

# Big Data Analytics: UNIT – 1

## INTRODUCTION TO BIG DATA ANALYTICS

Data Big Analytics involves examining large and complex datasets to uncover hidden patterns, correlations, and insights that can inform decision-making. Here's a brief overview:

### 1. What is Big Data?

Big Data refers to datasets that are so large or complex that traditional data-processing software is inadequate to handle them. It is often characterized by the "three Vs":

- **Volume:** The amount of data.
- **Velocity:** The speed at which data is generated and processed.
- **Variety:** The different types of data (structured, unstructured, semi-structured).

### 2. Key Components of Big Data Analytics

- **Data Collection:** Gathering data from various sources, such as social media, sensors, transactional systems, etc.
- **Data Storage:** Storing large volumes of data in scalable storage solutions like Hadoop Distributed File System (HDFS) or cloud-based storage.
- **Data Processing:** Using tools and technologies to process and analyze data. Common frameworks include Apache Hadoop and Apache Spark.
- **Data Analysis:** Applying statistical, machine learning, and data mining techniques to extract meaningful insights.
- **Data Visualization:** Presenting data insights through graphs, charts, and dashboards to make them understandable.

### 3. Tools and Technologies

- **Hadoop:** An open-source framework for distributed storage and processing.
- **Spark:** A fast, in-memory data processing engine.
- **NoSQL Databases:** Such as MongoDB and Cassandra, designed for handling unstructured data.
- **Data Warehouses:** Like Amazon Redshift and Google BigQuery, for large-scale data analysis.

### 4. Applications

- **Business Intelligence:** Enhancing decision-making through insights from customer data, sales, and operations.
- **Healthcare:** Analyzing patient data for better diagnosis and treatment.
- **Finance:** Fraud detection and risk management.
- **Retail:** Personalized marketing and inventory management.

### 5. Challenges

- **Data Privacy and Security:** Ensuring that data is handled securely and ethically.
- **Data Quality:** Maintaining accuracy and consistency in data.
- **Scalability:** Efficiently managing and processing growing data volumes.

### 6. Future Trends

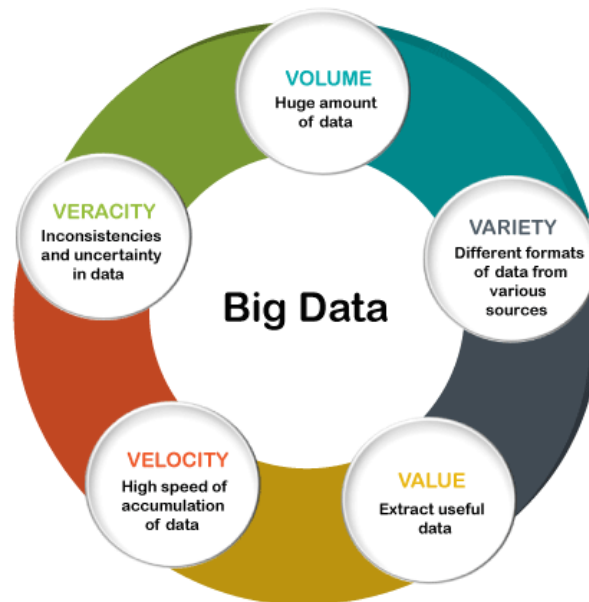
## Big Data Analytics: UNIT – 1

- **Artificial Intelligence:** Integrating AI to automate and enhance analytics processes.
- **Real-time Analytics:** Providing immediate insights and responses.
- **Edge Computing:** Processing data closer to where it is generated.

Big Data Analytics is a dynamic field that combines data science, statistics, and technology to drive informed decisions and innovations.

### CHARACTERISTICS OF BIG DATA

The characteristics of Big Data are often summarized by the "Three Vs," but as the field has evolved, more dimensions have been added to fully describe the nature of Big Data. Here's a detailed look:



#### 1. Volume

- **Definition:** Refers to the sheer amount of data being generated and collected. The scale can range from terabytes to petabytes and beyond.
- **Example:** Social media platforms generate massive volumes of user-generated content daily.

#### 2. Velocity

- **Definition:** The speed at which data is generated, processed, and analyzed. This includes the frequency of data creation and the rate at which it must be processed to be useful.
- **Example:** Streaming data from sensors in real-time for applications like autonomous vehicles.

#### 3. Variety

- **Definition:** The different types and formats of data. This can include structured data (e.g., databases), unstructured data (e.g., text, images), and semi-structured data (e.g., JSON, XML).
- **Example:** Combining data from customer transactions, social media posts, and email communications.

## Big Data Analytics: UNIT – 1

### 4. Veracity

- **Definition:** The trustworthiness and accuracy of the data. It involves addressing issues of data quality and reliability.
- **Example:** Ensuring that data from different sources is accurate and consistent for making reliable business decisions.

### 5. Value

- **Definition:** The usefulness and insights that can be derived from data. Not all data has inherent value; it must be processed and analyzed to provide actionable insights.
- **Example:** Analyzing customer purchase data to identify trends and improve marketing strategies.

### 6. Variability

- **Definition:** The inconsistency of data over time. Data can change in its format, meaning, or quality, making it challenging to handle.
- **Example:** Seasonal variations in sales data that affect forecasting.

### 7. Complexity

- **Definition:** The complexity of managing and integrating different types of data from various sources. It involves dealing with interconnections and dependencies between datasets.
- **Example:** Integrating data from multiple departments (e.g., sales, finance, customer service) to get a comprehensive view of company performance.

### 8. Actionability

- **Definition:** The ability to turn data insights into actionable strategies. This involves not just analyzing data but applying the results to drive decisions and actions.
- **Example:** Using predictive analytics to anticipate customer needs and optimize supply chain operations.

### 9. Data Integration

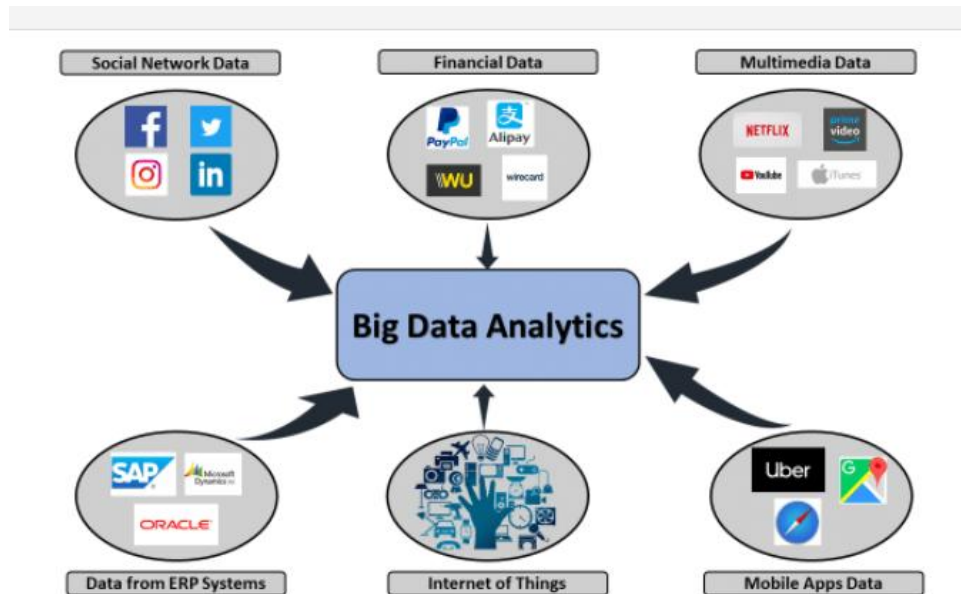
- **Definition:** Combining data from disparate sources into a unified view. This often involves overcoming challenges related to data formats, schemas, and sources.
- **Example:** Merging customer data from CRM systems with transaction data from sales platforms.

Understanding these characteristics helps in designing systems and strategies to effectively manage and leverage Big Data.

## SOURCES OF BIG DATA

Big Data comes from a variety of sources, each contributing different types and volumes of data. Here's a breakdown of common sources:

## Big Data Analytics: UNIT – 1



### 1. Social Media

- **Description:** Platforms like Facebook, Twitter, Instagram, LinkedIn, and others generate vast amounts of unstructured data including posts, comments, likes, shares, and multimedia content.
- **Examples:** User-generated content, sentiment analysis, and trend tracking.

