#### BASICS OF HADOOP - DATA FORMAT

Hadoop is an open-source framework for processing, storing, and analyzing large volumes of data in a distributed computing environment. It provides a reliable, scalable, and distributed computing system for big data.

## **Key Components:**

- Hadoop Distributed File System (HDFS): HDFS is the storage system of Hadoop, designed to store very large files across multiple machines.
- MapReduce: MapReduce is a programming model for processing and generating large datasets that can be parallelized across a distributed cluster of computers.
- YARN (Yet Another Resource Negotiator): YARN is the resource management layer of Hadoop, responsible for managing and monitoring resources in a cluster.

### **Advantages of Hadoop:**

- Scalability: Hadoop can handle and process vast amounts of data by distributing it across a cluster of machines.
- **Fault Tolerance:** Hadoop is fault-tolerant, meaning it can recover from failures, ensuring that data processing is not disrupted.
- Cost-Effective: It allows businesses to store and process large datasets cost-effectively, as it can run on commodity hardware.

#### Here are some basics:

#### 1. Data Storage in Hadoop:

 Hadoop uses the Hadoop Distributed File System (HDFS) to store data across multiple machines in a distributed fashion.
Data is divided into blocks (typically 128 MB or 256 MB in size), and each block is replicated across several nodes in the cluster for fault tolerance.

#### 2. Data Formats:

- Hadoop can work with various data formats, but some common ones include:
  - **Text**: Data is stored in plain text files, such as CSV or TSV.
  - **SequenceFile**: A binary file format optimized for Hadoop, suitable for storing key-value pairs.
  - Avro: A data serialization system that supports schema evolution. It's often used for complex data structures.
  - **Parquet**: A columnar storage format that is highly optimized for analytics workloads. It's efficient for both reading and writing.

#### 3. **Data Ingestion**:

Before analyzing data, you need to ingest it into Hadoop.
You can use tools like Apache Flume, Apache Sqoop, or simply copy data into HDFS using Hadoop commands.

### 4. Data Processing:

- Hadoop primarily processes data using a batch processing model. It uses a programming model called MapReduce to distribute the processing tasks across the cluster. You write MapReduce jobs to specify how data should be processed.
- In addition to MapReduce, Hadoop ecosystem also includes higher-level processing frameworks like Apache Spark, Apache Hive, and Apache Pig, which provide more userfriendly abstractions for data analysis.

## 5. Data Analysis:

Once data is processed, you can analyze it to gain insights.
This may involve running SQL-like queries (with Hive), machine learning algorithms (with Mahout or Spark MLlib), or custom data processing logic.

# 6. **Data Output**:

 After analysis, you can store the results back into Hadoop, or you can export them to other systems for reporting or further analysis.

## 7. Data Compression:

• Hadoop allows data compression to reduce storage requirements and improve processing speed. Common compression formats include Gzip, Snappy, and LZO.

### 8. Data Schema:

 When working with structured data, it's important to define a schema. Some formats like Avro and Parquet have built-in schema support. In other cases, you may need to maintain the schema separately.

## 9. Data Partitioning and Shuffling:

 During data processing, Hadoop can partition data into smaller chunks and shuffle it across nodes to optimize the processing pipeline.

## 10. Data Security and Access Control:

 Hadoop provides security mechanisms to control access to data and cluster resources. This includes authentication, authorization, and encryption.

### Step-by-step installation of Hadoop on a single Ubuntu machine.

Installing Hadoop on a single-node cluster is a common way to set up Hadoop for learning and development purposes. In this guide, I'll walk you through the step-by-step installation of Hadoop on a single Ubuntu

machine.

## **Prerequisites:**

- A clean installation of Ubuntu.
- · Java installed on your system.

Let's proceed with the installation:

#### **Step 1: Download Hadoop**

- 1. Visit the Apache Hadoop website (<a href="https://hadoop.apache.org">https://hadoop.apache.org</a>) and choose the Hadoop version you want to install. Replace X.Y.Z with the version number you choose.
- 2. Download the Hadoop distribution using wget or your web browser. For example:

bash

wget https://archive.apache.org/dist/hadoop/common/hadoop-X.Y.Z/hadoop-X.Y.Z.tar.gz

**Step 2: Extract Hadoop** 3. Extract the downloaded Hadoop tarball to your desired directory (e.g., /usr/local/):

bash

sudo tar -xzvf hadoop-X.Y.Z.tar.gz -C /usr/local/

**Step 3: Configure Environment Variables** 4. Edit your ~/.bashrc file to set up environment variables. Replace X.Y.Z with your Hadoop version:

bash

export HADOOP\_HOME=/usr/local/hadoop-X.Y.Z export PATH=\$PATH:\$HADOOP\_HOME/bin:\$HADOOP\_HOME/sbin

Apply these changes to your current shell:

bash

source ~/.bashrc

**Step 4: Edit Hadoop Configuration Files** 5. Navigate to the Hadoop configuration directory:

bash

cd \$HADOOP\_HOME/etc/hadoop

6. Edit the hadoop-env.sh file to specify the Java home directory. Add

the following line to the file, pointing to your Java installation:

bash export JAVA\_HOME=/usr/lib/jvm/default-java

7. Configure Hadoop's core-site.xml by editing it and adding the following XML snippet. This sets the Hadoop Distributed File System (HDFS) data directory:

