

## **What is Hadoop and its features → Hadoop:**

☞ Hadoop is a powerful, scalable, and distributed computing framework designed to handle massive datasets across a cluster of commodity hardware. It offers a cost-effective solution for storing, processing, and analyzing structured and unstructured data. At its core, Hadoop empowers organizations to extract valuable insights and make data-driven decisions with ease.

☞ It is an Open-Source Data Management with scale out storage and distributed processing.

☞ In the era of big data, organizations face the challenge of effectively managing and analyzing vast volumes of information. Hadoop, an open-source framework, has emerged as a game-changer in the field of big data processing.

### **→ Features of Hadoop:**

☞ Distributed Storage and Processing: Hadoop's core component, the Hadoop Distributed File System (HDFS), distributes data across multiple machines in a cluster. This distributed storage model ensures high availability, fault tolerance, and scalability. Furthermore, Hadoop's MapReduce programming model allows for parallel processing of data, enabling faster and efficient analysis.

☞ Scalability: One of the key features of Hadoop is its ability to scale horizontally. Organizations can easily add more machines to the cluster as data volumes grow, without any disruption to the existing infrastructure. This scalability makes Hadoop an ideal choice for accommodating the ever-increasing demands of big data.

☞ Fault Tolerance: Hadoop's fault tolerance mechanism ensures that data remains available even in the presence of hardware or software failures. HDFS achieves fault tolerance by replicating data across multiple nodes in the cluster. If a node fails, the data can be retrieved from other replicas, ensuring data integrity and uninterrupted processing.

☞ Data Locality: Hadoop leverages the concept of data locality, which means that the processing tasks are performed on the nodes where the data resides. By minimizing data transfer across the network, Hadoop maximizes performance and reduces network congestion, leading to faster data processing.

☞ Ecosystem of Tools: Hadoop offers a rich ecosystem of tools and frameworks that extend its capabilities. These tools enhance Hadoop's versatility and make it accessible to a wide range of users with varying skill sets.

## **What are the Limitations of Hadoop ?**

- Hadoop has played a pivotal role in the storage and processing of massive datasets across various industries. However, like any technology, Hadoop is not without its limitations.
- Complexity and Learning Curve: Hadoop is a complex framework comprising several components such as HDFS (Hadoop Distributed File System), YARN (Yet Another Resource Negotiator), and MapReduce. Setting up and configuring Hadoop clusters requires specialized knowledge and expertise. The learning curve for Hadoop can be steep, making it challenging for organizations with limited resources to adopt and implement effectively.
- Scalability and Performance: Hadoop's distributed architecture allows it to scale horizontally by adding more nodes to the cluster. However, there are certain limitations to its scalability. Hadoop's performance can be affected when dealing with small files or workloads with high metadata overhead. Additionally, real-time processing and low-latency requirements may not be well-suited for Hadoop's batch-oriented processing model.
- Data Recovery and Reliability: Hadoop's fault tolerance mechanisms, such as data replication across multiple nodes, ensure data durability and availability. However, data recovery from hardware failures can be time-consuming and resource-intensive. The process of replicating data across nodes introduces additional overhead and increases storage requirements, which might not be ideal for organizations with limited resources or tight budgets.
- Complex Data Processing: While Hadoop's MapReduce paradigm allows for distributed processing of large datasets, it may not be the best fit for all types of data processing tasks. Complex analytical operations that require multiple iterations or iterative algorithms might not perform optimally in a MapReduce environment. This limitation has led to the development of alternative data processing frameworks such as Apache Spark, which provide faster and more flexible processing models.
- Storage Overhead: Hadoop's default storage system, HDFS, introduces storage overhead due to data replication and the need to maintain multiple copies of data across the cluster. This replication factor, while providing fault tolerance, can significantly impact storage costs, especially for organizations dealing with petabytes or exabytes of data. The overhead can become prohibitive for scenarios where storage efficiency is a primary concern.

## What are the difference between RDBMS vs Hadoop ?

→ Relational Database Management Systems (RDBMS) and Hadoop are two fundamental technologies with significant differences when it comes to data management and analysis. Here's a concise overview of the key distinctions between RDBMS and Hadoop:

### → Data Structure and Processing Model:

☞ RDBMS: Designed for structured data with predefined schemas and follows a fixed schema approach. It utilizes SQL for data manipulation and retrieval, emphasizing ACID properties for transactional integrity.

☞ Hadoop: Tailored for unstructured and semi-structured data, employing a distributed file system (HDFS) and a processing framework called MapReduce. Hadoop utilizes a schema-on-read approach, providing flexibility for diverse data formats without enforcing a predefined schema.

### → Scalability and Performance:

☞ RDBMS: Vertically scalable, meaning it can handle increasing workloads by adding more resources to a single server. Well-suited for transactional processing and provides high performance for structured data retrieval.

☞ Hadoop: Horizontally scalable, allowing it to scale by adding more nodes to a cluster. Ideal for handling massive amounts of data and distributed processing, offering high throughput for large-scale data analysis and batch processing.

### → Data Storage and Processing Paradigm:

☞ RDBMS: Utilizes a row-based storage approach, storing data in tables as rows. Optimized for transactional processing, supporting complex joins, aggregations, and relational operations.

