

40 Backend Interview Questions with Answers

A complete guide for Recruiters, Hiring managers and Candidates

This document covers the most important Backend interview questions across fresher, intermediate, and expert levels.

HOW TO USE THIS GUIDE

This guide is built for **structured, competency-based Backend interviewing**. Each question includes:

- **The Question:** Ready to ask directly
- **What a Strong Answer Covers:** Key elements expected
- **Strong Answer Example:** What a top candidate sounds like
- **Weak Answer Example:** What bluffing/low-prep sounds like
- **Recruiter Evaluation Cue:** What to listen for
- **Score (1–5):** Use the scale below

Scoring Scale

	Label	What It Means
5	Exceptional	Field-ready, structured thinking, strong judgment
4	Strong	Good practical understanding, minor gaps
3	Competent	Basic understanding, limited field depth
2	Developing	Surface-level, generic answers
1	Not Ready	Incorrect / no clarity

Hire Threshold:

Candidates should average ≥ 3.5 across all questions for a conditional offer. A score of ≥ 4.0 on role-critical questions is strongly preferred.

PART 1: BACKEND INTERVIEW QUESTIONS FOR FRESHERS (Q1–Q15)

Section 1: Core Logic, Data Structures & Algorithms

Q1. What is the difference between an Array and a Linked List, and when would you choose one over the other?

Strong Answer: An array stores elements in contiguous memory locations, allowing $O(1)$ constant time access via index, but resizing or inserting elements in the middle is expensive $O(n)$. A linked list uses nodes pointing to the next element, making insertions and deletions fast $O(1)$ if you have the reference, but searching takes sequential linear time $O(n)$. I would choose an array for fixed-size lookup-heavy data, and a linked list for dynamic, insertion-heavy data.

Weak Answer: An array is just a list of items inside square brackets, and a linked list is a more advanced structure used for complex data routing when you don't know how big your data is.

Recruiter Cue: Tests **Foundational Data Structure Knowledge**.

Q2. How does a Hash Table/HashMap work, and what happens during a "hash collision"?

Strong Answer: A hash table uses a hash function to compute an index into an array of buckets, from which the desired value can be found. A collision occurs when two distinct keys produce the same hash value. It is usually resolved using "Chaining" (storing colliding elements in a linked list or binary tree at that index) or "Open Addressing" (searching for the next empty slot in the array).

Weak Answer: It maps keys to values instantly. Collisions happen when you enter duplicate keys by accident, which causes the program to throw an error or overwrite the previous value.

Recruiter Cue: Tests **Memory Mapping & Collision Resolution**.

Q3. Can you explain the difference between Recursion and Iteration, and what is a risk associated with deep recursion?

Strong Answer: Iteration uses a loop structural framework (like `for` or `while`) to repeat a block of code, executing within a single stack frame. Recursion happens when a function calls itself, creating a new stack frame for each call. The primary risk of deep recursion is a **Stack Overflow**, which occurs if the call stack runs out of memory because the recursion went too deep or missed its base case.

Weak Answer: Recursion is cleaner code and loops are faster. The risk of recursion is that it can run forever and freeze your computer if you forget to stop it.

Recruiter Cue: Tests **Memory Stack Mechanics**.

Q4. What is the purpose of "Big O Notation," and what does it mean if an algorithm has a time complexity of $O(n^2)$ versus $O(\log n)$?

Strong Answer: Big O notation is used to describe the upper bound performance or execution time of an algorithm as the input size (n) grows. $O(n^2)$ represents quadratic time complexity, meaning execution time grows proportionally to the square of the input (like nested loops; inefficient for large datasets). $O(\log n)$ represents logarithmic time, where the problem size is halved at each step (like Binary Search; highly efficient and scalable).

Weak Answer: Big O tells you exactly how many seconds a program will take to run on a standard server. $O(n^2)$ is bad because it takes twice as long, and $O(\log n)$ is good because it records logs efficiently.

Recruiter Cue: Tests **Algorithmic Efficiency Evaluation**.

Q5. When implementing a sorting algorithm, what is the trade-off between Bubble Sort and Quick Sort?

