

1. System Architecture

This defines the fundamental structure and organization of a software system.¹ The choices made here impact scalability, maintainability, and deployment.²

Monolithic vs. Microservices

Feature	Monolithic Architecture	Microservices Architecture
Structure	A single, unified codebase where all components (UI, business logic, data access) are tightly coupled and run as a single service.	The application is broken down into a collection of small, independent services, each responsible for a specific business capability.
Deployment	The entire application must be redeployed for any change, no matter how small.	Services can be deployed independently, allowing for faster and more flexible updates.
Scalability	You must scale the entire application, even if only one small part is a bottleneck. This can be inefficient.	You can scale individual services based on their specific needs, leading to more efficient resource utilization.
Technology	Constrained to a single technology stack (e.g., all Java, all Python).	Each service can use the best technology stack for its specific job (polyglot architecture).
Complexity	Simple to develop and test initially. Becomes complex and difficult to manage as the application grows.	More complex to set up initially due to the distributed nature. Simpler to understand and manage individual services.

Fault Isolation	A failure in one component can bring down the entire application.	A failure in one service can be isolated and is less likely to impact the rest of the system (if designed correctly).
------------------------	---	---

Modular Design

Modular design is a principle that applies to both monoliths and microservices. It involves breaking a system down into smaller, independent, and interchangeable modules.³

- **In a Monolith:** This means organizing the code into distinct packages or modules with well-defined boundaries (e.g., a "user management" module, a "payment" module).⁴ This improves maintainability even if it's a single deployable unit.
- **In Microservices:** The modular boundary is the service itself.⁵ Each microservice is a self-contained module.⁶

Key Architectural Qualities

- **Scalability:** The system's ability to handle a growing amount of work.⁷ This is often achieved through **horizontal scaling** (adding more machines/services), which is a natural fit for microservices.⁸
- **Availability:** The percentage of time the system is operational and accessible. It's often expressed in "nines" (e.g., 99.9% is "three nines"). High availability is achieved through redundancy, failover mechanisms, and eliminating single points of failure.⁹ Microservices can improve availability because a single service failure doesn't have to cause a total system outage.¹⁰
- **Consistency:** In distributed systems, this ensures that all nodes or services see the same data at the same time.¹¹ There is often a trade-off between consistency and availability, as defined by the **CAP Theorem** (Consistency, Availability, Partition Tolerance).¹² Systems may choose **Strong Consistency** (all replicas have the same data at all times) or **Eventual Consistency** (replicas will become consistent over time).¹³

2. Infrastructure

This refers to the core hardware and software components required to run the application.

- **Databases:** The heart of most applications, responsible for persistent data storage.¹⁴ The choice between SQL (for structured data with transactional needs) and NoSQL (for unstructured data and high scalability) is a critical architectural decision.
- **Storage:** This includes block storage for server disks (like Amazon EBS), object storage for files, images, and backups (like Amazon S3 or Google Cloud Storage), and file storage for shared access.
- **CPU (Compute):** The processing power required to run your application logic. The amount of CPU needed depends on the computational intensity of your tasks.
- **Bandwidth (Network):** The amount of data that can be transferred to and from your servers.¹⁵ This is crucial for user-facing applications and data-heavy internal communication.

Appropriate Scaling Methods

- **Vertical Scaling (Scaling Up):** Increasing the resources of a single server (more CPU, more RAM).¹⁶ It's simple but has a hard limit and can be expensive.
- **Horizontal Scaling (Scaling Out):** Adding more servers to distribute the load.¹⁷ This is the standard for modern, scalable applications and is enabled by load balancers.¹⁸
- **Auto-Scaling:** Automatically adjusting the number of servers based on real-time metrics like CPU utilization or request count.¹⁹ This is a cost-effective way to handle variable traffic.

