

arbitrarily chosen; and in real-world systems, the number of virtual nodes is much larger. Instead of using  $s0$ , we have  $s0_0$ ,  $s0_1$ , and  $s0_2$  to represent *server 0* on the ring. Similarly,  $s1_0$ ,  $s1_1$ , and  $s1_2$  represent *server 1* on the ring. With virtual nodes, each server is responsible for multiple partitions. Partitions (edges) with label  $s0$  are managed by *server 0*. On the other hand, partitions with label  $s1$  are managed by *server 1*.

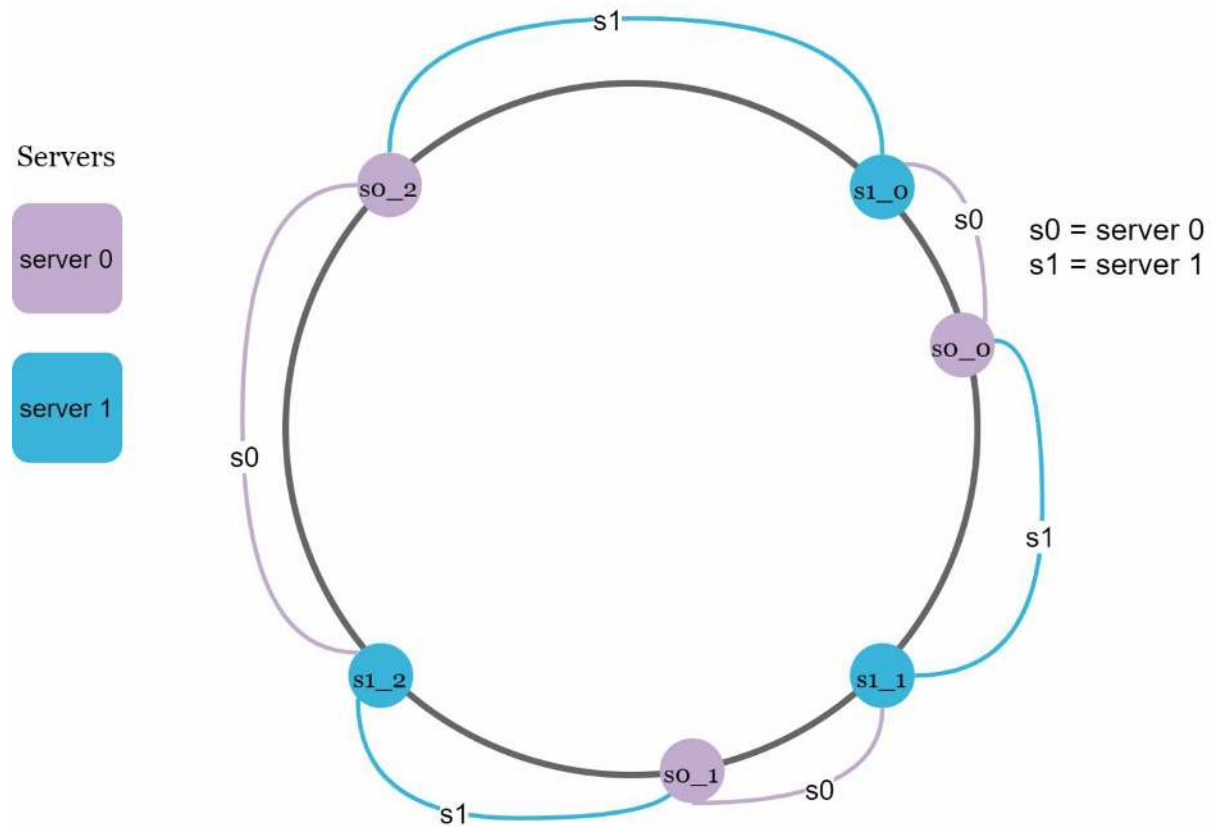


Figure 5-12

To find which server a key is stored on, we go clockwise from the key's location and find the first virtual node encountered on the ring. In Figure 5-13, to find out which server  $k0$  is stored on, we go clockwise from  $k0$ 's location and find virtual node  $s1_1$ , which refers to *server 1*.

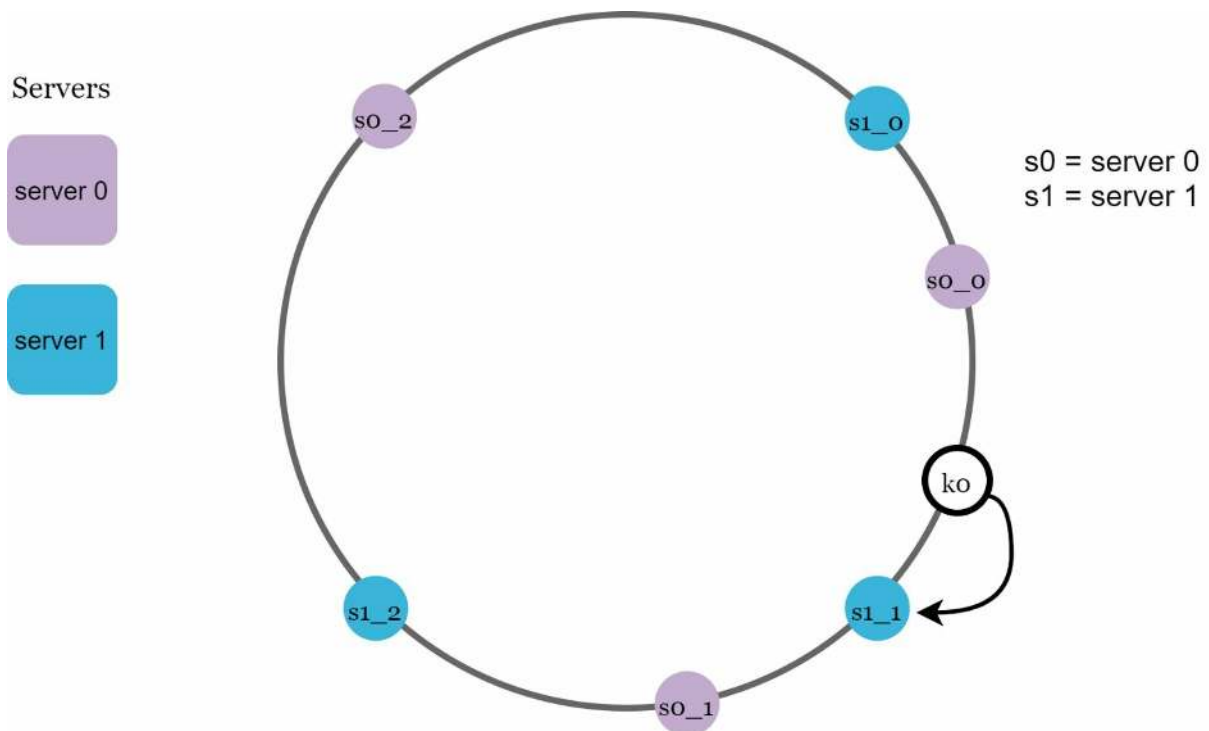


Figure 5-13

As the number of virtual nodes increases, the distribution of keys becomes more balanced. This is because the standard deviation gets smaller with more virtual nodes, leading to balanced data distribution. Standard deviation measures how data are spread out. The outcome of an experiment carried out by online research [2] shows that with one or two hundred virtual nodes, the standard deviation is between 5% (200 virtual nodes) and 10% (100 virtual nodes) of the mean. The standard deviation will be smaller when we increase the number of virtual nodes. However, more spaces are needed to store data about virtual nodes. This is a tradeoff, and we can tune the number of virtual nodes to fit our system requirements.

