-making

② Integration of AI and Machine Learning (ML)Appn in Big Data

→ AI models predict outcomes & provide actionable insights.

→ ML automates tasks like anomaly detection, data classification

ex → AI-driven chatbots, recommendation engines

③ Block Chain for data security

→ Why Block chain matters in Big data

→ Decentralized and immutable ledgers ensure data integrity.

→ Enables secure sharing of sensitive data across organizations

ex → Tracking supply chain data securely using Block chain

## ④ Realtime and Streaming Analytics

→ Processing data in real-time for instant insights.

use → Fraud detection in Banking

→ Personalized content delivery on streaming platforms

## ⑤ Cloud-Native Big Data Solutions

→ Amazon Redshift, Google Big Query

Scalability, cost-efficiency & on-demand resources

## Ethical Considerations in Big Data

### ① Data Privacy & Governance

→ importance of user consent & transparency.

### ② Bias in Analytics

→ Algorithms trained on biased data can perpetuate inequality

→ Need for diverse data sets & fair AI practices

### ③ Cybersecurity Risks

- As data grows, vulnerabilities inc.
- Need for robust encryption & secure systems.

## Future Directions for Big Data

### ① Quantum Computing

- Potential
- Solves complex Big Data problems exponentially faster
- App<sup>w</sup> in cryptography, climate modeling & genomics

### ② Democratization of Big Data Tools

- Make Big Data accessible to non-technical users

ex:- User-friendly platforms like Tableau, Power BI.

### ③ Autonomous Data Platforms

- Platforms that manage & analyze data with minimal human intervention,

ex Oracle's Autonomous Data Warehouse

## ④ Rise of Synthetic Data

→ Artificially generated data used to train AI and ML models

adv → Overcomes privacy concerns, fills gaps in real-world data.

## Hadoop Ecosystem

Apache Hadoop is an open source framework intended to make interaction with big data easier.

Hadoop has made its place in the industries & companies that need to work on large data sets which are sensitive and needs efficient handling.

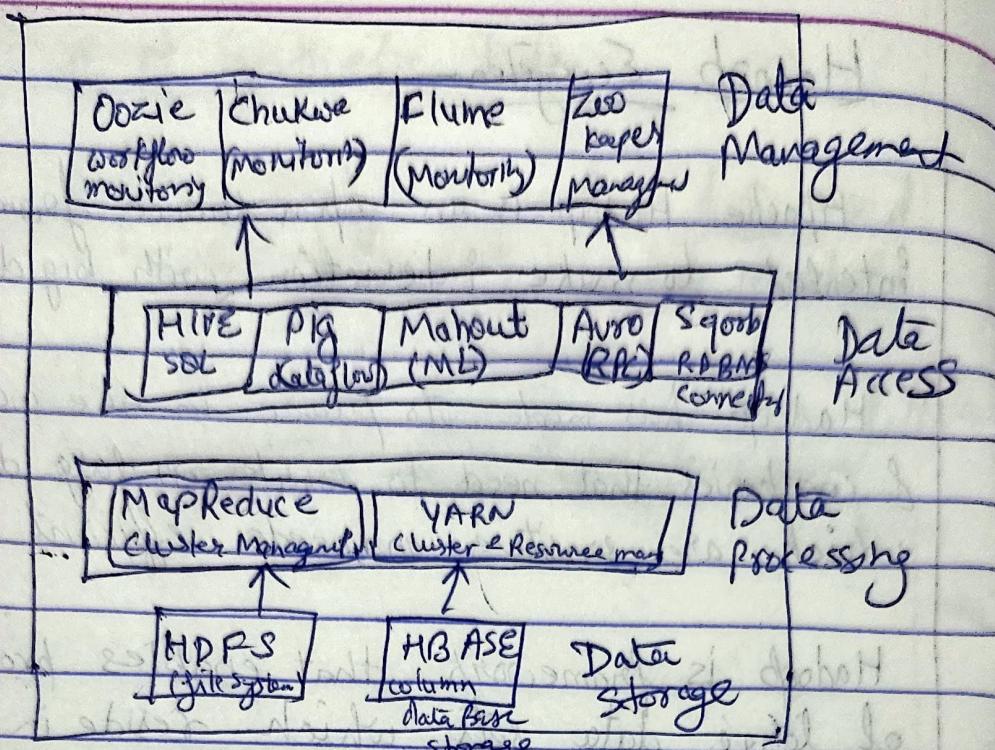
Hadoop is framework that enables processing of large data sets which reside in the form of clusters.

Being a framework, Hadoop is made up of several modules that are supported by a large ecosystem of technologies.

Hadoop ecosystem is a platform or a suite which provides various services to solve the big data problems.

There are 4 major elements of Hadoop i.e.

- ① HDFS
- ② Map Reduce
- ③ YARN
- ④ Hadoop common Utilities.



## Hadoop Ecosystem

HDFS → Hadoop Distributed File System.

HBASE → No SQL database

YARN → Yet Another Resource Negotiator

MapReduce → Programming based Data processing.

HIVE, PIG → Query Based processing of data services.

Mahout, Spark, MLlib:- machine learning algorithm libraries

~~Spark~~,

AVRO → Uses RPC to serialize data in  
remote procedural calls

Oozie : Job Scheduling

Zookeeper : → Managing clusters

Solr, Lucene : → Searching and Indexing

### HDFS

→ Primary or major component of Hadoop ecosystem and is responsible for storing large sets of structured or unstr. data across various nodes and thereby maintaining the metadata in the form of log files

HDFS consists of two core components i.e.

1. Name node → Prime node contains metadata required for resources than the data stores
2. Data node → actual data

HDFS maintains all the coordination b/w the clusters and hardware.

YARN → Yet Another Resource Negotiator

It is the one who helps to manage the resources across the clusters.

It performs scheduling & resource allocation for the hadoop System.

It consists of three major components:-

- ① Resource Manager → for allocating resources for the app<sup>n</sup>
- ② Nodes Manager → allocation of resources like CPU, memory, bandwidth per machine
- ③ App<sup>n</sup> Manager

→ works as an interface b/w resource manager & node manager. & performs negotiations as per the requirement of the two.

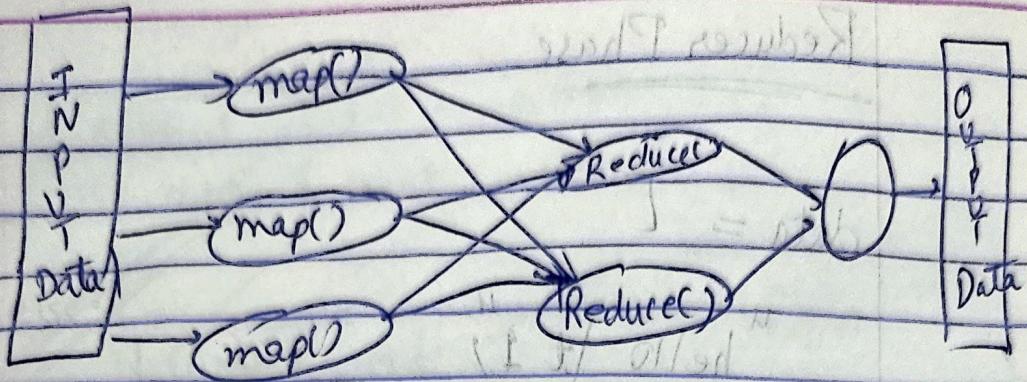
### ③ Map Reduce:

→ It is used to make it possible to carry over the processing logic & helps to write app<sup>n</sup> which transforms big data set into a manageable one.