3. Security

Security is not an afterthought; it must be designed into the system from the beginning.²⁰

- **JWT (JSON Web Token):** A standard for creating access tokens. When a user logs in, the server generates a JWT containing user details (the "payload") and signs it with a secret key.²¹ The token is sent to the client, which includes it in the header of subsequent requests.²² The server can then verify the token's signature to authenticate the user without needing to look up session data in a database, making it ideal for stateless, distributed systems.²³
- **IP Whitelisting:** A security measure where you create a list of approved IP addresses that are allowed to access your system or specific endpoints.²⁴ This is commonly used

to restrict access to sensitive internal admin panels or APIs to only office or VPN IP
addresses.²⁵

- **RSA Encryption:** An **asymmetric encryption** algorithm that uses a key pair: a **public key** (which can be shared with anyone) and a **private key** (which must be kept

secret).²⁶

- **Encryption:** Data encrypted with the public key can *only* be decrypted with the corresponding private key.²⁷ This is used to securely transmit data. A client can encrypt data with a server's public key, knowing that only the server can read it.
- **Digital Signatures:** Data signed with the private key can be verified by anyone with the public key. This is used to verify the identity of the sender (e.g., JWT signatures).

- **PII (Personally Identifiable Information) Handling:** PII is any data that can be used to identify a specific individual (e.g., name, address, email, phone number, national ID).²⁸

Handling PII requires strict security measures:

- **Encryption at Rest:** Encrypting the data when it is stored in the database or on disk.²⁹
 - **Encryption in Transit:** Using TLS/SSL to encrypt data as it travels over the network.³⁰
 - **Access Control:** Limiting who can access PII data based on their role.³¹
- **Data Protection Compliance:** Adhering to legal frameworks for data protection, such as:
 - **GDPR (General Data Protection Regulation):** EU law governing data protection and privacy for all individual citizens of the European Union and the European Economic Area.³²
 - **CCPA (California Consumer Privacy Act):** A state statute intended to enhance privacy rights and consumer protection for residents of California.³³

4. Observability

Observability is about understanding the internal state of your system by observing its external outputs.³⁴ It is often described by three pillars: logs, metrics, and traces.³⁵

- **Tools:**

- **Prometheus:** A powerful open-source tool for collecting and storing time-series data (**metrics**).³⁶ It periodically "scrapes" metrics from configured endpoints on your services.³⁷
 - **Grafana:** An open-source visualization tool.³⁸ It connects to data sources like Prometheus to create dashboards, graphs, and alerts, allowing you to visualize your system's health and performance.³⁹
 - **OpenTelemetry:** An open-source observability framework that provides a standardized way to generate and collect telemetry data (metrics, logs, and traces) from your applications.⁴⁰ It aims to unify how you instrument your code.
 - **Trace/Span IDs for Debugging:** In a microservices architecture, a single user request might travel through dozens of services. Distributed tracing is essential for debugging.⁴¹
 - **Trace ID:** A unique identifier assigned to an entire request as it enters the system.⁴² This ID is passed along to every service that the request touches.
 - **Span ID:** A unique identifier for a single unit of work or operation within a service (e.g., a database call or an API call to another service).⁴³
 - By logging the **Trace ID** and **Span ID** in every service, you can aggregate all the logs for a single request and visualize the entire call graph, making it easy to pinpoint which service is causing an error or latency.⁴⁴
-

5. Architecture Migration

Moving from a monolith to a microservices architecture is a complex process that must be done carefully.⁴⁵

Strategies for Transitioning from Monolith to Microservices

- **Strangler Fig Pattern:** This is the most common and recommended approach. Instead of a "big bang" rewrite, you gradually "strangle" the monolith by replacing its functionality piece by piece with new microservices.⁴⁶
 - An **API Gateway** or proxy is placed in front of the monolith.⁴⁷
 - You identify a module to extract (e.g., user profiles).⁴⁸
 - You build a new microservice for user profiles.

- You update the API Gateway to route all calls for user profiles to the new microservice instead of the monolith.
 - Repeat this process for other modules until the monolith is either gone or has shrunk to a manageable core.
 - **Handling Database Calls:** This is often the hardest part.
 - **Shared Database:** Initially, the new microservice might still read from and write to the monolith's database. This is a temporary anti-pattern but can ease the transition.
 - **Database per Service:** The ideal state is for each microservice to own its own data. This requires carefully migrating data from the monolith's database to a new database for the service and refactoring the monolith to call the new service's API for that data instead of accessing the table directly.
- 49

6. System Estimations (Back-of-the-Envelope Calculations)

These are quick, approximate calculations to estimate system requirements. They are crucial in system design interviews to demonstrate your ability to reason about scale.

Example: Estimating Bandwidth and Storage for a Simple Image Service

Assumptions:

- 1 million Daily Active Users (DAU).
- Each user uploads 1 image per day on average.
- Each user views 20 images per day on average.
- Average image size: 200 KB.

