

System Design 101: 15 Essential Concepts Every Engineer Should Master

Having encountered many system design interviews at top tech companies, I've identified the core concepts that separate good engineers from great ones. These aren't just theoretical frameworks — they're battle-tested approaches that solve real-world scaling problems.

1. Load Balancing

Load balancers are the traffic directors of your architecture, ensuring no single server bears too much burden.

Key considerations:

- **Layer 4 vs Layer 7 load balancing:** Network-level balancing is faster but application-level allows for smarter routing based on content
- **Balancing algorithms:** Round robin is simple but doesn't account for server capacity; least connections and least response time are more sophisticated
- **Session persistence:** Critical for stateful applications to maintain user context

Real-world implementation:

AWS Elastic Load Balancer configuration:

- Health check path: /health
- Health check interval: 30 seconds
- Healthy threshold: 2
- Unhealthy threshold: 3
- Timeout: 5 seconds
- Algorithm: Least outstanding requests

Companies like Netflix use sophisticated load balancing to handle millions of concurrent streams while maintaining sub-100ms response times.

2. Caching Strategies

Caching dramatically reduces database load and improves response times, but comes with complexity.

Implementation patterns:

- **Cache-aside (lazy loading):** Load data into cache only when needed
- **Write-through:** Update cache when database is updated

- **Write-behind (write-back):** Batch cache updates to the database
- **Refresh-ahead:** Predict and refresh items before expiration

Example Redis caching implementation:

```
// Cache-aside pattern
public Data getDataById(String id) {
    // Try to get from cache first
    String cacheKey = "data:" + id;
    Data data = redisTemplate.opsForValue().get(cacheKey);

    if (data == null) {
        // Cache miss - get from database
        data = repository.findById(id);

        // Store in cache with expiration
        redisTemplate.opsForValue().set(cacheKey, data, 1, TimeUnit.HOURS);
    }

    return data;
}
```

The TTL (Time To Live) strategy is crucial — too short and you lose cache benefits, too long and you serve stale data.

3. Database Sharding

When your database grows beyond what a single server can handle, sharding divides it into smaller, more manageable pieces.

Sharding strategies:

- **Horizontal sharding:** Split rows across multiple servers
- **Vertical sharding:** Split columns across multiple servers
- **Directory-based sharding:** Use a lookup service to determine which shard contains data

Practical example of horizontal sharding by user ID:

```
-- Users with IDs 1-1,000,000 go to Shard 1
CREATE TABLE users_shard_1 (
    user_id INT PRIMARY KEY,
    username VARCHAR(50),
    email VARCHAR(100),
    ... other columns
);

-- Users with IDs 1,000,001-2,000,000 go to Shard 2
CREATE TABLE users_shard_2 (
    user_id INT PRIMARY KEY,
    username VARCHAR(50),
    email VARCHAR(100),
    ... other columns
);
```

Instagram famously scaled to billions of users by implementing a carefully designed sharding strategy across thousands of database servers.

4. Consistent Hashing

Traditional hash-based distribution falls apart when servers are added or removed. Consistent hashing elegantly solves this problem.

How it works:

1. Map both servers and data to positions on a conceptual ring
2. For each data item, move clockwise around the ring to find the first server
3. When adding/removing servers, only a fraction of data needs to be remapped

```

public class ConsistentHash<T> {
    private final HashFunction hashFunction;
    private final int numberOfReplicas;
    private final SortedMap<Integer, T> circle = new TreeMap<>();
    public ConsistentHash(HashFunction hashFunction, int numberOfReplicas,
Collection<T> nodes) {
        this.hashFunction = hashFunction;
        this.numberOfReplicas = numberOfReplicas;
        for (T node : nodes) {
            add(node);
        }
    }
    public void add(T node) {
        for (int i = 0; i < numberOfReplicas; i++) {
            circle.put(hashFunction.hash(node.toString() + i), node);
        }
    }
    public T get(Object key) {
        if (circle.isEmpty()) {
            return null;
        }
        int hash = hashFunction.hash(key);
        if (!circle.containsKey(hash)) {
            SortedMap<Integer, T> tailMap = circle.tailMap(hash);
            hash = tailMap.isEmpty() ? circle.firstKey() :
tailMap.firstKey();
        }
        return circle.get(hash);
    }
}

```

Distributed caching systems like Memcached and Redis Cluster use consistent hashing to minimize cache invalidation during scaling operations.

5. CAP Theorem

The CAP theorem states you can only guarantee two of three properties in a distributed system:

- **Consistency:** All nodes see the same data at the same time
- **Availability:** The system remains operational even if nodes fail
- **Partition Tolerance:** The system continues to function despite network partitions

Real-world trade-offs:

- **CP systems** (MongoDB, HBase): Prioritize consistency over availability
- **AP systems** (Cassandra, DynamoDB): Prioritize availability over consistency

- **CA systems:** Don't exist in practical distributed systems due to inevitable network partitions

Companies must choose based on business requirements — financial services typically favor consistency, while e-commerce may prioritize availability.