☞ Hadoop: Relies on a distributed file system (HDFS) for data storage, distributing data across multiple nodes in a cluster. Data processing in Hadoop follows a batch-oriented approach (MapReduce), with parallel processing capabilities. Additionally, Hadoop supports alternative processing frameworks like Apache Spark for real-time and interactive processing.

### → Use Cases and Data Types:

☞ RDBMS: Excels in handling structured data, commonly used in transactional systems such as banking, e-commerce, and inventory management. Well-suited for scenarios requiring data consistency, strict schema enforcement, and complex queries.

☞ Hadoop: Specifically designed for large volumes of unstructured and semi-structured data, including log files, social media data, sensor data, and text documents. Widely used in data warehousing, data lakes, and big data analytics, providing scalability and parallel processing for massive datasets.

## **What is Big Data & Types of Big Data & their Significance ?**

→ Big Data, the massive volume of data generated at a rapid pace, holds immense potential for organizations. To harness its power, it's important to understand its types and significance. In this post, we explore the types of Big Data and their importance in driving insights and decision-making.

### → Types of Big Data:

☞ Structured Data: Structured data refers to organized and formatted information that fits neatly into predefined categories. It is typically stored in traditional databases and can be easily queried using structured query language (SQL).

☞ Unstructured Data: Unstructured data, on the other hand, does not have a predefined format or organization. It includes textual data, social media posts, images, videos, emails, sensor data, and more.

☞ Semi-Structured Data: Semi-structured data is a hybrid form that combines elements of both structured and unstructured data. It contains some organizational properties but does not adhere strictly to a predefined schema. Examples of semi-structured data include XML files, JSON documents, and log files. → Significance of Big Data Types:

☞ Enhanced Decision-making: Structured data enables organizations to analyze historical trends, identify customer preferences, and optimize operations. Unstructured data opens avenues for sentiment analysis, social media monitoring, and understanding customer sentiment. Semi-structured data helps in tracking machine-generated logs and IoT device data. The combination of these data types empowers decision-makers with comprehensive insights for strategic planning and informed actions.

☞ Improved Customer Understanding: By analyzing structured and unstructured data, businesses can uncover customer preferences, sentiment, and behavior patterns. The ability to harness data from multiple sources allows organizations to create a holistic view of their customers and deliver targeted, relevant interactions.

☞ Innovation and Competitive Advantage: The different types of Big Data hold the key to unlocking innovation and gaining a competitive edge. By analyzing structured, unstructured, and semi-structured data, organizations can identify market trends, detect emerging patterns, and anticipate customer needs. Big Data fuels innovation by providing the necessary insights to drive strategic initiatives and foster continuous improvement.

## **What are the 5 V's of Big Data ?**

→ Big Data is more than just a massive volume of information. It encompasses several key characteristics that define its complexity and value. These characteristics are commonly known as the 5 V's of Big Data: Volume, Velocity, Variety, Veracity, and Value. In this post, we will delve into each of these V's and explore their significance in unlocking the potential of Big Data.

### **→ Volume:**

☞ Volume refers to the sheer amount of data generated and collected. With the exponential growth of digital interactions, organizations are accumulating vast volumes of data like never before.

☞ This includes structured and unstructured data from diverse sources such as social media, IoT devices, and transaction records.

☞ The significance of volume lies in the potential insights hidden within this massive dataset.

☞ Analyzing large volumes of data allows organizations to identify patterns, trends, and correlations that can drive strategic decision-making and predictive analytics.

### **→ Velocity:**

☞ Velocity represents the speed at which data is generated and processed. In today's interconnected world, data is produced in real-time or near real-time.

☞ Social media updates, online transactions, sensor data, and streaming services are just a few examples of data sources with high velocity.

☞ The significance of velocity lies in the need for organizations to capture and process data quickly to extract timely insights.

☞ Real-time analytics and data processing enable businesses to respond rapidly to changing conditions, detect anomalies, and seize time-sensitive opportunities.

### **→ Variety:**

☞ Variety refers to the diverse types and formats of data available. Big Data encompasses structured, unstructured, and semi-structured data from multiple sources.

☞ This includes text, images, audio, video, social media posts, emails, and more.

☞ The significance of variety lies in the potential for organizations to gain a comprehensive understanding of their operations and customers.

☞ By analyzing and integrating various data types, organizations can derive deeper insights, uncover hidden relationships, and make informed decisions based on a holistic view of their data ecosystem.

## → Veracity

☞ Veracity refers to the reliability, accuracy, and trustworthiness of data. With the increasing volume and variety of data, ensuring data quality becomes a significant challenge. Data can be incomplete, inconsistent, or even intentionally misleading.

☞ The significance of veracity lies in the need for organizations to establish data governance and quality control processes.

☞ Validating and verifying data integrity is crucial to ensure that analysis and decision-making are based on accurate and reliable information.

## → Value:

☞ Value is perhaps the most important V of Big Data. It represents the ultimate goal of leveraging Big Data: extracting meaningful insights and generating tangible value for organizations.

☞ The significance of value lies in the ability to transform data into actionable information.

☞ By applying advanced analytics techniques, organizations can uncover patterns, identify trends, predict outcomes, and make data-driven decisions that drive efficiency, innovation, and competitive advantage.