Calculations:

- **Storage Estimation (per day):**
 - $1,000,000 \text{ users/day} * 1 \text{ image/user} * 200 \text{ KB/image} = 200,000,000 \text{ KB/day} = 200 \text{ GB/day}$
 - **Storage for 1 year:** $200 \text{ GB/day} * 365 \text{ days} = 73,000 \text{ GB} = 73 \text{ TB}$
- **Bandwidth Estimation (Egress/Outbound Traffic per day):**
 - $1,000,000 \text{ users/day} * 20 \text{ images/user} * 200 \text{ KB/image} = 4,000,000,000 \text{ KB/day} = 4,000 \text{ GB/day}$
 - **Bandwidth per second (average):** $4,000 \text{ GB} / (24 \text{ hours} * 3600 \text{ sec/hour}) \approx 46 \text{ MB/s}$
 - You would need to design for peak traffic, which might be 2-3x the average, so roughly $\sim 100-150 \text{ MB/s}$.
- **Cost Estimation (using AWS S3 and CloudFront as an example):**
 - **S3 Storage Cost:** Look up the current price for S3 Standard (e.g., $\sim \$0.023 \text{ per GB/month}$).⁵⁰ For 73 TB, this would be a significant cost.

- **CDN/Bandwidth Cost:** Look up AWS CloudFront data transfer out pricing (e.g., ~\$0.085 per GB).⁵¹ $4,000 \text{ GB/day} * 30 \text{ days} * \$0.085/\text{GB} = \$10,200/\text{month}$.

These rough numbers help you make informed decisions about technology choices (e.g., the need for a CDN is immediately obvious) and infrastructure capacity.

7. Database Scaling

As your application grows, the database often becomes the bottleneck.⁵²

- **Replication:**
 - This involves creating copies (replicas) of your database. A common pattern is **Primary-Replica (or Master-Slave) Replication**.⁵³
 - All **writes** go to the primary database.
 - The primary database then replicates the changes to one or more replica databases.⁵⁴
 - All **reads** are directed to the replicas.⁵⁵
 - This is an excellent way to scale for **read-heavy loads**, as you can add more replicas to handle increased read traffic. It also provides high availability; if the primary fails, a replica can be promoted to become the new primary.⁵⁶
- **Sharding (Partitioning):**
 - Replication helps with read loads but not write loads, as all writes still go to one primary server.⁵⁷ Sharding solves this by partitioning the data across multiple databases (shards).⁵⁸
 - Each shard contains a subset of the data.⁵⁹ For example, you could shard by `user_id`, where users 1-1,000,000 are on Shard A, users 1,000,001-2,000,000 are on Shard B, and so on.⁶⁰
 - When a request comes in, the application logic (or a proxy layer) uses a **shard key** (e.g., `user_id`) to determine which shard contains the relevant data.
 - Sharding allows for massive write scalability but adds significant complexity to the application.⁶¹

8. Interview Prep Focus Areas

For system design interviews, focus on demonstrating a structured thought process.

- **Load Balancers:** Understand their role in distributing traffic, enabling horizontal scaling, and improving availability. Know the difference between Layer 4 (TCP) and Layer 7 (HTTP) load balancing.
- **API Gateways:** Understand that they are more than just load balancers. They act as a single entry point for clients, handling concerns like authentication (JWT validation), rate limiting, routing to different microservices, and response aggregation.⁶²
- **Real-World Problem Solving:**
 1. **Clarify Requirements:** Always start by asking questions to understand the functional and non-functional requirements (e.g., "How many users?", "What is the expected latency?", "Is consistency critical?").
 2. **Make High-Level Decisions:** Start with a simple architecture (e.g., a load balancer, a few web servers, a database).
 3. **Identify Bottlenecks:** As you discuss scale, identify potential bottlenecks (e.g., "The database will be read-heavy").
 4. **Introduce Solutions:** Propose solutions for these bottlenecks (e.g., "We can add a caching layer," "We can use database replication").
 5. **Be Pragmatic:** Discuss trade-offs. No solution is perfect. Explaining *why* you chose one approach over another (e.g., "I chose eventual consistency to improve availability") shows a deep understanding.
 6. **Do Back-of-the-Envelope Calculations:** Use estimations to justify your design choices.