Strong Answer: Bubble Sort is simple to understand and implement with $O(1)$ auxiliary space complexity, but its average and worst-case time complexity is $O(n^2)$, making it unusable for large data. Quick Sort uses a divide-and-conquer strategy, yielding a much faster average-case time complexity of $O(n \log n)$, making it vastly superior for real-world production data despite a worst-case of $O(n^2)$ if the pivot is chosen poorly.

Weak Answer: Bubble Sort moves items one by one like bubbles in water, and Quick Sort is just a built-in library function that sorts arrays instantly using optimized machine code.

Recruiter Cue: Tests **Optimization Awareness**.

Q6. What are the key differences between a Stack and a Queue data structure, and can you name a real-world computing example for one?

Strong Answer: A Stack follows LIFO (Last In, First Out) mechanics, where elements are pushed and popped from the same end; a real-world example is the "Undo" history in a text editor or the browser back-button stack. A Queue follows FIFO (First In, First Out), where elements enter at the rear and leave from the front; an example is a print job queue or a messaging queue handling background tasks.

Weak Answer: Stacks store data vertically in memory, and Queues store them horizontally. They are both used to hold lists of elements until the backend is ready to process them.

Recruiter Cue: Tests **Execution Flow Design**.

Section 2: Databases, APIs & System Fundamentals

Q7. What is the fundamental difference between a Relational (SQL) and a Non-Relational (NoSQL) database?

Strong Answer: SQL databases are relational, table-based, have strict, predefined schemas, and support complex joins and ACID compliance (ideal for structured data like financial transactions). NoSQL databases (like document or key-value stores) are non-relational, have dynamic schemas for unstructured data, and scale horizontally across multiple servers easily (ideal for rapid scaling, real-time data, or content management).

Weak Answer: SQL is old and uses tables with rows, while NoSQL is modern, much faster, doesn't use schemas at all, and is written entirely in JSON format.

Recruiter Cue: Tests **Data Architecture Foundations**.

Q8. What are "ACID properties" in a database transaction? Can you explain at least two?

Strong Answer: ACID ensures database reliability. **Atomicity** means the entire transaction succeeds or completely rolls back; it's all-or-nothing. **Consistency** ensures data transitions from one valid state to another, maintaining constraints. **Isolation** ensures concurrent transactions don't interfere with each other. **Durability guarantees** committed data survives system crashes.

Weak Answer: It stands for Acid, Consistency, Inside, and Database. It basically means your database is protected from data leaks and corruption when multiple users log in at the same time.

Recruiter Cue: Tests **Data Integrity Awareness**.

Q9. What is a "Foreign Key" in SQL, and why is "Referential Integrity" important?

Strong Answer: A Foreign Key is a column or group of columns in one table that links to the Primary Key of another table. It enforces Referential Integrity by ensuring that the relationship between tables remains valid—meaning you cannot insert a row with a non-existent foreign key, nor can you delete a record if other tables still depend on its primary key.

Weak Answer: A foreign key is a key imported from an external database or API to help link accounts together without duplicating data fields.

Recruiter Cue: Tests **Relational Database Design**.

Q10. What are the main HTTP methods used in RESTful APIs, and what are their specific purposes?

Strong Answer: The primary methods are: **GET** for retrieving data; **POST** for creating new resources; **PUT** for replacing an existing resource completely; **PATCH** for making partial updates to a resource; and **DELETE** for removing data. GET and PUT/PATCH/DELETE should be idempotent, meaning multiple identical requests yield the same system state.

Weak Answer: They are the commands you type into your backend code to connect the front-end buttons to the database tables so information can move back and forth.

Recruiter Cue: Tests **API Design Conventions**.

Q11. What does it mean when an API returns a 404 status code versus a 500 status code?

Strong Answer: A **404 Not Found** is a **Client-Side Error**, indicating the server could connect but could not find the specific resource requested by the URL. A **500 Internal Server Error** is a **Server-Side Error**, indicating that the client's request was received, but the backend server encountered an unhandled exception or crash while processing it.

Weak Answer: 404 means the internet connection is broken or the website is down, and 500 means the database password was typed incorrectly in the config file.

Recruiter Cue: Tests **Debugging & Protocol Literacy**.

