

General Web Backend System Questions

1. What is the difference between horizontal scaling and vertical scaling?

In vertical scaling we increase the resources in the existing system like RAM and CPU in order to be able to handle more load, while in horizontal scaling we add more machines. When we don't need the load we can shut down other systems hence it is cost effective while in vertical scaling we can't shut down.

Relational databases are more like vertically scalable and NoSQL databases more like horizontal scalable.

2. What is the difference between API Gateway and Load Balancer?

An API gateway's primary function is to provide a unified interface for clients to access backend services. A load balancer's primary function is to distribute traffic across a group of servers.

An API gateway operates at the application layer, while a load balancer operates at the network layer.

3. What is the difference between Reverse Proxy and Forward Proxy?

A Forward Proxy sits between a client and the internet, masking the client's IP address. It's often used for anonymity, bypassing restrictions, or providing centralized internet access.

A Reverse Proxy sits between the frontend and backend servers, masking the server's IP address.

It's commonly used for load balancing, security, and caching. By distributing traffic across multiple servers, it improves performance and reliability. Additionally, it can protect servers from direct attacks by acting as a shield.

4. What is the difference between Microservices and Monolithic architecture?

Monolithic: Uses a single code base for the entire application, with all components tightly coupled. Monolithic applications are typically complex and difficult to modify, and scaling them can be inefficient.

Microservices: Uses multiple, independent code bases for each service, which are deployed on their own servers. Microservices are designed to be scalable, and each service can be updated, modified, or deployed independently.

Microservices are easier to understand and support because the complexity is moved from the application to the infrastructure. allow for independent teams to develop and manage their own services, while monolithic architectures typically involve large teams working together.

5. What is the difference between vertical and horizontal partitioning?

Vertical partitioning: Splits data by columns, so each partition contains a subset of the table's columns. This approach is good for optimizing access to columns that are often queried together.

Horizontal partitioning: Also known as sharding, this method splits data by rows, so each partition contains a subset of the table's rows. This approach is good for managing large datasets.

6. What is a Rate Limiter? How does it work?

A rate limiter is a tool that controls the number of requests a client can make to a server within a given time period. It's a cybersecurity technique that helps maintain system stability and security by: Limiting network traffic, Preventing users from exhausting system resources, Protecting against cyber attacks, and Ensuring fair usage of resources among multiple clients.

Rate limiting works by:

- a. Tracking IP addresses that requests are coming from
- b. Tracking how much time elapses between each request
- c. Temporarily blocking IP addresses that are making too many requests
- d. Using a sliding window algorithm to count the number of requests made within overlapping time windows

7. How does Single Sign-On (SSO) work?

Single Sign-On (SSO) is a user authentication method that allows users to access multiple applications with one set of credentials. It works by establishing a trust relationship between a service provider (SP) and an identity provider (IdP)

SSO implementation revolves around a central server. All applications trust this master server and use it to access login credentials. If your organization has multiple applications, they probably run on different servers, using different user directories and authentication mechanisms.

8. How does Apache Kafka work, and why is it so fast?

Apache Kafka is a messaging system that's fast and scalable because it uses a combination of techniques to process data efficiently:

Partitioning: Kafka uses a partitioned log model, which breaks up logs into segments called partitions. Each partition corresponds to a different subscriber, allowing for multiple subscribers to the same topic.

Sequential I/O: Kafka uses sequential I/O, which is faster than random I/O because the disk head can move in a straight line.

Zero copy: Kafka uses zero copy to save multiple data copies between the application context and kernel context.

Batch compression: Kafka uses batch compression, which is more effective as the data size increases.

Queues: Kafka uses queues instead of trees, which are better for a messaging system that requires many read and write operations.

Byte arrays: Kafka uses byte arrays to transfer data, which allows for compression using serializers and deserializers.

Kafka is also highly scalable because it can distribute data across multiple servers. It's designed to handle high throughput, low latency, and large volumes of data.

9. What is the difference between Kafka, ActiveMQ, and RabbitMQ?

Kafka: Designed for high-throughput and real-time stream processing, Kafka is a good choice for handling large amounts of data. It uses a distributed log-based model and partition-based design. Kafka is also scalable and can persist data on disk.

RabbitMQ: A good choice for low-latency messaging and flexible routing, RabbitMQ is versatile and developer-friendly. It uses a queue-based message model and supports a wide range of languages and legacy protocols. RabbitMQ is also good for complex message routing.

