## What Is Mapreduce ?

MapReduce is a programming model and processing paradigm designed to process and analyze large-scale datasets in a distributed computing environment.

It was popularized by Google and further developed by Doug Cutting and Mike Cafarella, leading to its integration into the Apache Hadoop framework.

MapReduce is a fundamental concept in the world of Big Data and is widely used for batch processing of large datasets.

## The Mapreduce Model Consists Of Two Main Steps

- 1. Map phase
- 2. Reduce phase.

## Map Phase

- Map Function: The Map Function is a step in the process that takes a big set of data and applies a custom "Map" function to each element of the set. The Map function works on the data and produces intermediate key-value pairs.
- 2. Key-Value Pairs: The information produced by the Map phase consists of intermediate sets of key-value pairs. The key represents the outcome of a certain action on the input data, while the value can be any related data.

Shuffling and Sorting:

After the Map phase, the MapReduce framework performs a "shuffling and sorting" step.

During this phase, the data is separated into smaller groups based on a common factor and organized in order. This makes sure that all the related information with the same factor is grouped together in one place, which is important for the next part.

## **Reduce Phase**

- Reduce Function: During the Reduce phase, the framework employs a custom "Reduce" function to handle each group of intermediate key-value pairs that share the same key. The Reduce function analyzes the values linked to each key and generates a collection of output values for that particular key.
- 2. The end result: The end result of the MapReduce task is the collection of key-value pairs created by the Reduce function.

MapReduce is a method created to efficiently deal with big sets of data by dividing the data processing across several computers in a group. Each computer independently handles a part of the data, and the outcomes are combined to get the final result. This parallel handling helps MapReduce work well and quickly with large data processing tasks.

MapReduce is particularly well-suited for batch processing tasks, such as data transformations, aggregations, and filtering. It has been widely used for tasks like log processing, web indexing, data extraction, and many other batch-oriented data processing jobs.

Apache Spark is a newer alternative to MapReduce that can process data faster by utilizing computer memory. This makes it a more suitable choice for tasks that require quick and

efficient data processing. However, MapReduce still holds significance in the realm of Big Data as it has paved the way for the development of various distributed data processing frameworks and technologies.