### Find affected keys

When a server is added or removed, a fraction of data needs to be redistributed. How can we find the affected range to redistribute the keys?

In Figure 5-14, *server 4* is added onto the ring. The affected range starts from *s4* (newly added node) and moves anticlockwise around the ring until a server is found (*s3*). Thus, keys located between *s3* and *s4* need to be redistributed to *s4*.

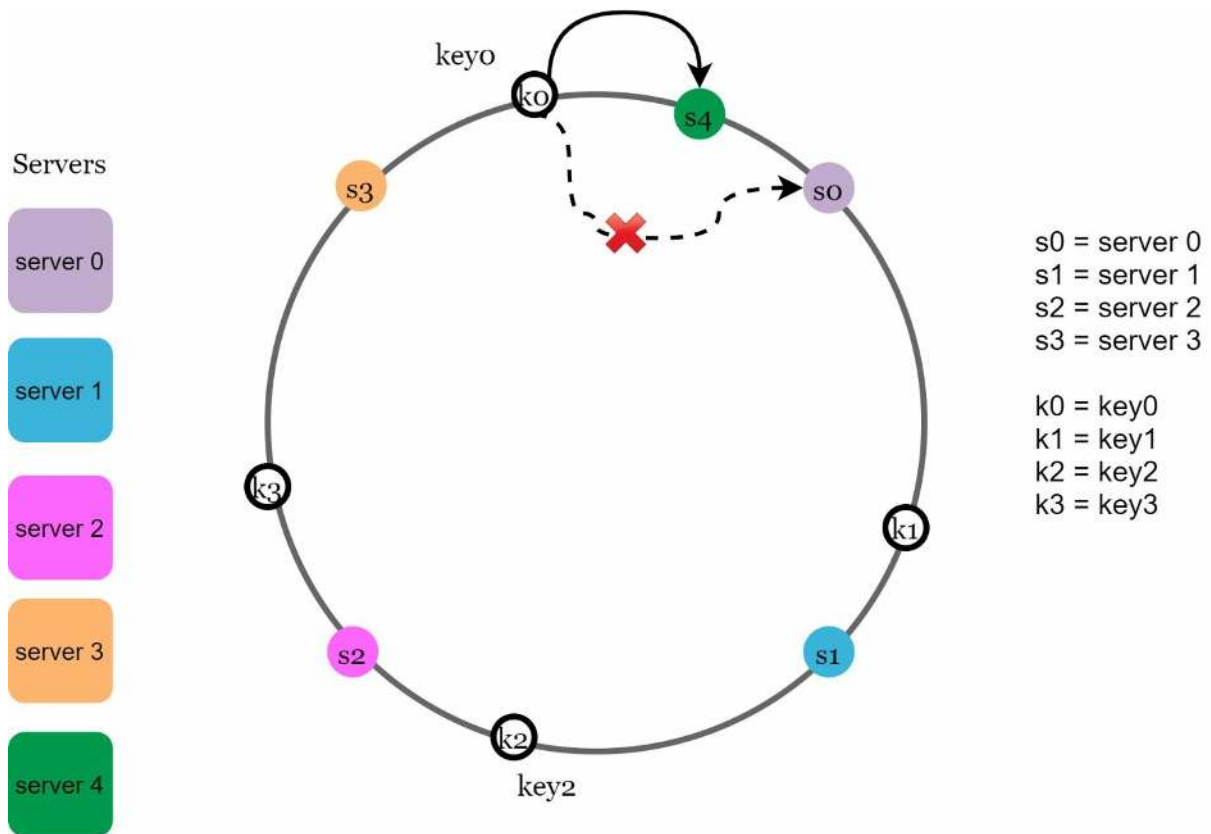


Figure 5-14

When a server ( $s1$ ) is removed as shown in Figure 5-15, the affected range starts from  $s1$  (removed node) and moves anticlockwise around the ring until a server is found ( $s0$ ). Thus, keys located between  $s0$  and  $s1$  must be redistributed to  $s2$ .

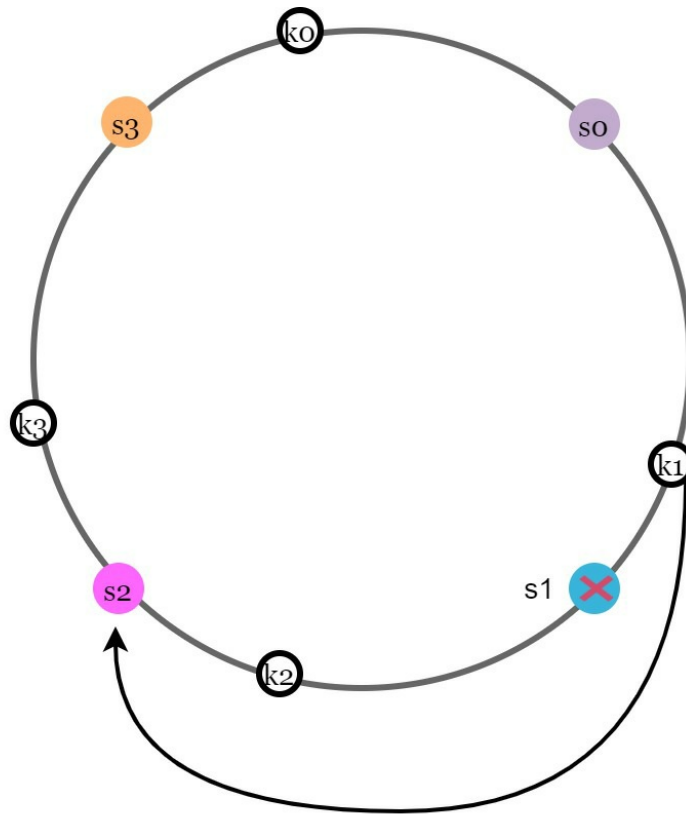
Servers

server 0

server 1

server 2

server 3



s0 = server 0  
s1 = server 1  
s2 = server 2  
s3 = server 3

k0 = key0  
k1 = key1  
k2 = key2  
k3 = key3

Figure 5-15

## Wrap up

In this chapter, we had an in-depth discussion about consistent hashing, including why it is needed and how it works. The benefits of consistent hashing include:

- Minimized keys are redistributed when servers are added or removed.
- It is easy to scale horizontally because data are more evenly distributed.
- Mitigate hotspot key problem. Excessive access to a specific shard could cause server overload. Imagine data for Katy Perry, Justin Bieber, and Lady Gaga all end up on the same shard. Consistent hashing helps to mitigate the problem by distributing the data more evenly.