☞ The value derived from Big Data lies in its potential to optimize processes, enhance customer experiences, and fuel strategic growth initiatives.

## **Explain Hadoop key components ?**

→ At the heart of Hadoop's power lies its key components, which work in tandem to provide a scalable and distributed computing environment. In this post, we will delve into the key components of Hadoop and explore their roles in handling big data.

→ Hadoop Distributed File System (HDFS):

HDFS is the primary storage system in Hadoop. It is designed to store large datasets across multiple machines in a distributed manner. The data is divided into blocks and replicated across different nodes for fault tolerance. HDFS ensures high throughput and provides fault tolerance, making it ideal for handling big data.

→ MapReduce:

MapReduce is a programming model and computational framework for processing large data sets in parallel. It divides the processing into two stages: map and reduce. The map phase processes input data and generates intermediate key-value pairs, while the reduce phase aggregates and summarizes the intermediate results to produce the final output. MapReduce allows for distributed and scalable processing of data across a Hadoop cluster.

→ YARN (Yet Another Resource Negotiator):

YARN is the resource management framework in Hadoop. It decouples the resource management and job scheduling functions from the MapReduce framework. YARN enables the processing of data with frameworks other than MapReduce, such as Apache Spark, Apache Hive, and Apache Flink. It efficiently allocates resources and manages the execution of applications in the cluster.

→ Hadoop Common:

Hadoop Common provides the libraries and utilities required by other Hadoop components. It includes the common utilities, such as input/output (I/O) operations, networking, and security, that are used across the Hadoop ecosystem. Hadoop Common ensures compatibility and standardization among different Hadoop distributions.

→ Hadoop Ecosystem:

Hadoop's strength lies in its rich ecosystem of tools and frameworks that extend its capabilities. Some prominent components of the Hadoop ecosystem include Apache Hive for data warehousing, Apache Pig for data flow scripting, Apache HBase for real-time read/write access to big data, Apache Spark for fast in-memory processing, and Apache ZooKeeper for distributed coordination.

## Explain the concept of InputFormat and OutputFormat in Hadoop ?

→ In the vast landscape of Hadoop, InputFormat and OutputFormat are crucial concepts that govern how data is ingested into and emitted from the Hadoop ecosystem. we will explore the concepts of InputFormat and OutputFormat in Hadoop and their significance in handling data input and output.

→ **InputFormat:** InputFormat is a class in Hadoop that defines how data is read and split into input records for processing by the MapReduce framework. It determines how data is divided into input splits, which are then assigned to different map tasks for parallel processing. InputFormat provides a set of methods and specifications for data processing, such as reading data from different file formats, handling compression, and defining record boundaries.

→ Hadoop offers various built-in InputFormats tailored to handle specific data sources, including:

☞ **TextInputFormat:** Processes plain text files, where each line is treated as a separate input record.

☞ **KeyValueInputFormat:** Handles key-value pair data, commonly used in scenarios like log files or key-value stores.

☞ **SequenceFileInputFormat:** Reads binary key-value pairs stored in a sequence file format.

☞ Developers can also create custom InputFormats to handle data from different sources or specialized file formats.

→ **OutputFormat:** OutputFormat, similar to InputFormat, is a class that defines how data is written or emitted from the MapReduce job. It specifies the format and destination for the output data generated by the reduce tasks. OutputFormat allows customization of the output behavior, such as controlling the file format, compression, and storage location of the output data.

→ Hadoop provides various built-in OutputFormats, including:

☞ **TextOutputFormat:** Writes the output data as plain text, with each key-value pair on a separate line.

☞ **KeyValueOutputFormat:** Stores key-value pairs in a simple text-based format.

☞ **SequenceFileOutputFormat:** Writes output data in a binary sequence file format.

☞ Similarly, developers can create custom OutputFormats to tailor the output data format based on their specific requirements.

→ **Customizing InputFormat and OutputFormat:**

☞ Hadoop's flexibility allows developers to define and use custom InputFormats and OutputFormats. This capability is particularly useful when working with specialized data sources, non-standard file formats, or complex data processing requirements.

☞ Developers can extend the InputFormat and OutputFormat classes, implementing their own logic for reading and writing data, respectively. This enables seamless integration of Hadoop with various data storage systems, databases, or APIs.

→ **Conclusion:**

InputFormat and OutputFormat are fundamental components in Hadoop's data processing workflow. InputFormat defines how data is read and split into input records, while OutputFormat specifies the format and destination for the output data generated by reduce tasks

## Explain Hadoop 1.x Architecture and its Daemons?

### ► Hadoop 1.x Architecture Overview:

- ☞ Hadoop 1.x, the earlier version of the Hadoop framework, laid the foundation for the distributed processing of big data. It introduced a scalable and fault-tolerant architecture that revolutionized data processing.
- ☞ The architecture of Hadoop 1.x consists of two core components: the Hadoop Distributed File System (HDFS) and the MapReduce processing framework. These components work together to store and process vast amounts of data across a cluster of commodity hardware.

### → HDFS:

- ☞ HDFS is a Hadoop Distributed FileSystem, where our BigData is stored using Commodity Hardware. It is designed to work with Large DataSets with default block size of 64MB.
- ☞ It follows a master-slave architecture, where the NameNode acts as the central coordinator and manages the file system metadata, while DataNodes store the actual data blocks.