### 2. IoT Devices

- **Description:** Internet of Things (IoT) devices generate data through sensors and connected devices that monitor various parameters.
- **Examples:** Smart thermostats, wearable fitness trackers, industrial sensors, and connected vehicles.

### 3. Transactional Data

- **Description:** Data generated from financial transactions, including sales, purchases, and payments.
- **Examples:** Point-of-sale (POS) systems, online transactions, and e-commerce platforms.

### 4. Log Files

- **Description:** Generated by systems and applications to record events, activities, and transactions.
- **Examples:** Web server logs, application logs, and system logs.

### 5. Sensors and Devices

- **Description:** Data collected from sensors embedded in various devices that monitor and record real-world conditions.
- **Examples:** Environmental sensors (temperature, humidity), GPS devices, and smart meters.

### 6. Multimedia Data

## Big Data Analytics: UNIT – 1

- **Description:** Includes images, videos, audio files, and other forms of multimedia content.
- **Examples:** Video surveillance footage, user-generated videos, and image repositories.

### 7. Enterprise Data

- **Description:** Data generated within organizations, encompassing various business processes and operations.
- **Examples:** Customer relation ip management (CRM) systems, enterprise resource planning (ERP) systems, and business databases.

### 8. Web Data

- **Description:** Data collected from web interactions and online activities.
- **Examples:** Clickstream data, web traffic logs, and user behavior analytics.

### 9. Public Data

- **Description:** Data made available to the public, often by governments, research institutions, or organizations.
- **Examples:** Government databases, open research data, and public records.

### 10. Communication Data

- **Description:** Data from various communication channels, including emails, messages, and call records.
- **Examples:** Email metadata, chat logs, and customer service interactions.

### 11. Health Data

- **Description:** Data from healthcare systems, including patient records, diagnostic results, and treatment data.
- **Examples:** Electronic health records (EHR), wearable health devices, and medical imaging data.

### 12. Data from Surveys and Research

- **Description:** Data collected from surveys, research studies, and other data-gathering methods.
- **Examples:** Market research surveys, academic studies, and customer feedback.

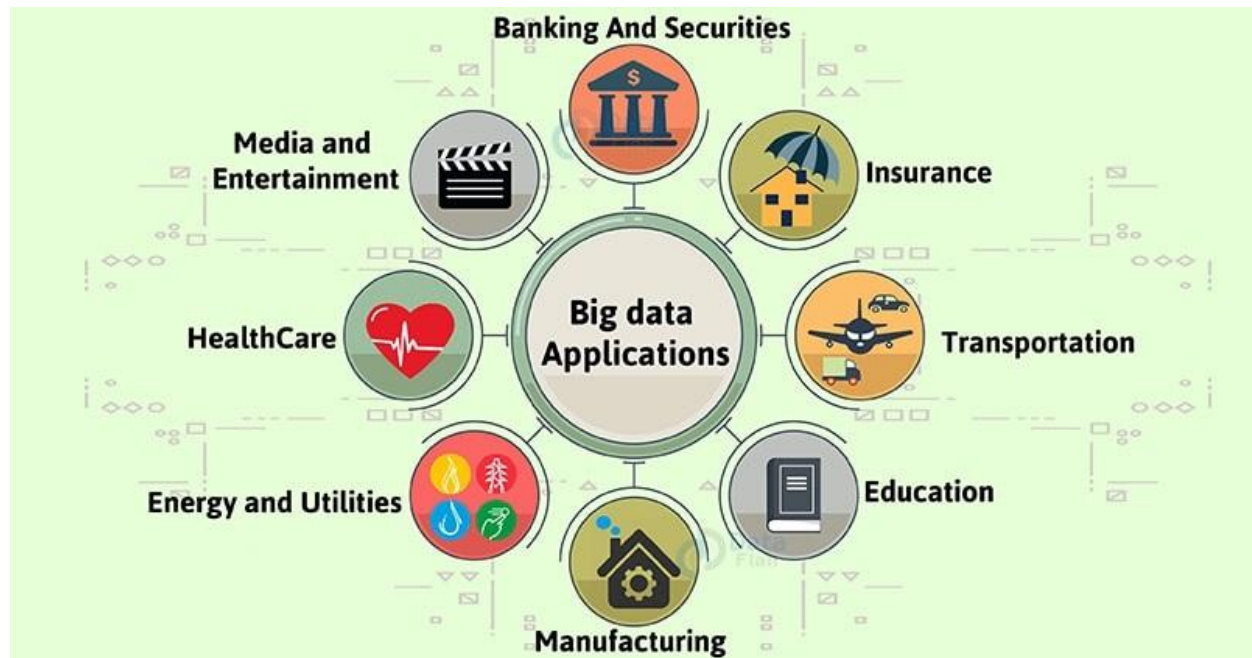
### 13. Geographic Data

- **Description:** Data related to geographical locations and spatial information.
- **Examples:** Geographic Information Systems (GIS) data, satellite imagery, and mapping data.

Each of these sources contributes to the overall landscape of Big Data, providing a diverse range of information that can be analyzed to generate insights and drive decision-making.

## Big Data Analytics: UNIT – 1

Big Data has a wide range of applications across various industries. Here's how it's used in different domains:



### 1. Healthcare

- **Predictive Analytics:** Forecast patient outcomes, disease outbreaks, and treatment responses.
- **Personalized Medicine:** Tailor treatments based on individual patient data and genetic information.
- **Operational Efficiency:** Optimize hospital operations, reduce costs, and improve patient care.

### 2. Finance

- **Fraud Detection:** Identify unusual patterns and anomalies in transactions to prevent fraud.
- **Risk Management:** Assess and mitigate financial risks through advanced analytics.
- **Algorithmic Trading:** Use complex algorithms and real-time data to make trading decisions.

### 3. Retail

- **Customer Personalization:** Deliver targeted marketing and personalized recommendations based on purchase history and behavior.
- **Inventory Management:** Optimize stock levels and supply chain operations using sales and demand data.
- **Customer Sentiment Analysis:** Monitor social media and reviews to gauge customer sentiment and improve service.

### 4. Manufacturing

## Big Data Analytics: UNIT – 1

- **Predictive Maintenance:** Monitor machinery and equipment to predict and prevent failures before they occur.
- **Quality Control:** Analyze production data to ensure quality and reduce defects.
- **Supply Chain Optimization:** Enhance logistics and inventory management with real-time data insights.

### 5. Transportation and Logistics

- **Route Optimization:** Improve delivery routes and reduce transportation costs using real-time traffic data.
- **Fleet Management:** Monitor vehicle performance and driver behavior to enhance fleet efficiency.
- **Demand Forecasting:** Predict transportation needs and adjust resources accordingly.

### 6. Telecommunications

- **Network Optimization:** Analyze network traffic and usage patterns to optimize performance and prevent outages.
- **Customer Churn Prediction:** Identify customers at risk of leaving and implement retention strategies.
- **Fraud Detection:** Detect and prevent fraudulent activities related to telecom services.

### 7. Government

- **Public Safety:** Analyze crime data to improve law enforcement strategies and public safety measures.
- **Urban Planning:** Use data to inform city planning, infrastructure development, and resource allocation.
- **Policy Making:** Base policy decisions on comprehensive data analysis to address social and economic issues.

### 8. Energy

- **Smart Grid Management:** Monitor and optimize energy distribution and consumption using data from smart meters and sensors.
- **Predictive Analytics:** Forecast energy demand and supply, and plan for future energy needs.
- **Maintenance:** Predict and prevent equipment failures in energy production and distribution.

### 9. Entertainment and Media

- **Content Recommendations:** Suggest movies, music, and shows based on user preferences and viewing history.
- **Audience Insights:** Analyze audience data to tailor content and marketing strategies.
- **Social Media Analytics:** Track trends and engagement to inform content creation and promotional activities.

### 10. Education

## Big Data Analytics: UNIT – 1

- **Student Performance:** Analyze student data to identify learning patterns and improve educational outcomes.
- **Adaptive Learning:** Personalize learning experiences based on individual student needs and progress.
- **Institutional Management:** Optimize administrative processes and resource allocation in educational institutions.

