



SNS COLLEGE OF ENGINEERING
Kurumbapalayam (Po), Coimbatore – 641 107
AN AUTONOMOUS INSTITUTION



**Accredited by NAAC-UGC with 'A' Grade, Accredited by NBA
Approved by AICTE & Affiliated to Anna University, Chennai.**

Hashing in Distributed Systems

What is Hashing?

Hashing is the process of transforming input data (of any size) into a fixed-size string of characters, which is typically a sequence of numbers and letters. The output is called a **hash value** or **hash code**. Hash functions have several important properties:

1. **Deterministic:** The same input will always produce the same output.
2. **Fixed Size:** The output hash value has a constant size, regardless of the size of the input.
3. **Fast Computation:** Hash functions are designed to compute the hash value quickly.
4. **Pre-image Resistance:** It should be computationally infeasible to reverse the hash function (i.e., to obtain the original input from the hash value).
5. **Collision Resistance:** It should be hard to find two different inputs that produce the same hash value.

Applications of Hashing in Distributed Systems

1. **Data Integrity and Verification:**
 - Hashing is used to verify the integrity of data. By generating a hash value for data before it is transmitted or stored, the system can later check whether the data has been altered by comparing the hash values.
2. **Distributed Hash Tables (DHTs):**
 - DHTs are a key component of many peer-to-peer networks and are used for efficient data retrieval. They use hashing to map keys to nodes in the network, allowing for distributed storage and retrieval of data.
 - Examples include **Chord**, **Kademlia**, and **Pastry**.
3. **Load Balancing:**
 - Hashing can help evenly distribute workloads across nodes in a distributed system. By hashing requests or data, the system can



SNS COLLEGE OF ENGINEERING
Kurumbapalayam (Po), Coimbatore – 641 107
AN AUTONOMOUS INSTITUTION



**Accredited by NAAC-UGC with 'A' Grade, Accredited by NBA
Approved by AICTE & Affiliated to Anna University, Chennai.**

determine which node should handle a specific request, ensuring that no single node becomes a bottleneck.

4. Sharding:

- In distributed databases, hashing is used to determine how data is partitioned across multiple nodes (shards). Each record can be assigned to a shard based on a hash of its key, ensuring an even distribution of data.

5. Caching:

- Hashing can be used to manage cache entries. By hashing the URL or query of a request, a distributed caching layer can quickly determine whether the requested data is already cached and retrieve it efficiently.

6. Blockchain and Cryptography:

- Hashing is essential in blockchain technology for ensuring data integrity and security. Each block in a blockchain contains a hash of the previous block, creating a secure and immutable chain of data.

Hashing is a critical technique used in distributed databases to efficiently manage and retrieve data across multiple nodes. By utilizing hash functions, distributed databases can achieve load balancing, data partitioning, and fast access to records. Below is a comprehensive overview of how hashing is applied in distributed databases, including its benefits, methodologies, and challenges.

1. Role of Hashing in Distributed Databases

A. Data Partitioning (Sharding)

- **Sharding** is the process of distributing data across multiple nodes to ensure that each node holds only a portion of the total dataset.
- **Hashing** is commonly used to determine which shard (or node) should store a particular record. This is achieved by applying a hash function to a key (e.g., user ID, product ID) to produce a hash value that corresponds to a specific shard.

Example:



SNS COLLEGE OF ENGINEERING
Kurumbapalayam (Po), Coimbatore – 641 107
AN AUTONOMOUS INSTITUTION



**Accredited by NAAC-UGC with 'A' Grade, Accredited by NBA
Approved by AICTE & Affiliated to Anna University, Chennai.**

- Consider a distributed database with 4 shards. If the hash function is designed to return values between 0 and 3, the hash value can be used to map records:
 - Record with key A -> $\text{hash}(A) \% 4 = 1$ (stored in shard 1)
 - Record with key B -> $\text{hash}(B) \% 4 = 2$ (stored in shard 2)

B. Load Balancing

- Hashing helps in distributing the data evenly across shards, ensuring that no single shard is overwhelmed with too much data or too many requests. This balances the workload and enhances performance.

C. Fast Data Retrieval

- When a query is made for a specific key, the same hash function can be used to quickly determine which shard contains the desired record. This allows for faster access compared to searching through all nodes.

