

The Principal Architect's Codex: Advanced System Design, Internals, and Operational Patterns for Senior and Staff Engineering Roles

1. Executive Summary: The Senior Engineering Evaluation Rubric

In the high-stakes environment of top-tier technology firms—spanning hyperscalers like Google, Meta, and Amazon, as well as high-velocity startups—the System Design Interview (SDI) serves as the primary mechanism for assessing an engineer's technical maturity, architectural intuition, and operational foresight. For candidates targeting Senior (L5/E5) and Staff (L6/E6) roles, the evaluation criteria shift dramatically from basic problem-solving to the demonstration of profound ownership and strategic depth. The expectation is no longer simply to construct a functional system but to exhibit the capability to lead complex engineering initiatives, navigate ambiguity, and proactively identify failure modes before they manifest in production.¹

The distinction between a mid-level engineer and a senior practitioner lies in the ability to "drive the discussion." While a mid-level candidate might wait for requirements to be fed to them, a senior candidate acts as a peer to the interviewer, challenging assumptions, clarifying vague constraints, and imposing a structured architectural vision upon the problem space. This proactive engagement is termed "driving the design," and it is the single most significant signal for senior-level competency. The candidate must demonstrate not only how to build a system but also how the system will evolve, fail, and be maintained over a multi-year horizon.³

This report provides an exhaustive technical deep dive into the topics required to achieve "over-preparedness" for these roles. It transcends the superficial application of design patterns, delving into the internal mechanics of database storage engines, the physics of networking protocols, the theoretical underpinnings of distributed consistency, and the practical realities of modern operational resilience. By synthesizing insights from production systems at companies like Discord, Stripe, and Netflix, this document establishes a knowledge base that enables candidates to justify every architectural decision with engineering rigor.⁵

1.1 The L5 vs. L6 Distinction

The expectations for Senior (L5) and Staff (L6) roles diverge significantly in scope and impact. An L5 engineer is expected to independently own major components, breaking down

high-impact problems into clear architectures while balancing simplicity with the complexity necessitated by scale. They must reason deeply about data consistency, sharding, and fault isolation, making decisions guided by Service Level Indicators (SLIs) and Service Level Objectives (SLOs) rather than intuition.¹

Conversely, an L6 Staff engineer operates at a broader organizational level. The interview assesses their ability to make technical decisions that shape the entire organization, anticipating "unknown unknowns" and managing the long-term ripple effects of architectural choices. A Staff candidate must demonstrate how they handle projects that do not go as planned, applying lessons learned to future designs, and influencing the roadmap through technical leadership. Red flags at this level include staying too shallow, playing it safe with competent but generic answers, and failing to demonstrate the independence required to drive the conversation without significant guidance.⁷

2. Structural Frameworks and Estimation Mechanics

To navigate the complexity of a 45-to-60-minute interview, candidates must employ a rigorous structural framework. The PEDALS framework—Process, Estimation, Design, Architecture, List, Scale—is a standard methodology that ensures comprehensive coverage, but senior candidates must execute it with superior efficiency and adaptability.⁹

2.1 The Art of Requirement Clarification

The initial phase of the interview, often termed the "Contract Negotiation," sets the trajectory for the entire session. Senior candidates must aggressively clarify functional and non-functional requirements. This involves asking precise questions regarding the scale of the system, such as Daily Active Users (DAU), read-to-write ratios, and specific consistency requirements. For instance, distinguishing between a system that requires strict serializability (e.g., a payment gateway) versus one that can tolerate eventual consistency (e.g., a social media feed) is critical. The candidate must also define the system's boundary conditions, asking questions about the expected behavior during peak loads or network partitions.²

2.2 Rigorous Back-of-the-Envelope Estimation

Estimation is not merely a mathematical exercise; it is a tool for justifying architectural decisions. A Senior engineer uses estimation to determine whether sharding is necessary, whether data can fit in memory, and what the bandwidth requirements are. For example, calculating the storage requirements for a system with 10 million DAU generating 1KB of data daily involves determining the total throughput and retention policy. If the calculation yields 500 QPS, the candidate should recognize that sharding is likely unnecessary for throughput, though it might be required for storage capacity or fault tolerance. Conversely, a calculation resulting in 100TB of storage necessitates a distributed storage solution.¹¹

Table 1: System Design Estimation Heuristics

Metric	Threshold	Architectural Implication
Throughput	< 1k QPS	Single instance DB likely sufficient (with replicas).
Throughput	> 10k QPS	Sharding or heavy caching layer required.
Latency	< 100ms	Requires in-memory caching (Redis/Memcached).
Storage	> 10TB	Sharding or distributed file system (S3/HDFS) required.
Availability	99.999%	Multi-region active-active deployment needed.

3. Advanced Networking and Traffic Management

At the senior level, the abstraction of "adding a load balancer" is insufficient. Candidates must demonstrate a nuanced understanding of the OSI model, packet flow, and the trade-offs involved in protocol selection.¹³

3.1 Layer 4 vs. Layer 7 Load Balancing Internals

The choice between Layer 4 (Transport) and Layer 7 (Application) load balancing represents a fundamental trade-off between performance and control. Layer 4 load balancers operate at the packet level, forwarding traffic based on IP addresses and ports without inspecting the payload. This approach, often utilizing Network Address Translation (NAT) or Direct Server Return (DSR), maintains a single TCP connection between the client and the backend server. It offers extremely high throughput and low latency because the load balancer does not need to buffer the request or perform complex parsing.¹⁴