Consistent hashing is widely used in real-world systems, including some notable ones:

- Partitioning component of Amazon's Dynamo database [3]
- Data partitioning across the cluster in Apache Cassandra [4]
- Discord chat application [5]
- Akamai content delivery network [6]
- Maglev network load balancer [7]

Congratulations on getting this far! Now give yourself a pat on the back. Good job!

## Reference materials

- [1] Consistent hashing: [https://en.wikipedia.org/wiki/Consistent\\_hashing](https://en.wikipedia.org/wiki/Consistent_hashing)
- [2] Consistent Hashing:  
<https://tom-e-white.com/2007/11/consistent-hashing.html>
- [3] Dynamo: Amazon's Highly Available Key-value Store:  
<https://www.allthingsdistributed.com/files/amazon-dynamo-sosp2007.pdf>
- [4] Cassandra - A Decentralized Structured Storage System:  
<http://www.cs.cornell.edu/Projects/ladis2009/papers/Lakshman-ladis2009.PDF>
- [5] How Discord Scaled Elixir to 5,000,000 Concurrent Users:  
<https://blog.discord.com/scaling-elixir-f9b8e1e7c29b>
- [6] CS168: The Modern Algorithmic Toolbox Lecture #1: Introduction and Consistent Hashing: <http://theory.stanford.edu/~tim/s16/l11.pdf>
- [7] Maglev: A Fast and Reliable Software Network Load Balancer:  
<https://static.googleusercontent.com/media/research.google.com/en//pubs/archive/44824.pdf>

## CHAPTER 6: DESIGN A KEY-VALUE STORE

A key-value store, also referred to as a key-value database, is a non-relational database. Each unique identifier is stored as a key with its associated value. This data pairing is known as a “key-value” pair.

In a key-value pair, the key must be unique, and the value associated with the key can be accessed through the key. Keys can be plain text or hashed values. For performance reasons, a short key works better. What do keys look like? Here are a few examples:

- Plain text key: “last\_logged\_in\_at”
- Hashed key: 253DDEC4

The value in a key-value pair can be strings, lists, objects, etc. The value is usually treated as an opaque object in key-value stores, such as Amazon dynamo [1], Memcached [2], Redis [3], etc.

Here is a data snippet in a key-value store:

Key	value
145	john
147	bob
160	Julia

Table 6-1

In this chapter, you are asked to design a key-value store that supports the following operations:

- put(key, value) // insert “value” associated with “key”
- get(key) // get “value” associated with “key”

## **Understand the problem and establish design scope**

There is no perfect design. Each design achieves a specific balance regarding the tradeoffs of the read, write, and memory usage. Another tradeoff has to be made was between consistency and availability. In this chapter, we design a key-value store that comprises of the following characteristics:

- The size of a key-value pair is small: less than 10 KB.
- Ability to store big data.
- High availability: The system responds quickly, even during failures.
- High scalability: The system can be scaled to support large data set.
- Automatic scaling: The addition/deletion of servers should be automatic based on traffic.
- Tunable consistency.
- Low latency.



## **Single server key-value store**

Developing a key-value store that resides in a single server is easy. An intuitive approach is to store key-value pairs in a hash table, which keeps everything in memory. Even though memory access is fast, fitting everything in memory may be impossible due to the space constraint. Two optimizations can be done to fit more data in a single server:

- Data compression
- Store only frequently used data in memory and the rest on disk

Even with these optimizations, a single server can reach its capacity very quickly. A distributed key-value store is required to support big data.

## Distributed key-value store

A distributed key-value store is also called a distributed hash table, which distributes key-value pairs across many servers. When designing a distributed system, it is important to understand CAP (**C**onsistency, **A**vailability, **P**artition Tolerance) theorem.

### CAP theorem

CAP theorem states it is impossible for a distributed system to simultaneously provide more than two of these three guarantees: consistency, availability, and partition tolerance. Let us establish a few definitions.

**Consistency:** consistency means all clients see the same data at the same time no matter which node they connect to.

**Availability:** availability means any client which requests data gets a response even if some of the nodes are down.

**Partition Tolerance:** a partition indicates a communication break between two nodes. Partition tolerance means the system continues to operate despite network partitions.

CAP theorem states that one of the three properties must be sacrificed to support 2 of the 3 properties as shown in Figure 6-1.

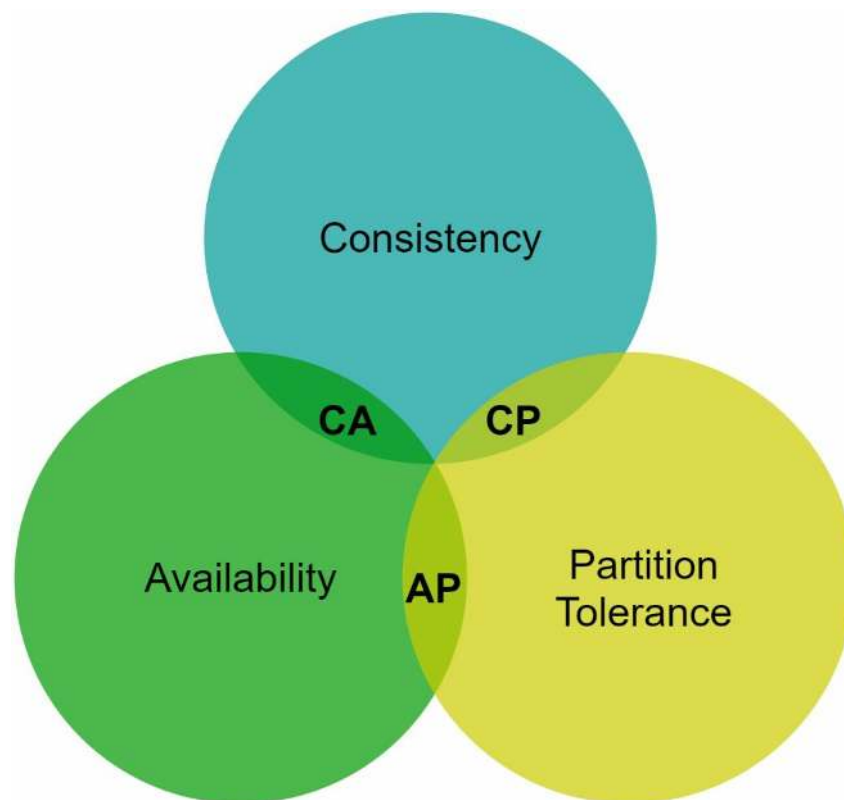


Figure 6-1

Nowadays, key-value stores are classified based on the two CAP characteristics they support:

**CP (consistency and partition tolerance) systems:** a CP key-value store supports consistency and partition tolerance while sacrificing availability.

**AP (availability and partition tolerance) systems:** an AP key-value store supports availability and partition tolerance while sacrificing consistency.