ActiveMQ: A versatile option with support for multiple protocols and enterprise features, ActiveMQ is a good choice for a robust solution. It supports protocols like OpenWire, STOMP, AMQP, MQTT, and HTTP web service protocols. ActiveMQ also supports SSL for encryption.

RabbitMQ excels in low-latency messaging and flexible routing, making it ideal for certain use cases. Kafka stands out for its high throughput, real-time stream processing capabilities, and scalability. ActiveMQ offers versatility with support for multiple protocols and advanced enterprise features.

10. What is the difference between JWT, OAuth, and SAML?

11. What is a Content Delivery Network (CDN), and how does it improve performance?

CDNs rely on a process called “caching” that temporarily stores copies of files in data centers across the globe, allowing you to access internet content from a server near you. Content delivered from a server closest to you reduces page load times and results in a faster, high-performance web experience.

12. Explain the CAP theorem. How do you deal with its trade-offs in distributed systems?

CAP theorem states that a distributed system can only guarantee two of the three properties (consistency, availability, and partition tolerance) at the same time. Therefore, when designing and implementing a distributed system, trade-offs must be made between these properties.

System designers often make dynamic adjustments based on the specific needs and circumstances of their systems. They can also use consensus algorithms like Raft or Paxos to maintain consistency and coordination among distributed nodes.

13. What is caching, and what are the different types of caches (e.g., in-memory, distributed)?

Caching is a technique that stores data in a temporary storage area, called a cache, to reduce the time and resources needed to fetch data from the primary storage. There are several types of caches, including:

Cache memory: There are two types of cache memory: primary and secondary. Primary cache memory is located on the CPU, while secondary cache memory is located on a separate chip close to the CPU.

Web server caching: Stores data for reuse, such as when a user visits a page for the first time and then requests the same page again.

Browser caching: Stores a copy or subset of data in the web browser to make cached websites load and look up faster.

Application caching: Allows an application to store data that is expensive to create.

CDN caching: Content delivery networks (CDNs) store cached data in multiple proxy servers that are geographically distributed.

Disk cache: Stores data that has been recently read from or written to a disk, such as a hard disk drive or solid-state drive.

Filesystem caching: Useful for caching static files or relatively static database queries.

Persistent cache: Data is not lost upon a system restart or crash.

Distributed Redis cache: Redis is a popular in-memory data store that can be used as a high-performance distributed cache.

14. What are the differences between strong consistency and eventual consistency?

Strong consistency means the latest data is returned, but, due to internal consistency methods, it may result with higher latency or delay. With eventual consistency, results are less consistent early on, but they are provided much faster with low latency.

Eventual consistency: The weakest consistency level, with no guarantees. It offers the best performance and lowest latency. It's ideal when an application doesn't need ordering guarantees.

Strong consistency: A common data consistency model for distributed databases. It ensures that writes to a record are applied in a specific order, and writes aren't skipped or reordered.

Sequential consistency: A higher level of consistency where all writes are globally ordered. It doesn't matter which thread made the write or which data item was written to.

Bounded staleness: A consistency level between session consistency and strong consistency. It allows reads to lag behind writes for a specific amount of time.

Replication: The process of copying data from one database or server to another. It's usually done to increase fault-tolerance, availability, and consistency between computing resources.

15. What is the difference between synchronous and asynchronous communication in distributed systems?

In distributed systems, synchronous and asynchronous communication differ primarily in how components interact with respect to time and dependencies. In synchronous communication, the sender waits for a response from the receiver before it can proceed with further operations. This means that the sender is "blocked" until it receives the necessary reply, which creates a real-time interaction where both parties must be available simultaneously.

Asynchronous communication allows the sender to continue its operations immediately after sending a message, without waiting for an immediate response. The sender and receiver can function independently, with the sender moving on to other tasks while the receiver processes the message whenever it becomes available.

16. What is the role of a message broker, and how do message queues improve system architecture?

A message broker plays a key role in distributed systems by acting as an intermediary that facilitates the exchange of messages between different services or components. Its primary function is to receive messages from a sender (producer) and route them to the appropriate receiver (consumer), ensuring that communication between services remains decoupled and reliable. By handling message storage, delivery, and potentially transformations or routing logic, a message broker allows different parts of a system to interact without needing to know each other's details or be online simultaneously.