It has two functions (1) Map() (2) Reduce()

A MapReduce job usually splits the input data-set into independent chunks which are processed by the map tasks in a completely parallel manner.

The framework sorts the outputs of the maps, which are then input to the reduce tasks.



## # WordCount

Mapper

import sys

myline = myline.strip()

→ Remove white space  
either side

words = myline.split()

→ Break the line  
into words

for myword in words:

print "%s %s" % (myword, 1)

OR

print f'{word} 1'

ex

Hello! How are you? How is your brother?

O/P

Hello -1  
! -1

How → 1

is → 1

How → 1

your → 1

are → 1

brother → 1

you → 1

? → 1

? → 1

writing

## Reducer Phase

data = [

"hello \t 1"

"! \t 1",

"How \t 1"

"are \t 1"

"you \t 1"

"? \t 1"

"How \t 1"

"is \t 1"

"your \t 1"

"bother \t 1"

"? \t 1"

]

~~current word = None~~

~~word-count = {}~~

~~current count = 0~~

~~result = {} + {format(" ", )}~~

for myline in data:  
myline = myline.strip()

try:

word, count = myline.split(' |t', 1)

count = int(count)

→ convert count to int

except ValueError:

continue

check if the current word  
matches the previous one

Page No.

Date: 11

if current-word == word:

    current-count += count

else:

    if current-word:

        result.append((current-word, current-count))

    current-word = word

    current-count = count

if current-word:

    result.append((current-word, current-count))

result

if word in word-counts:

    word-counts[word] += count

else:

    word-counts[word] = count

result = list(word-counts.items())

print(result)

(4)

## Hadoop Common Utilities

Hadoop common utilities are nothing but our java library & java files or javascripts that we need for all the components present in a Hadoop cluster.

like HDFS, YARN & MapReduce for running the cluster.

## Replication in HDFS & its importance

Replication ensures the availability of the data.

The no. of times you make a copy of that particular thing can be expressed as its replication factor.

default replication factor = 3

## Rack Awareness:

The rack is nothing but just the physical collection of nodes in our Hadoop cluster.

A large Hadoop cluster consists of so many racks with the help of this racks Information Namenode chooses the closest Datanode to achieve the max performance while performing the read/write info. which reduces the network traffic.

# We divide the whole file into blocks because of the increased seek time on the whole file.

## # HeartBeat:

If it is the signal that Datanode continuously sends to namenode. If namenode doesn't receive heartbeat from a datanode then it will be considered dead.

## Features of HDFS

- ① Distributed Storage
- ② Block reduces seek time
- ③ Data is highly available
- ④ High fault tolerance

## Limitations of HDFS

### ① Low latency data access

App which require low-latency access to data will not work well with HDFS

### ② Small file problem:-

Having lots of small files will result in lots of seeks & lots of movement from one datanode to another datanode to retrieve each small file

## Feature

### ① Scalable Storage

Horizontal scalability → support for large file

## Distributed Architecture

### ② Data Locality

HDFS moves computation tasks to the nodes where data resides, rather than transferring large dataset over the network

It Reduces data transfer latency & network congestion

### ③ Resilience

HDFS ensures high availability & reliability of data through redundancy & replication.

→ Replication factor

→ Automatic Recovery → In case of node failure the system automatically re-replicates blocks to maintain the replication factor.

### ④ Fault Tolerance

HDFS is designed to handle h/w & s/w failures without affecting data integrity or availability.

→ Heart Beat & Block Reports:-

The namenodes continuously monitors datanodes through heart beat & block reports.

→ Failure Mechanism : If a datanode fails, the namenode renounces tasks & re-replicates lost blocks.

## HDFS Administration

HDFS administration involves managing the Hadoop HDFS for efficient storage & processing of big data.

### Tasks of administration

#### ① Cluster Monitoring

→ Use tools like Ambari, Cloudera Manager or HDFS command-line utilities to monitor node health, disk usage & block distribution.

#### ② Replication Management

→ Adjust replication factor to balance fault tolerance & storage efficiency.

Use "hdfs dfsadmin -setrep" to configure replication levels

### ③ NameNode Operations

Monitor & manage the NameNode.  
→ Regularly back up NameNode metadata.

### ④ Balancing cluster Load

A tool is used to distribute data evenly across DataNodes.  
The tool is "Balances tool".

### ⑤ Disk Usage Management

Promiscuously clean up temporary & unused files to optimize storage.

### Setting up Hadoop Cluster

There are two ways in which we can set it up:  
→ Single Node Cluster

### Multi-Node Cluster

### Single Node Cluster

It has one DataNode running and setting up all the NameNode, DataNode, Resource Manager & Node Manager on a single machine.

This is used for studying & testing process.

## Multi Node Cluster

Here more than one DataNode running & each DataNode is running on different machines.

## Installation steps on a Single Node Cluster

- ① JAVA - Java JDK Install
- ② HADOOP - Hadoop package (download)

Extract Hadoop at C:\Hadoop

③ Setting up the Hadoop-HOME Variable

(Use windows environment variable setting for Hadoop path setting)

Go to environment properties or  
edit environment variables

- ④ Set JAVA\_HOME Variable

- ⑤ Set Hadoop & Java bin directory path

# Hadoop Configuration

We need to modify Six files

① core-site.xml

② Mapred-site.xml

③ Hdfs-site.xml

④ Yarn-site.xml

⑤ Hadoop-env.cmd

⑥ Create two folders  
datanode & namenode

## ① Core-site.xml

It contains the core configuration settings for Hadoop, primarily specifying fundamental parameters related to the HDFS & general behaviour of the cluster.