Q12. What is the purpose of an Index in a database, and what is the trade-off of creating too many indexes?

Strong Answer: An index speeds up data retrieval (**SELECT** queries) by creating a pointer structure (often a B-Tree) so the database doesn't have to perform a full-table scan. The trade-off is that every index slows down write operations (**INSERT**, **UPDATE**, **DELETE**) because the database must update the index structure every time the data changes, and it consumes extra disk space.

Weak Answer: An index numbers the rows in your table from 1 to N so that the backend knows exactly how many records exist in the system.

Recruiter Cue: Tests **Performance Trade-Off Evaluation**.

Q13. Explain the concept of Object-Relational Mapping (ORM). What is an advantage and a disadvantage of using it?

Strong Answer: An ORM is a tool that allows developers to interact with a database using their object-oriented programming language instead of writing raw SQL. An advantage is **faster development speed and code abstraction/readability**. A disadvantage is that it can generate **inefficient raw SQL queries under the hood** (like the N+1 query problem), which degrades performance compared to handwritten, optimized queries.

Weak Answer: An ORM converts a SQL database into a NoSQL database dynamically so that frontend languages can read the data easily without translation libraries.

Recruiter Cue: Tests **Tooling Abstraction Awareness**.

Q14. What is the difference between Authentication and Authorization?

Strong Answer: Authentication (AuthN) is verifying *who* a user is (e.g., via passwords, biometrics, or token validation). **Authorization (AuthZ)** is verifying *what* that authenticated user is allowed to do (e.g., determining if a user has "Admin" rights to delete a resource versus "Viewer" rights to only read it).

Weak Answer: They are the same thing, but authentication happens on the front end when logging in, and authorization happens on the back end inside the database server.

Recruiter Cue: Tests **Basic Security Literacy**.

Q15. Why is it a bad practice to store user passwords in plain text in a database, and how should they be stored instead?

Strong Answer: Storing passwords in plain text means a database breach exposes every user's credential instantly. Instead, passwords must be **hashed using a strong cryptographic algorithm (like bcrypt, Argon2, or PBKDF2) along with a unique "Salt"** (random data appended to the password before hashing). The salt prevents attackers from using precomputed tables (Rainbow Tables) to reverse-engineer cracked passwords.

Weak Answer: It's bad because anyone who logs into the database can see the passwords. They should be encrypted using standard base64 encoding so that they look like random letters.

Recruiter Cue: Tests **Data Security Fundamentals**.

PART 2: BACKEND INTERVIEW QUESTIONS FOR INTERMEDIATES

(Q16–Q25)

Section 1: Concurrency, Architecture & Performance

Q16. What is the difference between a Process and a Thread, and what is a "race condition"?

Strong Answer: A process is an independent executing program with its own isolated memory space allocated by the OS. A thread is a lightweight subunit of execution *inside* a process that shares that process's memory space with other sister threads. A race condition occurs when multiple threads concurrently read and write to a shared memory location without synchronization, causing the final system state to depend unpredictably on the exact execution timing or scheduling of the threads.

Weak Answer: A process is a backend program running on a server, and a thread is a line of code executed sequentially inside it. Race conditions happen when the server hardware runs out of CPU cores to execute tasks fast enough.

Recruiter Cue: Tests **Concurrency & OS-Level Fundamentals**.

Q17. Explain the "N+1 Query Problem" in the context of ORMs. How do you diagnose and fix it?

Strong Answer: It occurs when an application executes one initial query to fetch a list of parent records, and then executes $N+1$ subsequent individual queries to fetch related child records for each

parent (e.g., fetching 100 users, then running 100 queries to get each user's profile). It is diagnosed via database logs or profiling tools showing repetitive single lookups. It is fixed by using **Eager Loading** (e.g., using `JOIN FETCH`, `include`, or `select_related` depending on the framework) to fetch parent and child data in a single SQL operation.

Weak Answer: It means your query takes $N+1$ seconds to run because you forgot to create a primary key on the database table, which forces the database to search everything twice.

Recruiter Cue: Tests **Database Interaction Optimization**.

Q18. What is the difference between Synchronous (Blocking) and Asynchronous (Non-Blocking) I/O operations?