In contrast, Layer 7 load balancers operate at the application layer, terminating the TCP connection from the client and establishing a new connection to the backend server. This allows the load balancer to inspect HTTP headers, cookies, and URLs, enabling intelligent routing decisions such as directing traffic to specific microservices based on the URL path. However, this intelligence comes at the cost of latency and resource consumption, as the load

balancer must buffer the request and manage two separate TCP connections. Senior engineers must articulate these trade-offs, recommending L4 for high-throughput ingress points and L7 for complex microservices routing.¹⁴

3.2 Global Traffic Management: Anycast vs. Unicast

Routing users to the nearest data center is a critical component of global system design. Unicast routing relies on DNS-based load balancing, where a DNS server returns the IP address of the data center closest to the user. While simple, this approach suffers from DNS caching issues; if a data center fails, users with cached IP addresses may experience downtime until the Time-To-Live (TTL) expires. Additionally, DNS does not respect real-time network congestion or BGP routing changes.¹⁷

Anycast routing addresses these limitations by announcing the same IP address from multiple geographic locations via the Border Gateway Protocol (BGP). The internet's routing infrastructure automatically directs user traffic to the topologically nearest node. This provides immediate failover and lower latency, as traffic is routed to the optimal location without relying on client-side DNS behavior. However, Anycast can introduce "routing instability," where traffic fluctuates between data centers due to BGP path changes, potentially disrupting stateful connections. Senior engineers must discuss mitigation strategies, such as using consistent hashing at the edge to route users to the correct backend despite ingress shifts.¹⁹

3.3 The Evolution of Protocols: HTTP/3 and QUIC

The limitations of TCP and HTTP/2 have driven the adoption of newer protocols like HTTP/3 and QUIC. HTTP/2 introduced multiplexing to allow multiple streams over a single TCP connection, but it remains vulnerable to TCP Head-of-Line (HOL) blocking. If a single packet is lost, the operating system's TCP stack halts the delivery of all streams on that connection until the lost packet is retransmitted, negating the benefits of multiplexing on unreliable networks.²¹

QUIC, built on top of UDP, solves this by moving reliability and congestion control into user space. It manages streams independently, ensuring that packet loss in one stream does not block others. Furthermore, QUIC introduces the concept of Connection Migration, allowing a connection to survive changes in the client's IP address (e.g., switching from Wi-Fi to cellular) by using a persistent Connection ID. This capability is particularly valuable for mobile applications requiring seamless connectivity.²¹

3.4 TLS 1.3 and the 0-RTT Trade-off

TLS 1.3 significantly improves latency by reducing the handshake from two round-trips (2-RTT) to one (1-RTT). It also introduces the 0-RTT feature, which allows clients to send encrypted application data in the very first packet if they have a pre-shared key from a

previous session. While this drastically reduces latency for returning visitors, it introduces a security risk known as the Replay Attack. Because the 0-RTT data is not protected against replay, an attacker who intercepts the packet can resend it to the server. Senior engineers must identify this risk and mandate that 0-RTT be restricted to idempotent operations (like GET requests) to prevent duplicate transactions.²⁴

4. Database Internals and Storage Engine Mechanics

The selection of a database is often the most critical decision in a system design interview. Senior candidates must move beyond high-level comparisons and discuss the internal data structures that dictate performance characteristics.²⁶

4.1 B-Trees vs. LSM Trees: The Read/Write Dichotomy

The fundamental difference between relational databases (like PostgreSQL) and many NoSQL systems (like Cassandra) lies in their underlying storage engines.

- **B-Trees:** Used by engines like InnoDB, B-Trees organize data in sorted, fixed-size pages. This structure is optimized for reads, as lookups have a consistent $O(\log N)$ time complexity. However, writes (especially random updates) can be expensive, as they require locating and updating specific pages on disk, leading to random I/O and potential page fragmentation. This makes B-Trees ideal for read-heavy workloads requiring complex queries.²⁷
- **LSM Trees (Log-Structured Merge Trees):** Used by engines like RocksDB and LevelDB, LSM Trees are optimized for write-heavy workloads. Incoming writes are appended to an in-memory buffer (MemTable). When the buffer fills, it is flushed to disk as an immutable Sorted String Table (SSTable). This converts random writes into sequential writes, maximizing disk throughput. However, reads become more expensive, as the system must check the MemTable and multiple SSTables. To mitigate this, LSM engines use Bloom Filters to quickly rule out files that do not contain the requested key.²⁷

4.2 LSM Compaction Strategies and Write Stalls

A critical operational challenge in LSM-based systems is managing the accumulation of SSTables. Compaction processes run in the background to merge these tables and discard obsolete data (tombstones). If the write rate exceeds the compaction rate, the system may trigger a "Write Stall," effectively halting all writes to prevent disk saturation.

- **Leveled Compaction:** Used by RocksDB, this strategy aggressively merges files into levels. It minimizes space amplification and optimizes read performance but incurs high write amplification due to frequent re-writing of data. It is suitable for read-heavy workloads running on LSM engines.²⁹
- **Tiered Compaction:** Used by Cassandra, this strategy flushes files and leaves them until a threshold is reached. It minimizes write amplification, making it ideal for write-intensive

workloads, but suffers from high read amplification as queries must scan many files.²⁹

4.3 PostgreSQL MVCC and the Vacuum Problem

PostgreSQL uses Multi-Version Concurrency Control (MVCC) to handle concurrent transactions. Updates create new versions of rows (tuples) rather than overwriting them. This leaves behind "dead tuples" that must be cleaned up by the autovacuum process. If autovacuum cannot keep up—often due to long-running transactions holding open snapshots—the table bloats, slowing down queries. More critically, if the Transaction ID (XID) counter (a 32-bit integer) wraps around, the database essentially shuts down to prevent data corruption. To avoid this, PostgreSQL performs a "Vacuum Freeze" to mark old transaction IDs as frozen. Senior engineers must discuss tuning autovacuum_freeze_max_age and designing applications to avoid long transactions that block this critical maintenance process.³²

5. Distributed Systems Theory and Consistency Models

The ability to reason about consistency is a hallmark of a senior engineer. Candidates must distinguish between different types of consistency and understand the theoretical bounds of distributed systems.

