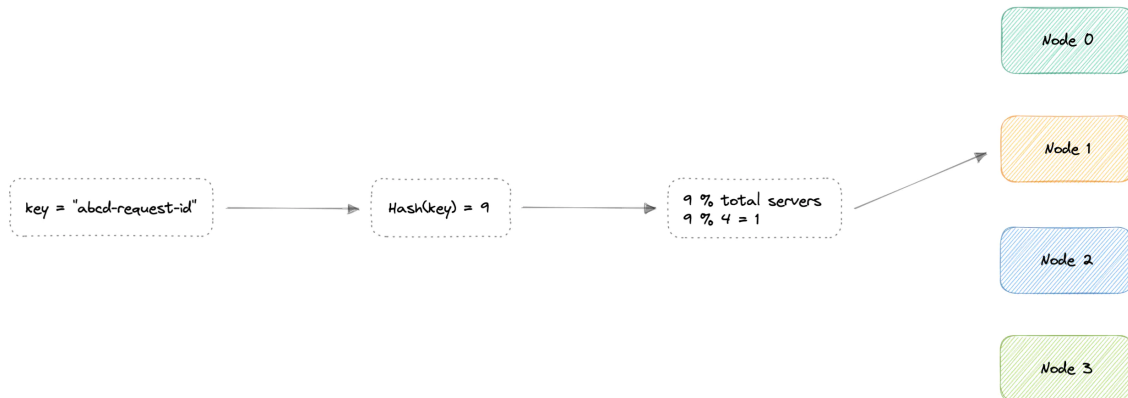


## Why do we need this?

In traditional hashing-based distribution methods, we use a hash function to hash our partition keys (i.e. request ID or IP). Then if we use the modulo against the total number of nodes (server or databases). This will give us the node where we want to route our request.



$\text{Hash}(\text{key}_1) \rightarrow H_1 \quad N = \text{Node}_0$   
 $\text{Hash}(\text{key}_2) \rightarrow H_2 \quad N = \text{Node}_1$   
 $\text{Hash}(\text{key}_3) \rightarrow H_3$

$N = \text{Node}_2 \dots \text{Hash}(\text{key}_n) \rightarrow H_n \quad N = \text{Node}_{n-1}$   
$$\begin{aligned} & \& \text{Hash}(\text{key}_1) \& \text{to } H_1 \& \bmod N = \text{Node}_0 \\ & \& \text{Hash}(\text{key}_2) \& \text{to } H_2 \& \bmod N = \text{Node}_1 \\ & \& \text{Hash}(\text{key}_3) \& \text{to } H_3 \& \bmod N = \text{Node}_2 \\ & \& \dots \& \text{Hash}(\text{key}_n) \& \text{to } H_n \& \bmod N = \text{Node}_{n-1} \end{aligned}$$
  
 $\text{Hash}(\text{key}_1) \rightarrow H_1 \bmod N = \text{Node}_0$   
 $\text{Hash}(\text{key}_2) \rightarrow H_2 \bmod N = \text{Node}_1$   
 $\text{Hash}(\text{key}_3) \rightarrow H_3 \bmod N = \text{Node}_2 \dots \text{Hash}(\text{key}_n) \rightarrow H_n \bmod N = \text{Node}_{n-1}$

Where,

**key**: Request ID or IP.

**H**: Hash function result.

**N**: Total number of nodes.

**Node**: The node where the request will be routed.

The problem with this is if we add or remove a node, it will cause **N** to change, meaning our mapping strategy will break as the same requests will now map to a different server. As a consequence, the majority of requests will need to be redistributed which is very inefficient.