Common examples of message brokers include **RabbitMQ**, **Apache Kafka**, and **ActiveMQ**.

Message broker introduces a queue-based mechanism where messages are stored until they are processed by consumers, message queues improve system architecture by enabling decoupling, asynchronous processing, scalability, and reliability. This results in more flexible, robust, and maintainable distributed systems.

17. What is sharding (or partitioning) in databases, and how does it improve performance?

Database sharding, also known as horizontal partitioning, is a technique that splits large databases into smaller parts, or shards, to improve performance and scalability:

Sharding distributes data across multiple database servers, or machines, so that each server can store and process a smaller amount of data. The shards are individual partitions that can be accessed, managed, and configured separately.

Reducing query scan time: Sharding can significantly improve query performance and reduce response times.

Providing fault tolerance: Sharding provides built-in fault tolerance.

Sharding is useful when you have large datasets that can be logically separated, such as by date, time, or category. It can also be used when all data still fits on a single server, but you want to improve performance.

18. What are the benefits of using a database index, and what are the trade-offs?

A database index is a data structure that improves the speed of data retrieval operations on a database table by providing quick access to rows. While indexes offer significant performance benefits, they also come with certain trade-offs, like Increased Storage Requirements. Every time data is inserted, updated, or deleted in an indexed column, the index must also be updated to reflect these changes. This can slow down write operations (INSERT, UPDATE, DELETE).

19. What is database replication, and what are the different types (e.g., master-slave, master-master)?

Database replication is the process of copying data from one database to another, which can be done in a variety of ways. Some common types of database replication include:

Master-slave replication: A primary database server (the master) is the source of data, and other database servers (the slaves) maintain copies of the master's data. The master handles writes, while the slaves handle reads. This type of replication can be synchronous or asynchronous, depending on when changes are propagated.

Master-master replication: Multiple databases can act as both sources and targets of replication. This type of replication is often used in collaborative platforms, such as Google Docs, where users can edit documents in real time.

Masterless replication: Data is replicated across multiple nodes that are all equal, so no single node can bring down the entire cluster.

Database replication can help ensure that data is always available to users and that the system can handle the load. For example, social media platforms like Facebook and Instagram use master-slave replication to distribute data across multiple servers.

20. What is the purpose of a load balancer, and what types are there (e.g., round-robin, least connections)?

A load balancer's purpose is to distribute network traffic across multiple servers to prevent any one server from becoming overloaded. Load balancers can be hardware or software devices.

There are several types of load balancing algorithms, including:

Round robin: A simple algorithm that distributes requests in a rotating order, ensuring that each server receives an equal share of requests.

Least connections: Directs traffic to the server with the fewest connections to clients. Weighted least connections: Similar to least connections, but allows administrators to assign different weights to each server based on their capacity.

Weighted response time: Considers the number of connections each server has open and its average response time to determine where to send traffic.

IP hash: Assigns a client's IP address to a fixed server.

21. How does a distributed transaction work, and what challenges does it introduce?

A distributed transaction is a series of operations that take place across multiple databases or systems, and it's a key part of distributed computing. Distributed transactions are complex and require coordination to ensure that data is consistent and accurate across the network.

Here's how distributed transactions work:

Coordination: A coordinator ensures that all participants in the transaction either commit or abort the transaction.

Atomicity: All or none of the operations in the transaction are performed. If an error occurs, all changes are undone and the system is returned to a consistent state.

Two-phase commit (2PC): A protocol that involves two phases:

1. **Prepare phase:** Each participant prepares to commit the transaction and votes either "yes" or "no".
2. **Commit phase:** If all participants vote "yes", the transaction is committed. If any participant votes "no", the transaction is aborted and changes are rolled back

Distributed transactions can introduce **challenges, such as:**

Network latency: Coordinating across different systems can introduce network latency, which can delay transaction processing.

Blocking issues: Protocols like 2PC can be vulnerable to blocking issues, which can leave transactions in limbo if a system component fails.

22. What is eventual consistency, and when is it acceptable in a system?

Eventual consistency is a data consistency model that ensures that updates to a distributed database will eventually be reflected across all nodes. This means that if no new updates are made to a data item, all reads of that item will eventually return the last updated value.

Eventual consistency is considered a "weak" consistency scheme. It helps ensure high availability of the data, but at the risk of some data being stale or out of date. It can also

be a problem for application developers because it puts them at risk of putting data into an inconsistent state.