5.1 Linearizability vs. Serializability vs. Strict Serializability

These terms are frequently conflated but represent distinct guarantees.

- **Linearizability:** A guarantee about single operations on single objects. It ensures that if operation A completes before operation B starts, then B must see the value written by A. It implies real-time recency.³⁴
- **Serializability:** A guarantee about transactions (multiple operations). It ensures that the execution of concurrent transactions produces the same result as some serial execution of those transactions. It provides isolation but does not guarantee real-time ordering.³⁴
- **Strict Serializability:** The combination of both. It ensures that transactions execute essentially in a serial order that corresponds to real-time. This is the strongest guarantee, provided by systems like Google Spanner.³⁵

5.2 Google Spanner and TrueTime

Spanner achieves strict serializability on a global scale through the use of TrueTime, a time API backed by atomic clocks and GPS receivers. TrueTime exposes time not as a single point, but as an interval [earliest, latest] with a known error bound. When a transaction commits, Spanner enforces a "Commit Wait" period, ensuring that the transaction does not finish until the uncertainty interval has passed. This guarantees that if one transaction finishes before another starts, the timestamp of the first is strictly less than the second, preserving external

consistency without requiring communication between distant data centers.³⁸

5.3 Distributed Consensus: Paxos, Raft, and Zab

Replicated state machines rely on consensus algorithms to agree on values.

- **Paxos:** The theoretical standard for consensus. It is leaderless (or uses an implicit leader) and is known for its complexity and difficulty of implementation.
- **Raft:** Designed for understandability, Raft uses a strong leader model where log entries flow only from the leader to followers. It explicitly handles leader election and log replication, making it easier to verify and implement.⁴⁰
- **Zab (ZooKeeper Atomic Broadcast):** Designed specifically for primary-backup systems like ZooKeeper. It guarantees a total ordering of updates (FIFO), which is critical for sequencing operations in distributed coordination services.⁴²

5.4 Conflict Resolution: CRDTs vs. Operational Transformation

For collaborative applications, dealing with concurrent edits is a primary challenge.

- **Operational Transformation (OT):** Used by Google Docs, OT requires a central server to order and transform operations (e.g., adjusting the index of an insert based on a concurrent delete). It provides strong consistency but introduces a single point of failure and complexity.⁴⁴
- **CRDTs (Conflict-free Replicated Data Types):** Used by Figma and Redis, CRDTs allow independent updates on different nodes that can be mathematically merged to a consistent state. They enable peer-to-peer synchronization and offline editing but often incur higher memory overhead due to the need to store tombstones for deleted items.⁴⁴

6. Distributed Transactions and Data Integrity

In microservices architectures, the breakdown of the monolithic database necessitates new patterns for managing transactions across service boundaries.

6.1 Two-Phase Commit (2PC) vs. Sagas

- **Two-Phase Commit (2PC):** A protocol where a coordinator asks all participants to "Prepare" and then "Commit." While it guarantees strong consistency, it is a blocking protocol; if the coordinator fails, participants may hold locks indefinitely, reducing system availability.⁴⁷
- **Sagas:** A sequence of local transactions where each step updates data within a single service. If a step fails, the system executes "compensating transactions" to undo the changes made by previous steps. Sagas favor availability over isolation, as the system may be in a temporarily inconsistent state during the execution of the sequence. They can be implemented via choreography (events) or orchestration (a central coordinator).⁴⁸

6.2 Distributed ID Generation

Generating unique identifiers in a distributed system is a classic design problem.

- **UUID:** 128-bit random identifiers. While they require no coordination, their randomness causes fragmentation in B-Tree indexes, leading to poor write performance.⁵⁰
- **Snowflake:** Developed by Twitter, this 64-bit ID combines a timestamp, a machine ID, and a sequence number. It is roughly time-ordered (k-ordered), making it efficient for DB indexing. However, it requires coordination (like ZooKeeper) to manage machine IDs.⁵⁰
- **ULID / KUUID:** These formats combine the sortability of Snowflake with the decentralized nature of UUIDs. They start with a timestamp component for locality and end with random bits for uniqueness, offering a "best of both worlds" solution for modern distributed systems.⁵⁰

7. High-Throughput Messaging and Streaming Internals

Understanding the internal architecture of messaging systems allows engineers to choose the right tool for high-scale data ingestion.

7.1 Kafka Internals: Zero-Copy and Sequential I/O

Kafka's high throughput is achieved by bypassing the JVM heap and leveraging the OS kernel.