Strong Answer: In synchronous I/O, the executing thread is blocked/paused while waiting for an external resource (like a database or a file system read) to complete, wasting CPU cycles. In asynchronous I/O, the thread initiates the I/O operation and immediately returns to handle other computational tasks. Once the I/O operation finishes, the thread is notified via a callback, promise, or event loop to process the result, allowing the system to handle significantly higher concurrent connections.

Weak Answer: Synchronous code runs fast because everything happens at the exact same time on the server, while asynchronous code runs slowly in the background whenever the CPU has free cycles.

Recruiter Cue: Tests **I/O Management & Scaling Foundations**.

Q19. When and why would you introduce an asynchronous message broker (like RabbitMQ or Apache Kafka) into a backend architecture?

Strong Answer: I would introduce a message broker to **decouple microservices** and handle intensive background tasks asynchronously. For example, when a user registers, instead of blocking the HTTP response while sending a welcome email, generating a PDF invoice, and updating analytics, the backend pushes an event message to the broker and returns a 201 success status immediately. Worker services consume these messages independently, improving system resilience, response time, and spike tolerance.

Weak Answer: You use a message broker when your main SQL database is too slow to handle write queries, so you store data in Kafka temporarily until the server recovers.

Recruiter Cue: Tests **Distributed Systems & Event-Driven Architecture**.

Q20. What is a "Connection Pool" (e.g., for a database), and why is it used instead of creating a new connection for every request?

Strong Answer: A connection pool is a cache of pre-established, persistent database connections maintained by the backend. Creating a new TCP and TLS database connection for every incoming

HTTP request introduces massive network overhead and CPU utilization. A connection pool allows threads to instantly borrow an active connection, execute a query, and return it back to the pool, dramatically reducing latency and protecting the database from running out of available file descriptors.

Weak Answer: It's a tool that merges all separate database tables into a single virtual pool so that you don't have to write complex SQL join statements in your code.

Recruiter Cue: Tests **Resource Management Optimization**.

Section 2: Data Engineering & System Design

Q21. Explain database "Sharding" versus "Horizontal Partitioning." When would you resort to sharding?

Strong Answer: Horizontal Partitioning splits rows of a single table across multiple separate tables *within the same database instance* based on a key (like partitioning by year). Sharding goes a step further by breaking up a massive dataset and distributing the resulting partitions across **entirely separate database servers/hardware instances**. I would resort to sharding when a single database server hits physical hardware bottlenecks (CPU, disk I/O, storage limits) and vertical scaling is no longer financially or technically viable.

Weak Answer: Partitioning is for SQL databases, and sharding is for NoSQL databases when you want to duplicate your data across multiple countries to prevent downtime.

Recruiter Cue: Tests **High-Scale Data Strategy**.

Q22. What are the trade-offs between using JWT (JSON Web Tokens) versus Server-Side Sessions for user authentication?

Strong Answer: JWTs are stateless; user data is encoded in a cryptographically signed token stored on the client, meaning the backend doesn't need to query a database to validate it on every request—making it highly scalable for microservices. However, JWTs are **difficult to revoke** before expiration without introducing a centralized blacklist database. Server-side sessions are stateful and easy to revoke instantly, but they require a centralized session store (like Redis) which adds operational complexity and scaling overhead.

Weak Answer: JWT is a secure framework used for modern APIs because it cannot be hacked, whereas server sessions are an older method that only works for desktop web applications.

Recruiter Cue: Tests **Architectural Security Trade-offs**.

Q23. Describe a scenario where a database deadlock occurs. How does a database engine resolve it, and how can you prevent it in your code?

Strong Answer: A deadlock occurs when Transaction A locks Row 1 and waits for Row 2, while Transaction B simultaneously locks Row 2 and waits for Row 1; both block each other indefinitely. Database engines resolve this via **Deadlock Detection algorithms** that actively kill one of the transactions (rolling it back) to free the locks. To prevent deadlocks in code, ensure all concurrent transactions always access and lock resources in the exact same sequential order, keep transactions as brief as possible, and implement optimistic locking.

Weak Answer: A deadlock happens when two users try to register with the exact same email address at the same second, causing the database to crash. It is prevented by increasing the server's RAM.