⇒ Specify the Default File System

↳ Defines the URI of the default Hadoop file system

(Uniform Resource Identifier)

property : fs.defaultFS → URI of the default file system

# hadoop.tmp.dir → Base directory for temporary files

# io.file.buffer.size → size of the buffer for reading & writing files

# hadoop.proxyuser.<user>.hosts → Define proxy host for a specific user

# hadoop.proxyuser.<user>.groups → Define groups allowed to use proxy for a specific user

# fs.s3a.access.key

→ Access key for Amazon S3 storage if using S3 as a file system

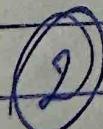
# fs.s3a.secret.key

→ Secret key for Amazon S3 storage

~~Advantage of core-site.xml~~



① acts as central configuration file that influences how hadoop interacts with underlying file system & network



② Plays a crucial role in setting up & running a Hadoop cluster effectively

③ Ensures compatibility & optimal performance for distributed system

## ② hdfs-site.xml

contains configurations specific to HDFS

# dfs.replication → no. of datablock replicas

# dfs.namenode.name.dir

(path for NameNode storage)

# dfs.datanode.data.dir

→ path for Data Node ~~for data~~ storage

## ③ mapred-site.xml

→ Contains MapReduce framework-specific configurations

# mapreduce.framework.name

→ /framework name  
(e.g. yarn)

# mapreduce.job.tracker

(Job Tracker hostname if part in a classic MapReduce Setup)

## ④ Yarn-site.xml

Contains configuration for YARN which manages cluster resources

# yarn.resourcemanager.hostname

(Resource Manager Hostname)

# yarn.nodemanager.resource.memory-mb

(Memory allocated for Node Manager)

# yarn.nodemanager.vmem-check-enabled

(Virtual memory checker)

⑤

hadoop-env.sh

A shell script to configure env. variables for hadoop such as java path

# JAVA\_HOME (path to Java installed)

# HADOOP\_HOME (path to the Hadoop installation)

# Data Loading into Hadoop

HDFS is designed to store large volume of data reliably & efficiently.

Loading data into HDFS & other Hadoop components can be achieved through various methods & tools, each suited for different types of data sources & use cases.

## Data loading Methods

- ① HDFS Commands
- ② Apache Sqoop
- ③ Apache Flume
- ④ Apache Kafka
- ⑤ Apache NIFI
- ⑥ DistCP (distributed Copy)
- ⑦ Using Hadoop APIs

### HDFS command

CLI → command line interface

Uses cases → ① Small to medium datasets

↳ ideal for loading data that can be managed via (CLI) without the need for automation

- ② Useful for quick one-time data uploads or downloads

~~Delta~~

- ① Easy to use → basic command-line operation
- ② No additional tool required  
direct interaction with HDFS

## Limitations

Manual process

↳ not suitable for large-scale or automated data ingestion

Handling of massive data is troublesome

②

## Apache Sqoop

↳ It is a tool designed for efficiently transferring bulk data from hadoop & structured datastores such as relational databases

Import data → RDBMS to HDFS

Export data → HDFS to RDBMS

Use

- ① Migrating data from enterprise databases to Hadoop for analytics
- ② Part of Extract, Transform, Load workflow to prepare data for processing.

~~Advantages~~ ✓ High performance → parallel data transfer utilizing Map Reduce,

- ① Seamlessly integrates with Hadoop ecosystem (Hive, HDFS)
- ② Automatically generates Hadoop schemas based on database tables

### Limitations

- ① Limited to structured
- ② Dependency on RDBMS

### Apache Flume

It is a distributed, reliable & available service for collecting, aggregating and moving large amounts of data log or event data from various sources to HDFS

- # Designed to handle high throughput data ingestion!
- # Ensures data is not lost during transfer with built-in fault tolerance
- # Flexible through customizable sources, channels & sinks

### Uses

- ① Log Aggregation → Collecting logs from multiple apps into HDFS for centralized storage & analysis
- ② Real time Analytics → Streaming event data for near

real-time processing with frameworks like Apache Storm or Spark Streaming.

### Adv

- ① Real time data ingestion
- ② Fault tolerance
- ③ Easy to use

### Limitations

- ① Customizing beyond std setups can be challenging
- ② High throughput ingestion may require significant cluster resources

### Apache Kafka

It is a distributed event streaming platform capable of handling large volumes of data in real-time.

While Kafka itself is not a data loader, it serves as a robust ingestion pipeline that can transport data to Hadoop components like HDFS, HIVE or spark.

- # High Throughput & low latency → millions of messages per sec
- # Durability & fault tolerance
- F Scalability → Horizontally Scalable

## ~~Advantages~~ Advantages of Real-Time Processing

① Versatility → can handle diverse data sources & sinks

### ③ Scalability

Challenges

① Operational complexity → Kafka cluster managing requires expertise.

② Very careful configurations are required to ensure minimal delays.

## Map/Reduce Concepts

### Map

Takes input data and transforms it into key-value pairs.

→ Processes each record independently.

→ Runs in parallel across multiple machines.

### Reduce

→ Takes the output from Map phase

→ Aggregates/combines values with the same key.

→ Produce final o/p.

→ Also runs in parallel for different keys.

Input → Map → Shuffle/Sort → Reduce → O/P.

# Manging Job Execution using Yarn

YARN involves several steps & components that work together to efficiently allocate resources & execute app's in a distributed environment.

Detailed Overview of How YARN manages job execution

## ① Job Submission

client Interaction :- The process begins when a client submits an appn to the YARN resource manager (RM)

This submission includes the appn code, configuration and resource requirements.

## ②

## Resource Manager (RM)

Resource Allocation ! The RM is the master daemon responsible for managing resources across the cluster.

→ It decides how to allocate resources to various app's based on available resources and scheduling policies.

Application Master (AM) Launch ! The RM allocates a container for the Appn Master (AM) on one of the

nodes.

→ The AM is responsible for managing the lifecycle of the appn.

### ③ Application Master (AM)

Resource Negotiation: The AM negotiates with the RM for additional resources (containers) needed to execute the appn tasks.

Task Coordination: The AM coordinates the execution of tasks, monitors their progress and handles task failures.

### ④ Node Manager (NM)

Container Management - Each node in the cluster runs a Node Manager, which is responsible for launching and monitoring containers on that node.

Resource Monitoring - The NM reports resource usage and container status back to the RM.

## ⑤ Container Execution

Task Execution: Containers are the units of resource allocation in YARN. They encapsulate environment, in which app<sup>n</sup> tasks run, including CPU, memory & other resources.

Isolation: Containers provide isolation between different tasks, ensuring that they do not interfere with each other.

## ⑥ Scheduling and Resource Management

Schedulers: → like FIFO manages how resources are allocated to apps.

Queue Management: Apps are submitted to queues, which can be configured with different resource capacities & priorities.

## ⑦ Monitoring and Logging

Web UI: YARN provides a web-based user interface for monitoring the status of app<sup>n</sup>, resource usage & cluster health.

Logs: App<sup>n</sup> logs are collected & can be accessed for debugging & performance analysis.

## ⑧ Fault Tolerance

Failure Handling → YARN is designed to handle failures gracefully. If a container or node fails, the AM can request new resources to restart the failed tasks.

AM Recovery → In case of AM failure, YARN can restart the RM to continue managing the app.

## ⑨ Security

Authentication & Authorization → YARN supports Kerberos-based authentication and provides mechanism for authorizing access to resources and apps.

Access Control → Queue-level access control lists (ACLs) can be configured to manage user permissions.

## ⑩ Optimization and Best Practices

Resource Configuration → Properly configure resource allocations for container optimization performance.

Queue Configuration → Set up queues to reflect organizational priorities & adjust configurations needed to improve efficiency.