- **Sequential I/O:** Kafka relies on the fact that sequential disk writes are significantly faster than random writes. It appends messages to an immutable log file.
- **Zero-Copy:** Kafka uses the sendfile system call to transfer data directly from the disk cache to the network socket. This avoids copying data into the application's user space, reducing CPU context switches and memory overhead. This mechanism allows Kafka to saturate network bandwidth even with modest CPU resources.⁵³
- **Exactly-Once Semantics (EOS):** Kafka supports EOS through idempotent producers (using sequence numbers to de-duplicate sends) and transactional writes that span multiple partitions, ensuring that a set of messages is either fully committed or fully aborted.⁵⁵

7.2 Broker Comparison: Kafka vs. RabbitMQ vs. Pulsar

- **Kafka:** A "dumb broker, smart consumer" model. It stores a durable log, and consumers manage their own offsets. It is optimized for replayability and massive throughput.⁵⁷
- **RabbitMQ:** A "smart broker, dumb consumer" model. It offers complex routing capabilities via exchanges but deletes messages once acknowledged. It is better suited for complex task routing rather than event streaming.⁵⁷
- **Apache Pulsar:** Distinguishes itself by separating compute (brokers) from storage (BookKeeper). This allows for independent scaling; if storage is full, you can add storage

nodes without rebalancing the brokers. It also supports tiered storage, offloading old data to S3, making it ideal for infinite retention use cases.⁶⁰

8. Modern Architecture Trends: Vector Databases and GenAI

The rise of Generative AI has introduced new components into the system design interview, specifically around vector search and LLM serving.

8.1 Vector Indexing: The HNSW Algorithm

Vector databases must perform similarity searches over billions of high-dimensional vectors. Brute-force search (KNN) is too slow ($\$O(N)$). The industry standard for Approximate Nearest Neighbor (ANN) search is HNSW (Hierarchical Navigable Small World). HNSW constructs a multi-layered graph where the top layers allow for long "skips" across the data space, and lower layers provide fine-grained navigation. This structure allows for $\$O(\log N)$ search complexity, though it consumes significant memory to store the graph connections.⁶²

8.2 LLM Serving Optimizations

Serving Large Language Models (LLMs) presents unique challenges due to memory bandwidth constraints.

- **PagedAttention (vLLM):** Traditional serving systems suffer from memory fragmentation because the size of the Key-Value (KV) cache is unknown ahead of time. PagedAttention, inspired by OS virtual memory, breaks the KV cache into non-contiguous blocks. This reduces memory waste from ~60% to <4%, allowing for larger batch sizes and higher throughput.⁶⁴
- **Continuous Batching:** Instead of waiting for a batch of requests to finish (where the speed is limited by the longest output), continuous batching inserts new requests into the batch as soon as previous ones complete. This "in-flight" batching maximizes GPU utilization.⁶⁶
- **Speculative Decoding:** To reduce latency, a smaller "draft model" predicts the next few tokens, which are then verified in parallel by the larger target model. This leverages the memory bandwidth more effectively, generating multiple tokens per forward pass of the large model.⁶⁸

9. Resilience, Security, and Operational Excellence

Designing a system is only the first step; ensuring it survives in a hostile production environment is where senior engineers demonstrate their value.

9.1 Observability and Sampling Strategies

- **Head-Based Sampling:** The decision to keep a trace is made at the start of the request. While simple and low-overhead, it may miss interesting errors that occur in requests that were not sampled.⁷⁰
- **Tail-Based Sampling:** The system buffers all spans for a request and decides whether to keep them after the request completes. This allows for capturing 100% of errors and high-latency traces but requires significant memory and infrastructure to buffer the telemetry data.⁷²

9.2 Resilience Patterns: Bulkheads and Shuffle Sharding

- **Bulkhead Pattern:** Inspired by ship design, this pattern isolates resources (like thread pools) for different services. If one service fails, it does not exhaust the resources of the entire system, preventing cascading failures.⁷⁴
- **Shuffle Sharding:** Used by AWS, this technique assigns each tenant to a unique subset of resources (e.g., a specific combination of queues or servers). This dramatically reduces the "blast radius" of a noisy neighbor. If one tenant overwhelms their assigned resources, only the small subset of other tenants sharing those specific resources is affected, rather than the entire platform.⁷⁶

9.3 Rate Limiting and Load Shedding

Stripe's rate limiting architecture serves as a prime example of operational maturity. It employs a multi-tiered approach:

1. **Hard Rate Limits:** Applied to non-critical traffic to protect infrastructure.
2. **Load Shedding:** During severe congestion, the system drops low-priority traffic (like analytics) to preserve capacity for critical operations (like payments).
3. **Federation:** Clients sync token counts periodically rather than for every request, trading strict accuracy for reduced latency on the Redis backend.⁷⁸

9.4 Security: Zero Trust and OAuth flow

- **Zero Trust Architecture:** Assumes no implicit trust within the network. Every request must be authenticated and authorized, often using mTLS and short-lived certificates (SPIFFE/SPIRE). It relies on strong identity rather than network perimeters.⁸⁰
- **OAuth2 / OIDC:** In microservices, the API Gateway often acts as an OAuth2 Resource Server, validating JWTs issued by an Identity Provider (IdP). This centralizes authentication logic while allowing services to make authorization decisions based on token scopes.⁸²

10. Architectural Case Studies

Applying theoretical knowledge to real-world architectures solidifies a candidate's expertise.

10.1 WhatsApp: The Erlang Concurrency Model

WhatsApp leverages Erlang's lightweight processes to handle millions of concurrent connections on a single server. Each user connection is an isolated process, ensuring that a crash in one does not affect others. For push notifications, WhatsApp maintains a persistent connection to the push notification servers (APNS/FCM). To ensure message ordering and delivery, the push notification often acts merely as a "tickle" to wake the app, which then connects to the server to fetch the actual encrypted message payload.⁸⁴

10.2 Discord: Migration to ScyllaDB

Discord's migration from Cassandra to ScyllaDB was driven by the need to eliminate Garbage Collection pauses and improve tail latency. ScyllaDB, written in C++, uses a shard-per-core architecture that avoids the overhead of the JVM. Discord also implemented a specialized Data Service in Rust to coalesce requests, ensuring that "hot" channels (with thousands of simultaneous readers) do not overwhelm the database partitions.⁸⁶