### → MapReduce:

- ☞ MapReduce is a Distributed and Batch Processing Programming Model.
- ☞ MapReduce also uses Commodity Hardware to process “High Volume of Variety of Data at High Velocity” in a reliable and fault-tolerant manner.
- ☞ It comprises two main components: the JobTracker and the TaskTrackers.

### ► Hadoop 1.x Daemons:

Let's dive into the key daemons in Hadoop 1.x and understand their roles in the data processing pipeline:

### → NameNode:

- ☞ The NameNode is the centerpiece of the HDFS architecture. It manages the file system namespace and regulates access to files and directories. The NameNode keeps track of the metadata, such as file locations, permissions, and block information, in memory.
- ☞ There is a single instance of this process which runs on a cluster and that is on a master node.
- ☞ It used to store Metadata about DataNode like “How many blocks are stored in Data Nodes, Which DataNodes have data, DataNode Details, locations, timestamps etc.”
- ☞ This process reads all the metadata from a file named fsimage and keeps it in memory. After this process is started, it updates metadata for newly added or removed files in RAM.
- ☞ It periodically writes the changes in one file called edits as edit logs. This process is heart of HDFS, if it is down HDFS is not accessible any more.

→ **DataNode:**

- ⌚ DataNodes are responsible for storing and retrieving data in HDFS. They manage data blocks and respond to read and write requests from clients. DataNodes communicate with the NameNode to provide information about block locations and report their health status.
- ⌚ There are many instances of this process running on various slave nodes.
- ⌚ It is responsible for storing individual file blocks on the slave nodes in Hadoop cluster.
- ⌚ Based on the replication factor, a single block is replicated in multiple slave nodes (only if replication factor is > 1) to prevent data loss.
- ⌚ This process periodically sends heart beats to NameNode to make NameNode aware that slave process is up and running.

→ **JobTracker:**

- ⌚ The JobTracker is the master daemon responsible for job scheduling and resource management in MapReduce. It receives job submissions, breaks them into tasks, and assigns tasks to available TaskTrackers. The JobTracker monitors task progress and handles task failures and re-executions.
- ⌚ Only one instance of this process runs on a master node.
- ⌚ Any MapReduce job is submitted to Job Tracker first.
- ⌚ Job Tracker checks for the location of various file blocks used in MapReduce processing.
- ⌚ Then it initiates the separate tasks on various DataNodes where blocks are present by communicating with Task Tracker Daemons.

→ **TaskTracker:**

- ⌚ TaskTrackers are worker daemons that run on individual nodes in the cluster. They execute Map and Reduce tasks assigned by the JobTracker. TaskTrackers report progress and status updates to the JobTracker and handle task failures by reassigning them if necessary.
- ⌚ This process has multiple instances running on the slave nodes.
- ⌚ It receives the job information from Job Tracker and initiates a task on that slave node.
- ⌚ In most of the cases, Task Tracker initiates the task on the same node where the physical data block is present.
- ⌚ Same as DataNode daemon, this process also periodically sends heart beats to Job Tracker to make Job Tracker aware that slave process is running.

→ **SecondaryNameNode:**

- ⌚ The SecondaryNameNode assists the NameNode by performing periodic checkpoints of the file system metadata. It merges the edits log with the file system image to prevent a single point of failure. Despite its name, the SecondaryNameNode does not act as a backup for the NameNode.
- ⌚ For this also single instance of this process runs on a cluster. This process can run on a master node (for smaller clusters) or can run on a separate node (in larger clusters) depends on the size of the cluster.
- ⌚ One misinterpretation from name is “This is a backup Name Node” but it’s NOT.

☞ It manages the metadata for the NameNode. In the sense, it reads the information written in edit logs (by Name Node) and creates an updated file of current cluster metadata.

☞ Then it transfers that file back to Name Node so that fsimage file can be updated.

☞ So, whenever NameNode daemon is restarted it can always find updated information in fsimage file.

### **What are the Limitations of Hadoop 1.x Architecture?**

→ Hadoop 1.x, the earlier version of the Hadoop framework, revolutionized distributed data processing and became a prominent solution for handling big data. However, as the volume and complexity of data grew, certain limitations within Hadoop 1.x became apparent. In this post, we will discuss some of the key limitations of Hadoop 1.x

→ Scalability: One of the primary limitations of Hadoop 1.x was its scalability. In this version, the JobTracker served as a single point of failure and resource bottleneck. It controlled the entire cluster's resource allocation, scheduling, and monitoring, leading to limitations in handling large-scale data processing workloads efficiently. As the size of the cluster increased, the JobTracker's performance and capacity became a limiting factor.

→ Resource Management: Hadoop 1.x lacked efficient resource management capabilities. The JobTracker used a static slot-based allocation model, where resources were allocated in fixed-size slots. This approach did not adapt well to varying job requirements, resulting in underutilized resources or inefficient allocation in some cases. Additionally, there was no support for multi-tenancy, making it challenging to prioritize and allocate resources among different users or applications.