# Data Analysis & Visualization

Data analysis is the process of inspecting, cleaning, transforming and modeling data to extract useful information, draw conclusions & support decision-making.

## Types of Data Analysis

① Descriptive Analysis → Summarizing data using statistics like mean, median, mode, variance etc.

② Diagnostic Analysis → Identifying reasons behind past outcomes

③ Predictive Analysis → Using ML models to forecast future trends

④ Prescriptive Analysis → Suggesting actions based on insights from data.

## Data Visualization

It is the graphical representation of data using charts, graphs & plots to make complex data more understandable.

ex: Box Chart, Line Chart, Pie Chart, etc.

## How to use Data Analysis & Visualization in projects

- (1) Define the problem statement  
Identify what you want to analyze
- (2) Collect & Clean Data
- (3) Perform Exploratory Data Analysis (EDA)
  - Identify patterns & trends using summary statistics
  - Create basic visualizations to understand distributions
- (4) Apply the <sup>MC</sup> models
- (5) Visualize Results

## Real-life Use cases in Data Analysis

- (1) Health Care Analytics → Hospitals use data analysis to predict disease outbreaks & optimize patient care
- (2) Stock Market Analysis → Traders use predictive analytics to forecast stock prices.

# Different Sources of Big Data

## ① Social Media Types

FB, Insta, Twitter

Data Type  $\rightarrow$  User Interactions, posts, likes, Shares, comments, multimedia

Characteristics  $\rightarrow$  High Unstructured, high volume and real-time data that offer insights into consumer behaviour, trends & sentiments

## ② Machine-Generated Data

Log files, system logs, application logs, clickstreams

Sources) Sensors, applications, network devices

Characteristics  $\rightarrow$  Str. and Semi Str. data that can be used for monitoring performance, security analysis & user behavior tracking.

### Transactional Data

(3)

Sources → Point-of-sale systems, Online Transactions,  
electronic payments, e-commerce purchases.

Type of data → Purchase details, payment records, invoices,  
receipts.

char.

Structured Data → critical for financial analysis,  
inventory management & customer relationship  
management.

(4)

### Sensor Data / IoT

Sources → Wearables, smart home devices, industrial sensors,  
environment sensors, GPS devices.

Type of data → temperature readings, motion detection, heart rate,  
location data.

Char. → High-volume, time-stamped data that is often continuous  
& requires real-time processing.

(5)

## Web Data

Sources → Websites, blogs, forums, online news articles.

Type of Data → Webpage content, metadata, user reviews, testimonials.

Char. → Semi-str. or Unstr. data useful for market research, competitive analysis & sentiment analysis.

(6)

## Scientific Research Data

Physics, CS, astronomy etc.

Type of Data → Simulation Data, Experimental results etc.

Char.

Extremely large datasets that require specialized processing capabilities, often used in advanced research & development.

(7)

## Audio & Video Data

Sources Surveillance cameras, video conferencing, media podcasts, podcasts.

Type of data

→ Multimedia files, streaming content.

Char.

Unstr. data that requires large storage capacity & specialized processing for analysis like.

speech recognition and video analytics.

### (8) Text Data

Sources Emails, documents, PDFs, books, SMS messages.

Type of Data Unstructured text requiring NLP for analysis.

Characteristics Valuable for extracting insights on communication patterns, topic modeling and information retrieval.

### Introduction to Data Ingestion

It is a fundamental process in the realm of data management & analytics.

If involves the collection, transportation & preparation of data from various sources into a storage or processing system where it can be analyzed, visualized or used in apps.

Data ingestion refers to the process of obtaining and importing data for immediate use or storage in a database.

## Introduction to Data Ingestion

It is a fundamental process in the realm of data management & analytics.

It involves the collection, transportation & preparation of data from various sources into a storage or processing system where it can be analyzed, visualized or used in apps.

Data ingestion refers to the process of obtaining and importing data for immediate use or storage in a database.

③

The data can come from a multitude of sources like

- ① Databases → Relational & Non Relational DB
- ② Files ) CSV, JSON, XML
- ③ API ) → Data retrieved through REST APIs & APIs
- ④ Streaming Data ) Social media, logs
- ⑤ Third-party Sources ) External data sets, partner data

## Why Data Ingestion Important

① Data Driven Decision Making

② Scalability

③ Data Integration

④ Operational Efficiency

⑤ Real-time Analytics

## Types of Data Ingestion

It can be categorized based on freq & method of data transfer

① Batch Ingestion → Data is transferred & collected in large blocks at scheduled time intervals

② Real-time Ingestion → Data is ingested continuously & processed in real-time as it arrives.

③ Micro-Batch Ingestion

↳ A hybrid approach where data is collected in small batches at frequent intervals.

### Date Ingestion Methods

① ETL (Extract, Transform, Load)

② ELT (Extract, Load, Transform)

③ CDC (Change Data Capture)

④ Streaming Ingestion

ETL →

Extract → Retrieve data from various sources (databases, APIs, flat files etc.)

Transform → Clean the data, enrich it and structure it (like removing duplicates, converting formats).

Load → Insert the transformed data into a target system

target system can be → data warehouse, databases.

~ ETL can be used in BI business Intelligence to process the structured data.

Adv.

→ It ensures data quality & consistency before loading.

→ Good for structured, historical data analysis.

→ Works well for scheduled batch processing.

Disadv. → Slower for large datasets as transformation occurs before loading.

→ Not ideal for real-time apps.

Tools

→ Apache Nifi, AWS Glue

ELT

Extract → Data is retrieved from source systems.

Load → Data is loaded directly into a cloud data lake or warehouse

~~Ex.~~ Snowflake, Google Big Query.

Transform:

Data is transformed within the destination system using SQL or processing systems.

ELT can be used for Big Data & cloud-based analytics.

Storing raw data for later analysis, for large data volumes.

Adv → faster ingestion as no pre-transformation is required

→ leverages cloud-based storage for efficient processing.

→ More flexible for data scientists to explore raw data.

## Disadv

- Requires powerful storage & processing infrastructure.
- If data is not well-structured, may lead to inefficient queries.

Tools Snowflake, AWS Redshift

## CDC (Change Data Capture)

It is a tech. that captures changes  
(insert, update, delete) kind of changes

in a database & propagates them to another system in real-time or near-real time.

Detect Changes → Identify modified records in source tables

Capture Changes → Log the changes as events

Replicate changes → Apply the changes to the target system

## Types of CDC

- ① Log-based CDC → Reads database transaction logs  
It is efficient, minimal impact.
- ② Trigger-based CDC → Uses database triggers  
It can slow performance.
- ③ Timestamp-based CDC → Compares timestamp P  
It will require additional indexing

CDC is used for Real-time database replication.  
 → for incremental data processing.  
 → To synchronize data across systems

### Adv.

- Real or near-real time updates
- Reduces data transfer costs by sending only changed data
- Useful for event-driven architectures

### Disadv.

- Maintaining logs, time stamps can lead to some additional load on source data
- log-based CDC may require database specific configuration

### Tools

→ Oracle Golden Gate, AWS DMS

(Database Migration Service)

## Streaming Data Ingestion

It is the process of continuously collecting, processing & transferring data in real time as it is generating.