When it's acceptable: Eventual consistency is sufficient for some applications, especially for datasets that are not updated frequently.

23. What is the difference between a stateless and a stateful service?

Stateful: Retains data from one session to the next. Stateful applications save client session data on the server, which allows for faster processing and historical context. They require backing storage, like a database, to preserve user data. Stateful applications are good for personalized experiences, but they require careful data management and can increase resource consumption.

Stateless: Doesn't preserve data between sessions. Stateless applications rely on external entities, like databases or cache, to manage state. Each request must contain all the necessary information to be processed independently. Stateless applications are generally more scalable and fault-tolerant than stateful applications.

24. What are the main challenges of designing a highly available system?

Some challenges to designing a highly available system include:

Cost: Designing for high availability can be expensive and requires more resources like hardware, software, bandwidth, and energy.

Complexity: Adding too many layers of redundancy and failover mechanisms can make the system more difficult to manage and increase the chance of misconfigurations.

Single points of failure: Failing to identify and address single points of failure can undermine the entire high availability strategy.

Consistency and integrity: Replicating or distributing data across multiple locations can create issues with consistency and integrity in the presence of failures or errors.

Threats and attacks: Scaling can open up the system to threats and attacks such as DoS, data breaches, or malicious insiders.

Disaster recovery: Disaster recovery is a challenging problem in IT platforms, especially in cloud computing.

Testing and maintenance: Regular testing is essential to ensure HA systems work as expected, which requires time and resources.

25. How would you design a notification system (e.g., for email, SMS, push notifications)?

26. What is the difference between ACID and BASE in terms of database properties?

The main difference between ACID and BASE database models is how they prioritize consistency and availability:

ACID: Prioritizes consistency over availability. ACID stands for Atomicity, Consistency, Isolation, and Durability. ACID databases are strongly consistent, which means they guarantee that transactions are processed reliably and predictably. ACID databases are good for applications that require data consistency, such as banking and financial systems. However, they can be less flexible and may restrict access to applications during network or power outages.

BASE: Prioritizes availability over consistency. BASE databases are weakly consistent, which means they allow the system to remain operational even if some parts fail. BASE databases are good for applications that require high availability and scalability, such as **social media platforms and e-commerce websites**. However, they may sacrifice some data consistency for a short period of time.

27. What is a circuit breaker pattern, and when would you use it?

The circuit breaker pattern is a software development design pattern that prevents an application from repeatedly trying to execute an operation that's likely to fail:

When to use it: You can use the circuit breaker pattern to prevent repeated failures when a component, service, or resource is temporarily unavailable. This is especially useful in distributed and microservices architectures.

How it works: The circuit breaker pattern places an intermediary service between a caller and a target. The intermediary service monitors the target's conditions and reroutes traffic to another service if a hazard occurs.

The circuit breaker pattern is inspired by the physical circuit breaker for electric lines. When an electric line becomes overloaded, the circuit breaker trips and stops electricity from flowing through the circuit's wire.

You can use the circuit breaker pattern along with other patterns such as **retry**, **fallback**, and **timeout pattern** to make the system more fault tolerant.

28. How does distributed caching work, and what are common strategies for invalidation?

Distributed caching stores frequently accessed data across multiple nodes in a network to improve data access performance. Cache invalidation is the process of removing or updating cached data to maintain data consistency.

How distributed caching works: A distributed caching system consists of cache nodes, cache clients, and a load balancer:

1. **Cache nodes:** Servers that store and manage cached data.
2. **Cache clients:** Applications or services that interact with the cache nodes to read and write data.
3. **Load balancer:** Distributes the load across the cache nodes.

Common Cache invalidation strategies:

Time-to-Live (TTL): Cached data is automatically invalidated after a specified time.

Event-based invalidation: When a significant event occurs, such as a data update, the cache is invalidated for the relevant data.

First in, first out (FIFO): The oldest cache entry is invalidated first.

Purge: Specific cached content is removed for a particular object.

29. What is the role of a reverse proxy in system architecture?

A reverse proxy sits between client devices and a web server. It intercepts incoming client requests, forwards them to the web server, and returns the server's response to the client.

A reverse proxy is a server, app, or cloud service that improves the security, performance, and scalability of websites, cloud services, and content delivery networks (CDNs) such as **Nginx, Apache, and HAProxy**.