→ Limited Processing Models: Hadoop 1.x primarily supported the MapReduce processing model, which was highly suitable for batch processing but less effective for real-time or interactive applications. The lack of support for other data processing frameworks limited the flexibility and versatility of Hadoop 1.x in handling different types of workloads.

→ Job Dependency Management: Hadoop 1.x had limited support for handling job dependencies and complex workflows. It relied on manual scripting or external tools to manage the dependencies between multiple MapReduce jobs, making it cumbersome and error-prone. Coordinating and orchestrating complex data processing pipelines within Hadoop 1.x was a challenging task.

→ Inefficient Resource Utilization: Due to the static allocation model and lack of fine-grained resource management, Hadoop 1.x often suffered from inefficient resource utilization. Jobs with varying resource requirements could not efficiently utilize the available cluster resources, resulting in underutilization of compute power and longer job execution times.

## Explain Hadoop 2.x Architecture and its Daemons ? ► Hadoop 2.x Architecture Overview:

☞ Hadoop 2.x introduced significant architectural changes to address scalability, reliability, and flexibility requirements. In this post, we will delve into the architecture of Hadoop 2.x and explore the daemons that power its distributed processing capabilities.

☞ Hadoop 2.x architecture consists of several components working together to process and store big data. The key architectural components include Hadoop Distributed File System (HDFS), Yet Another Resource Negotiator (YARN), MapReduce

☞ Hadoop 2.x has some common and new APIs which can be easily integrated with any third party applications to work with Hadoop.

☞ It has new Java APIs and features in HDFS and MapReduce which are known as HDFS2 and MR2 respectively.

☞ New architecture has added the architectural features like HDFS High Availability and HDFS Federation.

☞ Hadoop 2.x is not using Job Tracker and Task Tracker daemons for resource management now on-wards but it is using YARN (Yet Another Resource Negotiator) for Resource Management.

### ► Hadoop 2.x Core Components:

→ Storage: Identify the machine where data needs to be stored & Store the data.

→ Processing: Identify the machine where processing needs to be done & actual execution.

### ► Hadoop 2.x Architecture:

Hadoop 2.x architecture consists of several components working together to process and store big data. The key architectural components include:

#### → Hadoop Distributed File System (HDFS):

☞ HDFS is the primary storage system of Hadoop. It is designed to handle large datasets by distributing them across multiple nodes in a cluster.

☞ The architecture follows a master-worker model, with a NameNode serving as the master and DataNodes as the workers.

☞ The NameNode manages metadata and tracks the location of data blocks, while DataNodes store the actual data.

#### → Yet Another Resource Negotiator (YARN):

☞ YARN is a major addition to Hadoop 2.x, serving as a resource management and job scheduling framework.

☞ It decouples resource management from MapReduce, enabling the support of various data processing engines.

☞ YARN consists of two key components: ResourceManager and NodeManager. The ResourceManager manages cluster resources and schedules applications, while the NodeManager runs on each node and manages resources at the individual node level.

→ MapReduce:

☞ MapReduce is a programming model and processing framework for distributed data processing in Hadoop.

☞ It divides large data sets into smaller chunks and processes them in parallel across the cluster.

☞ MapReduce consists of two phases: the map phase, where data is transformed into key-value pairs, and the reduce phase, where the results of the map phase are aggregated.

☞ In Hadoop 2.x, MapReduce runs as an application on top of YARN, utilizing its resource management capabilities.

► Hadoop 2.x Daemons:

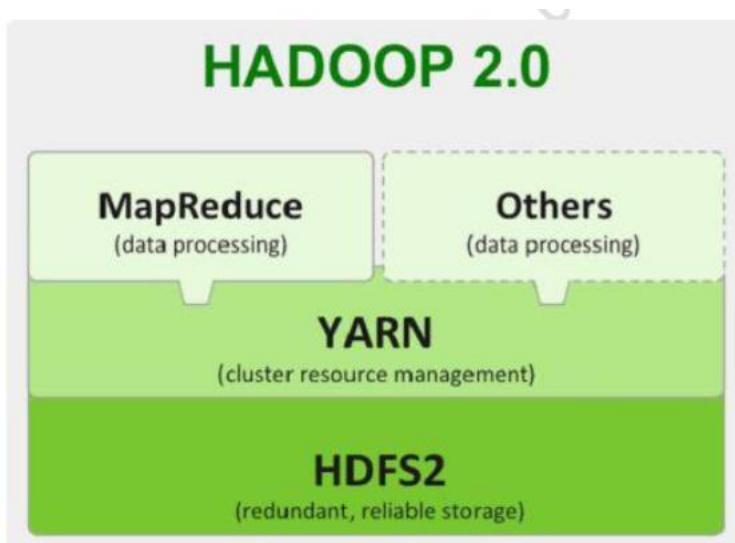
Let's dive into the key daemons in Hadoop 2.x and understand their roles in the data processing pipeline:

→ NameNode: The NameNode is the central metadata repository in HDFS. It keeps track of file system metadata, including file hierarchy, permissions, and block locations. The NameNode is a single point of failure, so it is typically configured in a high-availability mode with a standby NameNode ready to take over in case of failure.

→ DataNode: DataNodes are responsible for storing and retrieving data in HDFS. They communicate with the NameNode, report the status of blocks, and perform read/write operations on data blocks. DataNodes ensure data redundancy by replicating blocks across multiple nodes.

