# Scalable Chat System Design Documentation

## 1. Requirements and Assumptions

A. Functional Requirements:

  - Real-Time Messaging: Support one-to-one and group chat with near-instant message delivery.

  - Offline Messaging: Ensure reliable delivery of messages to users when they come online.

  - Presence and Status Updates: Provide live user status, online/offline indications, and typing indicators.

  - Message History: Persist chat history with efficient retrieval and searching capabilities.

  - Multimedia Support: Optionally, support file attachments (images, videos) with text messages.

  - Notifications: Deliver real-time push notifications for new messages and mentions.

B. Nonfunctional Requirements:

  - Low Latency: Target end-to-end message delivery and presence updates within 100-200 ms.

  - High Throughput & Scalability: The system must support billions of users with millions of concurrent connections.

  - Fault Tolerance and High Availability: Utilize replication and redundancy so that the system remains operational even if some components fail.

  - Consistency: Strong ordering for one-to-one messages and near real-time updates for group chats; eventual consistency acceptable for presence.

  - Security: All communication secured via TLS (HTTPS and secure WebSockets) and interservice calls secured with mutual TLS.

C. Assumptions:

  - Only a small fraction of the total registered users are online at any time.

  - Peak load can be many times higher than average traffic.

  - Reads (retrieving messages and presence) dominate writes (sending messages).

  - In one-to-one conversations, message order is critical; group chat ordering can be slightly relaxed for scalability.

## 2. High-Level Architecture and Component Responsibilities

A. Client Tier:

  - Client Applications (Mobile, Web, Desktop): Use HTTPS and secure WebSocket (over TLS) for persistent connections.

B. Global Access and Load Balancing:

  - Global DNS and Regional Load Balancers: Route user requests to the nearest data center to

reduce latency and distribute load evenly.

C. Gateway and Authentication Services:

  - API Gateway: Authenticates new connection requests, enforces rate limiting, and routes HTTP/WebSocket upgrade requests to appropriate chat servers.

    * Uses HTTPS externally and gRPC or REST over secured TCP internally.

  - Authentication Service: Verifies user credentials and issues secure tokens (e.g., JWT).

D. Core Messaging Services:

  - Chat (Messaging) Servers: Maintain persistent WebSocket connections, process incoming messages, and use asynchronous, non-blocking I/O to support millions of concurrent connections.

  - Message Broker Cluster: Implements a publish/subscribe model (using Apache Kafka, RabbitMQ, or similar) for decoupling message production from consumption and ensuring ordered delivery where needed.

  - Presence Service and Session Store: Maintain real-time user status and connection information (using an in-memory store like Redis).

  - Notification Service: Sends push notifications for offline users via mobile push platforms and in-app alerts.

E. Data Persistence Layer:

  - Message Storage: A distributed NoSQL database (e.g., Cassandra, DynamoDB) stores chat logs. Data is partitioned by conversation or user ID and replicated across nodes (replication factor of 3 or more) for durability.

  - User Data Store: Stores profiles and social connections; may use a graph database or a relational/NoSQL store optimized for fast queries.

  - Caching Layer: An in-memory cache (e.g., Redis cluster) holds recent messages and presence data to reduce load on persistent stores and speed up reads.

## 3. Detailed Workflow

A. Connection Establishment:

  - User logs in via HTTPS to the API Gateway, which authenticates using the Authentication Service.

  - The client establishes a secure WebSocket connection to a Chat Server, and the server registers the session in the Session Store, marking the user as online.

B. Sending a Message (Write Path):

# Scalable Chat System Design Documentation

- The client sends a message via the established WebSocket connection to the Chat Server.

- The Chat Server validates the message and publishes it to the Message Broker with a conversation ID for ordering.

- The Timeline Service (or relevant Chat Servers) subscribes to the broker topic, receives the message, and forwards it in real-time to the recipient(s).

- For offline users, the message is persisted in the Message Store for later delivery.

- Acknowledgements or read receipts may be sent back to the sender.

C. Retrieving Chat History (Read Path):

- Upon opening a conversation, the client requests recent messages from the Chat Server.

- The server checks the in-memory cache; on a cache miss, it retrieves older messages from the NoSQL database, supporting pagination.

D. Presence and Typing Indicators:

- Lightweight status messages (online, offline, typing) are sent directly over the WebSocket connection.

- These events are also recorded in the Presence Service so that contact lists can be updated in near real time.

## 4. Scalability, Fault Tolerance, and Data Partitioning

A. Horizontal Scalability:

- Chat Servers and Message Brokers can be scaled horizontally by adding more instances behind load balancers.

- The NoSQL database and in-memory caches are partitioned (using consistent hashing) to distribute data evenly across nodes.

B. Fault Tolerance and High Availability:

- Replication is employed in persistent stores (e.g., replication factor of 3) and cache clusters to ensure data availability even if individual nodes fail.

- Redundant service instances (for API Gateways, Chat Servers, etc.) ensure that failure in one region or component does not bring down the system.

C. Data Consistency and Message Ordering:

- One-to-one chats require strict ordering, achieved by partitioning by conversation ID in the Message Broker so that messages are processed sequentially.

- For group chats or presence updates, eventual consistency may be acceptable to optimize for

performance.

## 5. Protocols and External Infrastructure

A. Communication Protocols:

  - Client-to-Server: HTTPS for REST API calls and secure WebSocket connections over TLS for real-time messaging.

  - Interservice: gRPC or REST over secured TCP (with mutual TLS) is used for communications between microservices (API Gateway, Chat Servers, ID Generators, etc.).

B. External Infrastructure  Distributed ID Generator (if needed):

  - Although not required for every chat system, if unique message IDs are needed for strict ordering or tracking, an external Distributed ID Generator can be used.

  - We would use a microservice implementing the Twitter Snowflake algorithm.

  - Each generator node is assigned a unique machine ID, and upon request, it returns a 64-bit unique ID constructed from the current timestamp, machine ID, and a sequence number.

  - The Chat Servers or Messaging Service calls this external service via gRPC/REST and uses the unique ID for message sequencing.

C. Load Balancing and Autoscaling:

  - Global load balancing (via DNS) and regional load balancers distribute traffic.

  - Autoscaling policies ensure that as the number of concurrent connections and message volumes increase, additional Chat Servers, API Gateways, and backend services are deployed automatically.

## 6. Final Thoughts

This design for a scalable chat system is structured to support billions of users, millions of concurrent connections, and high volumes of real-time messaging. Key highlights include:

  - A multi-tier architecture that separates client access, connection management, messaging, and data persistence.

  - Use of global and regional load balancers, along with autoscaling, to manage traffic spikes efficiently.

  - A robust real-time messaging layer powered by WebSocket connections, a Message Broker for decoupling, and services that ensure message ordering where needed.

  - Data is persisted in a distributed NoSQL database with partitioning and replication, while an in-memory cache minimizes latency.

# Scalable Chat System Design Documentation

  - Secure protocols (HTTPS, TLS, gRPC with mutual TLS) protect both external and interservice communications.

Overall, the proposed architecture ensures high performance, fault tolerance, and scalability, making it well-suited for a global, real-time chat system.