

Key Strategies for High Concurrency & Low Latency

Handling high concurrency and low latency involves **using caching (Redis), asynchronous processing (message queues), database optimization (indexing, sharding, read replicas), and horizontal scaling**. Key techniques include connection pooling, non-blocking I/O (Node.js/Netty), and API gateways for throttling.

Key Strategies for High Concurrency & Low Latency

Caching Strategy: Implement fast, in-memory storage (e.g., Redis, Memcached) to cache frequently accessed data, which reduces database load and speeds up response times.

Asynchronous & Non-Blocking I/O: Utilize event-driven architectures (like Node.js) or asynchronous programming to handle thousands of requests without blocking threads.

Database Optimization:

Connection Pooling: Reuse database connections to eliminate the overhead of creating new ones.

Read/Write Splitting: Use master-slave replication to direct write traffic to the master and read traffic to replicas.

Sharding: Partition data horizontally to distribute load across multiple database instances.

System Architecture & Scaling:

Load Balancing: Distribute traffic across multiple backend instances.

Microservices: Break down monolithic applications into smaller, manageable services.

Message Queues: Use systems like Kafka or RabbitMQ to decouple tasks and process them asynchronously.

Code & Network Optimization:

Lightweight Data Formats: Use Protocol Buffers or MessagePack instead of JSON for smaller payloads.

Data Compression: Compress data to reduce network bandwidth and transfer time.

API Gateways: Use gateways for rate-limiting, authentication, and routing to protect backend services.

Monitoring & Performance Tuning:

Performance Monitoring: Use tools like Datadog, New Relic, or Zipkin to track latency and identify bottlenecks.

Profiling: Regularly analyze code to optimize critical paths.

Revision #1

Created 8 March 2026 12:45:46 by AI Channel

Updated 8 March 2026 12:46:39 by AI Channel