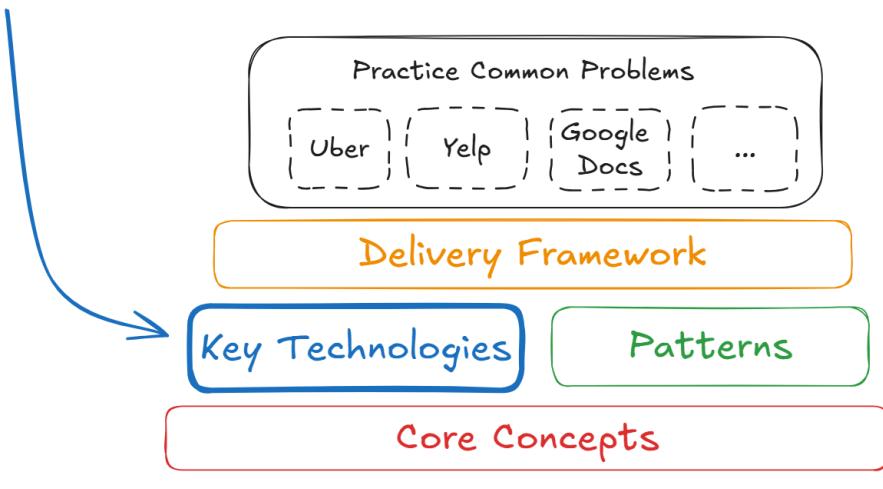


You are here



★ Understanding the Purpose of “Key Technologies” in System Design

System design interviews are not about memorizing fancy technologies.

They are about **assembling the right building blocks** to solve a problem reliably and scalably.

The passage essentially says:

You don't need to know *every* technology, but you must know at least *one good option* in each important category (database, cache, queue, load balancer, search index, etc.).

If you lack basic knowledge of one category (e.g., queues or caches), then:

- you will be **unable to scale the system**,
- and you will be unable to explain **why** you chose something,
- which becomes a **major red flag for an SDE or senior SDE interview**.

Why You Need Familiar “Building Blocks”

Think of system design like building a house:

- You do **not** need to know every brick company.
- But you **MUST** know that bricks are used for walls, cement binds them, steel rods add strength, etc.

Similarly, in system design:

- You must know **at least one database**
- **at least one cache**
- **at least one queue**
- **at least one search index**
- **at least one object storage**
- **at least one stream-processing tool**

The interviewer will NOT judge you for *which* tool you pick.

But they WILL judge you if:

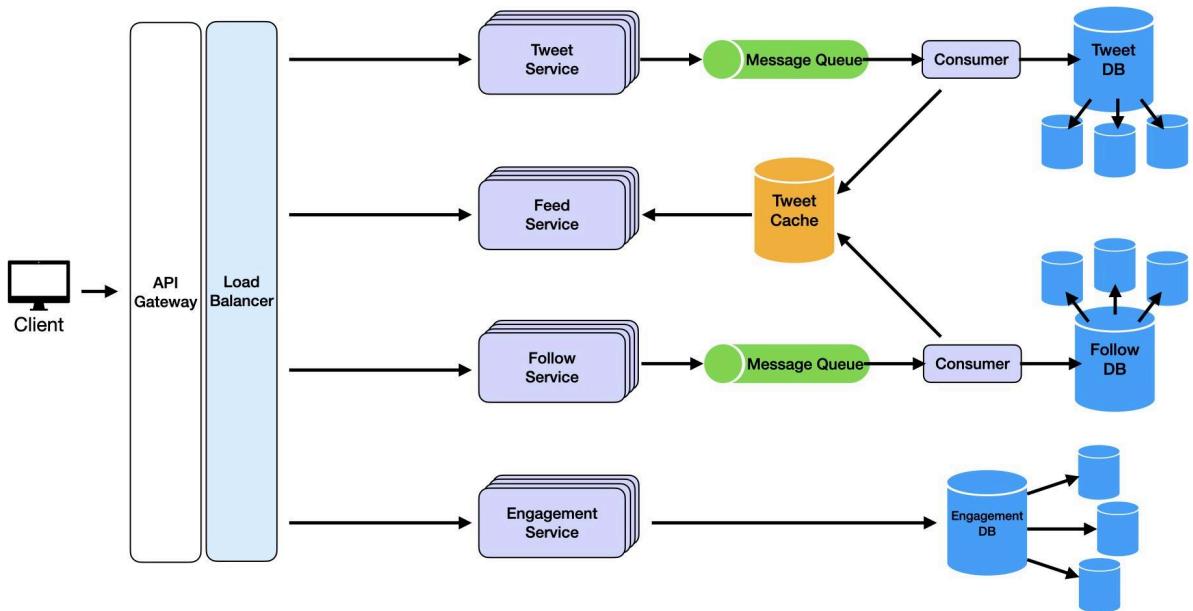
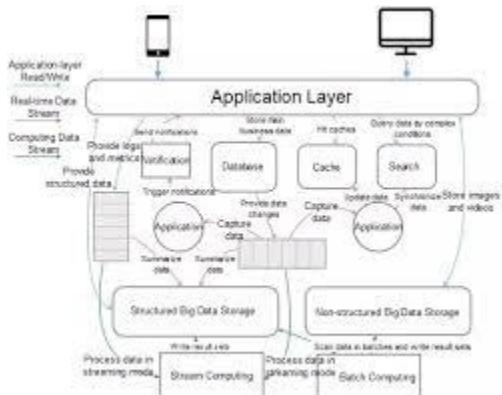
- you have **no tool** for a particular need
 - or you cannot explain **what the chosen tool does**
-