→ ResourceManager: The ResourceManager is the central authority for resource management in YARN. It receives resource requests from applications and coordinates the allocation of resources across the cluster. The ResourceManager also monitors the health of NodeManagers and restarts failed applications.

→ NodeManager: NodeManagers run on individual nodes in the cluster and manage resources at the node level. They monitor resource usage, report status to the ResourceManager, and launch and manage application containers.



## What is "map" and what is "reducer" in Hadoop ?

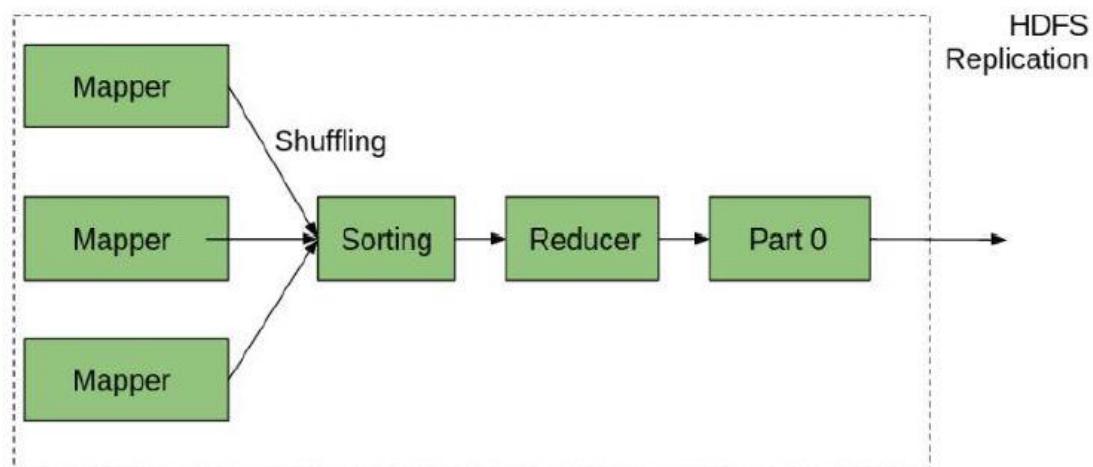
→ In the realm of Hadoop and its MapReduce framework, the terms "map" and "reducer" are integral concepts that define the processing logic for distributed data analysis. Understanding these components is crucial for harnessing the power of Hadoop and performing efficient and scalable data processing. In this post, we will explore what "map" and "reducer" mean in the context of Hadoop and their roles in the MapReduce paradigm.

### → What is the "Map" Function?

- ⌚ The "map" function is a fundamental component of the MapReduce framework in Hadoop.
- ⌚ It is responsible for processing input data in parallel across a distributed cluster. The map function takes individual input records and transforms them into a set of intermediate key-value pairs.
- ⌚ These key-value pairs serve as the input for the subsequent "reduce" phase.
- ⌚ The map function is highly flexible and can be customized based on specific data processing requirements.
- ⌚ It allows developers to write code that performs data transformations, filters, calculations, or any other necessary operations on the input data.
- ⌚ The map function can operate independently on each input record, making it suitable for parallel execution across multiple machines in the cluster.

### → What is the "Reducer" Function?

- ⌚ The "reducer" function is the second crucial component of the MapReduce framework.
- ⌚ It processes the intermediate key-value pairs generated by the map function and performs aggregations, summarizations, or any other required computations to produce the final output.
- ⌚ The reducer function takes a key and its associated set of values and produces a reduced set of output key-value pairs.
- ⌚ The reducer function works on a per-key basis, allowing it to group and process related data efficiently.
- ⌚ It receives all the values associated with a particular key, enabling operations such as data aggregation, statistical analysis, or any custom logic required for data processing.
- ⌚ Similar to the map function, the reducer function can be customized to suit the specific data processing needs of the application.



## Explain MapReduce Workflow in Hadoop ?

→ The MapReduce Workflow: In a Hadoop MapReduce job, the overall workflow can be summarized as follows:

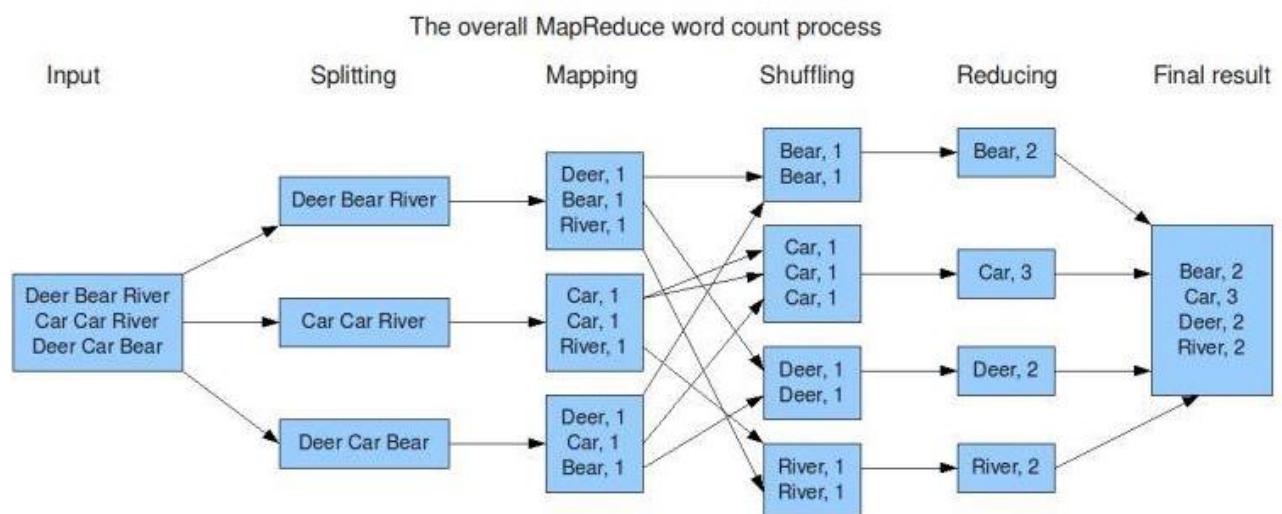
☞ Input Splitting: The input data is divided into smaller chunks called input splits, which are processed independently by different map tasks.