8. Configure Hadoop's hdfs-site.xml by editing it and adding the following XML snippet. This sets the HDFS data and metadata directories:

**Step 5: Format the HDFS Filesystem** 9. Before starting Hadoop services, you need to format the HDFS filesystem. Run the following command:

bash hdfs namenode -format

**Step 6: Start Hadoop Services** 10. Start the Hadoop services using the following command:

bash start-all.sh **Step 7: Verify Hadoop Installation** 11. Check the running Hadoop processes using the jps command:bashjsp

You should see a list of Java processes running, including NameNode, DataNode, ResourceManager, and NodeManager.

**Step 8: Access Hadoop Web UI** 12. Open a web browser and access the Hadoop Web UI at <a href="http://localhost:50070/">http://localhost:50070/</a> (for HDFS) and <a href="http://localhost:8088/">http://localhost:8088/</a> (for YARN ResourceManager).

You have successfully installed Hadoop on a single-node cluster. You can now use it for learning and experimenting with Hadoop and MapReduce.

## DATA FORMAT - ANALYZING DATA WITH HADOOP

## **Data Formats in Hadoop:**

- **Text Files:** Simple plain text files, where each line represents a record.
- **Sequence Files:** Binary files containing serialized key/value pairs.
- **Avro:** A data serialization system that provides rich data structures in a compact binary format.
- Parquet: A columnar storage file format optimized for use with Hadoop.

### **Analyzing Data with Hadoop:**

- MapReduce Programming Model: Data analysis tasks in Hadoop are accomplished using the MapReduce programming model, where data is processed in two stages: the Map stage processes and sorts the data, and the Reduce stage performs summary operations.
- **Hive:** Hive is a data warehousing and SQL-like query language for Hadoop. It allows users to query and manage large datasets stored in Hadoop HDFS.
- **Pig:** Pig is a high-level platform and scripting language built on top of Hadoop, used for creating MapReduce programs for data analysis.

Analyzing data with Hadoop involves understanding the data format and structure, as well as using appropriate tools and techniques for processing and deriving insights from the data. Here are some key considerations when it comes to data format and analysis with Hadoop:

#### 1. Data Format:

• Structured Data: If your data is structured, meaning it follows a fixed schema, you can use formats like Avro, Parquet, or ORC.

- These columnar storage formats are efficient for large-scale data analysis and support schema evolution.
- Semi-Structured Data: Data in JSON or XML format falls into this category. Hadoop can handle semi-structured data, and tools like Hive and Pig can help you query and process it effectively.
- **Unstructured Data**: Text data, log files, and other unstructured data can be processed using Hadoop as well. However, processing unstructured data often requires more complex parsing and natural language processing (NLP) techniques.

## 2. Data Ingestion:

 Before you can analyze data with Hadoop, you need to ingest it into the Hadoop Distributed File System (HDFS) or another storage system compatible with Hadoop. Tools like Apache Flume or Apache Sqoop can help with data ingestion.

## 3. Data Processing:

 Hadoop primarily uses the MapReduce framework for batch data processing. You write MapReduce jobs to specify how data should be processed. However, there are also high- level processing frameworks like Apache Spark and Apache Flink that provide more user-

friendly abstractions and real-time processing capabilities.

## 4. Data Analysis:

- For SQL-like querying of structured data, you can use Apache Hive, which provides a SQL interface to Hadoop. Hive queries get translated into MapReduce or Tez jobs.
- Apache Pig is a scripting language specifically designed for data processing in Hadoop. It's useful for ETL (Extract, Transform, Load) tasks.
- For advanced analytics and machine learning, you can use Apache Spark, which provides MLlib for machine learning tasks, and GraphX for graph processing.

### 5. Data Storage and Compression:

· Hadoop provides various storage formats optimized for analytics (e.g., Parquet, ORC) and supports data compression to reduce storage requirements and improve processing speed.

## 6. Data Partitioning and Shuffling:

 Hadoop can automatically partition data into smaller chunks and shuffle it across nodes to optimize the processing pipeline.

### 7. Data Security and Access Control:

• Hadoop offers mechanisms for securing data and controlling access through authentication, authorization, and encryption.

### 8. Data Visualization:

To make sense of the analyzed data, you can use data visualization tools like Apache Zeppelin or integrate Hadoop with business intelligence tools like Tableau or Power BI.

# **9. Performance Tuning:**

· Hadoop cluster performance can be optimized through configuration settings and resource allocation. Understanding how to fine-tune these parameters is essential for efficient data analysis.

# 10. Monitoring and Maintenance:

• Regularly monitor the health and performance of your Hadoop cluster using tools like Ambari or Cloudera Manager. Perform routine maintenance tasks to ensure smooth operation.

Analyzing data with Hadoop involves a combination of selecting the right data format, processing tools, and techniques to derive meaningful insights from your data. Depending on your specific use case, you may need to choose different formats and tools to suit your needs.