**CA (consistency and availability) systems:** a CA key-value store supports consistency and

availability while sacrificing partition tolerance. Since network failure is unavoidable, a distributed system must tolerate network partition. Thus, a CA system cannot exist in real-world applications.

What you read above is mostly the definition part. To make it easier to understand, let us take a look at some concrete examples. In distributed systems, data is usually replicated multiple times. Assume data are replicated on three replica nodes,  $n1$ ,  $n2$  and  $n3$  as shown in Figure 6-2.

### **Ideal situation**

In the ideal world, network partition never occurs. Data written to  $n1$  is automatically replicated to  $n2$  and  $n3$ . Both consistency and availability are achieved.

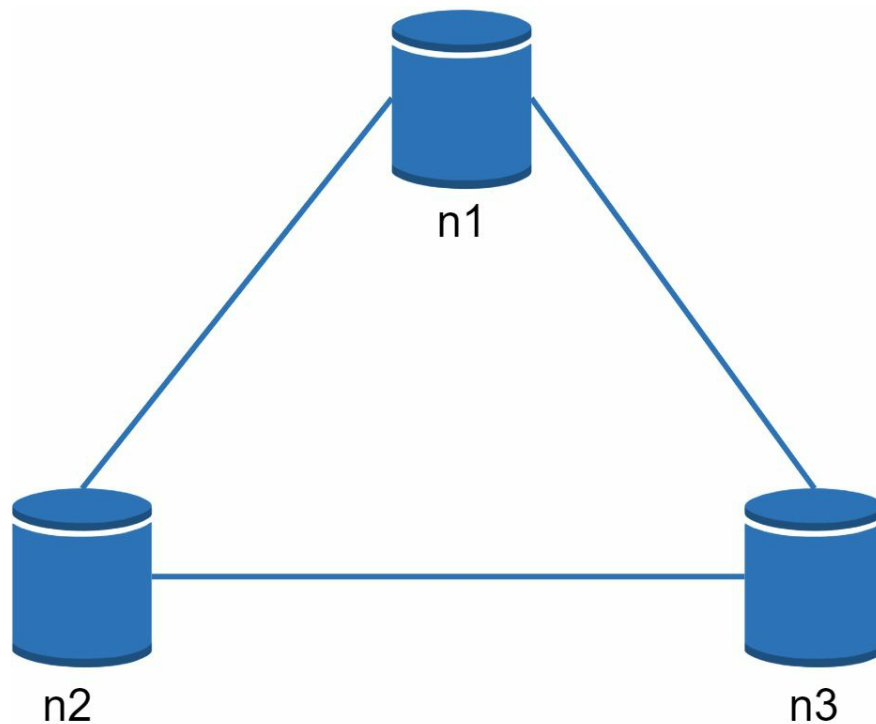


Figure 6-2

### **Real-world distributed systems**

In a distributed system, partitions cannot be avoided, and when a partition occurs, we must choose between consistency and availability. In Figure 6-3,  $n3$  goes down and cannot communicate with  $n1$  and  $n2$ . If clients write data to  $n1$  or  $n2$ , data cannot be propagated to  $n3$ . If data is written to  $n3$  but not propagated to  $n1$  and  $n2$  yet,  $n1$  and  $n2$  would have stale data.

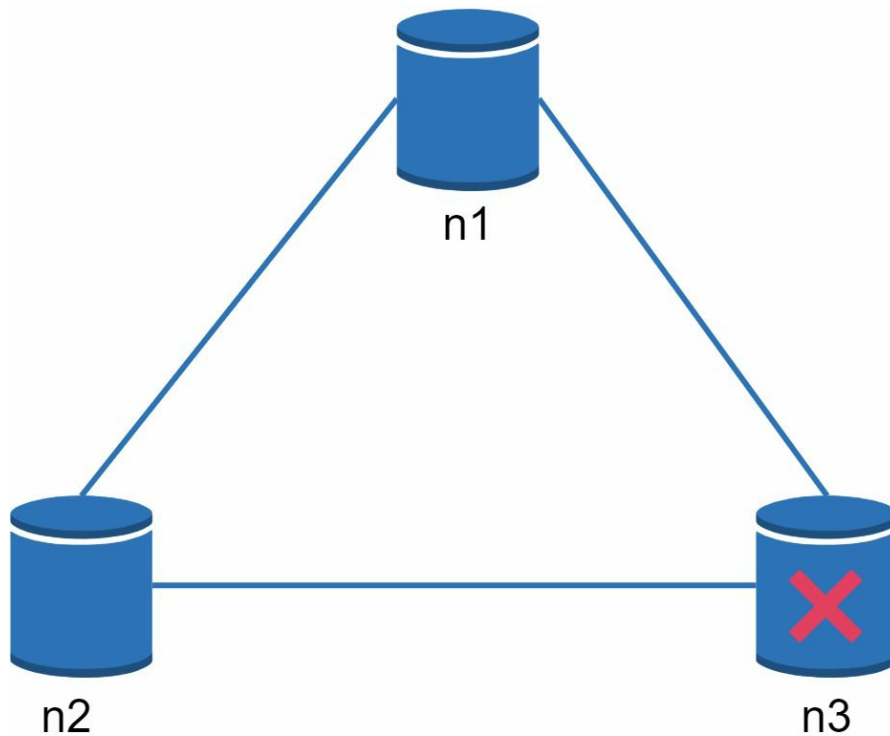


Figure 6-3

If we choose consistency over availability (CP system), we must block all write operations to *n1* and *n2* to avoid data inconsistency among these three servers, which makes the system unavailable. Bank systems usually have extremely high consistent requirements. For example, it is crucial for a bank system to display the most up-to-date balance info. If inconsistency occurs due to a network partition, the bank system returns an error before the inconsistency is resolved.

However, if we choose availability over consistency (AP system), the system keeps accepting reads, even though it might return stale data. For writes, *n1* and *n2* will keep accepting writes, and data will be synced to *n3* when the network partition is resolved.

Choosing the right CAP guarantees that fit your use case is an important step in building a distributed key-value store. You can discuss this with your interviewer and design the system accordingly.

### System components

In this section, we will discuss the following core components and techniques used to build a key-value store:

- Data partition
- Data replication
- Consistency
- Inconsistency resolution
- Handling failures
- System architecture diagram
- Write path
- Read path

The content below is largely based on three popular key-value store systems: Dynamo [4], Cassandra [5], and BigTable [6].

## Data partition

For large applications, it is infeasible to fit the complete data set in a single server. The simplest way to accomplish this is to split the data into smaller partitions and store them in multiple servers. There are two challenges while partitioning the data:

- Distribute data across multiple servers evenly.
- Minimize data movement when nodes are added or removed.

Consistent hashing discussed in Chapter 5 is a great technique to solve these problems. Let us revisit how consistent hashing works at a high-level.

- First, servers are placed on a hash ring. In Figure 6-4, eight servers, represented by  $s0$ ,  $s1$ , ...,  $s7$ , are placed on the hash ring.
- Next, a key is hashed onto the same ring, and it is stored on the first server encountered while moving in the clockwise direction. For instance,  $key0$  is stored in  $s1$  using this logic.