☞ Mapping: Each map task applies the map function to its assigned input split, generating intermediate key-value pairs.

☞ Shuffling and Sorting: The intermediate key-value pairs are sorted by key and grouped based on key-value pairs with the same key. This step prepares the data for the reduce phase.

☞ Reducing: Each reduce task processes a subset of the intermediate data, applying the reducer function to produce the final output.

☞ Output Generation: The final output is written to the desired output location, such as Hadoop Distributed File System (HDFS) or any other specified storage.



## What is the purpose of the MapReduce shuffle phase ?

→ The MapReduce shuffle phase is a crucial step in the Hadoop framework that occurs between the map and reduce phases. Although often overlooked, the shuffle phase serves a vital purpose in the distributed data processing paradigm.

The Objective & Significance of the Shuffle Phase:

☞ The MapReduce shuffle phase redistributes and sorts the intermediate key-value pairs generated by the map tasks.

☞ Its objective is to ensure that all values associated with the same key end up on the same reduce task.

☞ The shuffle phase minimizes data transfer by transmitting only the necessary data to each reduce task.

☞ It facilitates data locality by scheduling reduce tasks on nodes where the corresponding map output data is stored.

☞ The shuffle phase enables parallelism in reduce tasks by providing them with a sorted and grouped subset of the intermediate data.

☞ It can also include data compression to reduce network bandwidth utilization. ☞ The shuffle phase plays a vital role in optimizing data transfer, reducing network congestion, and enhancing the efficiency of Hadoop jobs.

## What is heartbeat in HDFS ?

→ In HDFS, a heartbeat refers to a periodic signal sent by certain entities within the system to convey their presence and operational state. It serves as a mechanism for nodes to report their status, availability, and functionality to the coordinating components. The heartbeat mechanism allows the HDFS cluster to monitor the health of its nodes and promptly respond to any failures or issues that may arise.

Significance of Heartbeat in HDFS:

⌚ Node Health Monitoring: Heartbeats enable the HDFS cluster to keep track of the health and availability of its constituent nodes. By regularly sending heartbeats, nodes inform the cluster about their operational state, including their ability to handle requests, availability of resources, and overall health. This information is vital for effective cluster management, fault detection, and recovery.

⌚ Failure Detection and Recovery: Heartbeats facilitate the detection of node failures or network partitions within the HDFS cluster. If a node fails to send a heartbeat within a specified time interval, it is considered unresponsive or failed. This triggers appropriate recovery mechanisms, such as reassigning the failed node's responsibilities to other healthy nodes or initiating data replication to maintain data redundancy and availability.

⌚ Load Balancing: Heartbeats also play a role in load balancing within the HDFS cluster. By monitoring the resource availability and utilization reported through heartbeats, the cluster can make informed decisions regarding task distribution and data placement. It allows the system to balance the workload across the available nodes, optimizing resource utilization and improving overall performance.

⌚ Cluster Coordination: Heartbeats establish a communication channel between the various components of the HDFS cluster. They enable coordination and synchronization between components such as the NameNode (metadata manager) and DataNodes (storage nodes). Heartbeats facilitate the exchange of important information, such as block reports, updates, and cluster-wide configurations, ensuring consistent and up-to-date cluster operation.

## What happens when a data node fails in hadoop ?

### → Data Node Failure in Hadoop:

☞ A data node failure refers to the situation when a node responsible for storing and serving data in a Hadoop cluster becomes unresponsive or fails. This can occur due to various reasons, including hardware issues, network problems, software failures, or unexpected shutdowns. When a data node fails, it affects the availability and accessibility of the data stored on that node.

### → Handling Data Node Failure:

Hadoop employs several mechanisms to handle data node failures and maintain the integrity of the distributed file system. The following steps outline the typical process when a data node fails:

☞ Heartbeat and Failure Detection: Hadoop utilizes a heartbeat mechanism where data nodes periodically send signals (heartbeats) to indicate their operational state. When a data node fails to send heartbeats within a specified time window, it is considered unresponsive or failed. The Hadoop cluster detects this failure through heartbeat monitoring.