Recruiter Cue: Tests **Concurrency Control & Transaction Safety**.

Q24. What is a "Cache Avalanche," and how do you protect a backend system against it?

Strong Answer: A cache avalanche occurs when a large volume of cached items expire at the exact same time, or if the caching server crashes under heavy load. This forces all subsequent incoming requests to hit the primary relational database simultaneously, leading to database failure and system downtime. It is prevented by **adding random jitter/offsets to cache expiration times** so keys expire gradually, using a circuit breaker pattern, and setting up a high-availability cluster for the cache layer.

Weak Answer: It happens when a user spams the refresh button on a website, causing the cache storage to overflow and delete random files from the server.

Recruiter Cue: Tests **Distributed Caching Strategies**.

Q25. What is Idempotency in API development, and why is it critical when designing a payment processing endpoint?

Strong Answer: An API endpoint is idempotent if making multiple identical requests yields the exact same system state without unintended side effects. In a payment system, if a network timeout occurs after a request leaves the client, the client may retry the request. If the endpoint is not idempotent, the user will be charged twice. It is built by requiring a unique **Idempotency Key** in the request headers; the backend checks a distributed cache (like Redis) for that key before processing, ensuring it returns the original transaction status if a duplicate request arrives.

Weak Answer: Idempotency means the API automatically encrypts the credit card details before sending them over the network so that they can't be stolen by an external attacker.

Recruiter Cue: Tests **API Fault-Tolerance & Business Logic Rigor**.

PART 3: BACKEND INTERVIEW QUESTIONS FOR FRESHERS (Q26–Q40)

Section 1: Distributed Systems, Scaling & Consistency

Q26. Explain the CAP Theorem. How do you design a system to handle a network partition, and what are the trade-offs between a CP and an AP architecture?

Strong Answer: The CAP Theorem states that a distributed system can guarantee at most two out of three properties: Consistency, Availability, and Partition Tolerance. Since network partitions (\$P\$) are inevitable in distributed systems, you must choose between Consistency (\$C\$) or Availability (\$A\$). A **CP system** (e.g., Etcd, Redis Cluster) prioritizes atomic consistency; if a partition occurs, it rejects writes on minority nodes to prevent stale data split-brain scenarios, sacrificing availability. An **AP system** (e.g., Cassandra, DynamoDB) prioritizes availability; nodes accept writes locally during a partition, sacrificing immediate consistency and relying on **Eventual Consistency** mechanisms like vector clocks or Conflict-Free Replicated Data Types (CRDTs) to resolve conflicts later.

Weak Answer: CAP stands for Data Consistency, App Availability, and Packet Protection. In a partition, you just duplicate your servers across multiple cloud regions so that if one region goes offline, the other one automatically takes over the traffic.

Recruiter Cue: Tests **Distributed Systems Foundations**.

Q27. How does the Raft or Paxos Consensus Algorithm work at a high level, and why is consensus critical in a distributed database?

Strong Answer: Consensus algorithms ensure that a cluster of independent machines can agree on a shared state or sequence of events, even if some nodes fail. In **Raft**, the cluster elects a single Leader node through randomized election timers. All writes go through the Leader, which logs the change and replicates it to Follower nodes. The Leader only commits the write once a **Quorum** (majority of nodes) acknowledges receipt. Consensus is critical for distributed databases to prevent split-brain conditions, manage cluster membership shifts, and maintain a linearizable write log without single points of failure.

Weak Answer: It's an algorithm that runs a voting poll among all background microservices every few seconds to see which server has the lowest CPU usage, then makes that server the primary database.

Recruiter Cue: Tests **Distributed Coordination Systems**.

Q28. What is the Dual-Write Problem in microservices, and how do the Outbox Pattern or Change Data Capture (CDC) solve it?

Strong Answer: The Dual-Write Problem occurs when a microservice needs to update its local database *and* notify an external system (like publishing an event to Kafka) in a single business transaction. If the database write succeeds but the network call to Kafka fails, the system enters an inconsistent state. The **Outbox Pattern** solves this by writing both the application state change and an event record into a dedicated **Outbox table** *within the same atomic database transaction*. A separate transaction log miner or **CDC tool** (like Debezium) asynchronously polls that table or reads the database WAL (Write-Ahead Log) to reliably stream the events to Kafka without missing data or blocking the main runtime.