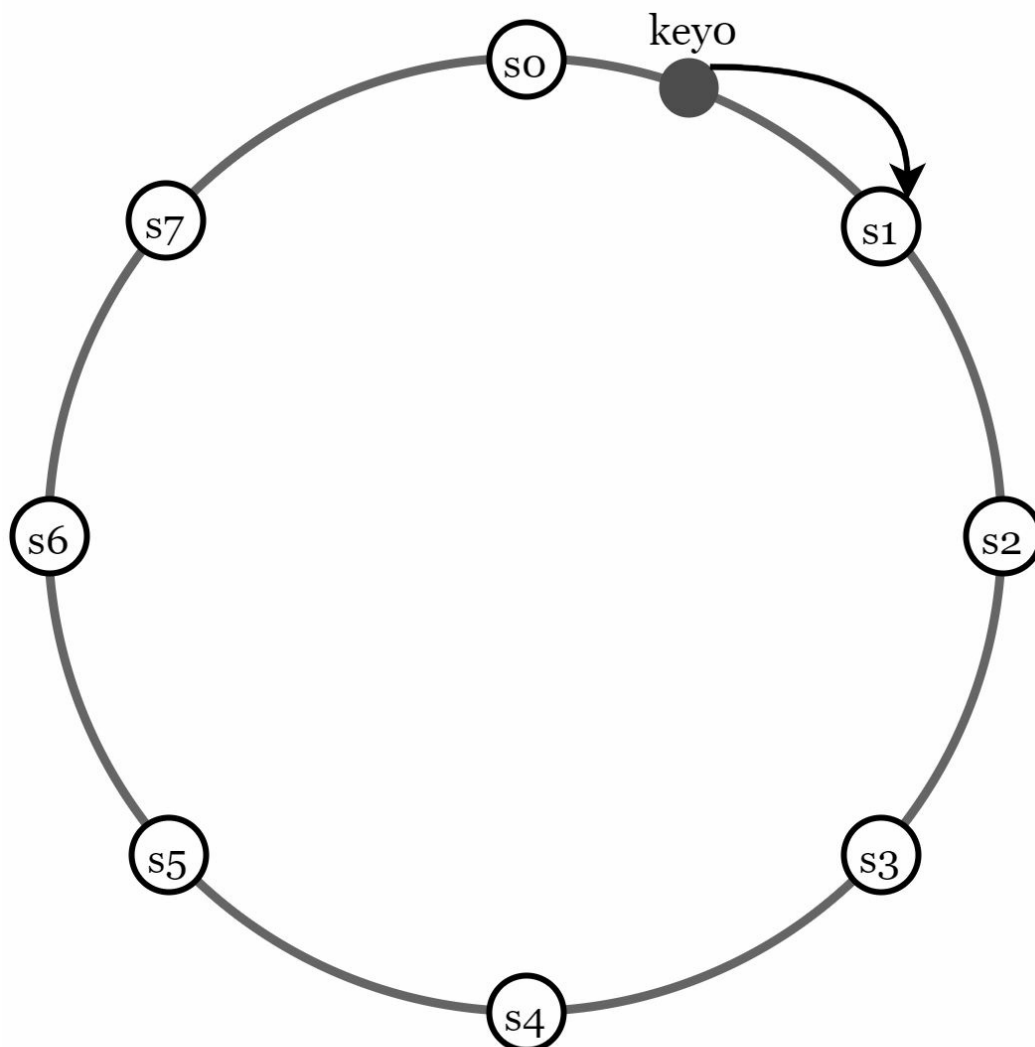


Figure 6-4

Using consistent hashing to partition data has the following advantages:

**Automatic scaling:** servers could be added and removed automatically depending on the

load.

**Heterogeneity:** the number of virtual nodes for a server is proportional to the server capacity. For example, servers with higher capacity are assigned with more virtual nodes.

### Data replication

To achieve high availability and reliability, data must be replicated asynchronously over  $N$  servers, where  $N$  is a configurable parameter. These  $N$  servers are chosen using the following logic: after a key is mapped to a position on the hash ring, walk clockwise from that position and choose the first  $N$  servers on the ring to store data copies. In Figure 6-5 ( $N = 3$ ), *key0* is replicated at *s1*, *s2*, and *s3*.

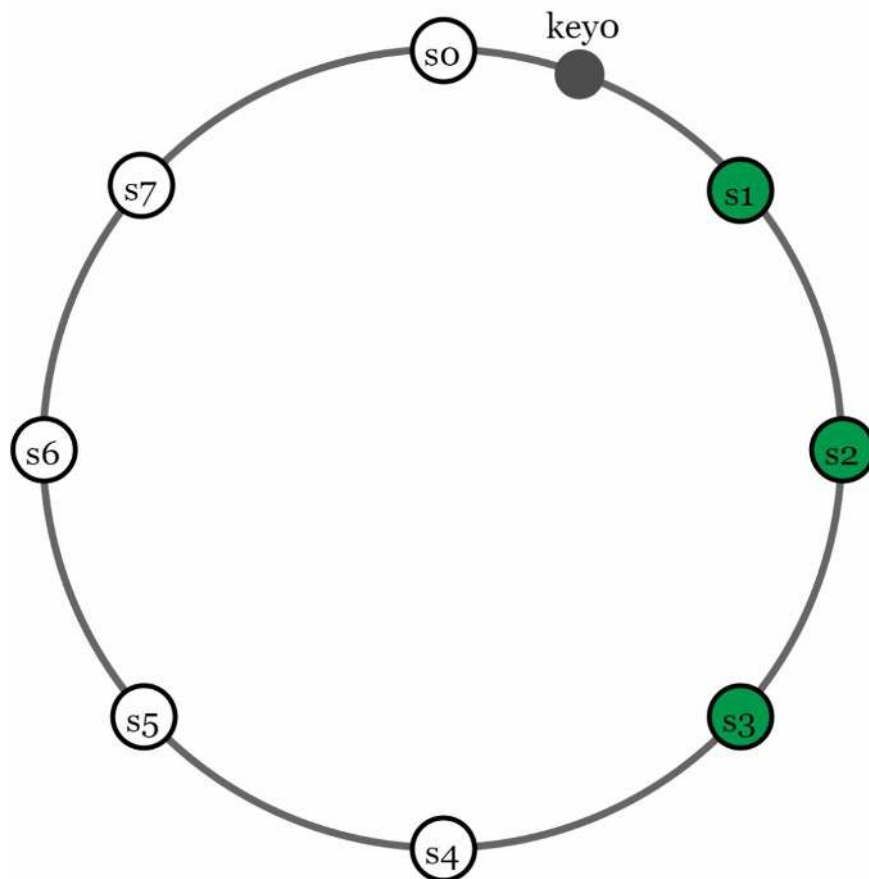


Figure 6-5

With virtual nodes, the first  $N$  nodes on the ring may be owned by fewer than  $N$  physical servers. To avoid this issue, we only choose unique servers while performing the clockwise walk logic.