10.3 Netflix: Open Connect CDN

Netflix built its own CDN, Open Connect, to deliver massive video files efficiently. The architecture splits the control plane (running on AWS for logic) from the data plane (Open Connect Appliances embedded in ISPs). By predictively pushing content to these appliances during off-peak hours based on regional popularity models, Netflix ensures that the vast majority of streaming traffic is served from the user's local ISP network, bypassing the congested internet backbone.⁸⁸

10.4 Facebook Haystack: Efficient Photo Storage

The "Needle in a Haystack" paper describes how Facebook solved the metadata bottleneck of traditional file systems. Storing billions of small photos as individual files would overwhelm the file system's metadata limits (inodes). Haystack aggregates millions of photos into massive 100GB files. It maintains an in-memory index of where each photo (needle) exists within these large files, allowing the server to retrieve a photo with a single disk seek.⁹⁰

10.5 Uber Ringpop: Application-Layer Sharding

Uber developed Ringpop to handle the real-time matching of riders and drivers. It uses a gossip protocol (SWIM) to maintain a consistent view of the cluster membership. Ringpop implements consistent hashing at the application layer, allowing any node to route a request (e.g., for a specific trip) to the correct node responsible for that trip's state. This decentralized approach eliminates the need for a central load balancer for internal service-to-service traffic.⁹²

11. Conclusion: The Bar-Raiser Mindset

The journey to Senior and Staff engineering roles is defined by the transition from implementation to ownership. It requires a mindset that embraces trade-offs, prioritizing simplicity and operability over complexity. The "over-prepared" candidate does not merely recite the definitions of CAP theorem or sharding strategies; they weave these concepts into a coherent narrative that addresses the specific constraints of the business problem. They demonstrate that they can navigate the "unknown unknowns," design for failure, and build systems that are not just theoretically sound but practically survivable in the face of real-world chaos. By mastering the internals of storage engines, the physics of networking, and the rigorous patterns of distributed coordination, a candidate proves they are ready to hold the pager for the systems that power the modern world.

Works cited

1. Google L5 System Design: A Complete Guide to Master Senior-Level Interviews, accessed on December 27, 2025,
<https://www.systemdesignhandbook.com/guides/google-l5-system-design/>
2. Google L5 Software Engineer 2025 Interview Guide: Senior-Level Coding, System Design & Leadership Mastery | Onsites.fyi, accessed on December 27, 2025,
<https://www.onsites.fyi/blog/article/google-L5-software-engineer-interview-questions>
3. Are We On Track? AI-Assisted Active and Passive Goal Reflection During Meetings - arXiv, accessed on December 27, 2025,
<https://arxiv.org/html/2504.01082v2>
4. How to pass system design interviews: drive the design - Formation, accessed on December 27, 2025,
<https://formation.dev/blog/how-to-pass-system-design-interviews-drive-the-design/>
5. The Complete System Design Interview Guide (2026 Edition), accessed on December 27, 2025,
<https://www.systemdesignhandbook.com/guides/system-design-interview/>
6. Top 50 System Design Interview Questions 2025, accessed on December 27, 2025,
<https://getsdeready.com/top-50-system-design-interview-questions-2025/>
7. How to distinguish between the performance of a strong L5 (senior) vs. an L6 (staff) interview performance? - Taro, accessed on December 27, 2025,
<https://www.jointaro.com/question/DT1Q2wrQy3CdILkqprpT/how-to-distinguish-between-the-performance-of-a-strong-l5-senior-vs-an-l6-staff-interview-performance/>
8. Expectations By Level | System Design Interview | AlgoMaster.io, accessed on December 27, 2025,
<https://algomaster.io/learn/system-design-interviews/expectations-by-level>
9. System Design Interview Guide: FAANG and Startups - Exponent, accessed on December 27, 2025,
<https://www.tryexponent.com/blog/system-design-interview-guide>
10. System Design Questions: How to Answer with the PEDALS Method™ by Lin and

- Patel (Official Video) - YouTube, accessed on December 27, 2025,
<https://www.youtube.com/watch?v=CzFAXOI-GLU>
11. Cracking Google Interview for Senior Software Engineer | by Vishal Saraogi - Medium, accessed on December 27, 2025,
<https://medium.com/@vsaraogi/cracking-google-interview-for-senior-software-engineer-79ec0b214fc5>
 12. A Framework For System Design Interviews - ByteByteGo, accessed on December 27, 2025,
<https://bytebytogo.com/courses/system-design-interview/a-framework-for-system-design-interviews>
 13. Networking Essentials - System Design in a Hurry - Hello Interview, accessed on December 27, 2025,
<https://www.hellointerview.com/learn/system-design/core-concepts/networking-essentials>
 14. Layer 4 vs Layer 7 Load Balancing | Glossary - A10 Networks, accessed on December 27, 2025,
<https://www.a10networks.com/glossary/how-do-layer-4-and-layer-7-load-balancing-differ/>
 15. Layer 4 Load Balancers: A Complete Deep Dive Guide | by Pawan Choudhary - Medium, accessed on December 27, 2025,
<https://medium.com/@pawanchoudhary276437/layer-4-load-balancers-a-complete-deep-dive-guide-03fa1053f4c9>
 16. System Design Series Part 5 — L4 vs L7 Load Balancers | by D Karthik Sainadh Reddy, accessed on December 27, 2025,
<https://medium.com/@karthiksainadhreddy/system-design-series-part-5-l4-vs-l7-load-balancers-0bc35a31861c>
 17. Load Balancing without Load Balancers - The Cloudflare Blog, accessed on December 27, 2025,
<https://blog.cloudflare.com/cloudflares-architecture-eliminating-single-p/>
 18. ELI5 : What is the difference between Unicast, Multicast and Anycast in telecommunications ? And who decide which "cast" it is ? : r/explainlikeimfive - Reddit, accessed on December 27, 2025,
https://www.reddit.com/r/explainlikeimfive/comments/1f1tlam/eli5_what_is_the_difference_between_unicast/
 19. Load-Balancing versus Anycast: A First Look at Operational Challenges - arXiv, accessed on December 27, 2025, <https://arxiv.org/html/2503.14351v1>
 20. Anycast vs Unicast: What a Global Business Needs to Know, accessed on December 27, 2025,
<https://www.anycast.com/resources/anycast-vs-unicast-what-a-global-business-needs-to-know>
 21. Networking | System Design Interview | AlgoMaster.io, accessed on December 27, 2025, <https://algomaster.io/learn/system-design-interviews/networking>
 22. [Frontend System Design] - Deep dive into HTTP [Part 1] - YouTube, accessed on December 27, 2025, <https://www.youtube.com/watch?v=gxF9fLo5XQw>
 23. HTTP/3 deep dive | Ably: Serious, serverless realtime infrastructure - Medium,