2. Hashing Techniques in Distributed Databases

A. Consistent Hashing

1. How Consistent Hashing Works

- **Hash Ring Structure:**
 - Consistent hashing uses a circular hash space (often referred to as a hash ring). Both nodes and data items (keys) are hashed to positions on this ring.
- **Node Placement:**
 - Each node in the system is assigned a position on the hash ring based on the hash of its identifier (e.g., IP address or node name). For instance, if a node hashes to position 150 on a 0-255 ring, it will be placed at that point.
- **Data Placement:**

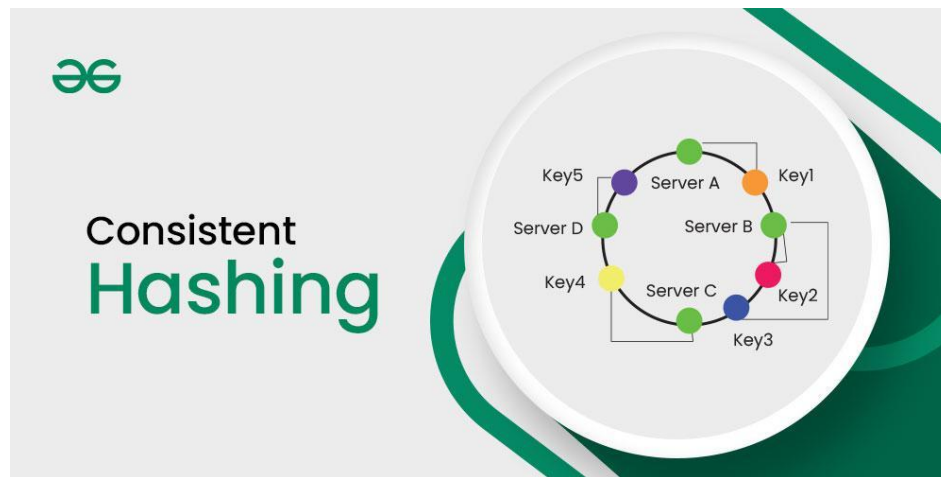


SNS COLLEGE OF ENGINEERING
Kurumbapalayam (Po), Coimbatore – 641 107
AN AUTONOMOUS INSTITUTION



**Accredited by NAAC-UGC with 'A' Grade, Accredited by NBA
Approved by AICTE & Affiliated to Anna University, Chennai.**

- When a data item is inserted, it is hashed to determine its position on the ring. The item is then stored on the first node that appears in a clockwise direction from its position. This allows for efficient storage of data items without needing to know the specific node beforehand.
- **Dynamic Node Management:**
 - When a new node is added or an existing node is removed, only a limited number of keys need to be redistributed, rather than all keys. This is a significant advantage over traditional hashing methods, where adding or removing a node often requires reshuffling the entire dataset.



Implementation of Consistent Hashing algorithm

1. **Choose a Hash Function:**
 - Select a hash function that produces a uniformly distributed range of hash values. Common choices include MD5, SHA-1, or SHA-256.
2. **Define the Hash Ring:**
 - Represent the range of hash values as a ring. This ring should cover the entire possible range of hash values and be evenly distributed.
3. **Assign Nodes to the Ring:**



SNS COLLEGE OF ENGINEERING
Kurumbapalayam (Po), Coimbatore – 641 107
AN AUTONOMOUS INSTITUTION



**Accredited by NAAC-UGC with 'A' Grade, Accredited by NBA
Approved by AICTE & Affiliated to Anna University, Chennai.**

- Assign each node in the system a position on the hash ring. This is typically done by hashing the node's identifier using the chosen hash function.
- 4. **Key Mapping:**
 - When a key needs to be stored or retrieved, hash the key using the chosen hash function to obtain a hash value.
 - Find the position on the hash ring where the hash value falls.
 - Walk clockwise on the ring to find the first node encountered. This node becomes the owner of the key.
- 5. **Node Additions:**
 - When a new node is added, compute its position on the hash ring using the hash function.
 - Identify the range of keys that will be owned by the new node. This typically involves finding the predecessor node on the ring.
 - Update the ring to include the new node and remap the affected keys to the new node.
- 6. **Node Removals:**
 - When a node is removed, identify its position on the hash ring.
 - Identify the range of keys that will be affected by the removal. This typically involves finding the successor node on the ring.
 - Update the ring to exclude the removed node and remap the affected keys to the successor node.
- 7. **Load Balancing:**
 - Periodically check the load on each node by monitoring the number of keys it owns.
 - If there is an imbalance, consider redistributing some keys to achieve a more even distribution.