### 11. E-Commerce

- **Customer Experience:** Enhance the online shopping experience through personalized recommendations and targeted promotions.
- **Price Optimization:** Adjust pricing strategies based on market trends, competitor pricing, and demand data.
- **Fraud Detection:** Monitor transactions and user behavior to identify and prevent fraudulent activities.

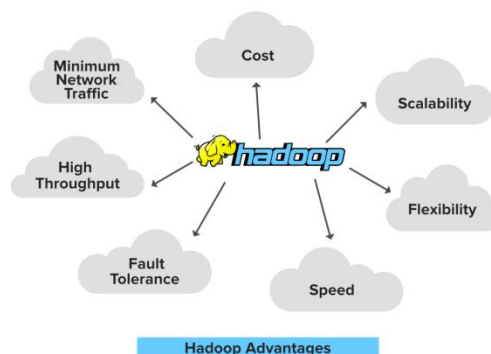
### 12. Research and Development

- **Scientific Research:** Analyze large datasets in fields like genomics, astronomy, and climate science to advance scientific knowledge.
- **Innovation:** Leverage data to drive innovation in product development and technological advancements.

Big Data enables organizations to make informed decisions, improve efficiency, and create new opportunities across diverse sectors. Its applications continue to evolve as data collection and analysis technologies advance.

## INTRODUCTION TO HADOOP

Hadoop is an open-source framework designed to handle and process large datasets in a distributed computing environment. It was developed by the Apache Software Foundation and is widely used for Big Data processing. Here's an introduction to its key components and features:



### 1. What is Hadoop?

Hadoop is a framework that allows for the distributed processing of large data sets across clusters of computers using simple programming models. It provides scalability, fault tolerance, and efficiency for managing and analyzing vast amounts of data.



# Big Data Analytics: UNIT – 1

## 2. Core Components of Hadoop

### a. Hadoop Distributed File System (HDFS)

- **Description:** A distributed file system designed to store large files across multiple machines. HDFS breaks down files into blocks and stores multiple copies of these blocks across different nodes in the cluster.
- **Features:**
  - **Fault Tolerance:** Data is replicated across multiple nodes to prevent data loss.
  - **Scalability:** Easily scales by adding more nodes to the cluster.
  - **High Throughput:** Optimized for large-scale data processing.

### b. MapReduce

- **Description:** A programming model for processing large datasets in parallel across a Hadoop cluster. It involves two main stages:
  - **Map:** Processes input data and converts it into a set of key-value pairs.
  - **Reduce:** Aggregates and processes the key-value pairs produced by the Map stage to produce the final output.
- **Features:**
  - **Parallel Processing:** Distributes tasks across multiple nodes to improve processing speed.
  - **Fault Tolerance:** Automatically recovers from failures during the processing.

### c. YARN (Yet Another Resource Negotiator)

- **Description:** A resource management layer for Hadoop that handles the allocation of resources and scheduling of tasks across the cluster.
- **Features:**
  - **Resource Management:** Manages and allocates resources to different applications.
  - **Job Scheduling:** Coordinates the execution of tasks and jobs.

### d. Hadoop Common

- **Description:** A set of shared utilities and libraries used by other Hadoop modules. It includes tools and APIs necessary for the functioning of the Hadoop ecosystem.
- **Features:**
  - **Support Libraries:** Provides essential components needed for Hadoop's various functionalities.

## 3. Ecosystem and Tools

Hadoop's ecosystem includes a variety of tools and projects that complement its core components:

- **Apache HBase:** A NoSQL database that runs on top of HDFS and provides real-time read/write access to large datasets.

## Big Data Analytics: UNIT – 1

- **Apache Hive:** A data warehousing tool that provides an SQL-like interface for querying and managing data stored in Hadoop.
- **Apache Pig:** A high-level scripting language used for processing and analyzing large datasets.
- **Apache Spark:** A fast, in-memory data processing engine that can work alongside Hadoop to provide real-time data processing capabilities.
- **Apache Flume:** A distributed service for collecting, aggregating, and moving large amounts of log data.
- **Apache Sqoop:** A tool for transferring data between Hadoop and relational databases.
- **Apache Oozie:** A workflow scheduler system that manages Hadoop jobs.

### 4. Key Features

- **Scalability:** Hadoop clusters can easily scale by adding more nodes to handle larger data volumes.
- **Fault Tolerance:** Designed to handle hardware failures by replicating data and redistributing tasks.
- **Cost-Effectiveness:** Runs on commodity hardware, making it a cost-effective solution for managing large datasets.

### 5. Use Cases

- **Data Storage and Processing:** Managing and analyzing massive volumes of structured and unstructured data.
- **Log Analysis:** Processing server logs for monitoring and troubleshooting.
- **Data Warehousing:** Storing and querying large datasets for business intelligence.

Hadoop provides a robust framework for managing and analyzing Big Data, and its ecosystem continues to expand with new tools and technologies to address various data processing needs.

## HADOOP COMPONENTS

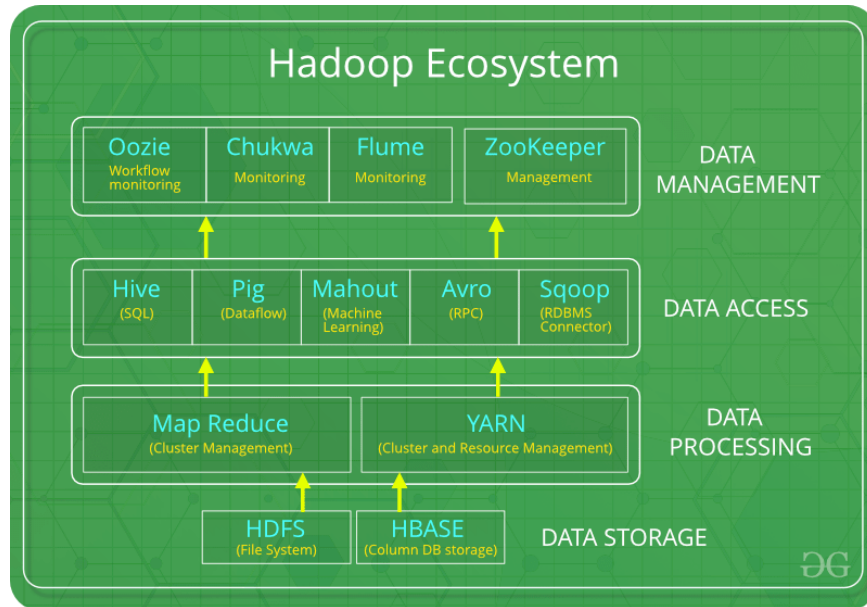
Hadoop has several key components that work together to provide a robust framework for processing and storing large datasets. Here's a detailed look at each component:

### 1. Hadoop Distributed File System (HDFS)

- **Function:** HDFS is the primary storage system in Hadoop, designed to store large files across multiple machines.
- **Key Features:**
  - **Block-Based Storage:** Files are split into fixed-size blocks (typically 128 MB or 256 MB) and distributed across a cluster of nodes.
  - **Replication:** Each block is replicated across multiple nodes to ensure fault tolerance and data reliability. The default replication factor is three.
  - **Fault Tolerance:** If a node fails, the data can be recovered from the replicated blocks on other nodes.

## Big Data Analytics: UNIT – 1

- **High Throughput:** Optimized for large-scale data access with high throughput rather than low latency.



### 2. MapReduce

- **Function:** MapReduce is a programming model used for processing and generating large datasets with a parallel, distributed algorithm.
- **Key Stages:**
  - **Map Stage:** Processes input data to produce a set of intermediate key-value pairs. Each mapper works on a portion of the data.
  - **Reduce Stage:** Aggregates and processes the intermediate key-value pairs produced by the mappers to produce the final output.
- **Key Features:**
  - **Parallel Processing:** Tasks are distributed across multiple nodes to improve processing efficiency.
  - **Fault Tolerance:** Automatically reassigns tasks if a node fails during processing.