Streaming Ingestion ensures: →

- ① Low Latency
- ② Scalability
- ③ Fault Tolerance
- ④ Event-Driven Processing → It works on individual data events as they arrive.

It can be used in

- Financial Transactions
- IoT Sensor Data
- Social Media Analytics → Trends
- Cyber Security

## Challenges in Streaming Ingestion

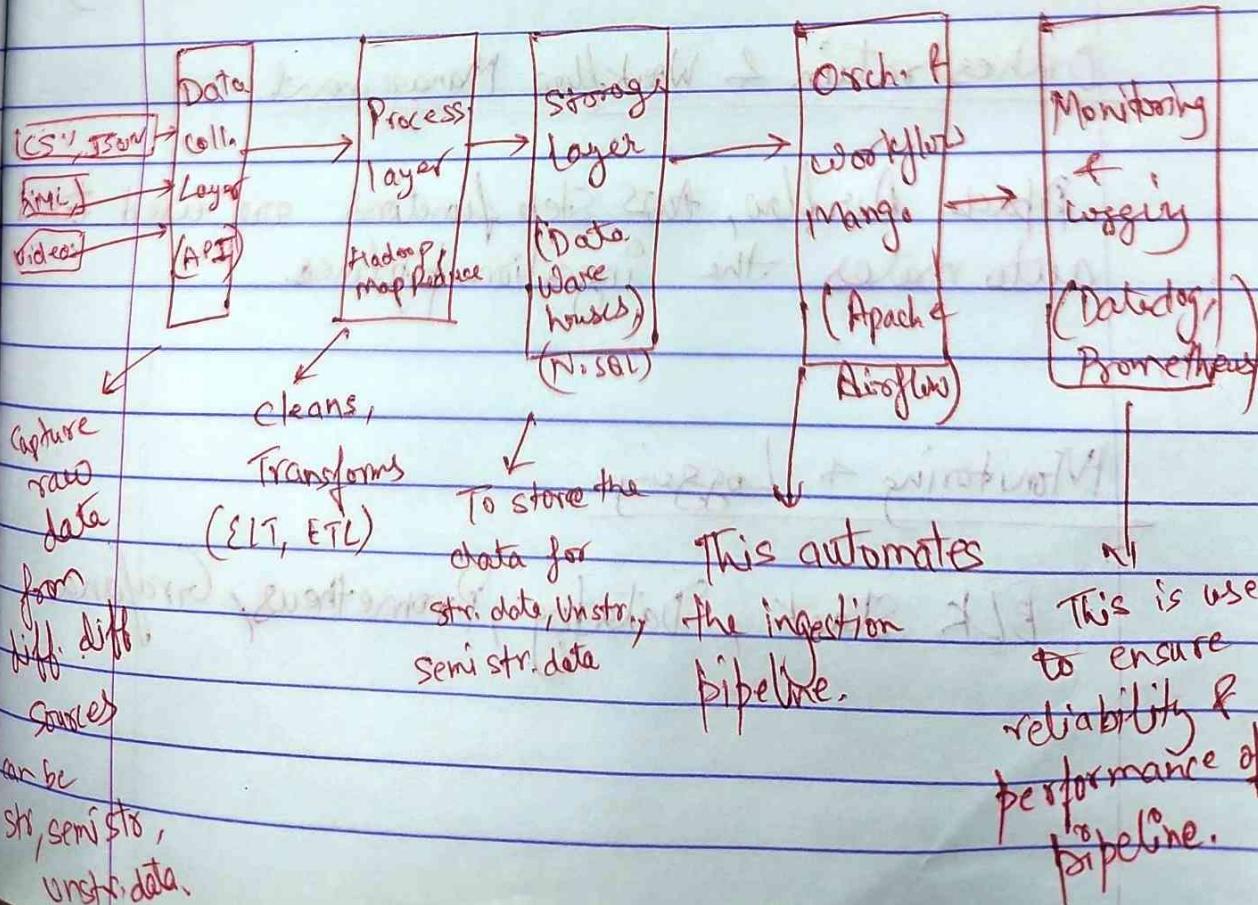
- Data Ordering & Consistency
- Scalability
- Latency vs Accuracy Trade-off

## Data Ingestion Pipeline

It is responsible for collecting, processing and storing large volumes of data efficiently.

### Key Components of a Data Ingestion Pipeline

1. Data Sources
2. Data collection layer
3. Processing Layer
4. Storage Layer
5. Orchestration & Workflow Management
6. Monitoring & Logging



Data Collection Examples → APIs, IoT sensors etc

Processing Layer → Apache Spark, Hadoop MapReduce, for Batch Tools.

Streaming Tools, Kafka Streams

Storage layer :-

Str. data → Data Warehouses like Google BigQuery.

Instr. data → Data Lakes like (AWS S3, Azure data Lake) for raw, unstructured data.

Semi Str. data → NoSQL Databases (MongoDB, Cassandra)

Orchestration & Workflow Management.

Apache Airflow, AWS Step functions are used to auto-mates the ingestion pipelines.

Monitoring & Logging

ELK Stack, Datadog, Prometheus, Grafana

Unit → 5

Page No.

Date:

Techniques for ingesting batch data from  
various sources.

## Sqoop Arch.

### Sqoop → Feature

- ① Helps to connect the result from SQL queries into HDFS
- ② Helps to load processed data into Hive or Hbase.
- ③ Can performs security operations of data (using kerberos).
- ④ can compress data.

### Main operations performed in Sqoop

- (1) Import
- (2) Export

Sqoop used the command-line interface (CLI) to process command of user.  
 (Can also use APIs)

Sqoop will only do import and export of data based on user command (Nothing else)  
 (no aggregation of data)

### Working

Import: Sqoop will first parses argument (from CLI) and send them to next stage where arguments are induced for Map only job.  
 (Only Mapper is used no Reducer).