Overall Structure of “Key Technologies” Section

This section of the guide is saying:

“We’ll cover the main technology categories you will use to solve 90% of system design problems.”

The categories generally include:



✓ Relational Database (Postgres, MySQL)

Stores structured, transactional data.

✓ NoSQL Databases (DynamoDB, MongoDB, Cassandra)

Stores massive-scale or flexible-schema data.

✓ Caches (Redis, Memcached)

Drastically reduce latency; avoid hitting the DB repeatedly.

✓ Message Queues (Kafka, SQS, RabbitMQ)

Decouple services, handle async processing, smooth traffic spikes.

✓ Object Storage (S3, GCS, Azure Blob)

Store images, videos, documents at massive scale.

✓ Search Index (ElasticSearch, OpenSearch, Solr)

Fast text search, ranking, filtering.

✓ Load Balancers, CDNs, Reverse Proxies

Distribute traffic, reduce latency, protect backend.

✓ Logging & Monitoring (Prometheus, Grafana, ELK)

Track system health.

These categories form your **toolbox**.

If you know at least ONE choice in each category, you can design almost any system.

⭐ Why “Breadth Before Depth” Matters

The guide is warning:

- A **mid-level candidate** only needs to know what each component *does* (e.g., “ElasticSearch is used for fast text search.”)
- A **senior candidate** must know how it works internally, at least conceptually (e.g., “ElasticSearch uses an inverted index → scaling involves sharding and replicas.”)

⚠ If a senior candidate cannot explain basics (like inverted index in ElasticSearch or partitioning in Kafka), it signals lack of experience. That's why the guide calls it a **yellow flag**.

So the rule is:

👉 Breadth first: Know what each building block does.

👉 Depth comes later: Know how each one works internally if you are senior level.

⭐ Why Interviewers Do NOT Care About Specific Tools

Interviewers don't care if you pick:

- Kafka or RabbitMQ
- Redis or Memcached
- ElasticSearch or OpenSearch

What matters is:

- ✓ **Do you know why that tool is needed?**
- ✓ **Do you know what problems it solves?**
- ✓ **Do you know its limitations and how it scales?**

Example:

You say:

“To handle spikes in write volume, I'll put a queue (Kafka) between the API and the order processor.”

The interviewer hears:

- You understand decoupling
- You understand buffering
- You understand backpressure

They don't care if you chose Kafka or RabbitMQ — they care that you thought of a **queue**.

⭐ Example of Good Breadth Answer (Mid-level)

“We’ll use Redis for caching, S3 for image storage, Kafka for queueing, Postgres for consistent transactional data, and ElasticSearch for fast search.”

This shows you know the essential components.

★ Example of Good Depth Answer (Senior)

“ElasticSearch uses an inverted index which makes text search extremely fast. For scaling, we can shard the index by user ID and use replicas for read throughput. But indexing is expensive, so we will buffer writes through Kafka.”

This shows both breadth and deep reasoning.

★ Summary in Simple Words

The passage is telling you:

- Learn the **main categories** of technologies.
- Know *one* tool per category deeply enough to use in a design.
- Don’t worry about memorizing 20 different tools.
- Interviewers expect **breadth**, not tool fetishism.
- Higher-level candidates need deeper internals knowledge.

1.) Core Database

1. Why “Core Database” Matters in System Design Interviews

In almost every system design problem—Twitter feed, e-commerce orders, Uber rides, Instagram posts—you need a place to store data.

Your *first responsibility* is to pick **one appropriate database** and explain why it fits your design.

Interviewers don't want you to waste time giving textbook comparisons like:

- “SQL is for relationships”
- “NoSQL is for scalability”

These are **oversimplifications** and indicate weak understanding.

Reality:

Both relational and NoSQL databases can handle relationships and scale extremely well *if designed properly*.

👉 So the interview trick is:

Talk about the **database you choose**, not about why you rejected the other.

2. How to Answer in an Interview (Important Tip)

Bad approach (yellow flag):

“I am choosing SQL because my data has relationships.”