6. Eventual Consistency

Perfect consistency across distributed systems is prohibitively expensive. Eventual consistency provides a practical middle ground.

Implementation approaches:

- **Conflict-free Replicated Data Types (CRDTs):** Mathematical structures that resolve conflicts automatically
- **Vector clocks:** Track the sequence of events to determine causality
- **Quorum-based systems:** Require a minimum number of nodes to agree before committing changes

```
// DynamoDB example configuration
{
  "TableName": "UserProfiles",
  "BillingMode": "PAY_PER_REQUEST",
  "AttributeDefinitions": [
    { "AttributeName": "UserId", "AttributeType": "S" }
  ],
  "KeySchema": [
    { "AttributeName": "UserId", "KeyType": "HASH" }
  ],
  "GlobalSecondaryIndexes": [...],
  "StreamSpecification": {
    "StreamEnabled": true,
    "StreamViewType": "NEW_AND_OLD_IMAGES"
  }
}
```

Amazon's shopping cart system famously uses eventual consistency — adding an item to your cart is prioritized for availability, while financial transactions require stronger consistency guarantees.

7. Message Queues & Event Streaming

Asynchronous communication via message queues and event streaming platforms decouples system components and improves resilience.

Common patterns:

- **Point-to-point queues:** Each message consumed by exactly one receiver

- **Publish-subscribe:** Messages broadcast to multiple consumers
- **Dead letter queues:** Handle failed message processing

Kafka configuration example:

```
bootstrap.servers: kafka-broker1:9092,kafka-broker2:9092
group.id: order-processing-group
enable.auto.commit: false
auto.offset.reset: earliest
max.poll.records: 500
```

LinkedIn processes over 7 trillion messages per day through Kafka, enabling real-time analytics and features like "Who viewed your profile."

8. Rate Limiting

Protect your services from traffic spikes and potential DoS attacks with rate limiting.

Implementation algorithms:

- **Token bucket:** Allows bursts up to a limit
- **Leaky bucket:** Processes requests at a constant rate
- **Fixed window:** Counts requests in fixed time periods
- **Sliding window:** Smooths boundaries between time periods

Example implementation with Redis:

```
public boolean allowRequest(String userId, int limit, int windowSeconds) {
    String key = "ratelimit:" + userId;
    Long currentCount = redisTemplate.opsForValue().increment(key, 1);

    // If this is the first request in this window, set expiration
    if (currentCount == 1) {
        redisTemplate.expire(key, windowSeconds, TimeUnit.SECONDS);
    }

    return currentCount <= limit;
}
```

GitHub's API uses rate limiting to ensure fair usage — authenticated requests are limited to 5,000 requests per hour while unauthenticated requests are limited to just 60.

9. Microservices Architecture

Breaking monoliths into microservices improves development velocity and system resilience but introduces new challenges.

Design principles:

- **Single responsibility:** Each service handles one business capability
- **Database per service:** Services own their data and expose it via APIs
- **Smart endpoints, dumb pipes:** Business logic in services, not in communication layer
- **Decentralized governance:** Teams choose their own technologies

```
// Example microservice architecture for an e-commerce platform
- User Service (Auth, Profiles)
- Catalog Service (Products, Categories)
- Inventory Service (Stock management)
- Order Service (Order processing)
- Payment Service (Payment handling)
- Notification Service (Emails, SMS)
- Analytics Service (Business intelligence)
```

Amazon's transition to microservices transformed their development process — they now deploy code every 11.7 seconds on average.

10. API Gateway Pattern

An API Gateway serves as the single entry point for all clients, handling cross-cutting concerns and routing.

Common features:

- **Request routing:** Direct requests to appropriate microservices
- **Authentication/Authorization:** Centralized security
- **Rate limiting:** Protect backend services
- **Response caching:** Improve performance
- **Request/response transformation:** Protocol translation

Example Kong gateway configuration:

```

services:
- name: user-service
  url: http://user-service:8080
  routes:
    - name: user-routes
      paths:
        - /api/users
  plugins:
    - name: rate-limiting
      config:
        minute: 60
    - name: jwt

```

Netflix's Zuul gateway handles billions of requests daily, protecting their microservices ecosystem from direct exposure to clients.

11. Circuit Breaker Pattern

Prevent cascading failures by breaking the circuit when a service is failing.