☞ Block Replication: HDFS ensures fault tolerance by employing block replication. When a data node fails, the blocks it hosted are no longer accessible. However, since each block is replicated across multiple data nodes, the system can recover the lost data from the replicas. Hadoop automatically detects the failed node's absence and initiates the replication process for the affected blocks.

☞ Replication and Recovery: The NameNode, the central metadata manager in HDFS, maintains information about the block locations and replication factor for each file. When a data node fails, the NameNode identifies the missing replicas and instructs other available data nodes to create new replicas to meet the desired replication factor. This process ensures data redundancy and prevents data loss.

☞ Rebalancing: To restore data distribution and optimize performance, Hadoop may trigger a rebalancing process after a data node failure. Rebalancing involves redistributing the blocks across the remaining healthy data nodes to achieve a more balanced utilization of resources.

☞ Data Node Replacement: In some cases, if a data node is permanently lost or deemed unrecoverable, Hadoop allows for the addition of a new data node to replace the failed one. The replacement node joins the cluster and goes through the necessary data replication and rebalancing processes to restore the desired data redundancy and distribution.

## **What is the relation between job and task in Hadoop ?**

→ Jobs and Tasks in Hadoop: →

Jobs:

☞ A job in Hadoop represents a high-level unit of work that describes a specific data processing operation. It defines the overall data flow, input data, processing logic, and output specifications. Jobs are submitted to the Hadoop cluster for execution. Examples of job types in Hadoop include MapReduce jobs, Spark jobs, or Hive queries.

→ Tasks: Tasks are the building blocks of jobs in Hadoop. They represent the individual computational units that perform the actual data processing work. In Hadoop, there are two main types of tasks: map tasks and reduce tasks.

☞ Map Tasks: Map tasks take input data and apply a map function defined within the job. They process small subsets of the input data in parallel across the cluster. Map tasks are responsible for transforming the input data into intermediate key-value pairs.

☞ Reduce Tasks: Reduce tasks receive the intermediate key-value pairs generated by the map tasks. They perform the reduce function defined within the job, which involves aggregating, summarizing, or further processing the intermediate data to produce the final output.

The Relationship between Jobs and Tasks:

→ Job Configuration: Jobs in Hadoop are composed of one or more tasks. The number of tasks created for a job depends on various factors, such as the size of the input data, the configuration settings, and the available resources in the cluster.

→ Task Distribution: The Hadoop framework and its underlying scheduler distribute tasks across the cluster's available compute resources. This distribution enables parallel processing, as different tasks can be executed simultaneously on separate nodes within the cluster.

→ Task Dependencies: Tasks within a job can have dependencies, particularly in MapReduce jobs. Map tasks are typically independent and can process their assigned data subsets concurrently. However, reduce tasks depend on the output of map tasks.

→ Fault Tolerance: The relationship between jobs and tasks is crucial for fault tolerance in Hadoop. If a task fails during execution due to node failures or other issues, the framework automatically detects the failure and attempts to reassign the failed task to another available node. This fault tolerance mechanism ensures that the job can continue execution without losing progress and helps maintain data consistency and reliability.

## What are the main differences between Hadoop 1 and Hadoop 2 (YARN) ?

→ Differences between Hadoop 1 and Hadoop 2 (YARN):

→ Architecture:

☞ Hadoop 1: In Hadoop 1, the architecture consists of two main components: the Hadoop Distributed File System (HDFS) for storage and the MapReduce framework for processing. The JobTracker is responsible for managing resources and scheduling tasks on the cluster.

☞ Hadoop 2 (YARN): With the introduction of YARN in Hadoop 2, the architecture underwent a significant change. YARN separates resource management and job scheduling capabilities, allowing for greater flexibility. It introduces the ResourceManager and NodeManager as key components.

→ Resource Management:

☞ Hadoop 1: In Hadoop 1, the JobTracker performs resource management by maintaining information about available resources, tracking task progress, and handling task assignments. However, this centralized approach limited scalability and caused a single point of failure.

☞ Hadoop 2 (YARN): YARN introduces a decentralized approach to resource management. The ResourceManager in YARN is responsible for allocating resources to different applications, tracking their resource usage, and managing the scheduling of tasks. This distributed resource management model improves scalability, fault tolerance, and allows for multi-tenancy.

→ Job Scheduling:

☞ Hadoop 1: In Hadoop 1, the JobTracker performs both resource management and job scheduling. It assigns map and reduce tasks to available task trackers based on data locality and other factors. However, this approach limits the types of applications that can run on the cluster.

☞ Hadoop 2 (YARN): YARN separates resource management from job scheduling, enabling a more diverse set of applications to run on the cluster. YARN introduces a pluggable scheduler framework that allows different scheduling policies to be employed based on application requirements.

→ Cluster Utilization: ☞ Hadoop 1: In Hadoop 1, the cluster's resources are primarily dedicated to running MapReduce jobs. The JobTracker assigns resources to jobs based on their priority and availability, with limited flexibility for resource sharing.

☞ Hadoop 2 (YARN): YARN enables better cluster utilization by allowing multiple applications to coexist and share resources efficiently. Applications running on YARN can dynamically request and release resources based on their needs, leading to higher resource utilization and improved overall cluster efficiency.