Weak Answer: It happens when two users press the save button on a form at the exact same millisecond, causing duplicate records in the database. It's solved by putting a unique constraint on the API endpoint layer.

Recruiter Cue: Tests **Distributed Transaction & Event Integrity**.

Q29. What is the difference between Optimistic and Pessimistic Locking, and how do you evaluate which one to use for a high-concurrency e-commerce inventory system?

Strong Answer: Pessimistic Locking explicitly locks rows at the database level (e.g., `SELECT ... FOR UPDATE`) upon reading, blocking all other transactions until the lock is released. It prevents conflicts entirely but severely limits throughput. **Optimistic Locking** assumes conflicts are rare; it reads data with a version number/timestamp and checks if that version changed before committing the write. If it did, the transaction rolls back and retries. For high-concurrency inventory (like a flash sale), pessimistic locking protects the data but can exhaust database connection pools due to thread blocking. I would implement optimistic locking or, better yet, **atomic increment queries** at the DB layer (`UPDATE inventory SET stock = stock - 1 WHERE id = X AND stock > 0`) combined with an asynchronous redis reservation queue to maximize throughput safely.

Weak Answer: Pessimistic locking is used for legacy relational databases to stop hackers, while optimistic locking is used in modern Node.js applications because async code handles multi-threading automatically.

Recruiter Cue: Tests **High-Concurrency Mitigation Strategy**.

Q30. Explain the "Thundering Herd Problem" (Cache Stampede). How do you architect a system to eliminate it when a high-traffic cache key expires?

Strong Answer: A Cache Stampede occurs when a highly popular cache key (e.g., homepage layout config) expires under massive concurrent load. Thousands of worker threads suddenly read a cache miss simultaneously and hit the primary database to recompute the data, instantly crashing the DB. To eliminate this, I use **Mutex/Locks on Cache Misses**: only the first thread that experiences a miss acquires a distributed lock (via Redis Redlock) to query the database, while other threads wait or serve slightly stale data. Alternatively, I implement **Probabilistic Early Expiration (XFetch)** or set up a background cron worker to proactively refresh the cache key *before* its hard TTL expires.

Weak Answer: It happens when a hacker runs a DDoS attack on your caching server. You stop it by setting up an auto-scaling group on your cloud provider to add more Redis instances automatically.

Recruiter Cue: Tests **Resilience Architecture & Traffic Control**.

Q31. How do you design an API to handle "Saga Pattern" orchestrations across multiple distributed microservices without relying on 2PC (Two-Phase Commit)?

Strong Answer: Two-Phase Commit is highly blocking and doesn't scale well in cloud microservices. Instead, the Saga Pattern manages distributed transactions using a sequence of local transactions. I choose between **Orchestration** (a centralized coordinator microservice commands participants what local transactions to run) or **Choreography** (services react to event messages published to Kafka asynchronously). Crucially, every step must have a corresponding **Compensating Transaction**—if step 3 (Payment) fails, the Saga emits reverse events to roll back step 2 (Inventory Allocation) and step 1 (Order Creation), ensuring eventual consistency.

Weak Answer: Saga is a microservice library that automatically connects all local databases together into a virtual cluster, so that if one database crashes, the entire workflow resets itself.

Recruiter Cue: Tests **Distributed Workflow & Orchestration Design**.

Q32. What is the "Slicing Window" or "Leaky Bucket" algorithm, and how would you build a distributed API rate limiter that scales across multiple geography nodes?

Strong Answer: A **Leaky Bucket** algorithm processes requests at a constant, smooth rate through a queue, dropping spikes. A **Sliding Window Counter** tracks request timestamps within a moving window, eliminating boundary spikes found in Fixed Window rate limiters. To scale this globally, I avoid centralized lookups on every request to prevent network latency. Instead, I implement a **hybrid architecture**: edge nodes (like Cloudflare Workers) run local token-bucket limiters, and periodically synchronize local counters asynchronously with a central, distributed Redis cluster using Lua scripts to ensure atomicity without blocking regional runtime speeds.