- accessed on December 27, 2025,
<https://medium.com/ably-realtime/http-3-deep-dive-9318f7d6834d>
24. TLS 1.2 vs. 1.3—Handshake, Performance, and Other Improvements - Catchpoint, accessed on December 27, 2025,
<https://www.catchpoint.com/http2-vs-http3/tls1-2-vs-1-3>
25. A Survey of TLS 1.3 O-RTT Usage - ETH Zürich, accessed on December 27, 2025,
https://ethz.ch/content/dam/ethz/special-interest/infk/inst-infsec/appliedcrypto/education/theses/masters-thesis_mihael-liskij.pdf
26. Data Structures Behind Databases - System Design School, accessed on December 27, 2025,
<https://systemdesignschool.io/fundamentals/data-structures-behind-databases>
27. Write throughput differences in B-tree vs LSM-tree based databases? - Reddit, accessed on December 27, 2025,
https://www.reddit.com/r/databasedevelopment/comments/187cp1g/write_throughput_differences_in_btree_vs_lsmtree/
28. Revisiting B+-tree vs. LSM-tree - USENIX, accessed on December 27, 2025,
<https://www.usenix.org/publications/loginonline/revisit-b-tree-vs-lsm-tree-upon-arrival-modern-storage-hardware-built>
29. LSM Tree Database Compaction Strategies: When to Use Size-Tiered, Leveled, or Time-Windowed | by Shivank Rastogi | Medium, accessed on December 27, 2025,
<https://medium.com/@rastogi.shivank16/lsm-tree-database-compaction-strategies-when-to-use-size-tiered-leveled-or-time-windowed-f40b5f839e3c>
30. LSM-Tree Compaction Visualization, accessed on December 27, 2025,
<https://disc-projects.bu.edu/compactionary/background.html>
31. Constructing and Analyzing the LSM Compaction Design Space - VLDB Endowment, accessed on December 27, 2025,
<https://vldb.org/pvldb/vol14/p2216-sarkar.pdf>
32. What is Freezing? - pganalyze, accessed on December 27, 2025,
<https://pganalyze.com/docs/vacuum-advisor/what-is-freezing>
33. How Postgres freezes tuples. Working with what we have | by Hussein Nasser - Medium, accessed on December 27, 2025,
<https://medium.com/@hnasr/how-postgres-freezes-tuples-4a9931261fc>
34. Linearizability in distributed systems - Eli Bendersky's website, accessed on December 27, 2025,
<https://eli.thegreenplace.net/2024/linearizability-in-distributed-systems/>
35. Understanding Linearizability vs Serializability - System Design School, accessed on December 27, 2025,
<https://systemdesignschool.io/blog/linearizability-vs-serializability>
36. Linearizability versus Serializability | Peter Bailis, accessed on December 27, 2025,
<http://www.bailis.org/blog/linearizability-versus-serializability/>
37. Linearizability And/Vs Serializability in Distributed Databases - Ajay Gupta - Medium, accessed on December 27, 2025,
<https://ajaygupta-spark.medium.com/linearizability-and-vs-serializability-in-distributed-databases-9da2462589d>
38. Spanner: Google's Globally-Distributed Database, accessed on December 27,

