

# **URL Shortener Deployment Considerations**

**By**  
**Gokul Balagopal**  
**PhD Student**  
**UT Dallas**

## Introduction

This document focuses on the key deployment considerations for scaling and ensuring high availability of a URL shortener service. This discussion is based on the simple URL shortener I have created using Python. The focus is on considering various strategies used to store and fetch the data in minimal time, while focusing on techniques to improve scalability and availability. It also explains the higher-level design aspects, such as horizontal scaling, the use of Twitter Snowflake for unique ID generation, and distributed deployments to ensure robust performance and continuous availability.

## Important Design Aspects

### *1. Database Management*

SQLite is used in this implementation for managing the database, for testing purposes. However, for production, we might want to consider using a more robust relational database like PostgreSQL, MySQL or other enterprise level solutions

### *2. Caching with Redis*

The use of Redis is for caching and it can significantly reduce the load on the database by storing frequently accessed data in memory.

### *3. URL Expiration*

This design allows for setting an expiration time for URLs. In my current design, when a new request is made to shorten a URL, the application checks if the long URL already exists in the database for that short URL. If it does exist, the application further checks if the existing short URL has expired. If the short URL has expired, then a new short URL would be generated, and if not, the existing short URL is returned.

### *4. Security Considerations*

The current design includes basic URL validation. For production, we might want to enforce stricter validation and also we may need to serve the application over HTTPS to ensure secure communication between the client and server.

### *5. Future Design Considerations:*

**Rate Limiting:** Implement rate limiting to protect the service from abuse (e.g., too many requests from a single IP).

**Custom Short URLs:** Allow users to specify custom short URLs rather than auto-generating them. This feature could be implemented with an additional check to ensure that the custom URL is unique (Ex: Similar to LinkedIn URLs)

## **For Higher level Design considering Scaling and Availability**

- Deploy the URL shortener application in a distributed environment. Configure load balancers to distribute traffic across multiple instances to avoid bottlenecks.
- Implement horizontal scaling by adding more nodes to handle increased load.
- Make unique and ordered ids using Twitter snowflake. IDs generated by Snowflake are time-based and roughly ordered, which is useful for sorting and analyzing data based on URL creation time. This can be distributed across DBs . Make sure there is clock synchronization to avoid duplication.
- Utilize a Base 62 hash function to convert long URLs into compact short URLs. Other hash functions like CRC32, MD5, or SHA-1 can be used, but they may lead to hash collisions. Hash collisions can be expensive to resolve as it requires querying the database to check for existing short URLs. Convert the raw hash value to Base 62 to create shorter and more user-friendly URLs. (Modern databases can handle long string so the URL length-based issue won't affect scalability much, but it may impact performance)
- We can use load balancers to distribute traffic across multiple servers. Load balancers can also detect if a server goes down and reroute traffic to healthy servers, avoiding service interruptions.
- We can store the URL mappings in a replicated database. Replication involves creating multiple copies of the database in different locations. If the primary database fails, the system can switch to a replica, minimizing downtime.
- We can implement database sharding to distribute the data across multiple databases or nodes. This improves performance and ensures that the failure of a single node doesn't affect the entire service. Sharding is especially useful when dealing with many URLs.
- Consider using distributed databases like Cassandra, which are designed to run across multiple nodes and provide high availability out of the box. These databases are designed to handle node failures and ensure that your data is always available.
- NoSQL databases like MongoDB offer built-in replication and sharding, making them suitable for high availability and scalability.
- Deploy the service across multiple geographical locations to ensure availability. If one server fails, others in different locations can take over, ensuring continuous service.

## **Conclusion**

By considering elements such as horizontal scaling, load balancing, distributed ID generation, and database sharding, the URL shortener system can be built to handle high traffic volumes with minimal downtime. With a focus on security, caching, and reliability, the deployment architecture can efficiently scale and provide seamless service, even in a production environment.