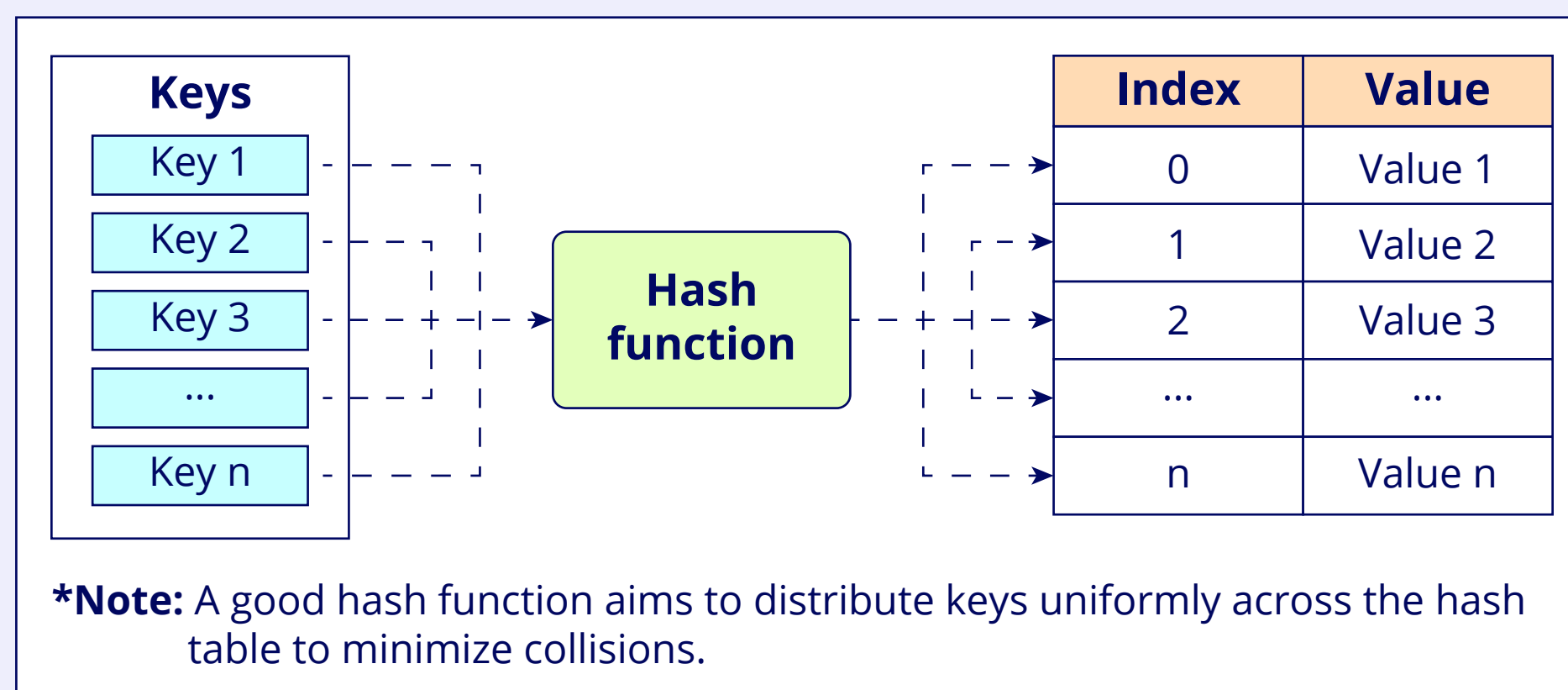


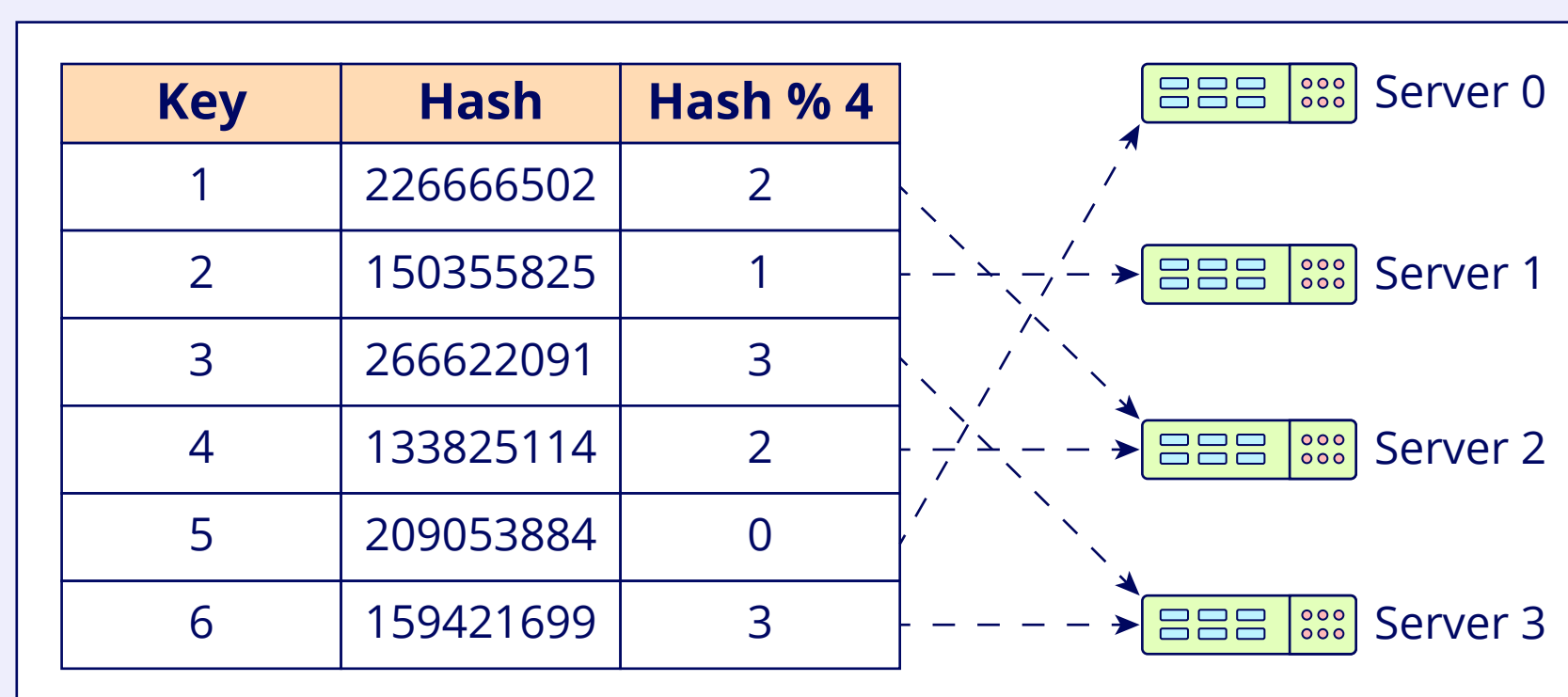
What is a Hash Function?

A hash function takes an input (key) and produces a fixed-size string of bytes (hash code or hash value). This hash code is used to index data in hash tables, enabling efficient data retrieval and lookup operations.



