

SNS COLLEGE OF TECHNOLOGY

Coimbatore-35 An Autonomous Institution

Accredited by NBA – AICTE and Accredited by NAAC – UGC with 'A++' Grade Approved by AICTE, New Delhi & Affiliated to Anna University, Chennai



DEPARTMENT OF MECHATRONICS ENGINEERING

UNIT 3 – HADOOP

VERTICAL SCALABILITY



VERTICAL SCALABILITY



Vertical scalability in Hadoop refers to the ability to increase the capacity of individual nodes within a cluster by adding more resources (such as CPU, RAM, or storage) to each node. This is in contrast to horizontal scalability, which involves adding more nodes to the cluster.

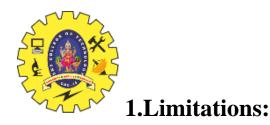
Vertical scalability can be beneficial in certain scenarios, especially when the existing hardware resources on individual nodes are not fully utilized or when specific tasks require more processing power or memory. However, it's important to note that vertical scalability has its limitations, and there are factors to consider when deciding whether to scale vertically in a Hadoop cluster:

1.Node Resource Upgrades:

1. Vertical scalability involves upgrading the resources (CPU, RAM, storage) of existing nodes within the cluster. This can be achieved by adding more powerful processors, increasing RAM capacity, or attaching additional storage to each node.

2.Use Cases for Vertical Scalability:

1. Vertical scalability is particularly useful when specific tasks or workloads within the Hadoop cluster require more computing power or memory. For example, if certain MapReduce jobs or data-intensive tasks are resource-intensive, upgrading individual nodes can improve the performance of those tasks.



VERTICAL SCALABILITY



1. There are practical limits to how much a node can be upgraded. For instance, there may be constraints on the maximum amount of RAM or number of CPU cores that a node can support. Eventually, these hardware limitations may be reached, and further scalability would require a combination of vertical and horizontal scaling.

2.Cost Considerations:

1. While upgrading individual nodes can enhance performance for specific tasks, it might not be as costeffective as horizontal scaling, especially when dealing with large-scale data processing. Adding more nodes (horizontal scaling) is often a more economical way to increase overall cluster capacity.

3.Impact on Cluster Configuration:

1. When upgrading the resources of individual nodes, it's crucial to consider the impact on the overall cluster configuration. The Hadoop ecosystem is designed to distribute and parallelize tasks across nodes, and a balanced cluster configuration is essential for optimal performance.



VERTICAL SCALABILITY



1.Dynamic Resource Allocation:

 Modern Hadoop distributions and resource management frameworks like YARN (Yet Another Resource Negotiator) provide dynamic resource allocation. This allows the cluster to adapt to changing workloads, making it more flexible in utilizing available resources.

2.Combining Vertical and Horizontal Scaling:

1. In some cases, a combination of vertical and horizontal scaling may be the most effective approach. For example, upgrading the resources of existing nodes while also adding more nodes to the cluster can provide a balanced and scalable solution.

In summary, vertical scalability in Hadoop involves upgrading the resources of individual nodes within a cluster to meet specific performance or workload requirements. While it can be useful in certain scenarios, careful consideration should be given to the practical limits of node upgrades, cost-effectiveness, and the impact on overall cluster performance and configuration. In many cases, a balanced approach that includes both vertical and horizontal scaling may be the most suitable solution.



INTRODUCTION TO HADOOP



How Hadoop Works:

1.Storage:

1. Data is stored in HDFS, which distributes the data across nodes in the cluster.

2.Processing:

1. MapReduce programming paradigm is used for distributed processing. Map tasks process data in parallel across the nodes, and reduce tasks aggregate the results.

3.Fault Tolerance:

1. Hadoop provides fault tolerance by replicating data across multiple nodes. If a node fails, tasks are redirected to other nodes with copies of the data.

Use Cases:

1.Batch Processing:

1. Hadoop is well-suited for batch processing of large volumes of data, making it a key tool for tasks like log analysis, data warehousing, and ETL (Extract, Transform, Load) processes.

2.Data Storage and Retrieval:

1. HDFS is used for storing vast amounts of data, and Hadoop provides mechanisms for efficient retrieval and processing.

3.Data Transformation:

1. Hadoop facilitates the transformation of raw data into a structured and usable format through its MapReduce paradigm. **4.Data Analysis and Exploration:**

1. Hadoop is commonly employed for data analysis tasks, enabling organizations to derive insights from large datasets.



INTRODUCTION TO HADOOP



Ecosystem:

- Hadoop has a rich ecosystem of related projects and tools that extend its capabilities. Some notable components include:
- •Apache Hive: Data warehousing and SQL-like query language for Hadoop.
- •Apache Pig: Platform for analyzing large datasets using a high-level scripting language.
- •Apache HBase: Distributed, scalable, and NoSQL database for Hadoop.
- •Apache Spark: In-memory data processing engine for faster and more flexible data processing.
- Challenges:
- •Complexity: Setting up and configuring a Hadoop cluster can be complex.
- •Skill Requirements: Users need to be familiar with Java and the Hadoop ecosystem tools.

•Data Movement Overhead: Data needs to be moved to the Hadoop cluster, which can incur additional overhead.

Evolution:

Hadoop has evolved over the years, with advancements such as the introduction of YARN, improvements in performance, and the integration of new tools like Apache Spark. As the big data landscape continues to evolve, Hadoop remains a foundational technology for many organizations dealing with large-scale data processing and analytics.