Once the Map receives arguments it then gives command of release of multiple mappers depending upon the number defined by the user (default 4) at the of CLI)

Each mapper task is assigned with respective part of data that is imported on the basis of key which is defined by the user in CLI)

To increase the efficiency of process Sqoop uses parallel processing technique

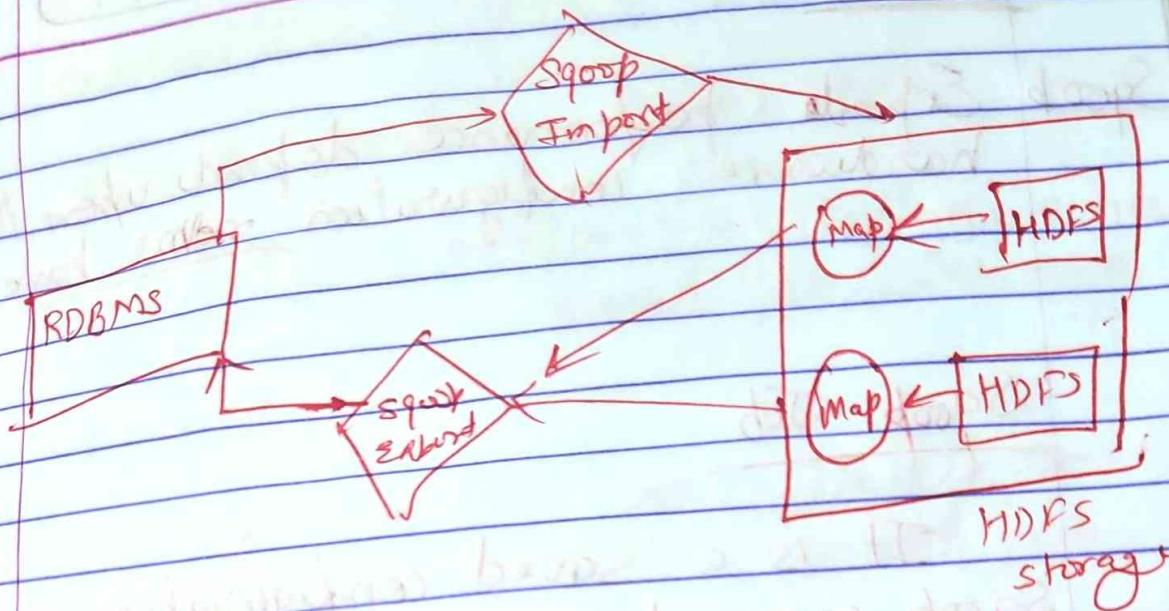
In this data is distributed equally among all mappers.

Then each mapper creates an individual connection with the database & then fetches individual part of data assigned by Sqoop.

Once the data is fetched then the data is been written in HDFS or HBase or Hive.

## Export

Export is done in the same way, Sqoop export tool which available performs the operation by allowing set of files from the Hadoop system back to RDBMS.



Adv

- Can use variety of structured data stores like Oracle, MySQL etc.
- Can perform ETL operations very fast + cost-effective way.
- parallel processing of data leads to overall faster processes.
- Supports fault Tolerance (multiple Map)

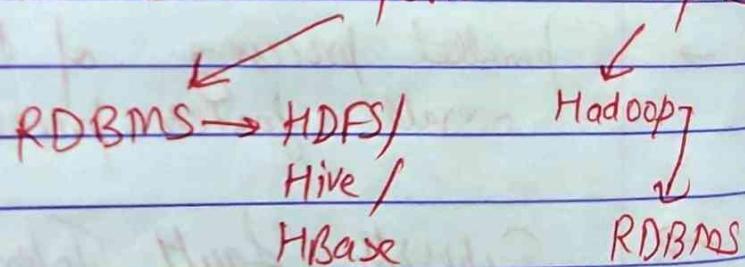
Sqoop Export performance depends upon the hardware configuration networks have.

### Sqoop Job

It is a saved configuration of a Sqoop command which can be used later to run or schedule without having to re-enter all the parameters each time.

It encapsulates all the necessary details for moving data between a relational database and a Hadoop ecosystem component.

It can be data imports or exports



It is reusable. It stores the parameters like database connection URL, credentials, table name, target directory in HDFS and can be other things.

It can be simply trigger the job whenever needed without recreating the

command line.

How to create a Sqoop job.

# --create <job-id> :> Define a new saved job with the specified job-id

# --delete <job-id> :> Delete a saved job

# --exec <job-id> :> Given a job defined with --create, run the saved job.

# --show <job-id> :> Show the parameters for a saved job.

# --list :> List all saved jobs.

--create my-job

\$ sqoop job --create my-job --import  
-connect jdbc:mysql://example.com/  
db --table my-table

To execute

\$ sqoop job --exec my-job

Sqoop jobs simplifies repetitive tasks ensures consistency and is especially useful in scheduled or production environment where tasks need to be automated.

## Incremental Import

The process of importing only the data that has changed since the last import operation.

Only newly added records need to be updated. It will reduce the time to completely import the whole dataset.

There are two modes of incremental import

- Append Mode
- Last modified Mode

Append Mode :-

Sqoop imports only the rows where the value in a designated column is greater than the last value

imported. This mode works best when new rows are added to the table.

Last modified Mode → Sqoop uses a timestamp column to track changes.

### File Handling in incremental Imports

To handle we have to specify ~~the~~ to sqoop where to write the imported in HDFS.

So, --target-dir option is used

For incremental imports, we must decide whether to write to the same target directory or to use a new directory for each run.

To write in same target directory we can use --append flag

It means adding the new files alongside the other files.

## Handling & generation of small files

As the volume of data from incremental changes can be relatively small, this could lead to many small files.

Large number of small files can lead to performance issues when processing or querying the data later.

→ To solve this issue, periodical consolidation of these files is necessary. This can be done using file merging or compaction job.

# When writing to the same target directory over multiple incremental runs, ensures that the files won't get overwritten accidentally.

In case of new target directory creation every time, we'll need a file handling process which merges or organizes these directories for processing.

## Best Practices

- 1) Merge Small files
- 2) Logging + Monitoring
  - ↳ monitor the target directories
- 3) Directory Organization
  - ↳ organize HDFS directories in a way that makes it easy to separate full imports from incremental imports.

## Introduction to message queuing systems

Message queues enable communication between various system components.

A message queue is a form of asynchronous service-to-service communication used in serverless & microservices.

Messages are stored on the queue until they are processed & deleted.

→ This is done in order to increase system reliability, scalability & flexibility.

This is achieved by adding a ~~message~~ message broker (like RabbitMQ or Amazon SQS) to pass messages ~~as~~ or data b/w processes or systems.

Messages are sent by producer & the message queue, where they are held until consumed by a consumer.

↓  
This decouples the producer from the consumer, allowing them to operate independently.

↳ means data can be sent to a queue without waiting for immediate, real-time processing

Message Queuing can be used for distributed computing, task scheduling & workflow management.