### 3. Yet Another Resource Negotiator (YARN)

- **Function:** YARN is the resource management layer of Hadoop that manages and schedules resources for various applications.
- **Key Components:**
  - **ResourceManager:** Manages resource allocation across the cluster and schedules jobs.
  - **NodeManager:** Monitors resource usage on individual nodes and reports to the ResourceManager.
  - **ApplicationMaster:** Manages the lifecycle of applications, including job scheduling and resource allocation.
- **Key Features:**

## Big Data Analytics: UNIT – 1

- **Resource Allocation:** Manages resources across different applications and users.
- **Job Scheduling:** Handles scheduling and execution of tasks.

### 4. Hadoop Common

- **Function:** Provides a set of shared utilities, libraries, and tools that support other Hadoop components.
- **Key Features:**
  - **Support Libraries:** Includes necessary libraries for Hadoop's operations, such as serialization, configuration, and security.
  - **Utilities:** Provides essential tools for managing and interacting with the Hadoop ecosystem.

### 5. Hadoop Ecosystem Tools

The Hadoop ecosystem includes various tools and projects that enhance the functionality of Hadoop:

#### a. Apache HBase

- **Function:** A NoSQL database that runs on top of HDFS, providing real-time read/write access to large datasets.
- **Key Features:** Supports random, real-time read/write access; designed for scalability and high performance.

#### b. Apache Hive

- **Function:** A data warehousing tool that provides an SQL-like query language (HiveQL) for querying and managing data stored in Hadoop.
- **Key Features:** Simplifies data querying and analysis for users familiar with SQL

## CONFIGURATION OF HADOOP

Configuring Hadoop involves setting up its various components to ensure efficient and optimized operation. The configuration typically involves editing configuration files to define cluster settings, resource allocation, and various operational parameters. Here's a step-by-step guide to the basic configuration of Hadoop:

### 1. Prerequisites

- **Java:** Hadoop requires Java to run. Ensure that the Java Development Kit (JDK) is installed and configured on all nodes.
- **S :** Secure Shell (S ) needs to be set up for communication between nodes without requiring a password.

### 2. Hadoop Installation Directories

- **Hadoop Home Directory:** This is the directory where Hadoop is installed.
- **Data Directories:** Directories where HDFS stores its data and metadata.

### 3. Key Configuration Files

Hadoop has several key configuration files located in the conf or etc/hadoop directory:

#### a. core-site.xml

## Big Data Analytics: UNIT – 1

- **Purpose:** Defines core configuration settings, including the default file system and I/O settings.
- **Key Properties:**

xml

```
<configuration>
  <property>
    <name>fs.defaultFS</name>
    <value>hdfs://namenode:9000</value>
  </property>
  <property>
    <name>hadoop.tmp.dir</name>
    <value>/app/hadoop/tmp</value>
  </property>
</configuration>
```

### b. hdfs-site.xml

- **Purpose:** Configures HDFS-specific settings, including replication factor and data directory paths.
- **Key Properties:**

xml

```
<configuration>
  <property>
    <name>dfs.replication</name>
    <value>3</value>
  </property>
  <property>
    <name>dfs.namenode.name.dir</name>
    <value>file:///var/hadoop/namenode</value>
  </property>
  <property>
    <name>dfs.datanode.data.dir</name>
    <value>file:///var/hadoop/datanode</value>
  </property>
</configuration>
```

### c. mapred-site.xml

- **Purpose:** Contains MapReduce framework-specific settings.
- **Key Properties:**

## Big Data Analytics: UNIT – 1

xml

```
<configuration>
  <property>
    <name>mapreduce.framework.name</name>
    <value>yarn</value>
  </property>
</configuration>
```

### d. yarn-site.xml

- **Purpose:** Configures YARN-specific settings, including resource management and scheduling.
- **Key Properties:**

xml

```
<configuration>
  <property>
    <name>yarn.resourcemanager.hostname</name>
    <value>resourcemanager</value>
  </property>
  <property>
    <name>yarn.nodemanager.aux-services</name>
    <value>mapreduce_uffle</value>
  </property>
</configuration>
```

### e. hadoop-env.

- **Purpose:** Sets environment variables for Hadoop, such as the Java home directory.
- **Key Properties:**

```
export JAVA_HOME=/path/to/java
export HADOOP_HOME=/path/to/hadoop
export HADOOP_CONF_DIR=${HADOOP_HOME}/etc/hadoop
export HADOOP_HEAPSIZE=1024
```

## 4. Configuring S for Hadoop

- **Generate S Key:**

```
s -keygen -t rsa
```

## Big Data Analytics: UNIT – 1

- **Copy S Key to Nodes:**

s-copy-id user@hostname

### 5. Formatting the Namenode

Before starting the Hadoop cluster, the HDFS namenode needs to be formatted:

hdfs namenode -format

### 6. Starting and Stopping Hadoop Services

- **Start HDFS:**

start-dfs.

- **Start YARN:**

start-yarn.

- **Stop HDFS:**

stop-dfs.

- **Stop YARN:**

stop-yarn.

### 7. Verifying the Hadoop Installation

- **Access HDFS Web UI:** Typically available at <http://namenode:50070>
- **Access YARN Web UI:** Typically available at <http://resourcemanager:8088>

### 8. Configuring Resource Management and Monitoring

- **Resource Allocation:** Set properties in yarn-site.xml for resource allocation and limits.
- **Monitoring Tools:** Integrate monitoring tools like Apache Ambari, Cloudera Manager, or third-party monitoring solutions.

### 9. Tuning Performance

- **Heap Sizes:** Adjust Java heap sizes for various Hadoop daemons in hadoop-env. .
- **Block Size:** Configure the HDFS block size in hdfs-site.xml based on workload requirements.
- **Replication Factor:** Set an appropriate replication factor in hdfs-site.xml to balance data redundancy and storage efficiency.

By properly configuring these components and files, you can optimize Hadoop for your specific environment and workload, ensuring efficient data processing and storage.

## THE DESIGN OF HDFS

# Big Data Analytics: UNIT – 1

Hadoop Distributed File System (HDFS) is designed to store and manage large datasets efficiently across multiple nodes in a distributed computing environment. Its design principles emphasize scalability, reliability, and fault tolerance. Here's an overview of the key design aspects of HDFS:

## 1. Architecture

### a. NameNode

- **Role:** The master server that manages the filesystem namespace and controls access to files by clients.
- **Responsibilities:**
  - Maintains the directory tree of all files and folders in the filesystem.
  - Manages metadata (e.g., filenames, permissions, and the mapping of files to data blocks).
  - Keeps track of the locations of data blocks on DataNodes.

### b. DataNodes

- **Role:** The worker nodes that store the actual data.
- **Responsibilities:**
  - Store and retrieve blocks as instructed by the NameNode.
  - Report the status of the data blocks to the NameNode periodically.
  - Perform block creation, deletion, and replication based on the NameNode's instructions.

### c. Secondary NameNode

- **Role:** A helper node that periodically merges the namespace image with the edit logs to prevent the NameNode from becoming a single point of failure.
- **Responsibilities:**
  - Takes snapshots of the NameNode's metadata.
  - Helps in minimizing the downtime of the NameNode during failure recovery by providing a recent checkpoint.

## 2. Data Storage and Replication

### a. Blocks

- **Definition:** Files in HDFS are split into large blocks (default size: 128 MB) and distributed across DataNodes.
- **Advantages:**
  - Facilitates large-scale processing by allowing the system to read and write data in parallel.
  - Simplifies data management and recovery.

### b. Replication

- **Definition:** Each block is replicated across multiple DataNodes to ensure fault tolerance and data availability (default replication factor: 3).



## Big Data Analytics: UNIT – 1

- **Replication Strategy:**
  - One replica on a node in the local rack.
  - Another replica on a node in a different (remote) rack.
  - A third replica on a different node in the same remote rack.

### 3. Fault Tolerance

#### a. Data Integrity

- **Checksums:** HDFS creates checksums for data blocks and verifies data integrity during read/write operations.
- **Automatic Recovery:** If a corrupted block is detected, HDFS retrieves another copy from a different DataNode.

#### b. Heartbeats and Block Reports

- **Heartbeats:** DataNodes send regular heartbeats to the NameNode to confirm their availability.
- **Block Reports:** DataNodes send block reports to the NameNode periodically, detailing the blocks they store.

#### c. Rebalancing

- **Load Balancing:** HDFS can redistribute blocks across DataNodes to balance the storage utilization and optimize performance.