- 2025, <https://research.google.com/archive/spanner-osdi2012.pdf>
- 39. Spanner: TrueTime and external consistency - Google Cloud Documentation, accessed on December 27, 2025,
<https://docs.cloud.google.com/spanner/docs/true-time-external-consistency>
 - 40. Consensus Protocols: Paxos vs Raft - System Design Interview Guide | bugfree.ai, accessed on December 27, 2025,
<https://bugfree.ai/knowledge-hub/consensus-protocols-paxos-vs-raft>
 - 41. Paxos vs. Raft vs. ZAB: A Comprehensive Dive into Distributed Consensus Protocols | by Remis Haroon | Medium, accessed on December 27, 2025,
<https://medium.com/@remisharoon/paxos-vs-raft-vs-zab-a-comprehensive-dive-into-distributed-consensus-protocols-6243a3f6539b>
 - 42. A Brief Analysis of Consensus Protocol: From Logical Clock to Raft - Alibaba Cloud, accessed on December 27, 2025,
https://www.alibabacloud.com/blog/a-brief-analysis-of-consensus-protocol-from-logical-clock-to-raft_594675
 - 43. Raft: A more understandable consensus algorithm that is equivalent to Paxos : r/compsci - Reddit, accessed on December 27, 2025,
https://www.reddit.com/r/compsci/comments/1c1rjx/raft_a_more_understandable_consensus_algorithm/
 - 44. Building Collaborative Interfaces: Operational Transforms vs. CRDTs - DEV Community, accessed on December 27, 2025,
<https://dev.to/puritanic/building-collaborative-interfaces-operational-transforms-vs-crdts-2obo>
 - 45. Google Docs Architecture: Real-Time Collaboration with OT vs. CRDTs - Sde Ray, accessed on December 27, 2025,
<https://sderay.com/google-docs-architecture-real-time-collaboration/>
 - 46. Differences between OT and CRDT - Codemedia, accessed on December 27, 2025,
https://codemedia.io/knowledge-hub/path/differences_between_ot_and_crdt
 - 47. Mastering Distributed Transactions: From 2PC to the Saga Pattern | by Aseem | Medium, accessed on December 27, 2025,
<https://medium.com/@aseem2372005/mastering-distributed-transactions-from-2pc-to-the-saga-pattern-690483d565c8>
 - 48. Difference Between Two-Phase Commit and Saga Pattern | Baeldung on Computer Science, accessed on December 27, 2025,
<https://www.baeldung.com/cs/two-phase-commit-vs-saga-pattern>
 - 49. Difference Between SAGA Pattern and 2-Phase Commit in Microservices - GeeksforGeeks, accessed on December 27, 2025,
<https://www.geeksforgeeks.org/system-design/difference-between-saga-pattern-and-2-phase-commit-in-microservices/>
 - 50. Generating Unique and Sortable IDs in Distributed Systems - Level Up Coding, accessed on December 27, 2025,
<https://levelup.gitconnected.com/generating-unique-and-sortable-ids-in-distributed-systems-e0bb7a008dab>
 - 51. What Are the Differences Between UUID, ULID, KSUID, and Snowflake IDs, and How Do I Choose? - Design Gurus, accessed on December 27, 2025,

<https://www.designgurus.io/course-play/grokking-scalable-systems-for-interviews/doc/what-are-the-differences-between-uuid-ulid-ksuid-and-snowflake-ids-and-how-do-i-choose>

52. Faster ULID generation — Ferroid - Medium, accessed on December 27, 2025,
<https://medium.com/@s0l0ist/faster-ulid-generation-ferroid-1dd65ee53e11>
53. The Zero Copy Principle With Apache Kafka | by Gautam Goswami - Medium, accessed on December 27, 2025,
<https://gautambangalore.medium.com/the-zero-copy-principle-with-apache-kafka-749dfd2ef2df>
54. Zero Copy. One Of Reason Behind Why Kafka So Fast. | by ANKIT SHEORAN - Medium, accessed on December 27, 2025,
<https://medium.com/@ankitsheoran127201/zero-copy-one-of-reason-behind-why-kafka-so-fast-7154c9d74b0a>
55. Exactly-once Semantics is Possible: Here's How Apache Kafka Does it - Confluent, accessed on December 27, 2025,
<https://www.confluent.io/blog/exactly-once-semantics-are-possible-heres-how-apache-kafka-does-it/>
56. What is Kafka Exactly Once Semantics - GitHub, accessed on December 27, 2025,
<https://github.com/AutoMQ/automq/wiki/What-is-Kafka-Exactly-Once-Semantics>
57. What's the Difference Between Kafka and RabbitMQ? - AWS, accessed on December 27, 2025,
<https://aws.amazon.com/compare/the-difference-between-rabbitmq-and-kafka/>
58. Kafka vs RabbitMQ: Choosing the Right Messaging System - Level Up Coding, accessed on December 27, 2025,
<https://levelup.gitconnected.com/kafka-vs-rabbitmq-choosing-the-right-messaging-system-10ea83e75567>
59. RabbitMQ vs Apache Kafka: Architectural and Conceptual Differences - Medium, accessed on December 27, 2025,
<https://medium.com/@steffankharmaaiarvi/rabbitmq-vs-apache-kafka-architectural-and-conceptual-differences-37f986a8d5f5>
60. The Battle of Message Queues: Apache Kafka vs Apache Pulsar | by Kunal Kumar Gupta, accessed on December 27, 2025,
<https://kunalkkg.medium.com/the-battle-of-message-queues-apache-kafka-vs-apache-pulsar-102a6100b023>
61. How is Apache Pulsar different from Apache Kafka? - Milvus, accessed on December 27, 2025,
<https://milvus.io/ai-quick-reference/how-is-apache-pulsar-different-from-apache-kafka>
62. Vector Database Basics: HNSW | Tiger Data, accessed on December 27, 2025,
<https://www.tigerdata.com/learn/vector-database-basics-hnsw>
63. Hierarchical Navigable Small Worlds (HNSW) - Pinecone, accessed on December 27, 2025, <https://www.pinecone.io/learn/series/faiss/hnsw/>
64. The Architecture Behind vLLM: How PagedAttention Improves Memory Utilization - Medium, accessed on December 27, 2025,
<https://medium.com/@mandeep0405/the-architecture-behind-vllm-how-paged>

[attention-improves-memory-utilization-2f9b25272110](#)

