



"Provide detailed explanations for the following system design terms, along with relevant sources (articles, blogs, etc.) for further reading:

- 1. System Design**
- 2. DNS Resolution**
- 3. Fault Tolerance**
- 4. Vertical Scaling**
- 5. Downtime**
- 6. Over-Optimization**
- 7. Horizontal Scaling**
- 8. Load Balancer**
- 9. Load Balancer Algorithms**
- 10. Health (Active, Up, and Running)**
- 11. ELB (Elastic Load Balancer)**
- 12. Microservices Architecture**
- 13. Auth Services**
- 14. API Services**
- 15. Routing (Server)**
- 16. API Gateway**
- 17. EC2 (Elastic Compute Cloud)**
- 18. Auth0**
- 19. Batch Processing**
- 20. Synchronous vs. Asynchronous Processing**
- 21. Rate Limiting**
- 22. Bottleneck**
- 23. Queue System**
- 24. SQS (Simple Queue Service)**
- 25. Pull (Short vs. Long Pull) and Push Systems**
- 26. Event-Driven Architecture**
- 27. Dead Letter Queue**
- 28. Fan-Out Architecture**
- 29. Overwhelm Rate Limiting (Leaky Bucket, Token Bucket)**
- 30. Database Primary Node, Analytic, Read, and Write Replicas**
- 31. CDN (Content Delivery Network), Anycast, and CloudFront**
- 32. Latency**

For each term, provide a clear and concise definition, and include relevant sources for deeper understanding. Ensure explanations are technical and accurate, focusing on system design principles and practices."

A thorough set of detailed explanations with relevant sources will be compiled for each of the 32 system design terms requested. Here is the first part of the detailed definitions along with references for further reading:

1. System Design

System design is the process of defining the architecture, components, modules, interfaces, and data flow of a complex software system to meet specific functional and non-functional requirements like scalability, reliability, and maintainability. It involves creating a blueprint that guides implementation and future growth, balancing trade-offs among performance,

cost, and complexity.

Source: GeeksforGeeks, [Swimm.io](#), [AlgoMaster.io](#), Snappify [\[1\]](#) [\[2\]](#) [\[3\]](#) [\[4\]](#)

2. DNS Resolution

DNS (Domain Name System) resolution is the process of translating human-readable domain names (e.g., [example.com](#)) into IP addresses that computers use to route network requests. It operates as a distributed hierarchical database system and an application layer protocol essential for internet functionality.

Source: GeeksforGeeks [\[5\]](#) [\[6\]](#)

3. Fault Tolerance

Fault tolerance is the ability of a system to continue operating correctly even when some components fail. This is achieved through redundancy, replication, error detection, and recovery mechanisms that ensure reliability and availability despite faults.

Source: [Swimm.io](#) [\[2\]](#)

4. Vertical Scaling

Vertical scaling (scaling up) means increasing the capacity of a single machine by adding more resources such as CPU, RAM, or storage to handle greater workload. It's simple but limited by hardware constraints and can require downtime during upgrades.

Source: GeeksforGeeks, Multiplayer, LinkedIn [\[7\]](#) [\[8\]](#) [\[9\]](#)

5. Downtime

Downtime is the period when a system or service is unavailable or non-functional. It impacts reliability and user experience and ideally is minimized through fault tolerance and failover strategies.

Source: General system reliability principles [\[2\]](#)

6. Over-Optimization

Over-optimization happens when a system or code is excessively tuned prematurely, leading to unnecessary complexity, maintainability challenges, and possibly degrading overall performance or scalability.

Source: General system design best practices [\[3\]](#)

7. Horizontal Scaling

Horizontal scaling (scaling out) increases capacity by adding more machines or nodes to a system, distributing the load across them. It enhances availability and fault tolerance better than vertical scaling but requires additional system management complexity.

Source: [Swimm.io](#), Multiplayer [\[7\]](#) [\[2\]](#)

8. Load Balancer

A load balancer is a device or software that distributes incoming network or application traffic evenly across multiple servers to optimize resource use, maximize throughput, minimize response time, and avoid overload on any single server.

Source: GeeksforGeeks [\[6\]](#) [\[5\]](#)

9. Load Balancer Algorithms

Load balancer algorithms determine how traffic is distributed and can include round-robin, least connections, IP hash, random, weighted round-robin, and others, each suited to different scenarios for balancing workload efficiently.

Source: General load balancing principles [\[5\]](#) [\[6\]](#)

10. Health (Active, Up, and Running)

Health of a system or component means it is actively functioning within expected parameters, responsive, and able to serve requests without errors. Monitoring health is essential for fault detection and triggering failover mechanisms.

Source: System monitoring principles^[5]

11. ELB (Elastic Load Balancer)

ELB is a managed load balancing service provided by AWS that automatically distributes incoming application traffic across multiple EC2 instances or containers, scales elastically, and integrates health checks for fault tolerance.

Source: AWS documentation general knowledge^[5]

12. Microservices Architecture

Microservices architecture is a design pattern where an application is built as a collection of loosely coupled services, each responsible for a specific business capability and independently deployable, enhancing scalability and maintainability.