### 4. High Throughput and Scalability

#### a. Large Block Size

- **Efficient Data Transfers:** The large block size minimizes the overhead of metadata management and optimizes data transfer rates.
- **Batch Processing:** Designed to support high-throughput data access suitable for batch processing jobs like MapReduce.

#### b. Horizontal Scaling

- **Scalability:** HDFS can scale horizontally by adding more DataNodes to the cluster, allowing it to handle growing datasets seamlessly.

### 5. File System Namespace

#### a. Hierarchical Structure

- **Directories and Files:** HDFS organizes files in a hierarchical directory structure similar to traditional file systems.
- **Metadata Storage:** Metadata such as file permissions, modification timestamps, and block locations are managed by the NameNode.

#### b. Write-Once, Read-Many

- **Design Philosophy:** HDFS follows a write-once, read-many access model, allowing for simple coherency semantics and optimized for batch processing workloads.

### 6. Data Access and Client Interaction

#### a. Client Interaction

## Big Data Analytics: UNIT – 1

- **File Operations:** Clients interact with the NameNode to perform file operations like create, delete, and rename.
- **Data Read/Write:** Clients communicate directly with DataNodes to read/write data blocks, reducing the load on the NameNode.

### b. API and Interfaces

- **File System API:** HDFS provides a POSIX-like file system interface for seamless integration with applications.
- **Streaming Data Access:** Supports high-throughput streaming data access, making it ideal for large-scale data processing applications.

## 7. Security and Access Control

### a. User Authentication

- **Kerberos:** HDFS can integrate with Kerberos for secure user authentication.
- **Access Control Lists (ACLs):** Provides granular access control for files and directories.

### b. Data Encryption

- **Encryption at Rest:** Supports encryption of data stored in HDFS to protect sensitive information.
- **Encryption in Transit:** Ensures data is encrypted during transfer between clients and DataNodes.

HDFS's design aims to provide a reliable, scalable, and high-performance distributed file system for Big Data applications. Its architecture and features make it well-suited for environments requiring large-scale data storage and processing.

## HDFS CONCEPTS

Hadoop Distributed File System (HDFS) is designed to handle large data sets reliably and efficiently by distributing data across multiple nodes in a cluster. Here are the key concepts and features of HDFS:

### 1. Blocks

- **Definition:** HDFS stores each file as a series of blocks. These blocks are the basic units of storage.
- **Block Size:** The default block size is typically 128 MB, but it can be configured.
- **Advantages:** Large block sizes minimize the overhead of metadata management and support efficient data streaming.

### 2. Replication

- **Definition:** HDFS replicates data blocks to ensure fault tolerance and high availability.
- **Replication Factor:** The default replication factor is 3, meaning each block is stored on three different nodes.
- **Replica Placement:** One replica is stored on a node in the local rack, another on a node in a different rack, and the third on a different node in the same remote rack.

# Big Data Analytics: UNIT – 1

## 3. NameNode and DataNode

- **NameNode:**
  - **Role:** The master node that manages the filesystem namespace and regulates access to files.
  - **Responsibilities:** Maintains metadata, including the directory structure, file permissions, and the mapping of blocks to DataNodes.
  - **Fault Tolerance:** The NameNode is a single point of failure (though this can be mitigated with a Secondary NameNode or High Availability setup).
- **DataNode:**
  - **Role:** The worker nodes that store actual data blocks.
  - **Responsibilities:** Store and retrieve blocks as instructed by the NameNode, perform block creation, deletion, and replication, and periodically report the status of blocks to the NameNode.

## 4. Secondary NameNode

- **Role:** A helper node that performs regular checkpoints of the NameNode's namespace metadata.
- **Responsibilities:** Merges the namespace image with the edit logs to produce a new checkpoint, which helps in minimizing the downtime of the NameNode during recovery.

## 5. High Availability (HA)

- **Concept:** Provides a failover mechanism to ensure the continuous availability of the NameNode.
- **Implementation:** Typically involves two NameNodes (Active and Standby) and a pair of JournalNodes for merging metadata updates.

## 6. Heartbeats and Block Reports

- **Heartbeats:** DataNodes send regular heartbeats to the NameNode to confirm their availability.
- **Block Reports:** DataNodes send periodic reports listing all blocks they store, allowing the NameNode to maintain an accurate mapping.

## 7. Data Integrity

- **Checksums:** HDFS calculates and stores checksums for data blocks to detect and recover from data corruption.
- **Verification:** During data retrieval, checksums are verified to ensure data integrity.

## 8. Data Pipelining

- **Concept:** When writing data, HDFS uses pipelining to distribute data blocks to multiple DataNodes in a sequential manner.
- **Process:** The client writes the first block to the first DataNode, which forwards it to the second DataNode, and so on.

## 9. Rack Awareness

## Big Data Analytics: UNIT – 1

- **Definition:** HDFS is rack-aware and uses the network topology to place replicas to improve data reliability and availability.
- **Rack Awareness:** Helps in optimizing data transfer rates and fault tolerance by placing replicas across different racks.

### 10. File System Namespace

- **Structure:** HDFS provides a hierarchical file system namespace similar to traditional file systems.
- **Operations:** Supports operations like create, delete, rename, and move files and directories.

### 11. Write Once, Read Many (WORM)

- **Concept:** HDFS follows a write-once, read-many model for files, meaning once a file is written, it cannot be modified.
- **Advantages:** Simplifies data coherency and supports efficient data streaming and batch processing.

### 12. Federation

- **Concept:** Allows multiple independent NameNodes to share the same DataNodes to scale the namespace horizontally.
- **Benefits:** Enhances scalability by partitioning the namespace across multiple NameNodes.

### 13. Snapshots

- **Definition:** HDFS supports snapshots, which are read-only point-in-time copies of the filesystem or parts of it.
- **Usage:** Useful for backup, recovery, and data archiving.

### 14. Quotas

- **Function:** HDFS can enforce quotas on the number of files and directories as well as storage space usage.
- **Purpose:** Helps in managing resources and preventing a single user or application from consuming excessive resources.

### 15. Data Encryption

- **Encryption at Rest:** Data stored in HDFS can be encrypted to protect it from unauthorized access.
- **Encryption in Transit:** Ensures that data is encrypted during transfer between clients and DataNodes to secure data communication.

HDFS is designed to handle large volumes of data with high reliability and performance, making it a crucial component of the Hadoop ecosystem for big data processing and storage.

## THE COMMAND LINE INTERPRETER

The command-line interpreter for HDFS, also known as the Hadoop shell, allows users to interact with the Hadoop Distributed File System (HDFS) using a set of commands. These commands are

## Big Data Analytics: UNIT – 1

similar to Unix shell commands and provide functionalities for managing files and directories in HDFS. Here's an overview of the key commands and their usage:

### Basic HDFS Commands

#### 1. hdfs dfs

- **Description:** The primary command to interact with HDFS.
- **Syntax:** hdfs dfs [options] <args>

### File and Directory Management

#### 2. ls

- **Description:** Lists the contents of a directory.
- **Syntax:** hdfs dfs -ls <path>
- **Example:**

```
hdfs dfs -ls /user/hadoop
```

#### 3. mkdir

- **Description:** Creates a directory in HDFS.
- **Syntax:** hdfs dfs -mkdir <path>
- **Example:**

```
hdfs dfs -mkdir /user/hadoop/newdir
```

#### 4. mkdir -p

- **Description:** Creates a directory and any necessary parent directories.
- **Syntax:** hdfs dfs -mkdir -p <path>
- **Example:**

```
hdfs dfs -mkdir -p /user/hadoop/newdir/subdir
```

#### 5. rmdir

- **Description:** Removes an empty directory.
- **Syntax:** hdfs dfs -rmdir <path>
- **Example:**

```
hdfs dfs -rmdir /user/hadoop/newdir
```

### File Operations

#### 6. put

- **Description:** Copies files from the local filesystem to HDFS.

