# 1. Introduction to System Design Interviews

System design interviews assess your ability to design scalable, maintainable, and efficient systems. They test how well you can handle real-world engineering challenges that involve distributed systems, scalability, reliability, and trade-offs.

## What Interviewers Evaluate:

- Structured thinking and clarity of communication
- Handling ambiguity through questions
- Awareness of scalability, fault tolerance, and performance
- Ability to justify trade-offs (cost, latency, consistency)
- Familiarity with system components and data flow

The goal isn't a perfect design — it's to show your thought process, how you make informed engineering decisions, and how you adapt under constraints.

---

# 2. Step-by-Step Framework

## 1 Clarify Requirements

Ask clarifying questions:

- What are functional vs non-functional requirements?
- What's the scale (users, requests/sec, data size)?
- What are the latency and availability expectations?

## 2 Define Scope

Avoid over-engineering. Focus on core use cases and ignore optional features.

## 3 Sketch a High-Level Architecture

Start with big blocks:
Client → Load Balancer → Application Servers → Database → Cache → CDN (if needed)

## 4 Identify Core Components

Define how each part fits in:

- API Gateway: Entry point for clients
- Application Layer: Contains business logic
- Database Layer: Manages persistent data
- Cache: Improves response times
- Queue: Handles asynchronous tasks

## 5 Design Data Model

Define entities, relationships, and query patterns. Choose between SQL or NoSQL based on use case.

## 6 Handle Bottlenecks and Scaling

Think about:

- Vertical vs horizontal scaling
- Database replication/sharding
- Load balancing and caching

## 7 Discuss Trade-offs

Explain pros/cons of design decisions:

- SQL vs NoSQL
- Strong vs eventual consistency
- Monolith vs microservices

## 8 Wrap-Up

End with a summary:

- How your design meets requirements
- Bottlenecks and future improvements
- Monitoring and observability plans

# 3. Core System Components

## Client Layer

User interface or external service sending requests (web, mobile, API).

## Load Balancer

Distributes incoming traffic across multiple servers to prevent overload.

## Application Servers

Process requests, enforce business logic, and communicate with backend systems.

## Database Layer

Stores persistent data. Choose:

- SQL (PostgreSQL, MySQL) for transactions
- NoSQL (MongoDB, Cassandra) for flexibility or high write throughput

## Cache

Stores frequently accessed data in memory (e.g., Redis, Memcached) to reduce latency and DB load.

## CDN

Delivers static assets (images, videos) quickly from edge locations.

## Queue System

Handles asynchronous or delayed tasks (Kafka, RabbitMQ, AWS SQS).

# 4. Backend-Focused Design Principles

## Microservices

Break the system into small, independent services that communicate through REST, gRPC, or message queues.

## Database Sharding

Split large datasets across multiple shards (horizontal partitioning) to handle high traffic.

## Replication

Maintain multiple copies of data for fault tolerance and faster reads.

## Caching

Reduce latency and improve performance by storing computed or frequently accessed results.

## Message Queues

Enable asynchronous processing:

- Kafka for event streaming
- RabbitMQ for reliable delivery

## Monitoring & Logging

Use:

- Prometheus / Datadog for metrics (latency, error rate, uptime)
- ELK stack (Elasticsearch, Logstash, Kibana) for centralized logging

# 5. Example Designs

## A. URL Shortener

Goal: Convert long URLs into short unique links.

Approach:

- Use hashing or base62 encoding to create short IDs
- Store mapping in a NoSQL database (high read/write)
- Cache popular URLs for instant lookup
- Handle collisions and rate limiting

Focus Areas: Hashing, caching, scalability

## B. Chat Application

Goal: Enable real-time communication between users.

Approach:

- Use WebSocket connections for real-time message delivery
- Store messages in a persistent store (Cassandra, DynamoDB)
- Use message queues to handle delivery reliability
- Partition users or chat rooms for scalability

Focus Areas: Real-time updates, message persistence, partitioning

## C. Notification System

Goal: Deliver notifications (push, SMS, email) reliably.

Approach:

- Producers publish events to Kafka
- Consumers process messages and deliver notifications
- Use retry and dead-letter queues for failures
- Track delivery metrics and retries

Focus Areas: Event streaming, asynchronous design, reliability

## 6. System Design Checklist

- ✅ Clarify scope and assumptions
- ✅ Identify main system components
- ✅ Choose appropriate data stores
- ✅ Handle scaling (load balancing, caching, sharding)
- ✅ Discuss trade-offs (latency, consistency, cost)
- ✅ Include monitoring, logging, and alerting
- ✅ Summarize your design and future improvements

---

## 7. Pro Tips

- Think aloud — your reasoning matters more than final architecture
- Prioritize bottlenecks over completeness
- Mention real metrics: latency goals, throughput, data size
- Always include failure handling and monitoring
- End strong: summarize trade-offs and how your design evolves at scale

---

## 📘 Conclusion

System design is less about drawing diagrams and more about communicating trade-offs. A great candidate shows structured thinking, flexibility, and awareness of real-world challenges.

"The best system design answers sound like engineering conversations, not rehearsed scripts."