Weak Answer: Rate limiting is done by checking the user's IP address against an array stored in global memory on the server, then returning a 403 error if they hit the endpoint too many times.

Recruiter Cue: Tests **Edge Computing & Global Scalability**.

Section 2: Database Internals, Infrastructure & Security Operations

Q33. What is the fundamental internal architectural difference between an LSM-Tree (Log-Structured Merge-Tree) and a B-Tree database engine, and how does that affect their read/write profiles?

Strong Answer: **B-Trees** (used in PostgreSQL, MySQL) maintain a balanced, page-structured hierarchy on disk. They offer fast, predictable read performance ($O(\log n)$) because data is modified in place, but writes are heavy due to page splits and random disk I/O. **LSM-Trees** (used in Cassandra, RocksDB) append all writes directly to an in-memory structural framework (MemTable) and a sequential commit log. Once full, the MemTable is flushed to disk as immutable SSTables. Background threads periodically run **Compaction** to merge files. This architecture makes LSM-Trees exceptionally fast for write-intensive systems, at the expense of slower read lookups, which require scanning multiple SSTables (mitigated via Bloom Filters).

Weak Answer: B-Trees are used for simple index sorting algorithms, while LSM-Trees are modern cloud databases that store data as a giant flat file on an SSD to bypass indexing overhead completely.

Recruiter Cue: Tests **Database Engine Storage Internals**.

Q34. Describe how "Symmetric Encryption" differs from "Asymmetric Encryption," and explain how both are leveraged in a standard TLS/SSL handshake for a secure backend connection.

Strong Answer: Symmetric encryption uses a single shared key for both encryption and decryption (highly efficient, fast). Asymmetric encryption uses a mathematically linked Public/Private key pair; anyone can encrypt with the public key, but only the holder of the private key can decrypt. During a **TLS handshake**, asymmetric encryption is used initially for **authentication and key exchange**—the client validates the server's SSL certificate and securely shares a pre-master secret. Once the identity is verified, both sides derive a temporary **Symmetric Session Key** used to encrypt all actual application data payload transmission, optimizing performance without sacrificing transport security.

Weak Answer: Symmetric encryption is for internal database records, and asymmetric is for external API tokens. TLS handshakes just use asymmetric keys because they are impossible to brute-force over public networks.

Recruiter Cue: Tests **Cryptographic Transport Protocols**.

Q35. How do you handle the "N+1 Problem" when designing a high-throughput GraphQL backend compared to a traditional REST API?

Strong Answer: In REST, N+1 occurs on database joins, which we fix with eager loading. In GraphQL, N+1 is uniquely problematic because the resolver function for a nested field executes independently for *every single parent node* returned by the top-level query. To resolve this, I implement the **DataLoader Pattern** (batching and caching). DataLoader intercepts individual resolver calls during the same execution tick, batches the requested IDs into a single array query (e.g., `SELECT * WHERE id IN (...)`), and caches the resulting promises, ensuring the backend executes exactly one query for the relation regardless of the nesting depth.

Weak Answer: GraphQL handles N+1 automatically because the client specifies exactly what data fields it wants, meaning the backend server never queries data that isn't requested in the payload query.

Recruiter Cue: Tests **Graph Execution Loop Architecture**.

Q36. What is "Database Sharding Topology Rebalancing," and how does Consistent Hashing minimize data movement when adding a new database shard node to a cluster?

Strong Answer: In standard modular hashing ($\text{hash}(key) \pmod N$), changing the number of shard servers (N) completely alters the destination shard for almost every single key in the database, requiring a catastrophic cluster-wide data re-indexing migration. **Consistent Hashing** maps

both keys and database nodes onto a virtual circular ring structure. A key is routed to the next closest node clockwise on the ring. When a new shard node is added, it only intercepts a fraction of the keys from its immediate neighbor on the ring. This minimizes data movement, requiring only $\frac{K}{N}$ keys to be migrated (where K is total keys, N is shard count), preserving cluster availability during scaling.

Weak Answer: Rebalancing means copying all database rows into a temporary CSV file, adding the new server hardware, and running an automated parsing script to split the rows equally between the new nodes.