Security: A reverse proxy protects the origin server's identity by masking its IP address.

Performance: A reverse proxy can distribute requests across multiple servers to ensure no single server is overwhelmed.

Scalability: A reverse proxy can provide seamless scaling options.

Load balancing: A reverse proxy can decide where and how to route HTTP sessions.

Caching: A reverse proxy can store cached content for faster retrieval.

SSL encryption: A reverse proxy can offload SSL/TLS processing.

Live activity monitoring and logging: A reverse proxy can monitor all the requests that get passed through it

30. What is a heartbeat mechanism in distributed systems, and why is it important?

A heartbeat mechanism in a distributed system is a periodic communication process that monitors the health of nodes, networks, and interfaces. It's a common technique in mission critical systems that helps ensure high availability and fault tolerance.

Here's how a heartbeat mechanism works:

1. A node sends a heartbeat, which is a small packet of data, at regular intervals to another node or a monitoring system.
2. The monitoring system checks for the heartbeat within a specified time frame.
3. If the monitoring system doesn't receive a heartbeat, it assumes the node is down or experiencing issues.
4. The system can then take corrective actions, such as restarting the service, alerting administrators, or rerouting traffic.

Heartbeats are important because they help prevent cluster partitioning and detect network or system failures.

31. How do you handle traffic spikes in high-demand systems?

Handling traffic spikes in high-demand systems requires strategies that ensure scalability and performance. Auto-scaling adjusts server capacity automatically based on demand, while load balancing distributes traffic across multiple servers to prevent overloading.

Caching (e.g., using CDNs and in-memory caches like Redis) reduces database load by serving frequently requested data quickly. Rate limiting and throttling control the number of requests a client can make, protecting against overloads or abuse.

Message queues (e.g., RabbitMQ) allow tasks to be processed asynchronously, preventing system slowdowns during spikes. Database scaling through read replicas, sharding, or connection pooling ensures the database can handle increased traffic efficiently. Using a microservices architecture enables independent scaling of services,

while monitoring and alerting systems help detect and respond to spikes in real time. Together, these techniques ensure resilience and performance during traffic surges.

32. What is a dead letter queue (DLQ), and when would you use it?

A dead letter queue (DLQ) is a message queue that stores messages that a software system cannot process or deliver.

DLQs prevent the source queue from overflowing with unprocessed messages.

Messages are moved to the DLQ after a set number of unsuccessful attempts to process them. The messages are then stored for analysis and no longer affect the activity of applications.

Messages can be removed from the DLQ and inspected. An application might let a user correct issues and resubmit the message.

33. How would you design a system to handle eventual consistency in a large distributed setup?

Handling eventual consistency in distributed systems ensures that while updates aren't immediately reflected everywhere, the system eventually becomes consistent. Use databases like **Cassandra or DynamoDB**, designed for eventual consistency, with data sharding and replication across nodes for high availability.

To resolve conflicts, use **timestamps** (last-write-wins) or vector clocks to track concurrent changes. Implement quorum-based reads and writes, where only a subset of replicas needs to acknowledge requests, balancing availability and consistency.

Asynchronous replication ensures updates propagate across nodes in the background without blocking, keeping systems responsive. For event-driven systems, use message brokers (e.g., Kafka) and event sourcing to propagate and store changes.

Client-side strategies like caching and retry mechanisms help manage temporary inconsistencies. Monitoring tools like Prometheus ensure visibility into replication lag and system health, ensuring eventual consistency is maintained while optimizing for performance and availability.

34. What is the role of an API gateway in microservices architecture?

An API gateway in a microservices architecture acts as a mediator between client applications and backend services. It's a software layer that performs several tasks, including:

Routing: Forwards requests to the appropriate microservice based on URL, content, or business rules

Security: Applies authentication, access control, and threat protection

Traffic control: Performs load balancing, caching, and rate limiting

Orchestration: Discovers services and handles failures like retries and circuit breaking

Observability: Logs, monitors, and traces

Transformation: Translates protocols and formats responses

Authentication and authorization: Verifies credentials like API keys and implements authorization policies

Protocol translation: Translates protocols between clients and microservices

Dynamic routing: Routes the request to the appropriate service instance or load balancer

An API gateway can help:

1. Decrease complexity
2. Enhance security and manageability
3. Provide a unified entry point for all API calls
4. Aggregate data and resources and deliver them to the user in a unified manner
5. Allow clients to create a single call request data
6. Enable protocol translations automatically
7. Reduce the number of request/response calls

35. What are idempotent operations, and why are they important in distributed systems?

Idempotent operations are operations that produce the same result when performed multiple times. They are important in distributed systems because they help ensure consistency and reliability when dealing with issues like network failures, request retries, or duplicated requests.