Nodes in the same data center often fail at the same time due to power outages, network issues, natural disasters, etc. For better reliability, replicas are placed in distinct data centers, and data centers are connected through high-speed networks.

## Consistency

Since data is replicated at multiple nodes, it must be synchronized across replicas. Quorum consensus can guarantee consistency for both read and write operations. Let us establish a few definitions first.

$N$  = The number of replicas

$W$  = A write quorum of size  $W$ . For a write operation to be considered as successful, write operation must be acknowledged from  $W$  replicas.

$R$  = A read quorum of size  $R$ . For a read operation to be considered as successful, read operation must wait for responses from at least  $R$  replicas.

Consider the following example shown in Figure 6-6 with  $N = 3$ .

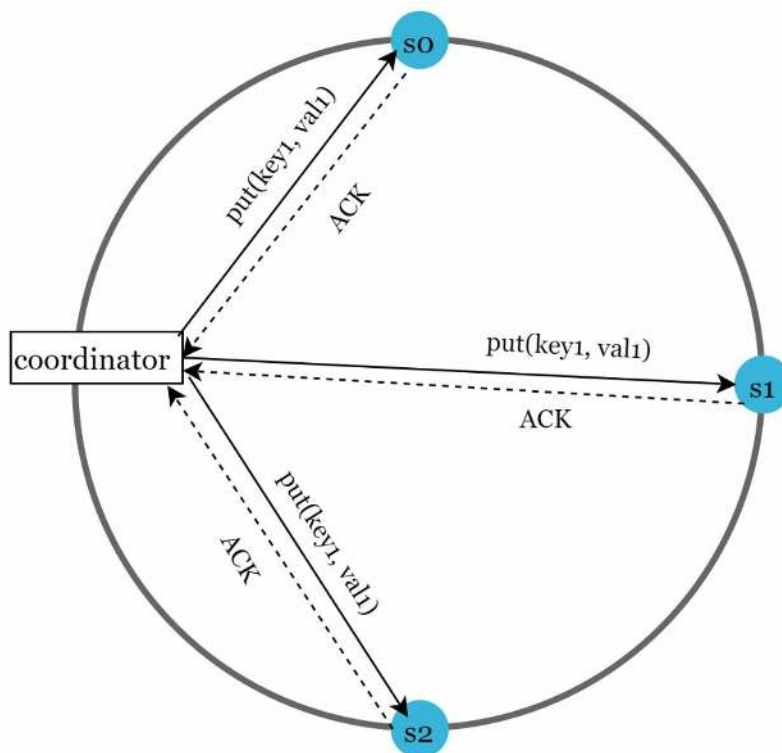


Figure 6-6 (ACK = acknowledgement)

$W = 1$  does not mean data is written on one server. For instance, with the configuration in Figure 6-6, data is replicated at  $s0$ ,  $s1$ , and  $s2$ .  $W = 1$  means that the coordinator must receive at least one acknowledgment before the write operation is considered as successful. For instance, if we get an acknowledgment from  $s1$ , we no longer need to wait for acknowledgements from  $s0$  and  $s2$ . A coordinator acts as a proxy between the client and the nodes.

The configuration of  $W$ ,  $R$  and  $N$  is a typical tradeoff between latency and consistency. If  $W = 1$  or  $R = 1$ , an operation is returned quickly because a coordinator only needs to wait for a response from any of the replicas. If  $W$  or  $R > 1$ , the system offers better consistency; however, the query will be slower because the coordinator must wait for the response from the slowest replica.

If  $W + R > N$ , strong consistency is guaranteed because there must be at least one

overlapping node that has the latest data to ensure consistency.

How to configure  $N$ ,  $W$ , and  $R$  to fit our use cases? Here are some of the possible setups:

If  $R = 1$  and  $W = N$ , the system is optimized for a fast read.

If  $W = 1$  and  $R = N$ , the system is optimized for fast write.

If  $W + R > N$ , strong consistency is guaranteed (Usually  $N = 3$ ,  $W = R = 2$ ).

If  $W + R \leq N$ , strong consistency is not guaranteed.

Depending on the requirement, we can tune the values of  $W$ ,  $R$ ,  $N$  to achieve the desired level of consistency.

### **Consistency models**

Consistency model is other important factor to consider when designing a key-value store. A consistency model defines the degree of data consistency, and a wide spectrum of possible consistency models exist:

- Strong consistency: any read operation returns a value corresponding to the result of the most updated write data item. A client never sees out-of-date data.
- Weak consistency: subsequent read operations may not see the most updated value.
- Eventual consistency: this is a specific form of weak consistency. Given enough time, all updates are propagated, and all replicas are consistent.

Strong consistency is usually achieved by forcing a replica not to accept new reads/writes until every replica has agreed on current write. This approach is not ideal for highly available systems because it could block new operations. Dynamo and Cassandra adopt eventual consistency, which is our recommended consistency model for our key-value store. From concurrent writes, eventual consistency allows inconsistent values to enter the system and force the client to read the values to reconcile. The next section explains how reconciliation works with versioning.



## Inconsistency resolution: versioning

Replication gives high availability but causes inconsistencies among replicas. Versioning and vector locks are used to solve inconsistency problems. Versioning means treating each data modification as a new immutable version of data. Before we talk about versioning, let us use an example to explain how inconsistency happens:

As shown in Figure 6-7, both replica nodes *n1* and *n2* have the same value. Let us call this value the original value. *Server 1* and *server 2* get the same value for `get("name")` operation.

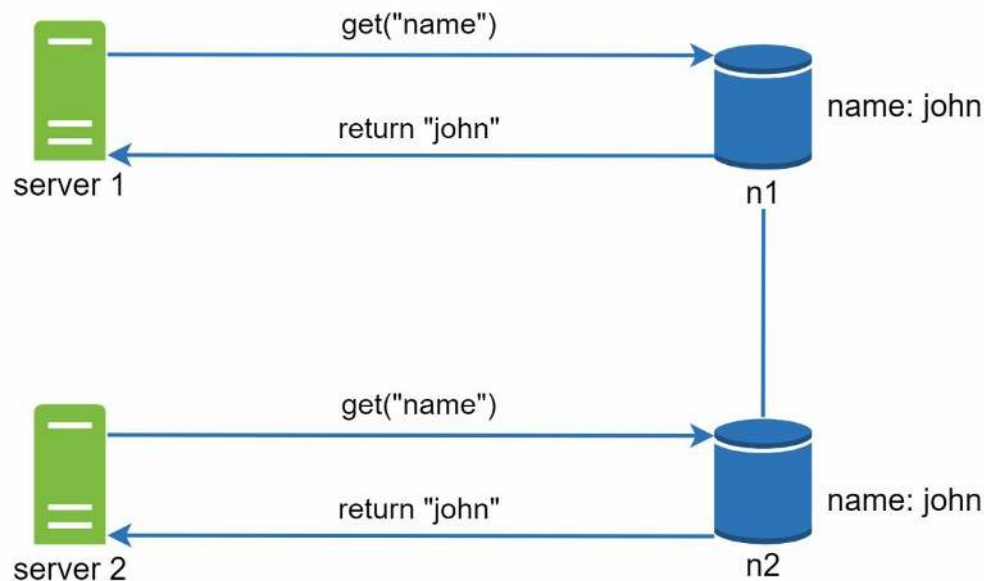


Figure 6-7

Next, *server 1* changes the name to “johnSanFrancisco”, and *server 2* changes the name to “johnNewYork” as shown in Figure 6-8. These two changes are performed simultaneously. Now, we have conflicting values, called versions *v1* and *v2*.