Distributed Hash Tables

Often, the size of a hash table becomes very large; this means it has to be split into several parts. Each part of the hash table is stored on a different server, as shown below:



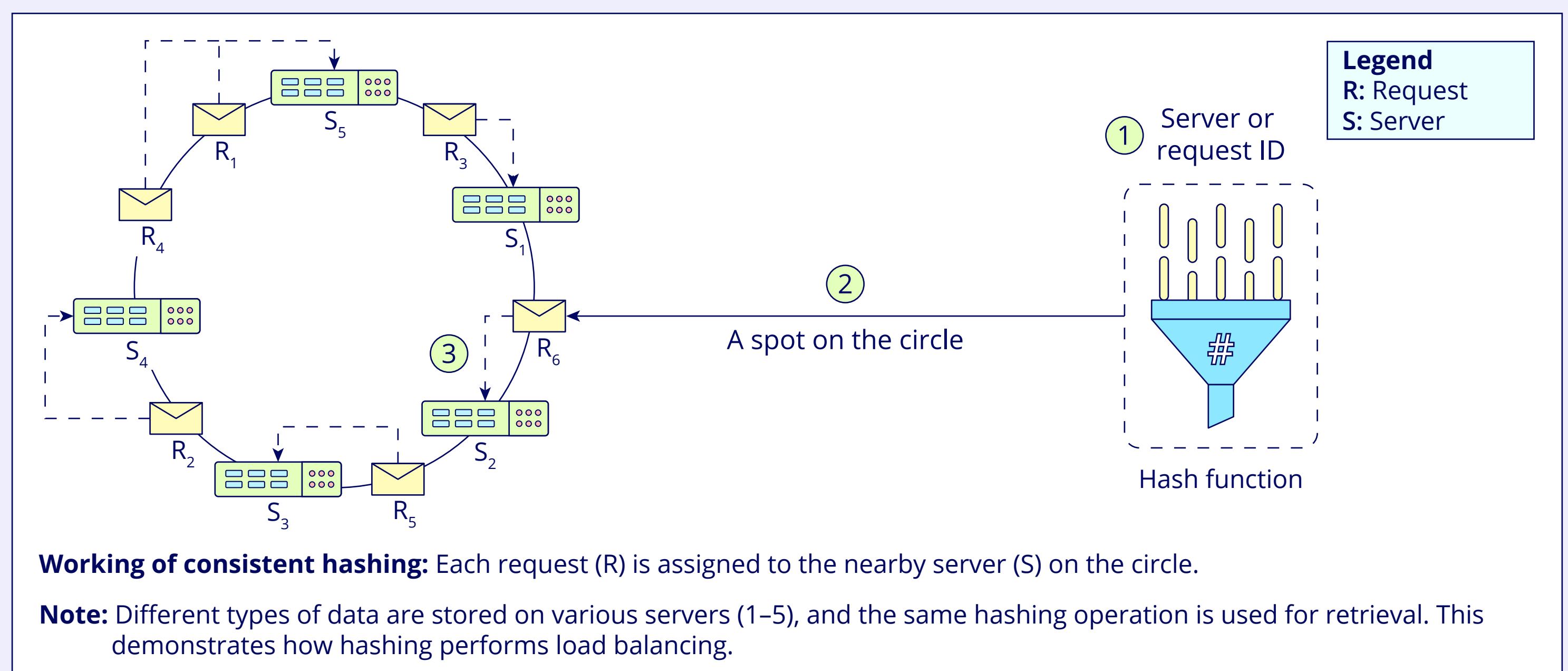
Problem!

If the number of servers is reduced in the cluster, the modulo must be recalculated accordingly for each value (an expensive operation!).