36. How does OAuth 2.0 work, and how is it used in authorization?

Resource owner: This is the user that is granting third-party access to their data.

Client: This is the third-party application that is requesting access to the resource owner's data. When the resource owner grants access, the client gets an access token that can be used to request the resources within the granted scope.

Authorization server: The authorization server acts as an intermediary between the client, resource owner, and resource server by issuing access tokens to the client after the resource owner successfully authorizes the request.

Resource server: This is the server that hosts the protected resources. It is responsible for accepting and responding to requests to access protected resources using an access token.

Here is a step-by-step breakdown of everything that happens under the hood:

1. The client (the meal planning app) asks the user for access to their resources on the resource server of the fitness app.
2. The user grants access to the client through the authorization server by logging in to the fitness app with their credentials. These credentials are not shared with the client. Instead, an authorization code is generated and shared with the client.
3. The client uses this authorization code to request an access token from an endpoint that is provided by the authorization server.
4. The authorization server generates and returns an access token, which the client can use to access the user's resources on the resource server.

5. The client sends the access token to the resource server to request access to the user's resources.
6. The resource server validates the access token with the authorization server. If the token is valid, it grants the client access to the user's resources.

37. How do you prevent distributed denial of service (DDoS) attacks in a system?

To prevent Distributed Denial of Service (DDoS) attacks, a combination of techniques is used to mitigate the risk of overwhelming system traffic.

1. **Traffic Filtering and Rate Limiting:** Implement rate limiting to restrict the number of requests a client can make, and filter traffic by blocking known malicious IPs or allowing trusted IPs.
2. **Web Application Firewalls (WAF):** WAFs (e.g., Cloudflare WAF, AWS WAF) block malicious traffic by filtering HTTP requests for attack patterns, protecting applications from DDoS attempts.
3. **Content Delivery Networks (CDNs):** CDNs (e.g., Cloudflare, Akamai) distribute traffic across global servers, absorbing excess load and preventing volumetric DDoS attacks from reaching the origin server.
4. **Auto-Scaling and Load Balancing:** Auto-scaling in cloud environments and load balancers distribute traffic across multiple servers, ensuring high availability even during spikes caused by DDoS attacks.
5. **DDoS Protection Services:** Services like AWS Shield, Google Cloud Armor, or Cloudflare provide automatic protection at both the network and application layers, mitigating large-scale attacks.
6. **Anomaly Detection:** Real-time monitoring tools (e.g., Prometheus, Datadog) detect unusual traffic spikes early, enabling quick mitigation.

These strategies together help absorb, distribute, and block malicious traffic, ensuring system resilience.

38. What is a quorum in distributed systems, and why is it important for consensus?

In a distributed system, a quorum is the minimum number of nodes that must agree on a value before it's considered valid. Quorums are important for consensus because they ensure that a distributed system can function reliably and make consistent decisions, even when nodes fail or networks partition.

The most common type of quorum is majority-based, where more than half of the nodes must agree on a decision. For example, in a system with five nodes, at least three must agree for a decision to be made.

Quorums are particularly important in distributed databases, replication, and clustering.

39. How would you design a real-time analytics platform?

Real-time analytics architecture requires integrating data from multiple sources, such as IoT sensors, social media, and transactional systems, in real time. This can be a complex process that requires expertise in data management and integration.

A scalable streaming platform is a key component of a real-time analytics platform. Some options include Apache Kafka, Redpanda, AWS Kinesis, Azure Stream, and Google Dataflow.

We can use scrapers, collectors, agents, and listeners to capture live streaming data and store it in a database

40. What is a Service Mesh, and how does it help in managing microservices?

A service mesh is a software layer that manages communication between microservices in an application. It helps with microservices management.

This layer is composed of containerized microservices. As applications scale and the number of microservices increases, it becomes challenging to monitor the performance of the services.

A service mesh works by using lightweight containers to create a transparent infrastructure layer over the application. This layer is made up of proxies, called "sidecars", that run alongside each service. The sidecars route requests between microservices, and the control plane coordinates their behavior.