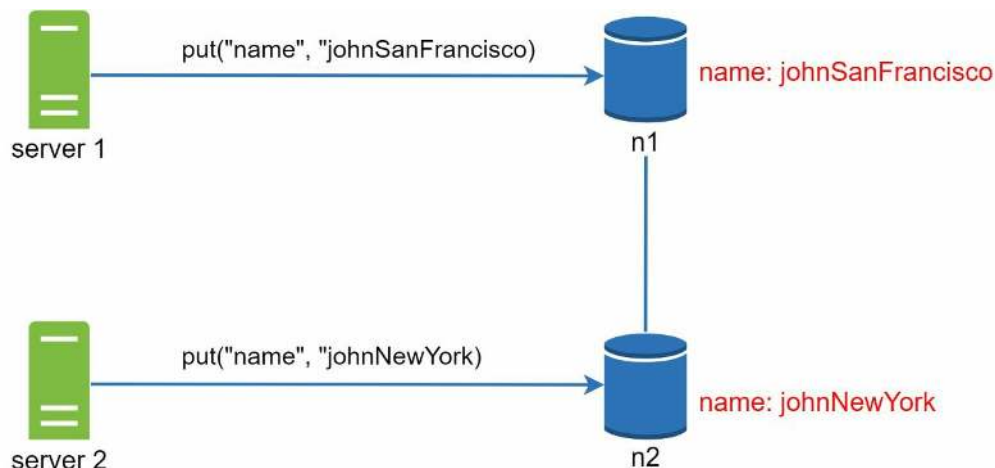


Figure 6-8

In this example, the original value could be ignored because the modifications were based on it. However, there is no clear way to resolve the conflict of the last two versions. To resolve this issue, we need a versioning system that can detect conflicts and reconcile conflicts. A vector clock is a common technique to solve this problem. Let us examine how vector clocks work.

A vector clock is a  $[server, version]$  pair associated with a data item. It can be used to check if one version precedes, succeeds, or in conflict with others.

Assume a vector clock is represented by  $D([S1, v1], [S2, v2], \dots, [Sn, vn])$ , where  $D$  is a data item,  $v1$  is a version counter, and  $s1$  is a server number, etc. If data item  $D$  is written to server  $Si$ , the system must perform one of the following tasks.

- Increment  $vi$  if  $[Si, vi]$  exists.
- Otherwise, create a new entry  $[Si, 1]$ .

The above abstract logic is explained with a concrete example as shown in Figure 6-9.

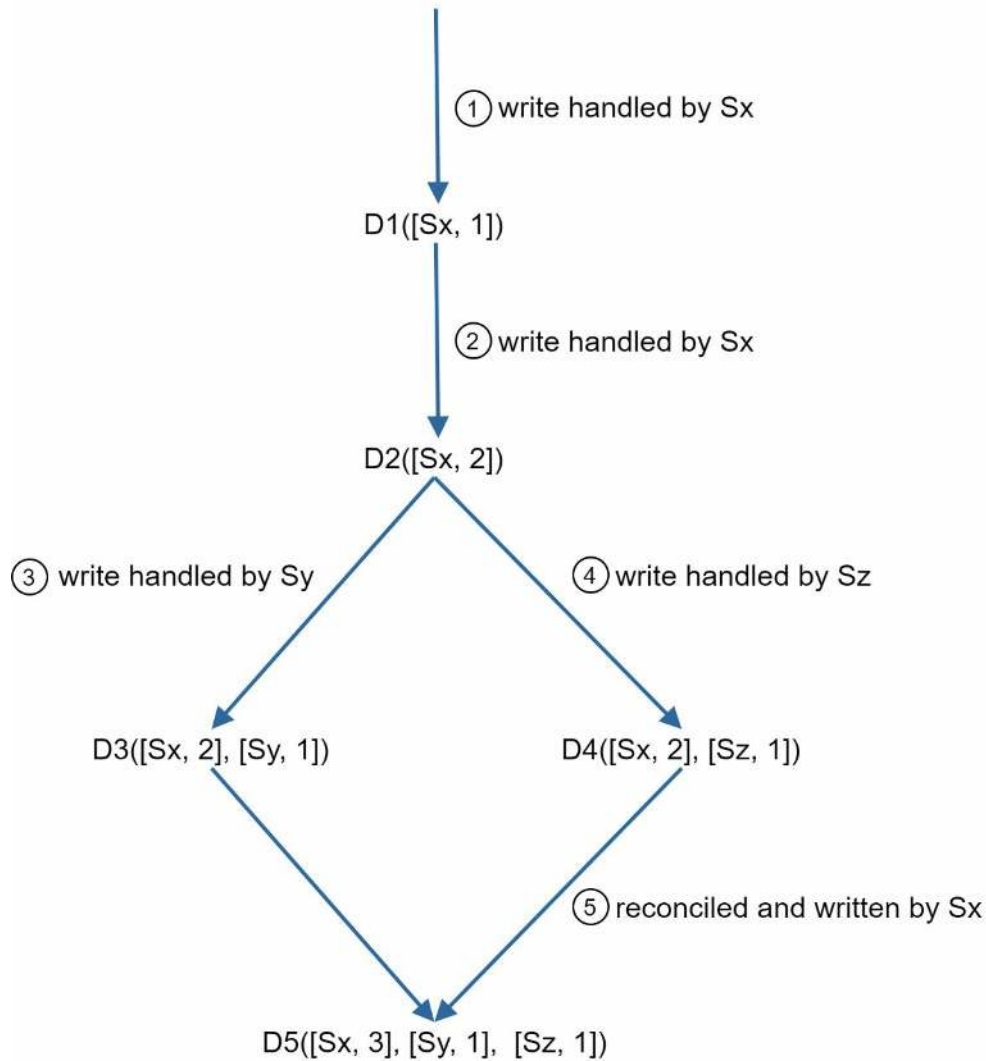


Figure 6-9 (source: [4])

1. A client writes a data item  $D1$  to the system, and the write is handled by server  $Sx$ , which now has the vector clock  $D1([Sx, 1])$ .
2. Another client reads the latest  $D1$ , updates it to  $D2$ , and writes it back.  $D2$  descends from  $D1$  so it overwrites  $D1$ . Assume the write is handled by the same server  $Sx$ , which now has vector clock  $D2([Sx, 2])$ .
3. Another client reads the latest  $D2$ , updates it to  $D3$ , and writes it back. Assume the write is handled by server  $Sy$ , which now has vector clock  $D3([Sx, 2], [Sy, 1])$ .
4. Another client reads the latest  $D2$ , updates it to  $D4$ , and writes it back. Assume the write

is handled by server  $S_z$ , which now has  $D4([S_x, 2], [S_z, 1])$ .