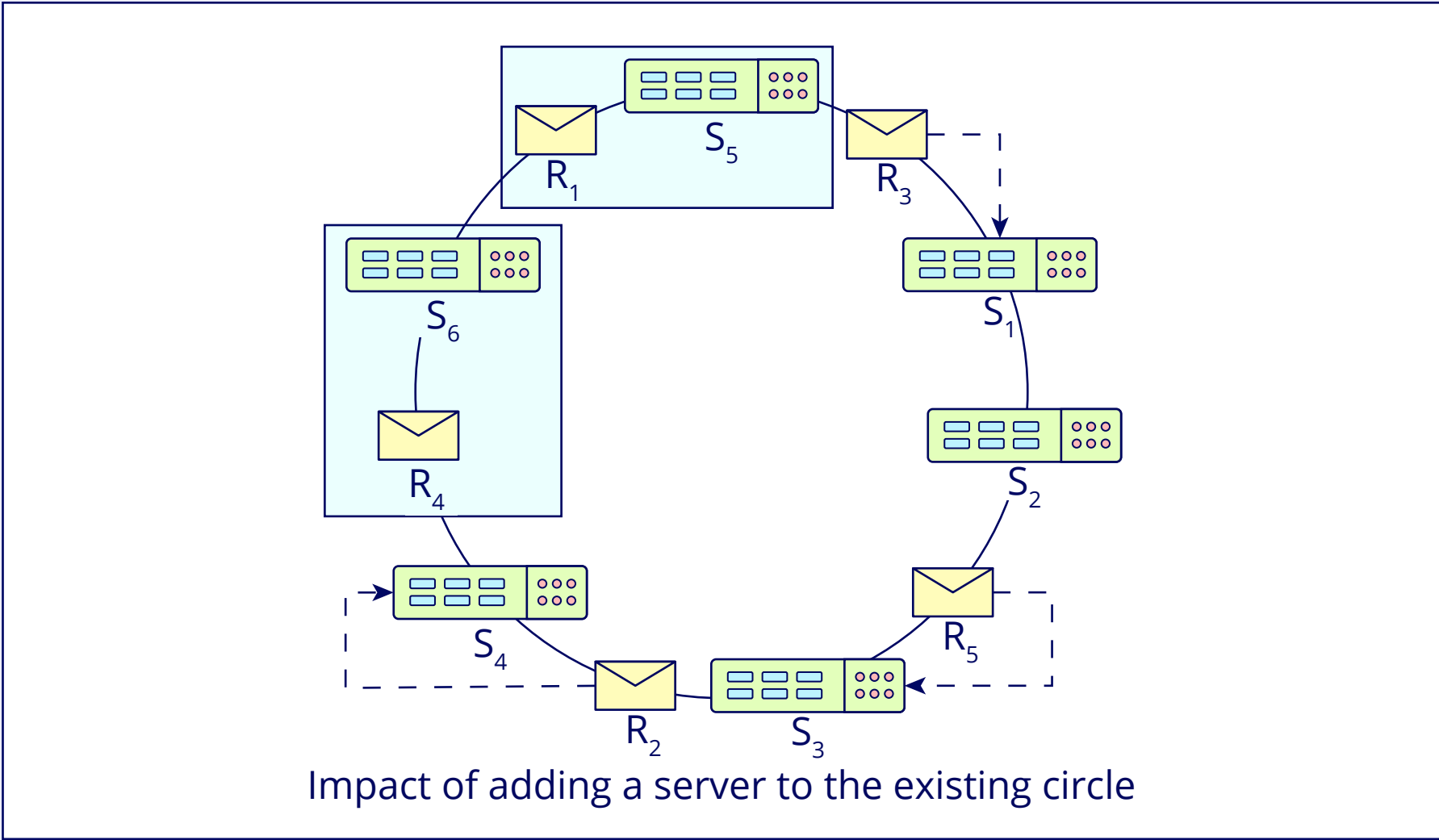
Solution: Consistent hashing

What is Consistent Hashing?

Consistent hashing is an effective technique for distributing the workload among a set of servers efficiently placed on an imaginary circle. It minimizes the number of keys that need to be remapped when a server is added or removed from the cluster.