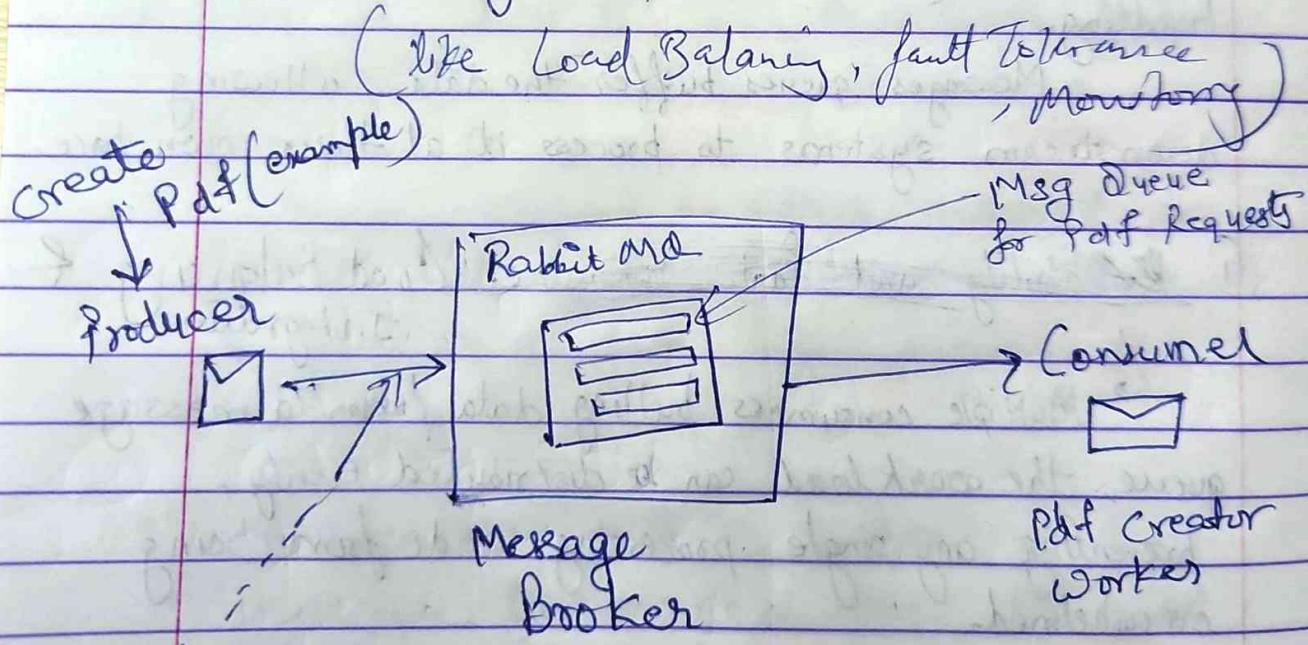
### Basic Features

- ① Decoupling Producers & Consumers
- ② Handling High Throughput & Scalability.
- ③ Asynchronous Processing  
↳ data ingestion often requires asynchronous handling.
- ④ Reliability and Fault Tolerance, Load Balancing & Integration.  
↳ Multiple consumers pulling data from a message queue, the workload can be distributed evenly, preventing any single processing node from being overwhelmed.
- ⑤ Reliability & Fault Tolerance.

These systems typically offer features such as message persistence, automatic retries, & fault tolerance.

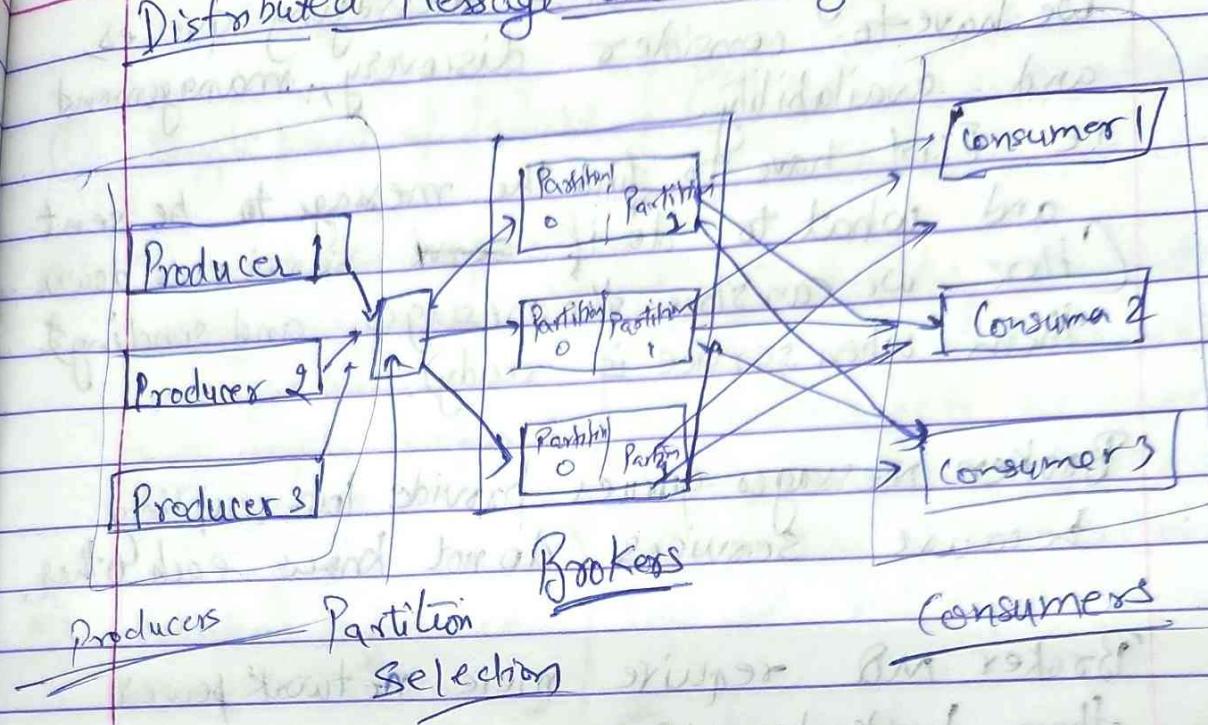
## Arch. for Message Driven Systems

- ① Producer layer
- ② Messaging Layer (Queue or Broker)
- ③ Consumer Layer
- ④ Data Processing & Storage
- ⑤ Auxiliary Components & Considerations



Message  
including  
user information

# Distributed Message Queue System



## Broker v/s Broker Less messaging

Broker messaging arch. means we have a centralized broker acting like a mediator.  
ex: RabbitMQ, IBM MQ

In Brokerless messaging arch., no broker exists and participants communicate directly with themselves.

ex: ZeroMQ, NanoMsg

In this separate broker process does not exist.

When using Brokerless messaging queues, we have to consider discovery, management and availability.

i.e. First have to find the message to be sent and what to do if service is down. (Here we can store the messages and sending them when service is ready).

Broker messages queues provide low coupling because services do not know each other.

Broker MQ require more network power than brokerless MQ.

### Adv. of Broker

- ① Loose coupling
- ② Scalability → horizontal scaling possible
- ③ Reliability → message persistence & guaranteed delivery
- ④ Centralized Control
- ⑤ Message Transformation & routing → Brokers have powerful routing capabilities.

### Adv. of Brokerless

- ① Simplicity
- ② Improved Performance without the need to route messages through a broker, brokerless messaging can offer lower latency & higher throughput, especially for local comm. b/w services.

Disadv. of Broker

- ① Single point of failure
- ② Increased complexity: setting & configuring a message broker can be more complex compared to brokerless.
- ③ Potential Performance impact: Additional overhead of routing messages

Disadv. of Brokerless

- ① Lack of centralized control.
- ② Scalability limitations: may not scale as well as brokers because responsibilities for message distribution & load balancing lies with individual services
- ③ No built-in guaranteed delivery

NoteMessage Routing

It is the process of determining how messages are directed or forwarded from one point to another within a messaging system.  
 (It is based on a set of rules, ~~or~~ criteria or configurations)

## Exchanges and Exchange Types

Exchange is a critical component responsible for routing messages to one or more queues based on defined rules.