## Big Data Analytics: UNIT – 1

- **Syntax:** `hdfs dfs -put <local_src> <dest>`
- **Example:**

```
hdfs dfs -put localfile.txt /user/hadoop/
```

### 7. copyFromLocal

- **Description:** Copies files from the local filesystem to HDFS (same as put).
- **Syntax:** `hdfs dfs -copyFromLocal <local_src> <dest>`
- **Example:**

```
hdfs dfs -copyFromLocal localfile.txt /user/hadoop/
```

### 8. get

- **Description:** Copies files from HDFS to the local filesystem.
- **Syntax:** `hdfs dfs -get <src> <local_dest>`
- **Example:**

```
hdfs dfs -get /user/hadoop/file.txt localfile.txt
```

### 9. copyToLocal

- **Description:** Copies files from HDFS to the local filesystem (same as get).
- **Syntax:** `hdfs dfs -copyToLocal <src> <local_dest>`
- **Example:**

```
hdfs dfs -copyToLocal /user/hadoop/file.txt localfile.txt
```

### 10. moveFromLocal

- **Description:** Moves files from the local filesystem to HDFS.
- **Syntax:** `hdfs dfs -moveFromLocal <local_src> <dest>`
- **Example:**

```
hdfs dfs -moveFromLocal localfile.txt /user/hadoop/
```

### 11. moveToLocal

- **Description:** Moves files from HDFS to the local filesystem.
- **Syntax:** `hdfs dfs -moveToLocal <src> <local_dest>`
- **Example:**

## Big Data Analytics: UNIT – 1

```
hdfs dfs -moveToLocal /user/hadoop/file.txt localfile.txt
```

### 12. cat

- **Description:** Displays the contents of a file.
- **Syntax:** hdfs dfs -cat <path>
- **Example:**

```
hdfs dfs -cat /user/hadoop/file.txt
```

### 13. rm

- **Description:** Deletes files from HDFS.
- **Syntax:** hdfs dfs -rm <path>
- **Example:**

```
hdfs dfs -rm /user/hadoop/file.txt
```

### 14. rm -r

- **Description:** Recursively deletes a directory and its contents.
- **Syntax:** hdfs dfs -rm -r <path>
- **Example:**

```
hdfs dfs -rm -r /user/hadoop/dir
```

## File and Directory Information

### 15. du

- **Description:** Displays disk usage of files and directories.
- **Syntax:** hdfs dfs -du <path>
- **Example:**

```
hdfs dfs -du /user/hadoop/
```

### 16. df

- **Description:** Displays the filesystem disk space usage.
- **Syntax:** hdfs dfs -df
- **Example:**

```
hdfs dfs -df
```

## Big Data Analytics: UNIT – 1

### 17. stat

- **Description:** Prints statistics about a file or directory.
- **Syntax:** `hdfs dfs -stat <path>`
- **Example:**

```
hdfs dfs -stat /user/hadoop/file.txt
```

### 18. count

- **Description:** Counts the number of directories, files, and bytes under a given path.
- **Syntax:** `hdfs dfs -count <path>`
- **Example:**

```
hdfs dfs -count /user/hadoop/
```

### Permissions and Owner ip

### 19. chown

- **Description:** Changes the owner of a file or directory.
- **Syntax:** `hdfs dfs -chown <owner>[:<group>] <path>`
- **Example:**

```
hdfs dfs -chown hadoop:supergroup /user/hadoop/
```

### 20. chmod

- **Description:** Changes the permissions of a file or directory.
- **Syntax:** `hdfs dfs -chmod <mode> <path>`
- **Example:**

```
hdfs dfs -chmod 755 /user/hadoop/
```

### 21. chgrp

- **Description:** Changes the group association of a file or directory.
- **Syntax:** `hdfs dfs -chgrp <group> <path>`
- **Example:**

```
hdfs dfs -chgrp supergroup /user/hadoop/
```

### Advanced Commands

### 22. fsck



## Big Data Analytics: UNIT – 1

- **Description:** Runs a filesystem check.
- **Syntax:** `hdfs fsck <path>`
- **Example:**

```
hdfs fsck /user/hadoop/
```

### 23. setrep

- **Description:** Changes the replication factor of files.
- **Syntax:** `hdfs dfs -setrep <replication> <path>`
- **Example:**

```
hdfs dfs -setrep 2 /user/hadoop/file.txt
```

### 24. tail

- **Description:** Displays the last kilobyte of the file.
- **Syntax:** `hdfs dfs -tail <path>`
- **Example:**

```
hdfs dfs -tail /user/hadoop/file.txt
```

### Summary

The HDFS command-line interpreter provides a powerful set of commands for managing files and directories in the Hadoop Distributed File System. These commands cover a wide range of operations, from basic file manipulations to more advanced administrative tasks, allowing users to effectively interact with and manage their data within HDFS.

## BASIC FILE SYSTEM OPERATIONS

Basic file system operations in HDFS (Hadoop Distributed File System) involve managing files and directories, such as creating, listing, copying, moving, deleting, and viewing files. These operations can be performed using the `hdfs dfs` command-line interface, which provides a set of commands similar to Unix file system commands. Here are the basic file system operations in HDFS:

### 1. Listing Files and Directories

#### ls

- **Description:** Lists the contents of a directory.
- **Syntax:** `hdfs dfs -ls <path>`
- **Example:**

## Big Data Analytics: UNIT – 1

```
hdfs dfs -ls /user/hadoop
```

### ls -R

- **Description:** Recursively lists the contents of a directory.
- **Syntax:** hdfs dfs -ls -R <path>
- **Example:**

```
hdfs dfs -ls -R /user/hadoop
```

## 2. Creating Directories

### mkdir

- **Description:** Creates a directory in HDFS.
- **Syntax:** hdfs dfs -mkdir <path>
- **Example:**

```
hdfs dfs -mkdir /user/hadoop/newdir
```

### mkdir -p

- **Description:** Creates a directory and any necessary parent directories.
- **Syntax:** hdfs dfs -mkdir -p <path>
- **Example:**

```
hdfs dfs -mkdir -p /user/hadoop/newdir/subdir
```

## 3. Copying Files

### put

- **Description:** Copies files from the local filesystem to HDFS.
- **Syntax:** hdfs dfs -put <local\_src> <dest>
- **Example:**

```
hdfs dfs -put localfile.txt /user/hadoop/
```

### copyFromLocal

- **Description:** Copies files from the local filesystem to HDFS (same as put).
- **Syntax:** hdfs dfs -copyFromLocal <local\_src> <dest>
- **Example:**

```
hdfs dfs -copyFromLocal localfile.txt /user/hadoop/
```

## Big Data Analytics: UNIT – 1

### get

- **Description:** Copies files from HDFS to the local filesystem.
- **Syntax:** `hdfs dfs -get <src> <local_dest>`
- **Example:**

```
hdfs dfs -get /user/hadoop/file.txt localfile.txt
```

### copyToLocal

- **Description:** Copies files from HDFS to the local filesystem (same as get).
- **Syntax:** `hdfs dfs -copyToLocal <src> <local_dest>`
- **Example:**

```
hdfs dfs -copyToLocal /user/hadoop/file.txt localfile.txt
```

## 4. Moving Files

### moveFromLocal

- **Description:** Moves files from the local filesystem to HDFS.
- **Syntax:** `hdfs dfs -moveFromLocal <local_src> <dest>`
- **Example:**

```
hdfs dfs -moveFromLocal localfile.txt /user/hadoop/
```

### moveToLocal

- **Description:** Moves files from HDFS to the local filesystem.
- **Syntax:** `hdfs dfs -moveToLocal <src> <local_dest>`
- **Example:**

```
hdfs dfs -moveToLocal /user/hadoop/file.txt localfile.txt
```

### mv

- **Description:** Moves files or directories within HDFS.
- **Syntax:** `hdfs dfs -mv <src> <dest>`
- **Example:**

```
hdfs dfs -mv /user/hadoop/file.txt /user/hadoop/backup/
```

## 5. Deleting Files and Directories

### rm

## Big Data Analytics: UNIT – 1

- **Description:** Deletes files from HDFS.
- **Syntax:** `hdfs dfs -rm <path>`
- **Example:**

```
hdfs dfs -rm /user/hadoop/file.txt
```

### **rm -r**

- **Description:** Recursively deletes a directory and its contents.
- **Syntax:** `hdfs dfs -rm -r <path>`
- **Example:**

```
hdfs dfs -rm -r /user/hadoop/dir
```

## 6. Viewing File Contents

### **cat**

- **Description:** Displays the contents of a file.
- **Syntax:** `hdfs dfs -cat <path>`
- **Example:**

```
hdfs dfs -cat /user/hadoop/file.txt
```

### **tail**

- **Description:** Displays the last kilobyte of the file.
- **Syntax:** `hdfs dfs -tail <path>`
- **Example:**

```
hdfs dfs -tail /user/hadoop/file.txt
```

## 7. Checking File and Directory Information

### **stat**

- **Description:** Prints statistics about a file or directory.
- **Syntax:** `hdfs dfs -stat <path>`
- **Example:**

```
hdfs dfs -stat /user/hadoop/file.txt
```

### **count**

- **Description:** Counts the number of directories, files, and bytes under a given path.