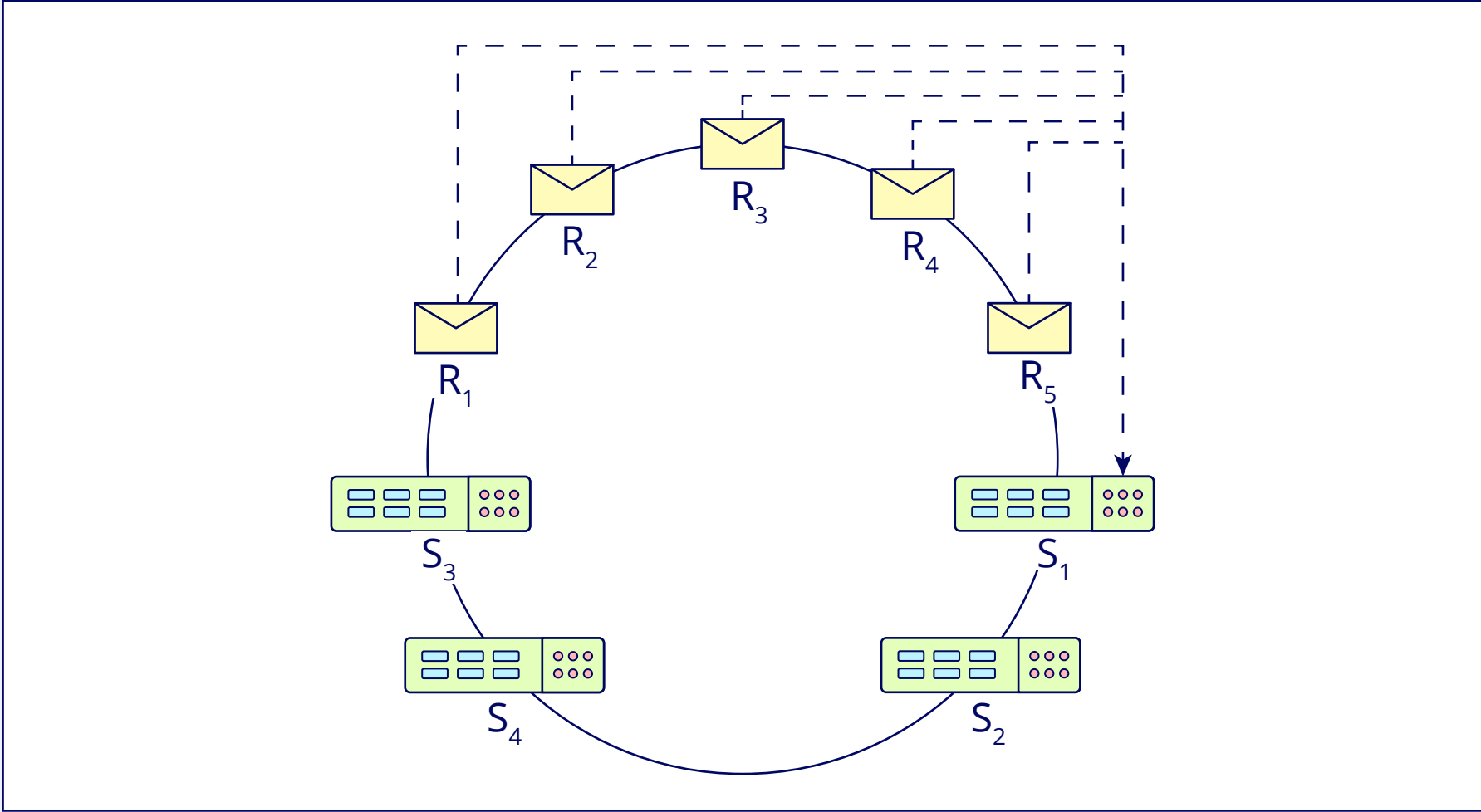


Adding new servers reduces the burden on other servers in the cluster using consistent hashing.



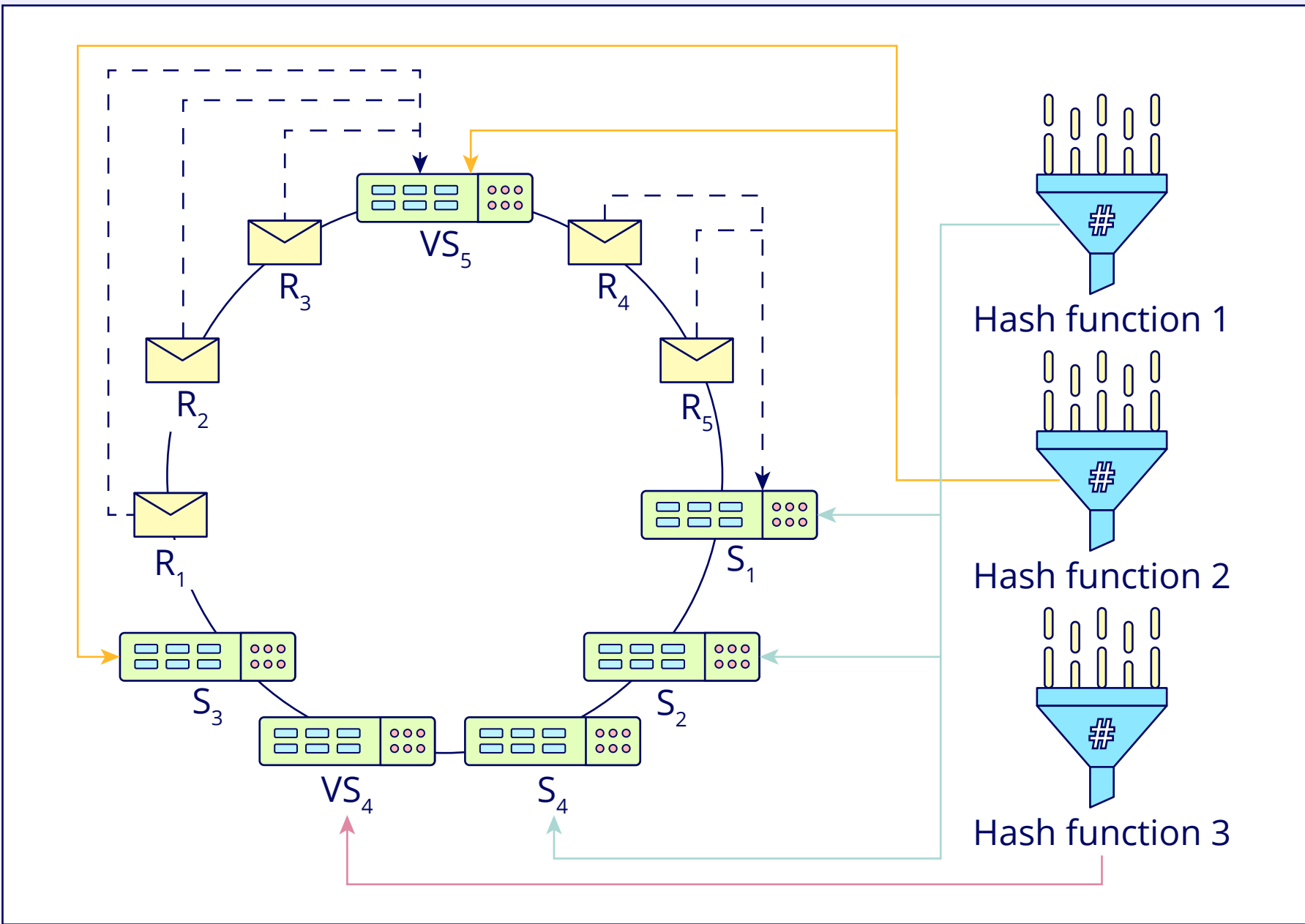
What is a Hotspot in Consistent Hashing?

As depicted below, most of the incoming requests lie between the S_3 and S_1 servers. Consequently, S_1 has to handle most of the requests compared to other servers, becoming a hotspot.



Hotspot Prevention: Use Multiple Hash Functions

To prevent hotspots, we use virtual servers (see VS_2 and VS_4) for balanced workload distribution. Virtual servers map each physical server to multiple locations using multiple hash functions.



Some important use cases of consistent hashing

- Web caching and CDNs
- Distributed databases
- Sharded counters
- Distributed file systems