Recruiter Cue: Tests **High-Scale System Elasticity**.

Q37. Explain the security vulnerability known as "SQL Injection (SQLi)" at a compilation level. How do Prepared Statements / Parameterized Queries completely nullify it?

Strong Answer: SQL Injection occurs when untrusted user input is directly concatenated into an SQL command string, causing the database engine to interpret the input data as executable code instructions. **Prepared Statements** nullify this by separating the query code compilation from the data insertion. The database engine pre-compiles the template query layout (e.g., `SELECT * FROM users WHERE email = ?`). When the parameters are passed later, the engine treats them strictly as literal scalar values within the execution context, completely preventing the query compiler from altering the execution tree even if the input contains SQL commands like `UNION` or `DROP TABLE`.

Weak Answer: SQLi happens when a hacker enters a long string into a text input. Prepared statements stop it by running a regex filter that strips out special characters like quotes or semicolons before saving.

Recruiter Cue: Tests **Database Engine Parser Security**.

Q38. How do you implement "Circuit Breaker Pattern" in a microservices network layer, and how does it prevent Cascading Failures across an enterprise backend?

Strong Answer: A Circuit Breaker wraps network calls to a downstream service and monitors failure rates. It operates in three states: **Closed** (normal traffic flow), **Open** (failures exceed a threshold; traffic is instantly rejected/short-circuited with a fast-fail error response to preserve resources), and **Half-Open** (a small amount of probe traffic is allowed through to check if the downstream service has recovered). This pattern prevents **Cascading Failures** because if a downstream payment gateway goes down, the upstream services stop waiting on network timeouts, freeing up thread pools and memory allocations that would otherwise choke the entire upstream infrastructure stack.

Weak Answer: A circuit breaker is a cloud infrastructure script that cuts off the server's electricity supply or power grid connection if the server temperature or CPU usage climbs dangerously high.

Recruiter Cue: Tests **Fault-Tolerant Network Topology**.

Q39. What is "Database Isolation Levels" (Read Uncommitted, Read Committed, Repeatable Read, Serializable), and what are the specific data anomalies each level prevents?

Strong Answer: Database isolation levels define the degree to which transaction integrity is visible to other concurrent transactions.

- **Read Uncommitted** prevents nothing; allows **Dirty Reads** (reading uncommitted changes that might roll back).
- **Read Committed** prevents Dirty Reads, but allows **Non-Repeatable Reads** (rereading the same row yields different data because another transaction committed an update).
- **Repeatable Read** prevents Dirty and Non-Repeatable reads, but can allow **Phantom Reads** (rereading a range yields new rows added by another concurrent transaction).
- **Serializable** provides absolute isolation by locking ranges or using MVCC (Multi-Version Concurrency Control), preventing all anomalies at the cost of high transaction abortion rates and lowest performance.

Weak Answer: Isolation levels dictate who has permission to view administrative tables inside the database dashboard based on their security clearance profile.

Recruiter Cue: Tests **Transaction Theory & Concurrency Mechanics**.

Q40. You need to upgrade a production database schema for a multi-million user system with zero downtime. Describe your deployment sequence strategy.

Strong Answer: I use the **Expand and Contract Pattern** split across multiple backward-compatible phases to avoid schema lock-ups. For example, if changing a column name:

1. **Expand Phase:** Deploy a migration to add the *new* column alongside the old one while keeping the old database schema intact.
2. **Dual-Write Phase:** Deploy backend code that writes to *both* columns but reads exclusively from the old column. Run an asynchronous background script to backfill data from the old column to the new column for historical rows.
3. **Read Pivot Phase:** Update the backend code to read from the new column.
4. **Contract Phase:** Drop the old column from the schema and remove the dual-write code pathway once performance metrics normalize.

Weak Answer: I run the migration script at 3:00 AM on a Sunday when user activity is at its lowest, and put up a temporary maintenance screen for 5 minutes while the database alters the table columns.

Recruiter Cue: Tests **High-Availability Zero-Downtime Operations**.

Standardize and scale hiring for backend roles with this checklist. [Talk to our experts today.](#)