65. Introduction to vLLM and PagedAttention | Runpod Blog, accessed on December 27, 2025, <https://www.runpod.io/blog/introduction-to-vllm-and-pagedattention>
66. How to Optimize Batch Processing for LLMs - Ghost, accessed on December 27, 2025,
<https://latitude-blog.ghost.io/blog/how-to-optimize-batch-processing-for-langs/>
67. Static, dynamic and continuous batching | LLM Inference Handbook - Bentoml, accessed on December 27, 2025,
<https://bentoml.com/llm/inference-optimization/static-dynamic-continuous-batching>
68. An Introduction to Speculative Decoding for Reducing Latency in AI Inference, accessed on December 27, 2025,
<https://developer.nvidia.com/blog/an-introduction-to-speculative-decoding-for-reducing-latency-in-ai-inference/>
69. Speculative Decoding: A technique that makes LLMs faster without sacrificing quality, accessed on December 27, 2025,
<https://medium.com/@itssujeeth/speculative-decoding-a-technique-that-makes-llms-faster-without-sacrificing-quality-a2e712b52866>
70. accessed on December 27, 2025,
<https://uptrace.dev/opentelemetry/sampling#:~:text=A%20disadvantage%20of%20head%2Dbased,collect%20more%20data%20than%20desired.>
71. Mastering Distributed tracing: data volume challenges, and Datadog's approach to efficient sampling, accessed on December 27, 2025,
<https://www.datadoghq.com/architecture/mastering-distributed-tracing-data-volume-challenges-and-datadogs-approach-to-efficient-sampling/>
72. The Complete Guide to Sampling in Distributed Tracing | Logz.io, accessed on December 27, 2025, <https://logz.io/learn/sampling-in-distributed-tracing-guide/>
73. Sampling in OpenTelemetry: A Beginner's Guide | Better Stack Community, accessed on December 27, 2025,
<https://betterstack.com/community/guides/observability/opentelemetry-sampling/>
74. Bulkhead pattern in Microservices: Spring Boot example | by S Sivaraman - Medium, accessed on December 27, 2025,
<https://medium.com/@sivaramansankar2019/bulkhead-pattern-in-microservices-spring-boot-example-d82db2c0cc5>
75. Bulkhead Pattern - GeeksforGeeks, accessed on December 27, 2025, <https://www.geeksforgeeks.org/system-design/bulkhead-pattern/>
76. Isolate tenant workflows using shuffle sharding | Grafana Loki documentation, accessed on December 27, 2025,
<https://grafana.com/docs/loki/latest/operations/shuffle-sharding/>
77. Handling billions of invocations – best practices from AWS Lambda | AWS Compute Blog, accessed on December 27, 2025,
<https://aws.amazon.com/blogs/compute/handling-billions-of-invocations-best-practices-from-aws-lambda/>
78. Scaling your API with rate limiters - Stripe, accessed on December 27, 2025,

- <https://stripe.com/blog/rate-limiters>
79. How Stripe Does Rate Limiting - Quastor, accessed on December 27, 2025,
<https://blog.quastor.org/p/stripe-rate-limiting>
80. What Is Zero Trust? Architecture, Principles, and Technology - Tigera.io, accessed on December 27, 2025, <https://www.tigera.io/learn/guides/zero-trust/>
81. Implementing Zero Trust Architecture in Microservices: An In-Depth Guide - Medium, accessed on December 27, 2025,
<https://medium.com/@platform.engineers/implementing-zero-trust-architecture-in-microservices-an-in-depth-guide-a66417447621>
82. Microservice Architecture | a real business-world example with API Gateway and OAuth2/OIDC Login - Medium, accessed on December 27, 2025,
<https://medium.com/@a.zagarella/microservice-architecture-a-real-business-world-example-with-api-gateway-and-oauth2-oidc-login-c77c31a957fb>
83. Authentication and authorization in a microservice architecture: Part 2, accessed on December 27, 2025,
<https://microservices.io/post/architecture/2025/05/28/microservices-authn-authz-part-2-authentication.html>
84. Understanding WhatsApp Architecture & System Design - TechDotBit, accessed on December 27, 2025,
<https://techdotbit.com/whatsapp-architecture-software-development-architecture/>
85. Unpacking WhatsApp's System Design: The Power Behind Your Messages - Medium, accessed on December 27, 2025,
<https://medium.com/@lovejot.singh/unpacking-whatsappss-system-design-the-power-behind-your-messages-316280b38f78>
86. How Discord Migrated Trillions of Messages from Cassandra to ScyllaDB, accessed on December 27, 2025,
<https://www.scylladb.com/tech-talk/how-discord-migrated-trillions-of-messages-from-cassandra-to-scylladb/>
87. How Discord Stores TRILLIONS of Messages | by Mohammed Vaghjipurwala | Medium, accessed on December 27, 2025,
<https://mohammedvaghjipurwala.medium.com/how-discord-stores-trillions-of-messages-27906ef2a6ff>
88. Netflix Architecture: A Deep Dive into Seamless Global Streaming - Talent500, accessed on December 27, 2025,
<https://talent500.com/blog/netflix-streaming-architecture-explained/>
89. Inside Netflix's Architecture: Encoding, Open Connect & Chaos Engineering - YouTube, accessed on December 27, 2025,
https://www.youtube.com/watch?v=FBwrKr_kiFI
90. Finding a needle in Haystack: Facebook's photo storage - Meta Research, accessed on December 27, 2025,
<https://research.facebook.com/publications/finding-a-needle-in-haystack-facebook-photo-storage/>
91. Needle in a haystack: efficient storage of billions of photos - Engineering at Meta - Facebook, accessed on December 27, 2025,

<https://engineering.fb.com/2009/04/30/core-infra/needle-in-a-haystack-efficient-storage-of-billions-of-photos/>

92. Architecture, Design, and Implementation — Ringpop 0.1.0 documentation, accessed on December 27, 2025,
https://ringpop.readthedocs.io/en/latest/architecture_design.html
93. How Ringpop from Uber Engineering Helps Distribute Your Application, accessed on December 27, 2025,
<https://www.uber.com/blog/ringpop-open-source-nodejs-library/>