5. When another client reads  $D3$  and  $D4$ , it discovers a conflict, which is caused by data item  $D2$  being modified by both  $S_y$  and  $S_z$ . The conflict is resolved by the client and updated data is sent to the server. Assume the write is handled by  $S_x$ , which now has  $D5([S_x, 3], [S_y, 1], [S_z, 1])$ . We will explain how to detect conflict shortly.

Using vector clocks, it is easy to tell that a version  $X$  is an ancestor (i.e. no conflict) of version  $Y$  if the version counters for each participant in the vector clock of  $Y$  is greater than or equal to the ones in version  $X$ . For example, the vector clock  $D([s0, 1], [s1, 1])$  is an ancestor of  $D([s0, 1], [s1, 2])$ . Therefore, no conflict is recorded.

Similarly, you can tell that a version  $X$  is a sibling (i.e., a conflict exists) of  $Y$  if there is any participant in  $Y$ 's vector clock who has a counter that is less than its corresponding counter in  $X$ . For example, the following two vector clocks indicate there is a conflict:  $D([s0, 1], [s1, 2])$  and  $D([s0, 2], [s1, 1])$ .

Even though vector clocks can resolve conflicts, there are two notable downsides. First, vector clocks add complexity to the client because it needs to implement conflict resolution logic.

Second, the  $[server: version]$  pairs in the vector clock could grow rapidly. To fix this problem, we set a threshold for the length, and if it exceeds the limit, the oldest pairs are removed. This can lead to inefficiencies in reconciliation because the descendant relationship cannot be determined accurately. However, based on Dynamo paper [4], Amazon has not yet encountered this problem in production; therefore, it is probably an acceptable solution for most companies.

## Handling failures

As with any large system at scale, failures are not only inevitable but common. Handling failure scenarios is very important. In this section, we first introduce techniques to detect failures. Then, we go over common failure resolution strategies.

### Failure detection

In a distributed system, it is insufficient to believe that a server is down because another server says so. Usually, it requires at least two independent sources of information to mark a server down.

As shown in Figure 6-10, all-to-all multicasting is a straightforward solution. However, this is inefficient when many servers are in the system.

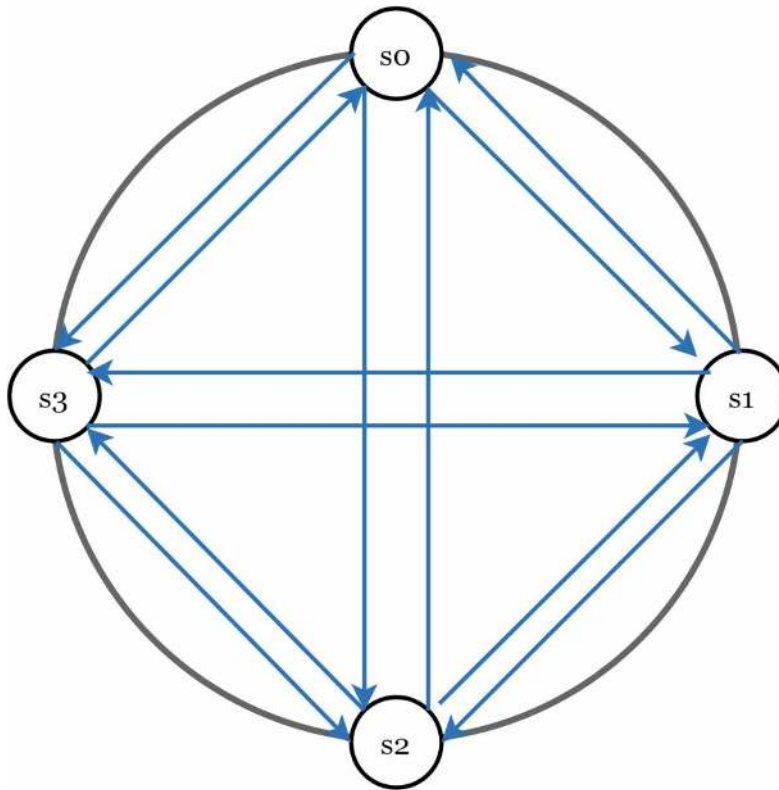


Figure 6-10

A better solution is to use decentralized failure detection methods like gossip protocol.

Gossip protocol works as follows:

- Each node maintains a node membership list, which contains member IDs and heartbeat counters.
- Each node periodically increments its heartbeat counter.
- Each node periodically sends heartbeats to a set of random nodes, which in turn propagate to another set of nodes.
- Once nodes receive heartbeats, membership list is updated to the latest info.
- If the heartbeat has not increased for more than predefined periods, the member is considered as offline.

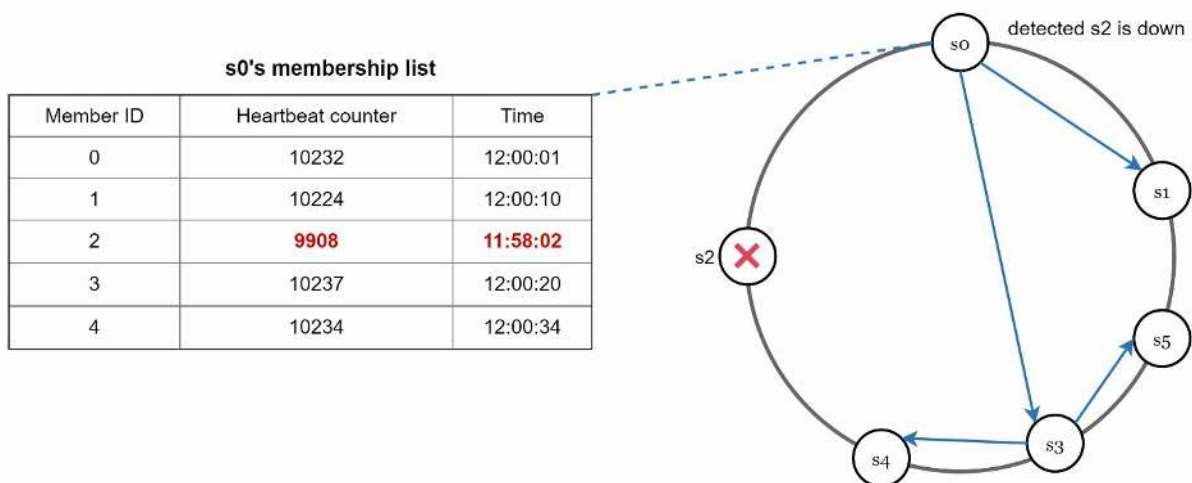


Figure 6-11

As shown in Figure 6-11: