

2. Scaling Model

Perform capacity planning for a chat application that supports 40,000 concurrent users.

Estimate the network, storage, and compute requirements.

Use simple, back-of-the-napkin math to justify the scaling model.

Deliverables:

Diagrams & Documentation: Present findings in a markdown file, spreadsheet,

Perform capacity planning for a chat application that supports 40,000 concurrent users.

Capacity planning for a chat application with 40,000 concurrent users requires a comprehensive assessment of various factors. Here's a breakdown of the process:

1. Define Application Requirements:

- **Message Rate:**
 - Estimate the average number of messages per user per minute or second.
 - Consider peak usage periods.
- **Message Size:**
 - Determine the average message size (text, images, files).
- **Features:**
 - Identify features that impact resource usage (e.g., group chats, file sharing, video calls, presence indicators).
- **Data Storage:**
 - Estimate the amount of data to be stored (messages, user profiles, files).
 - Consider data retention policies.
- **Latency Requirements:**
 - Define acceptable latency for message delivery and other operations.
- **Availability Requirements:**
 - Determine the desired uptime (e.g., 99.9%, 99.99%).

2. Estimate Resource Requirements:

- **Server Resources:**
 - **CPU:**
 - Message processing, real-time communication, and application logic.
 - Consider using load testing to determine CPU usage under peak load.
 - **Memory (RAM):**
 - User sessions, message buffers, and application data.
 - Estimate memory usage per user and multiply by the number of concurrent users.
 - **Network Bandwidth:**
 - Message traffic, file transfers, and real-time communication.

- Calculate bandwidth based on message rate and size.
- **Storage:**
 - Message storage, user profiles, and file storage.
 - Estimate storage requirements based on data volume and retention policies.
- **Database Resources:**
 - **Database Servers:**
 - Message storage, user profiles, and other application data.
 - Consider using a scalable database solution (e.g., NoSQL, cloud-based database).
 - **Database Storage:**
 - Estimate database storage requirements based on data volume and retention policies.
 - Database IOPS.
- **Real-time Communication Resources:**
 - **WebSockets/SignalR Servers:**
 - Handle real-time communication between users.
 - Estimate the number of connections and message traffic.
 - **Load Balancers:**
 - Distribute traffic across multiple servers.
 - Ensure high availability and scalability.

3. Capacity Planning Calculations (Example):

- **Assumptions:**
 - Average message rate: 5 messages per user per minute.
 - Average message size: 1 KB.
 - 40,000 concurrent users.
- **Calculations:**
 - Total messages per minute: $40,000 \text{ users} * 5 \text{ messages/user} = 200,000$ messages.
 - Total data per minute: $200,000 \text{ messages} * 1 \text{ KB/message} = 200 \text{ MB}$.
 - Total data per second: $200\text{MB}/60 = \sim 3.3\text{MB/s}$
 - Bandwidth requirements depend on the connection type and overhead.
- **Server Calculations:**
 - These calculations are highly dependant on the software used, and the architecture of the application. Load testing is required.
 - Example: if each user requires 1MB of ram, then 40GB of ram is required.
 - CPU usage testing is required.

4. Technology Choices:

- **Cloud Platforms:** AWS, Azure, GCP (provide scalability and elasticity).
- **Databases:**
 - NoSQL databases (e.g., Cassandra, MongoDB) for high-volume message storage.
 - Relational databases (e.g., PostgreSQL, MySQL) for user profiles and other structured data.
 - Cloud-based databases (e.g., Azure Cosmos DB, AWS DynamoDB) for scalability and availability.

- **Real-time Communication:**
 - WebSockets.
 - SignalR.
 - Socket.IO.
- **Load Balancing:**
 - Cloud-based load balancers (e.g., Azure Load Balancer, AWS ELB).
 - Nginx, HAProxy.
- **Caching:**
 - Redis, Memcached.

5. Load Testing:

- **Simulate Peak Load:**
 - Use load testing tools (e.g., JMeter, Locust, Gatling) to simulate 40,000 concurrent users and their message traffic.
- **Measure Performance:**
 - Monitor CPU usage, memory usage, network bandwidth, and latency during load testing.
- **Identify Bottlenecks:**
 - Identify performance bottlenecks and optimize application and infrastructure.
- **Scale Gradually:**
 - Increase the number of concurrent users gradually to determine the application's breaking point.

6. Monitoring and Scaling:

- **Real-time Monitoring:**
 - Implement real-time monitoring of application and infrastructure metrics.
- **Auto-Scaling:**
 - Configure auto-scaling to automatically adjust resources based on demand.
- **Alerting:**
 - Set up alerts for critical metrics (e.g., high CPU usage, high latency).

Important Considerations:

- **Architecture:** The application's architecture plays a significant role in its scalability. Consider using a microservices architecture for better scalability.
- **Optimization:** Optimize database queries, message processing, and network communication for performance.
- **Redundancy:** Implement redundancy to ensure high availability and fault tolerance.
- **Security:** Implement security measures to protect user data and prevent unauthorized access.
- **Geographic Distribution:** If your users are geographically dispersed, consider using a content delivery network (CDN) and deploying your application in multiple regions.

By following these steps, you can create a robust capacity plan for your chat application and ensure that it can handle 40,000 concurrent users.

Estimate the network, storage, and compute requirements.

Let's break down the estimation of network, storage, and compute requirements for a chat application supporting 40,000 concurrent users, building upon the previous capacity planning outline.

Assumptions (For Example Calculations):

- Average message rate: 5 messages per user per minute.
- Average message size: 1 KB.
- 40,000 concurrent users.
- Data retention: 1 year.
- Average user profile size: 5 KB.
- Average file upload per user per day: 1 file, 1 MB in size.
- Peak traffic is 3 times the average.

1. Network Requirements:

- **Average Message Traffic:**
 - Total messages per minute: $40,000 \text{ users} * 5 \text{ messages/user} = 200,000$ messages.
 - Total data per minute: $200,000 \text{ messages} * 1 \text{ KB/message} = 200 \text{ MB}$.
 - Total data per second: $200 \text{ MB} / 60 \text{ seconds} = \sim 3.33 \text{ MB/s}$.
 - Peak Traffic: $3.33 \text{ MB/s} * 3 = \sim 10 \text{ MB/s}$.
- **File Transfer Traffic:**
 - Total file uploads per day: $40,000 \text{ users} * 1 \text{ MB/user} = 40,000 \text{ MB}$ (40 GB).
 - This will be spread throughout the day, so it will not add the full amount to the peak traffic.
- **Real-time Communication Overhead:**
 - WebSockets/SignalR connections introduce overhead.
 - Estimate 20% overhead for control messages and keep-alive signals.
- **Total Network Bandwidth:**
 - Consider peak traffic, real-time overhead, and file transfer bursts.
 - A minimum of 1 Gbps (125 MB/s) bandwidth is recommended. This allows for headroom.
 - For a cloud environment, ensure that the virtual machines or container instances have sufficient network throughput.

2. Storage Requirements:

- **Message Storage:**
 - Total messages per year: $200,000 \text{ messages/minute} * 60 \text{ minutes/hour} * 24 \text{ hours/day} * 365 \text{ days/year} = \sim 105 \text{ billion messages}$.
 - Total message storage: $105 \text{ billion messages} * 1 \text{ KB/message} = \sim 105 \text{ TB}$.
- **User Profile Storage:**
 - Total user profile storage: $40,000 \text{ users} * 5 \text{ KB/user} = 200 \text{ MB}$. This is very small.
- **File Storage:**
 - Total file storage per day: 40 GB.
 - Total file storage per year: $40 \text{ GB/day} * 365 \text{ days/year} = \sim 14.6 \text{ TB}$.
- **Database Storage:**
 - Include indexes, logs, and other database overhead.

- Add a 20-30% buffer to message and user profile storage.
- **Total Storage:**
 - Approximately 120-130 TB (including buffer and overhead).
 - Consider object storage for files, and a scalable database for messages.

3. Compute Requirements:

- **Application Servers:**
 - CPU:
 - Message processing, real-time communication, and application logic.
 - Load testing is essential to determine precise CPU requirements.
 - Start with a cluster of medium-sized VMs or containers (e.g., 8 vCPUs each).
 - Use auto-scaling to adjust based on CPU utilization.
 - Memory:
 - User sessions, message buffers, and application data.
 - Estimate memory usage per user and multiply by the number of concurrent users.
 - Start with 16-32 GB RAM per server.
- **Database Servers:**
 - CPU:
 - Database operations (queries, indexing).
 - Use a database cluster with sufficient CPU power.
 - Memory:
 - Database caching and operations.
 - Allocate sufficient RAM for the database.
 - Database IOPS:
 - The database will be very write heavy. Plan for high IOPS.
- **Real-time Communication Servers (WebSockets/SignalR):**
 - CPU:
 - Connection management and message routing.
 - Use a cluster of servers with load balancing.
 - Memory:
 - Connection state and message buffers.
 - Allocate sufficient RAM for connection management.
- **Load Balancers:**
 - Distribute traffic across application and real-time communication servers.
 - Use cloud-based load balancers for scalability and availability.

Important Notes:

- **Load Testing:** Load testing is crucial to validate these estimates and identify performance bottlenecks.
- **Cloud Scalability:** Cloud platforms provide elasticity, allowing you to scale resources as needed.
- **Monitoring:** Implement comprehensive monitoring to track resource utilization and performance.
- **Optimization:** Optimize application code, database queries, and network communication for performance.
- **Redundancy:** Implement redundancy to ensure high availability and fault tolerance.
- **Caching:** Use caching to reduce database load and improve performance.

- **Microservices:** Consider using a microservices architecture for better scalability and maintainability.
- **Geographic Distribution:** If your users are geographically dispersed, consider using a CDN and deploying your application in multiple regions.

Using simple, back-of-the-napkin math to justify the scaling model

let's do some quick, back-of-the-napkin calculations to justify our scaling model for the chat app with 40,000 concurrent users.

1. Message Traffic (Network):

- **Users:** 40,000
- **Messages/User/Minute:** 5
- **Total Messages/Minute:** $40,000 * 5 = 200,000$
- **Message Size:** 1 KB
- **Data/Minute:** 200,000 KB = 200 MB
- **Data/Second:** 200 MB / 60 seconds ≈ 3.3 MB/s
- **Peak (3x Average):** $3.3 \text{ MB/s} * 3 \approx 10 \text{ MB/s}$

Justification:

- Even at peak, 10 MB/s isn't a huge amount. But, we need to account for overhead and bursts.
- A 1 Gbps (125 MB/s) connection gives us ample headroom, ensuring smooth performance even during spikes.

2. Storage:

- **Messages/Year:** 200,000 messages/minute * 60 * 24 * 365 ≈ 105 billion messages
- **Message Storage:** 105 billion KB ≈ 105 TB
- **File Uploads/Day:** 40,000 users * 1 MB ≈ 40 GB
- **File Storage/Year:** 40 GB/day * 365 days ≈ 14.6 TB

Justification:

- Message storage will dominate. 105 TB is a lot, so we need a scalable storage solution (like cloud object storage).
- File storage adds another 15 TB, so we're looking at roughly 120-130 TB total.

3. Compute (Simplified):

- **Users:** 40,000
- **"User Load Unit":** Let's say, for simplicity, each user requires a small, consistent amount of compute, a "User Load Unit" (ULU). This ULU represents a combination of CPU and RAM.
- **Server Capacity:** A modest server can handle, say, 1,000 ULUs.
- **Total ULUs:** 40,000 users * 1 ULU/user = 40,000 ULUs
- **Servers Needed:** 40,000 ULUs / 1,000 ULUs/server = 40 servers (roughly)

Justification:

- This is a very rough estimate. The actual ULU value depends on the application's efficiency.
- This shows we'll need a cluster of servers, not just one.
- We need to add more servers for redundancy and peak loads.
- We also need to have separate servers for the database, and real time communication.

Key points:

- **Network:** We need sufficient bandwidth to handle peak message traffic and bursts.
- **Storage:** We need a scalable storage solution to accommodate the large volume of messages and files.
- **Compute:** We need a cluster of servers to handle the load, with auto-scaling to adjust based on demand.
- **Load Testing is Key:** These are just rough estimates. Load testing is essential to validate and refine our scaling model.
- **Cloud Benefits:** Cloud platforms offer the elasticity and scalability we need to handle this scale.