First producer sends a message to an exchange and then exchange decides which queue the message should be delivered to.

It is done utilizing attributes like routing keys and binding rules.

### Routing Keys

A routing key is a string attached to a message by the producer when it is sent to an exchange.

It is like a label or tag that describes the message's content or category.

### Binding Rules

Binding rules are configurations that associate a queue with an exchange using a binding key or a matching pattern.

When a queue is bound to an exchange, the binding rule specifies which messages the queue is interested in receiving based on the routing keys.

- # ① Producer assigns a routing key to each message.
- ② The exchange uses this routing key to evaluate the binding rules set by queues.

③ If the routing key of the message matches a queue's binding rule, the message is routed to that queue.

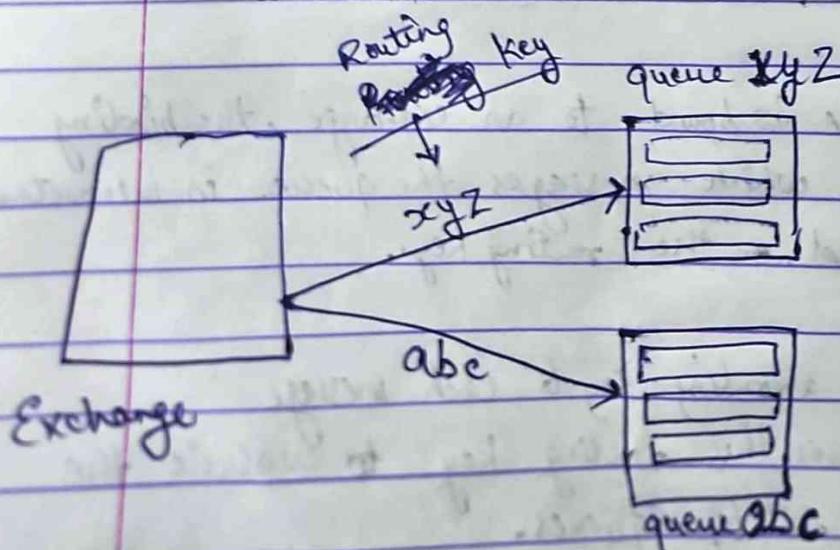
### Types of Exchanges

- ① Direct Exchange
- ② Topic Exchange
- ③ Fanout Exchange
- ④ Headers Exchange

#### ① Direct Exchange

A message goes to the queues whose binding key exactly matches the routing key of the message.

When you need a simple one-to-one routing.  
ex → sending ~~to~~ routing tasks to specific workers.



⑦ If the message routing key does not match any binding key, the message is discarded.

⑧ If more than one queue is bound to exchange with same bounding key, the direct exchange will broadcast the message to all matching queues.

## ② Topic Exchange

It is similar to direct exchange, but the routing is done acc. to the routing pattern.

It uses wildcards instead of using fixed routing keys.

↓  
(\* & #)

The routing key must consist of list of words delimited by a period "•".

"\*" is used to match a word in a specific position of the routing keys

or  
data.\*.\*.list

→ only match routing keys where the first word is "data" & fourth ~~the~~ word is "list".

"#" → it is used to match zero or more words

or  
data.list.# → means all routing keys starting with "data.list.",

### ③ Fanout Exchange

This routes a message to all queues that are bound to it regardless of routing keys or patterns.

It is useful when the same message needs to be sent to one or more queues with consumers who may process the same message in different ways.

④

### Headers Exchange

It uses the message header attributes for routing.

A Header Exchange routes messages on arguments containing headers.

i.e Header contains some arguments like (key,value) pair, format, type, x-match.

x-match is a special argument added in the binding b/w exchange & queue.

It specifies if all headers must match or just one.

x-match argument property can have two different values ① Any ② all.

# Default value is 'All' → i.e all header pairs must pass, while value of 'any' means at least one of the headers pairs must match.

## Distributed publish-subscribe messaging system

It is an architectural pattern in which messages produced by publishers are distributed among multiple consumers (subscribers).

In this model, publishers (producers) post messages to a central topic or channel without knowing which subscribers will receive them, and subscribers express interest by subscribing to channels or topics.

# Publishers and subscribers operate independently i.e decoupling.

↳ It improves scalability, flexibility.

# Topics and Channels

→ In many publish-subscribe systems, messages are organized by topics or channels.

A publisher sends a message to a specific topic, and all subscribers to that topic receive the messages.

# These systems are designed to run across clusters of machine to distribute the load, manage network partitions & provide high throughput & fault tolerance.

Components are horizontally scalable

# Modern distributed pub-sub systems often provide guarantees about message delivery and ordering.

e.g. Apache Kafka, Apache Pulsar

### Working

Pub-sub system has four key components

- ① Messages ② Topics/channel
- ③ Publisher ④ Subscribers

### Messages

A msg is communication data sent from sender to receiver.

Msg datatypes can be anything from strings to complex objects representing text, video, audio etc.

### Topics

The topic acts like an intermediary channel b/w senders and receivers.

It maintains a list of receivers who are interested in messages about that topic.

Topics serve as a logical grouping to organize messages by subject type, or any

other categorical attribute,

Topic is a channel identified by a unique name.

Every message has a topic associated with it.

### Subscribers

A ~~topic~~ subscriber is the message recipient. Subscribers have to register to topics of interest.

### Publishers

The publisher is the component that sends messages. It creates messages about a topic & sends them once only to all subscribers of that topic.

The publisher doesn't need to know who is using the information it is broadcasting & the subscribers don't need to know where the message comes from.

### Usage

Perform parallel asynchronous processing  
Balance workloads, Deliver app<sup>n</sup> & system alerts.

## Kafka Broker

It is a server that handles the storage, retrieval & distribution of messages in an Apache Kafka cluster.

It is responsible for managing topics, partitions, producers (publisher) & consumers (subscribers).

A Kafka cluster typically consists of multiple brokers to ensure fault tolerance, scalability & high availability.

Apache Zookeeper manages the metadata, leader elections & broker coordination.

It is process where Apache kafka dynamically selects a broker as the leader for a given partition.

The main task of the leader is to handle all read & write operations for that partition, while other brokers act as followers, maintaining replicas of the data.

If a leader broker fails, Zookeeper detects it & elects a new leader from the in-sync replicas (ISR).

## Key Functions of a Kafka Broker

- ① Message Storage :- Brokers store messages in topics, which are further divided into partitions.
- ② Message Distribution :- distribute messages across multiple partitions for scalability & fault tolerance.
- ③ Load Balancing :- Efficient distribution is done.
- ④ Replication :- Supports replication of data.
- ⑤ Leader & follower Roles :- One broker acts as a leader other acts as followers.
- ⑥ Handling Requests :- Brokers handle process read/write requests from producers & consumers.

## Producers & Consumers

Do it Yourself