We want to uniformly distribute requests among different nodes such that we should be able to add or remove nodes with minimal effort. Hence, we need a distribution scheme that does not depend directly on the number of nodes (or servers), so that,

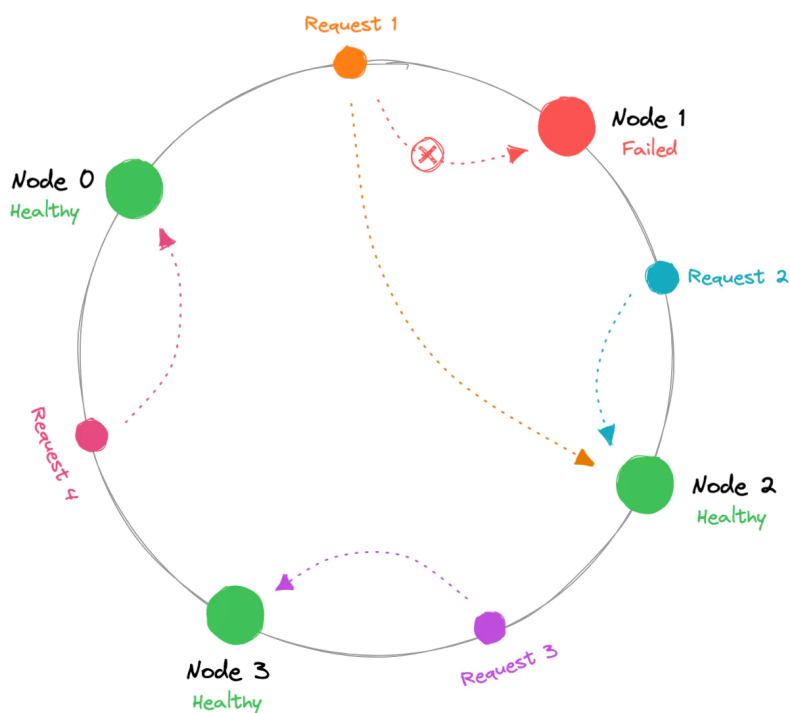
when adding or removing nodes, the number of keys that need to be relocated is minimized.

Consistent hashing solves this horizontal scalability problem by ensuring that every time we scale up or down, we do not have to re-arrange all the keys or touch all the servers.

Now that we understand the problem, let's discuss consistent hashing in detail.

## How does it work

Consistent Hashing is a distributed hashing scheme that operates independently of the number of nodes in a distributed hash table by assigning them a position on an abstract circle, or hash ring. This allows servers and objects to scale without affecting the overall system.



Using consistent hashing, only  $\frac{K}{N}$  data would require re-distributing.

$$R = \frac{K}{N} \times \frac{K}{N} = \frac{K}{N}$$

Where,

$R$ : Data that would require re-distribution.

$K$ : Number of partition keys.

$N$ : Number of nodes.

The output of the hash function is a range let's say  $0 \dots m-1$  which we can represent on our hash ring. We hash the requests and distribute them on the ring depending on what the output was. Similarly, we also hash the node and distribute them on the same ring as well.

$$\text{Hash}(\text{key}_1) = P_1 \text{Hash}(\text{key}_2) = P_2 \text{Hash}(\text{key}_3) = P_3 \dots \text{Hash}(\text{key}_n) = P_{m-1}$$

$$\text{Hash}(\text{key}_1) = P_1 \ \& \ \text{Hash}(\text{key}_2) = P_2 \ \& \ \text{Hash}(\text{key}_3) = P_3 \ \& \ \dots \ \& \ \text{Hash}(\text{key}_n) = P_{m-1}$$

Where,

**key**: Request/Node ID or IP.

**P**: Position on the hash ring.

**m**: Total range of the hash ring.

Now, when the request comes in we can simply route it to the closest node in a clockwise (can be counterclockwise as well) manner. This means that if a new node is added or removed, we can use the nearest node and only a *fraction* of the requests need to be re-routed.

In theory, consistent hashing should distribute the load evenly however it doesn't happen in practice. Usually, the load distribution is uneven and one server may end up handling the majority of the request becoming a *hotspot*, essentially a bottleneck for the system. We can fix this by adding extra nodes but that can be expensive.

