# Scalable URL Shortening Service Design Documentation

## 1. Problem Overview and Requirements

A. Functional Requirements:

  - URL Shortening: Convert a long URL into a unique, compact short URL.

  - Redirection: When a short URL is accessed, redirect the user swiftly to the original long URL.

  - Custom Aliasing (Optional): Allow users to specify custom aliases if desired.

  - Analytics (Optional): Track metrics such as click counts, geolocation, referrers, and timestamps.

B. Nonfunctional Requirements:

  - High Performance and Low Latency: Redirections must occur in milliseconds for a smooth experience.

  - Scalability: The service should support millions of URL creations and high read volumes per day.

  - Fault Tolerance and High Availability: The system must remain operational despite node failures or traffic spikes.

  - Data Consistency: Eventual consistency is acceptable as URLs are immutable post-creation.

  - Cost Efficiency: Resources must be managed efficiently with elastic scaling.

C. Assumptions and Estimates:

  - Write Load: Assume around 1 million URL shortenings per day (with peaks much higher than the average).

  - Read Load: Redirection requests are significantly higher than writes, making caching critical.

  - Data Volume: Over time, the database will accumulate billions of records, requiring partitioning and replication.

## 2. High-Level Architecture and Component Responsibilities

A. Client and Global Access Layer:

  - Client Devices: Mobile and web clients send HTTPS requests to create or access URLs. Real-time interactions (if needed) use WebSocket (over TLS).

  - Global Load Balancers: DNS-based and regional load balancers route incoming requests to the nearest available server clusters, reducing latency and evenly distributing load.

B. API Gateway and Authentication:

  - API Gateway: Authenticates incoming requests, enforces rate limiting, and routes them to the appropriate service. Autoscaling ensures handling of tens of thousands of requests per second using HTTPS externally and gRPC/REST (secured with mutual TLS) internally.

  - Authentication Service: Validates client credentials and issues secure tokens.

# Scalable URL Shortening Service Design Documentation

C. URL Shortening Service (Core Creation Logic):

  - Workflow:

    1. Validate the long URL and optionally check for duplicates (for idempotency).

    2. Generate a unique ID by calling an external Distributed ID Generator (using the Twitter Snowflake algorithm or similar).

     * The Distributed ID Generator is a dedicated microservice where each node is assigned a unique machine ID.

     * It creates a unique 64-bit ID using the current millisecond timestamp, the machine ID, and a sequence number.

     * The URL Shortening Service invokes this service via an internal API (using gRPC or REST over secured channels) and receives the unique ID.

    3. Convert the numeric ID into a Base62 string to produce a compact short code.

    4. Persist the mapping (short code  long URL) plus metadata (e.g., timestamp) in a NoSQL distributed database.

    5. Optionally update a Redis cache for rapid lookup during redirection.

    6. Return the short URL to the client.

D. URL Redirection Service:

  - Workflow:

    1. When a short URL is accessed, route the request through the API Gateway to the Redirection Service.

    2. First, attempt to retrieve the long URL from the Redis cache.

    3. On a cache miss, query the NoSQL database and update the cache.

    4. Issue an HTTP 301/302 redirect to send the client to the long URL.

  - Performance Goal: Achieve sub-100ms latency for the redirection process.

E. Data Persistence and Storage:

  - NoSQL Database: Stores URL mappings as records containing the short code, long URL, creation timestamp, and optional analytics data.

    * Data is partitioned by short code (or its hash) using consistent hashing.

    * Replication (e.g., factor of 3) ensures high availability and fault tolerance.

    * Scales horizontally by adding additional nodes.

  - Cache Layer: An in-memory store (e.g., Redis cluster) caches URL mappings to reduce database load and improve lookup speed.

# Scalable URL Shortening Service Design Documentation

- Optional Analytics Store: A separate layer for aggregating and analyzing click data, referrers, and geolocation information.

## 3. Detailed Data Flow

A. URL Shortening (Creation) Flow:

  1. A client sends a POST request (over HTTPS) with a long URL to the API Gateway.

  2. The API Gateway authenticates the request, enforces rate limiting, and routes it to the URL Shortening Service.

  3. The URL Shortening Service validates the URL, checks for duplicates, and calls the Distributed ID Generator to obtain a unique 64-bit ID.

  4. The unique ID is converted into a Base62 string to form the short URL code.

  5. The mapping is stored in the NoSQL database (with partitioning and replication) and optionally added to the Redis cache.

  6. The shortened URL is returned to the client.

B. URL Redirection Flow (Read Path):

  1. When a short URL is accessed, the request is routed via the API Gateway to the URL Redirection Service.

  2. The service first attempts to retrieve the long URL from the Redis cache.

  3. If not found in cache, it queries the NoSQL database, updates the cache, and fetches the long URL.

  4. An HTTP redirect (301/302) is issued to send the client to the long URL.

## 4. External Infrastructure, Protocols, and Traffic Handling

A. External Infrastructure  Distributed ID Generator:

  - We use a dedicated microservice implementing the Twitter Snowflake algorithm (or similar) to generate globally unique 64-bit IDs.

  - Each instance of the generator is assigned a unique machine ID.

  - When a new ID is requested, it combines the millisecond timestamp, the machine ID, and a sequence number to produce a unique value.

  - The URL Shortening Service invokes the generator via an internal API (using gRPC or REST over secured channels with mutual TLS).

B. Protocols:

# Scalable URL Shortening Service Design Documentation

   - Client-to-Server Communication: HTTPS is used for secure interactions between clients and the service.

   - Interservice Communication: gRPC or internal REST over secured TCP (with mutual TLS) ensures low-latency, secure communication between microservices.

C. Traffic Handling and Scalability:

   - Global and regional load balancers (using DNS-based routing and tools like HAProxy or AWS ELB) distribute incoming traffic efficiently.

   - Autoscaling ensures that services (API Gateway, URL Shortening Service, Distributed ID Generator) can handle tens of thousands of requests per second during peak periods.

   - The Redis cache significantly reduces database load by serving frequent redirection lookups.

   - The NoSQL database scales horizontally by adding nodes and uses partitioning and replication to manage billions of records over time.

## 5. Final Thoughts

This design for a URL shortening service (similar to Bitly) covers every critical aspect, including:

   - The URL Shortening Service uses a Distributed ID Generator based on the Twitter Snowflake algorithm to produce unique IDs, which are then Base62-encoded into compact short URLs.

   - Data persistence is achieved through a NoSQL database that partitions and replicates records to handle massive volumes.

   - An in-memory cache (Redis) ensures sub-100ms redirection response times.

   - Secure protocols and autoscaling strategies help manage high traffic efficiently while maintaining low latency and high availability.

This comprehensive approach results in a robust, scalable, and low-latency URL shortening service capable of supporting high volumes of traffic with excellent performance and reliability.