## Big Data Analytics: UNIT – 1

- **Syntax:** `hdfs dfs -count <path>`
- **Example:**

```
hdfs dfs -count /user/hadoop/
```

### **du**

- **Description:** Displays disk usage of files and directories.
- **Syntax:** `hdfs dfs -du <path>`
- **Example:**

```
hdfs dfs -du /user/hadoop/
```

### **df**

- **Description:** Displays the filesystem disk space usage.
- **Syntax:** `hdfs dfs -df`
- **Example:**

```
hdfs dfs -df
```

## 8. Changing Permissions and Owner ip

### **chown**

- **Description:** Changes the owner of a file or directory.
- **Syntax:** `hdfs dfs -chown <owner>[:<group>] <path>`
- **Example:**

```
hdfs dfs -chown hadoop:supergroup /user/hadoop/
```

### **chmod**

- **Description:** Changes the permissions of a file or directory.
- **Syntax:** `hdfs dfs -chmod <mode> <path>`
- **Example:**

```
hdfs dfs -chmod 755 /user/hadoop/
```

### **chgrp**

- **Description:** Changes the group association of a file or directory.
- **Syntax:** `hdfs dfs -chgrp <group> <path>`
- **Example:**

## Big Data Analytics: UNIT – 1

```
hdfs dfs -chgrp supergroup /user/hadoop/
```

These basic file system operations in HDFS provide a comprehensive set of tools for managing files and directories in a distributed environment, ensuring efficient data storage, access, and manipulation.

### EXPLAIN HADOOP FILE SYSTEM

The Hadoop File System (HDFS) is the primary storage system used by Hadoop applications. It is a distributed file system that provides high-throughput access to data and is designed to be highly fault-tolerant. HDFS is an integral part of the Apache Hadoop project and follows a master-slave architecture. Here's an in-depth explanation of HDFS:

#### Key Features of HDFS

1. **Scalability:**
  - Designed to store and manage large datasets, scaling to petabytes of data across thousands of nodes.
2. **Fault Tolerance:**
  - Data is replicated across multiple nodes to ensure reliability and fault tolerance. If one node fails, the data can still be accessed from another node with a replica.
3. **High Throughput:**
  - Optimized for large streaming read and write operations rather than random access, ensuring high throughput for data-intensive applications.
4. **Cost-Effectiveness:**
  - Runs on commodity hardware, reducing costs compared to traditional enterprise storage systems.

#### Architecture of HDFS

HDFS follows a master-slave architecture comprising the following key components:

1. **NameNode:**
  - **Role:** The master server that manages the file system namespace and regulates access to files.
  - **Responsibilities:** Maintains the directory tree of all files in the file system and tracks where across the cluster the file data is kept. It also records changes in the namespace or in the properties of the files and directories.
  - **Fault Tolerance:** The NameNode is a single point of failure, but high availability can be achieved using secondary NameNodes or standby NameNodes in an HA setup.
2. **DataNodes:**

## Big Data Analytics: UNIT – 1

- **Role:** Worker nodes that store and retrieve data blocks when requested by clients or the NameNode.
  - **Responsibilities:** Perform block creation, deletion, and replication upon instruction from the NameNode. They also send periodic heartbeats and block reports to the NameNode to confirm their health and report on the blocks they are storing.
3. **Secondary NameNode:**
- **Role:** Assists the NameNode by periodically merging the namespace image with the edit log to create a new image.
  - **Responsibilities:** Reduces the load on the NameNode and helps ensure that the NameNode restarts quickly by keeping the size of the edit log within bounds.

### HDFS Data Flow

1. **Data Storage:**
  - Data is stored in large blocks (default 128 MB) across the cluster.
  - Each block is replicated (default replication factor is 3) to different DataNodes for fault tolerance.
2. **Data Access:**
  - **Read Operation:** A client requests the NameNode for the locations of the blocks of a file. The NameNode responds with the list of DataNodes containing the blocks. The client then directly reads the blocks from the DataNodes.
  - **Write Operation:** A client requests the NameNode to create a file. The NameNode allocates DataNodes to host the replicas of the first block of the file. The client writes the data to the allocated DataNodes, which then replicate the blocks to other DataNodes.

### HDFS Operations

1. **File Creation:**
  - When a file is created, it is divided into blocks. These blocks are stored on DataNodes and replicated based on the replication factor. The metadata about the file and its block locations is stored in the NameNode.
2. **File Deletion:**
  - When a file is deleted, the NameNode removes the metadata associated with the file. The DataNodes later delete the blocks of the file during their regular block reporting process.
3. **File Reading:**
  - The client retrieves the block locations from the NameNode and then directly reads the blocks from the DataNodes.
4. **File Writing:**

## Big Data Analytics: UNIT – 1

- The client writes the blocks to a sequence of DataNodes. Once the primary DataNode receives the data, it writes it to its local storage and forwards the data to the next DataNode in the pipeline.

### Data Integrity and Reliability

#### 1. Replication:

- Ensures data availability even if some nodes fail. The replication factor can be configured, and the system can automatically replicate blocks to maintain the desired replication factor.

#### 2. Heartbeat and Block Reports:

- DataNodes send regular heartbeats and block reports to the NameNode to confirm their health and the status of the blocks they store. If a DataNode fails to send a heartbeat, it is considered failed, and its data blocks are replicated to other nodes.

#### 3. Rack Awareness:

- HDFS is rack-aware, meaning it understands the network topology and optimizes data placement for fault tolerance and network efficiency. Blocks are stored on different racks to ensure that data is available even if an entire rack fails.

### HDFS Commands

Some basic HDFS commands include:

- `hdfs dfs -ls /`: List files in the root directory.
- `hdfs dfs -mkdir /newdir`: Create a new directory.
- `hdfs dfs -put localfile.txt /newdir`: Copy a local file to HDFS.
- `hdfs dfs -get /newdir/localfile.txt .`: Copy a file from HDFS to the local filesystem.
- `hdfs dfs -rm /newdir/localfile.txt`: Delete a file in HDFS.

### Conclusion

HDFS is a robust and scalable distributed file system that plays a critical role in the Hadoop ecosystem. It is designed to handle large amounts of data efficiently, ensuring high throughput, fault tolerance, and reliability. By distributing data across multiple nodes and replicating it, HDFS ensures that data remains accessible even in the face of hardware failures.

## INTERFACES DATA FLOW

In the context of Hadoop and HDFS (Hadoop Distributed File System), the data flow interfaces refer to the various APIs and command-line interfaces that allow users and applications to interact with the Hadoop ecosystem. These interfaces enable data storage, retrieval, processing, and management. Here's an overview of the key interfaces and how data flows through them:

### 1. HDFS Client Interface

The HDFS client interface allows users to interact with the HDFS file system. It includes both command-line tools and APIs for programmatic access.



# Big Data Analytics: UNIT – 1

## Command-Line Interface (CLI)

- **Description:** Provides a set of commands for managing files and directories in HDFS.
- **Common Commands:**
  - `hdfs dfs -put localfile.txt /user/hadoop/:` Copies a local file to HDFS.
  - `hdfs dfs -get /user/hadoop/file.txt localfile.txt:` Retrieves a file from HDFS.
  - `hdfs dfs -ls /user/hadoop/:` Lists the contents of a directory in HDFS.

## API Interface

- **Description:** Provides programmatic access to HDFS for applications written in Java, Python, or other languages.
- **Key Classes and Methods (Java API):**
  - `FileSystem fs = FileSystem.get(conf);` Obtain a FileSystem object.
  - `fs.copyFromLocalFile(new Path("localfile.txt"), new Path("/user/hadoop/"));` Copy a file to HDFS.
  - `fs.copyToLocalFile(new Path("/user/hadoop/file.txt"), new Path("localfile.txt"));` Copy a file from HDFS.