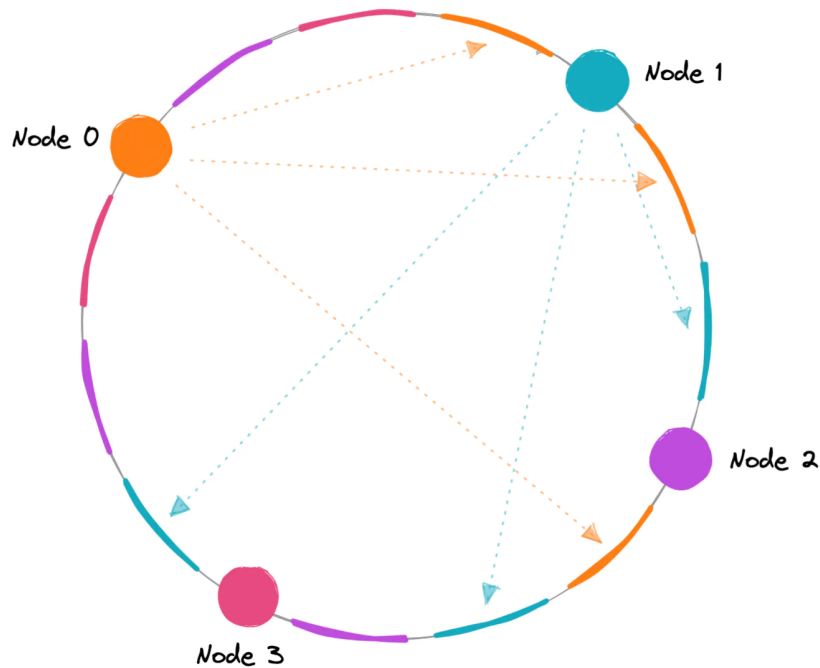
Let's see how we can address these issues.

## Virtual Nodes

In order to ensure a more evenly distributed load, we can introduce the idea of a virtual node, sometimes also referred to as a VNode.

Instead of assigning a single position to a node, the hash range is divided into multiple smaller ranges, and each physical node is assigned several of these smaller ranges. Each of these subranges is considered a VNode. Hence, virtual nodes are basically existing physical nodes mapped multiple times across the hash ring to minimize

changes to a node's assigned range.



For this, we can use  $k$  number of hash functions.

$$\begin{aligned} \text{Hash1}(\text{key}_1) &= P_1 \text{Hash2}(\text{key}_2) = P_2 \text{Hash3}(\text{key}_3) = P_3 \dots \text{Hash}_k(\text{key}_n) = P_{m-1} \\ \text{Hash1}(\text{key}_1) &= P_1 \ \& \ \text{Hash}_2(\text{key}_2) = P_2 \ \& \ \text{Hash}_3(\text{key}_3) = P_3 \ \& \ \dots \ \& \ \text{Hash}_k(\text{key}_n) \\ &= P_{m-1} \end{aligned}$$

Where,

$\text{key}$ : Request/Node ID or IP.

$k$ : Number of hash functions.

$P$ : Position on the hash ring.

$m$ : Total range of the hash ring.

As VNodes help spread the load more evenly across the physical nodes on the cluster by dividing the hash ranges into smaller subranges, this speeds up the re-balancing process after adding or removing nodes. This also helps us reduce the probability of hotspots.

## Data replication

To ensure high availability and durability, consistent hashing replicates each data item on multiple  $N$  nodes in the system where the value  $N$  is equivalent to the *replication factor*.

The replication factor is the number of nodes that will receive the copy of the same data. In eventually consistent systems, this is done asynchronously.

## Advantages

Let's look at some advantages of consistent hashing:

- Makes rapid scaling up and down more predictable.
- Facilitates partitioning and replication across nodes.
- Enables scalability and availability.
- Reduces hotspots.

## Disadvantages

Below are some disadvantages of consistent hashing:

- Increases complexity.
- Cascading failures.
- Load distribution can still be uneven.
- Key management can be expensive when nodes transiently fail.

## Examples

Let's look at some examples where consistent hashing is used:

- Data partitioning in [Apache Cassandra](#).
- Load distribution across multiple storage hosts in [Amazon DynamoDB](#).