States of a circuit breaker:

- **Closed:** Normal operation, requests flow through
- **Open:** Service failing, requests immediately rejected
- **Half-open:** Testing if service has recovered

```

// Spring Cloud Circuit Breaker example
@CircuitBreaker(name = "paymentService", fallbackMethod =
"paymentServiceFallback")
public PaymentResponse processPayment(PaymentRequest request) {
    return paymentServiceClient.processPayment(request);
}
public PaymentResponse paymentServiceFallback(PaymentRequest request,
Exception e) {
    return new PaymentResponse(PaymentStatus.PENDING, "Payment queued for
later processing");
}

```

Resilience4j configuration:

```

resilience4j.circuitbreaker:
  instances:
    paymentService:
      slidingWindowSize: 10
      failureRateThreshold: 50
      waitDurationInOpenState: 10000
      permittedNumberOfCallsInHalfOpenState: 5

```

Netflix's Hystrix pioneered this pattern at scale, dramatically improving their system resilience during partial outages.

12. Idempotency in Distributed Systems

Network failures can cause duplicate requests. Idempotent operations ensure the same outcome regardless of how many times they're executed.

Implementation strategies:

- **Idempotency keys:** Client-generated unique identifiers for operations
- **Conditional updates:** Only update if conditions are met (e.g., if-match headers)
- **Deduplication stores:** Track already processed requests

```
// REST API idempotent endpoint example
@PostMapping("/payments")
public ResponseEntity<PaymentResult> createPayment(
    @RequestHeader("Idempotency-Key") String idempotencyKey,
    @RequestBody PaymentRequest request) {

    // Check if we've seen this idempotency key before
    Optional<PaymentResult> existingResult =
        idempotencyStore.get(idempotencyKey);

    if (existingResult.isPresent()) {
        // Return the stored result from the previous execution
        return ResponseEntity.ok(existingResult.get());
    }

    // Process the payment
    PaymentResult result = paymentService.processPayment(request);

    // Store the result with the idempotency key
    idempotencyStore.store(idempotencyKey, result);

    return ResponseEntity.ok(result);
}
```

Stripe's payment API rigorously implements idempotency to ensure customers are never charged twice, even if network issues cause duplicate requests.

13. Data Partitioning vs. Replication

Understand when to partition data (for scale) versus replicate it (for availability and read performance).

Partitioning strategies:

- **Range-based:** Partition by value ranges (e.g., dates, alphabetical)
- **Hash-based:** Distribute evenly using hash functions

- **Directory-based:** Use a lookup service to find data location

Replication patterns:

- **Master-slave:** One writable master, multiple read-only replicas
- **Multi-master:** Multiple writable masters with conflict resolution
- **Quorum-based:** Majority of nodes must agree on changes

MongoDB configuration example:

```
# Sharded cluster with replication
sharding:
  clusterRole: shardsvr
replication:
  replSetName: rs0
net:
  port: 27017
storage:
  dbPath: /var/lib/mongodb
```

Google's Spanner database combines partitioning and replication with TrueTime to provide globally consistent distributed transactions.

14. System Monitoring & Observability

You can't manage what you can't measure. Comprehensive observability is essential for complex distributed systems.

The three pillars:

- **Metrics:** Quantitative measurements over time
- **Logging:** Structured event records
- **Tracing:** Tracking requests across services

Prometheus configuration example:

```
global:
  scrape_interval: 15s
  evaluation_interval: 15s
scrape_configs:
  - job_name: 'api-gateway'
    static_configs:
      - targets: ['api-gateway:9090']
  - job_name: 'user-service'
    static_configs:
      - targets: ['user-service:9090']
  - job_name: 'payment-service'
    static_configs:
      - targets: ['payment-service:9090']
```

Facebook's Canopy system processes over 100 trillion events daily, providing insights into system performance and user behavior.

15. Chaos Engineering

Proactively testing system resilience by deliberately injecting failures.

Key principles:

- **Start small:** Begin with non-critical systems
- **Contain blast radius:** Limit potential damage
- **Run in production:** Test real environments
- **Automate experiments:** Create reproducible tests

Example chaos experiment:

```
# Chaos Monkey configuration
chaos:
  enabled: true
  leashed: true # Safety mode
  schedule:
    cronExpression: "0 0 * * * ?" # Run daily at midnight
  strategy:
    type: "ShutdownInstance"
    probability: 0.005 # 0.5% chance per instance per day
  excludedInstances:
    - "prod-database-primary"
    - "payment-gateway-primary"
```

Netflix's Chaos Monkey randomly terminates EC2 instances in production, ensuring their systems are resilient to unexpected failures.

Conclusion

System design isn't just about theoretical knowledge — it's about making thoughtful trade-offs based on your specific requirements. No single architecture fits all problems. The best engineers deeply understand these fundamental concepts and know when and how to apply them.

Before your next system design interview or architecture review, consider how these patterns might apply to your specific challenge. What consistency model makes sense? How would you handle failure scenarios? What performance bottlenecks might emerge as you scale?

By mastering these 15 concepts, you'll be well-equipped to design resilient, scalable systems that stand the test of time and growth.