## 2. MapReduce Framework

The MapReduce framework is a core component of Hadoop for processing large datasets in a distributed manner.

### Data Flow

1. **Input Data Splits:**
  - Input data is split into chunks called input splits, each processed by a separate Map task.
2. **Map Phase:**
  - Each Map task processes an input split and produces key-value pairs.
3. **Shuffle and Sort Phase:**
  - The intermediate key-value pairs are shuffled and sorted to group values by keys.
4. **Reduce Phase:**
  - Each Reduce task processes the grouped key-value pairs to produce the final output.
5. **Output:**
  - The final output is written to HDFS.

### Interfaces

- **Mapper Interface:**
  - Processes input splits and generates key-value pairs.
  - ```
public class MyMapper extends Mapper<LongWritable, Text, Text, IntWritable> {  
    ...  
}
```
- **Reducer Interface:**
  - Processes grouped key-value pairs to produce the final output.

## Big Data Analytics: UNIT – 1

- `public class MyReducer extends Reducer<Text, IntWritable, Text, IntWritable> {  
... }`

### 3. YARN (Yet Another Resource Negotiator)

YARN manages resources and schedules tasks in a Hadoop cluster.

#### Data Flow

1. **Resource Request:**
  - Applications request resources (CPU, memory) from the ResourceManager.
2. **Resource Allocation:**
  - The ResourceManager allocates resources to ApplicationMasters.
3. **Task Execution:**
  - The ApplicationMaster negotiates with NodeManagers to launch tasks on nodes.
4. **Task Monitoring:**
  - NodeManagers monitor the execution of tasks and report the status to the ResourceManager.

#### Interfaces

- **ResourceManager:**
  - Manages cluster resources and schedules tasks.
- **NodeManager:**
  - Manages resources on a single node and launches tasks.
- **ApplicationMaster:**
  - Manages the execution of a specific application.

### 4. Hive and Pig

Hive and Pig are high-level interfaces for processing and analyzing large datasets in Hadoop.

#### Hive

- **Description:** A data warehouse system that provides SQL-like querying capabilities.
- **Data Flow:**
  1. **Query Submission:** Users submit queries using HiveQL.
  2. **Compilation:** Queries are compiled into MapReduce or Tez jobs.
  3. **Execution:** The jobs are executed on the Hadoop cluster.
  4. **Result Retrieval:** Results are retrieved and presented to the user.

#### Pig

- **Description:** A high-level platform for creating data analysis programs using a scripting language called Pig Latin.
- **Data Flow:**
  1. **Script Submission:** Users write Pig Latin scripts.
  2. **Compilation:** Scripts are compiled into MapReduce jobs.
  3. **Execution:** The jobs are executed on the Hadoop cluster.
  4. **Result Retrieval:** Results are retrieved and presented to the user.

# Big Data Analytics: UNIT – 1

## 5. HBase

HBase is a NoSQL database that runs on top of HDFS and provides real-time read/write access to large datasets.

### Data Flow

1. **Data Ingestion:**
  - Data is ingested into HBase through APIs or bulk load operations.
2. **Data Storage:**
  - Data is stored in HBase tables, which are divided into regions and distributed across the cluster.
3. **Data Access:**
  - Data is accessed through HBase APIs, providing random read/write capabilities.

### Interfaces

- **Java API:**
  - `HTable table = new HTable(conf, "mytable");`: Access an HBase table.
  - `Get get = new Get(Bytes.toBytes("row1"));`: Create a Get object to retrieve data.
  - `Result result = table.get(get);`: Retrieve data from HBase.

### Conclusion

The Hadoop ecosystem provides a variety of interfaces for data flow, ranging from low-level file system operations to high-level data processing frameworks. Each interface is designed to handle specific types of data flow and processing requirements, enabling efficient and scalable data management and analysis.

### PARALLEL COPYING WITH DISTCP

distcp (distributed copy) is a powerful tool provided by Hadoop for efficiently copying large datasets in parallel across different clusters within HDFS or from one HDFS to another. It leverages the MapReduce framework to distribute the copy operation across multiple nodes, ensuring high throughput and fault tolerance.

### Key Features of distcp

1. **Scalability:** Utilizes MapReduce to parallelize the copying process across many nodes.
2. **Fault Tolerance:** Provides retries and fault-tolerant mechanisms to handle node failures during the copy process.
3. **Efficiency:** Handles large data transfers efficiently by dividing the task into multiple smaller tasks.
4. **Resumable:** Supports resuming partially completed copy operations.
5. **Preserves Attributes:** Can preserve file attributes such as permissions, block size, replication factor, and checksums.

### How distcp Works

## Big Data Analytics: UNIT – 1

distcp operates by creating a MapReduce job where each map task is responsible for copying a portion of the data. The copied data can be within the same HDFS, between different HDFS clusters, or from HDFS to other Hadoop-compatible file systems.

### Basic Usage of distcp

The basic syntax for distcp is:

```
hadoop distcp [options] <source_path> <target_path>
```

### Common Options and Examples

#### 1. Basic Copy

- Copies data from source to destination.

```
hadoop distcp hdfs://source_cluster/path/to/data hdfs://destination_cluster/path/to/data
```

#### 2. Preserve Attributes

- Preserves file attributes like permissions, owner ip, and timestamps.

```
hadoop distcp -p hdfs://source_cluster/path/to/data hdfs://destination_cluster/path/to/data
```

#### 3. Overwrite Target Files

- Overwrites files at the destination if they already exist.

```
hadoop distcp -overwrite hdfs://source_cluster/path/to/data  
hdfs://destination_cluster/path/to/data
```

#### 4. Update

- Only copies files that are missing or have been modified at the destination.

```
hadoop distcp -update hdfs://source_cluster/path/to/data  
hdfs://destination_cluster/path/to/data
```

#### 5. Skip CRC

- Skips checksum verification during the copy process.

```
hadoop distcp -skipcrccheck hdfs://source_cluster/path/to/data  
hdfs://destination_cluster/path/to/data
```

#### 6. Bandwidth Control

## Big Data Analytics: UNIT – 1

- Limits the bandwidth used by the copy operation.

```
hadoop distcp -bandwidth <bandwidth_in_MB_per_sec> hdfs://source_cluster/path/to/data  
hdfs://destination_cluster/path/to/data
```

### Advanced Usage of distcp

#### Copying Large Data Sets

When copying large data sets, it's essential to manage resources efficiently to avoid overloading the network or the clusters involved. You can control the number of map tasks and the bandwidth used by distcp.

```
hadoop distcp -m 100 -bandwidth 50 hdfs://source_cluster/large_data  
hdfs://destination_cluster/large_data
```

- -m 100: Specifies 100 map tasks.
- -bandwidth 50: Limits the bandwidth to 50 MB per second.

#### Incremental Data Transfer

For incremental data transfer, where only new or modified files are copied, use the -update option.

```
hadoop distcp -update hdfs://source_cluster/data hdfs://destination_cluster/data
```

#### Resuming a Failed Copy

If a distcp operation fails, you can resume it using the -i (ignore failures) option.

```
hadoop distcp -i hdfs://source_cluster/data hdfs://destination_cluster/data
```

#### Internals of distcp

1. **Job Setup:**
  - distcp initializes a MapReduce job with each map task responsible for copying a subset of the data.
2. **Input Splits:**
  - The input paths are divided into splits. Each split is assigned to a map task.
3. **Copying Process:**
  - Each map task copies its assigned files from the source to the destination. The copying is done in parallel, leveraging the distributed nature of Hadoop.
4. **Monitoring and Fault Tolerance:**

## Big Data Analytics: UNIT – 1

- distcp monitors the progress of the copy operation. If a task fails, it retries based on the configured number of attempts.

### 5. Completion:

- Once all map tasks complete successfully, the distcp job is marked as successful.

### Conclusion

distcp is a robust tool for parallel copying of large datasets within HDFS or between different HDFS clusters. By leveraging Hadoop's distributed architecture, distcp ensures efficient, scalable, and fault-tolerant data transfers. It is particularly useful for migrating data between clusters, backing up data, or synchronizing data across multiple Hadoop environments.