Source: GeeksforGeeks, DesignGurus^{[6] [5]}

13. Auth Services

Authentication services manage user identity verification, allowing access to resources only to authenticated users, using methods like OAuth, JWT, and session tokens.

Source: General security and auth service principles^[6]

14. API Services

API services provide defined interfaces for external or internal clients to interact with software components, enabling communication and data exchange via well-defined protocols like REST or gRPC.

Source: System design best practices^[6]

15. Routing (Server)

Server routing is the process of directing incoming requests to appropriate backend services or resources based on URL paths, HTTP methods, or other criteria.

Source: Web server architecture basics^[6]

16. API Gateway

An API Gateway acts as a single entry point for all client requests in a microservices architecture, routing requests, enforcing policies like authentication and rate limiting, and aggregating responses.

Source: GeeksforGeeks, DesignGurus^{[5] [6]}

17. EC2 (Elastic Compute Cloud)

EC2 is an AWS cloud service offering resizable compute capacity (virtual servers) in the cloud, allowing users to run applications on scalable virtual machines.

Source: AWS general knowledge^[5]

18. Auth0

Auth0 is a cloud identity platform that provides authentication and authorization as a service, simplifying user identity management across applications.

Source: Auth0 company info and security blogs^[6]

19. Batch Processing

Batch processing involves executing a series of jobs or tasks on a large set of data without

manual intervention, often scheduled and used for offline or deferred processing.

Source: General system design concepts [\[10\]](#)

20. Synchronous vs. Asynchronous Processing

Synchronous processing waits for a task to complete before continuing, while asynchronous processing allows tasks to proceed independently, improving efficiency and responsiveness.

Source: General programming concepts [\[10\]](#)

21. Rate Limiting

Rate limiting controls the number of requests a user or service can make to a system within a time period to prevent overload and ensure fair usage.

Source: System design blogs [\[10\]](#) [\[6\]](#)

22. Bottleneck

A bottleneck is a point in the system where performance or throughput is limited, constraining overall system capacity.

Source: Performance tuning literature [\[10\]](#)

23. Queue System

A queue system facilitates message passing between components or services asynchronously, decoupling producer and consumer systems.

Source: Distributed system texts [\[6\]](#)

24. SQS (Simple Queue Service)

SQS is an AWS-managed message queuing service that enables decoupling and scaling of microservices, serverless applications, and distributed systems.

Source: AWS general knowledge [\[6\]](#)

25. Pull (Short vs. Long Poll) and Push Systems

Pull systems actively request data (short poll checks quickly, long poll waits for data); push systems send data to consumers as soon as it's available.

Source: Messaging system principles [\[6\]](#)

26. Event-Driven Architecture

Event-driven architecture uses events (notifications) to trigger communication and actions between loosely coupled components, allowing highly scalable reactive systems.

Source: System design patterns [\[10\]](#)

27. Dead Letter Queue

A dead letter queue is a queue that stores messages that could not be processed successfully for later analysis or reprocessing.

Source: Messaging systems [\[6\]](#)

28. Fan-Out Architecture

Fan-out architecture involves sending a message or request to multiple recipients or services in parallel to increase efficiency and responsiveness.

Source: Distributed systems literature [\[10\]](#)

29. Overwhelm Rate Limiting (Leaky Bucket, Token Bucket)

These are algorithms to control data flow and prevent system overwhelm—leaky bucket smooths bursts by a fixed outflow rate; token bucket allows bursting up to token capacity.

Source: Rate limiting and traffic shaping papers [\[6\]](#)

30. Database Primary Node, Analytic, Read, and Write Replicas

Primary node handles all writes; replicas are copies for read or analytics workloads, improving performance and availability.

Source: Database design principles^[6]

31. CDN (Content Delivery Network), Anycast, and CloudFront

CDNs distribute content geographically; Anycast routes user requests to nearest server; CloudFront is AWS's CDN offering.

Source: CDN technical docs^[6]

32. Latency

Latency is the time delay between a request and the corresponding response, critical for user experience and system performance.

Source: Performance metrics literature^[10]

Would a deeper dive with source-by-source extraction for topics like DNS Resolution, Load Balancer Algorithms, or others be preferable before proceeding with completing the entire list?

**

1. <https://www.geeksforgeeks.org/system-design/getting-started-with-system-design/>
2. <https://swimm.io/learn/system-design/system-design-complete-guide-with-patterns-examples-and-techniques>
3. <https://snappify.com/blog/system-design-components>
4. <https://algomaster.io/learn/system-design/what-is-system-design>
5. <https://www.geeksforgeeks.org/system-design/system-design-interview-bootcamp-guide/>
6. <https://www.designgurus.io/blog/system-design-interview-fundamentals>
7. <https://www.multiplayer.app/distributed-systems-architecture/system-design-primer>
8. <https://www.geeksforgeeks.org/system-design/system-design-horizontal-and-vertical-scaling/>
9. <https://www.linkedin.com/pulse/system-design-horizontal-scaling-vs-vertical-harsh-kumar-sharma-jadpf>
10. https://dev.to/ajmal_hasan/mastering-system-design-24-key-concepts-for-beginners-198b