Better approach (what interviewers love):

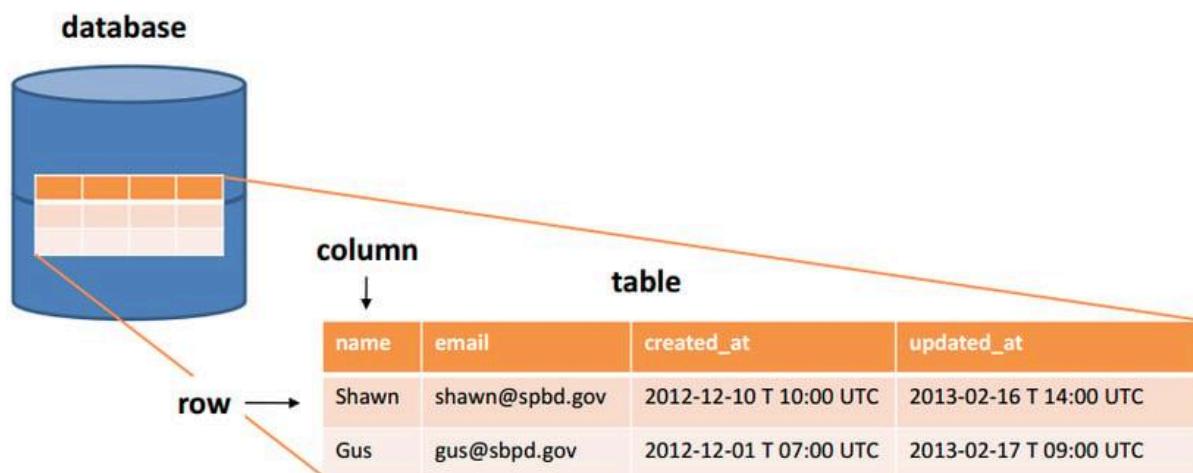
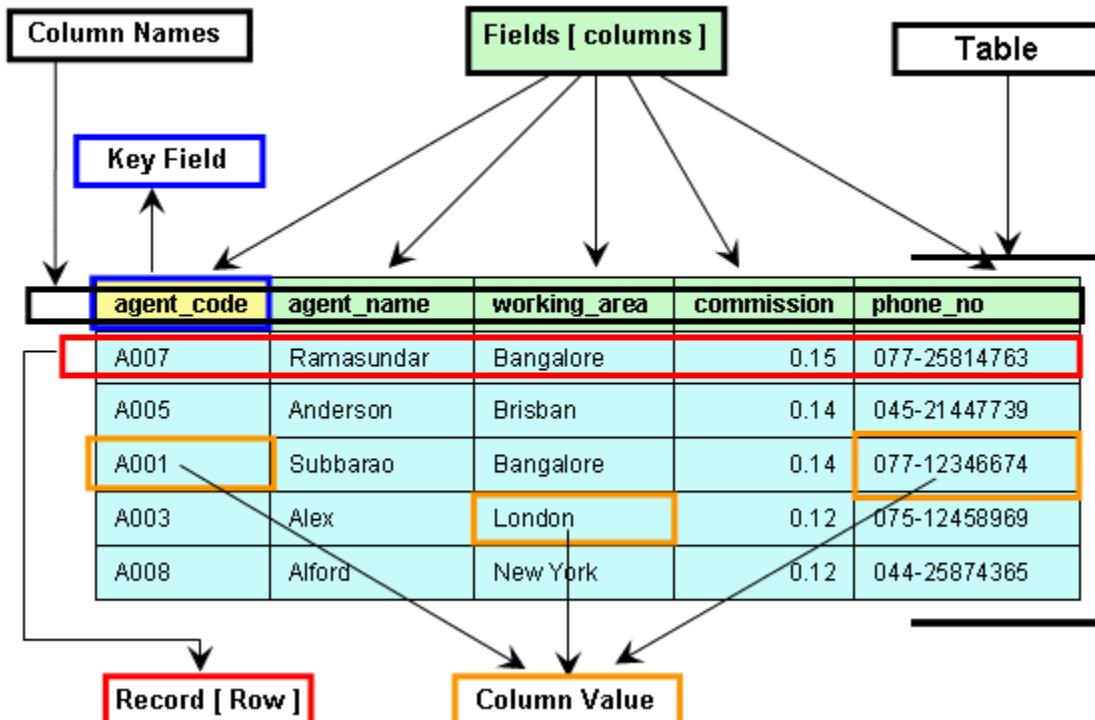
“I'll use Postgres since strong ACID guarantees help maintain data integrity for orders and payments.”

OR for NoSQL:

“I'll use DynamoDB here since the access pattern is simple key-value reads at very high scale.”

You justify based on **your design's needs**, not generic statements.

✓ 3. Relational Databases (RDBMS) – What They Are



A relational database stores data in **tables**.

Each table has:

- **Columns** → fields (name, email, price, etc.)
- **Rows** → each entry/record (one user, one order)

Examples:

id	name	email
1	Raj	raj@gmail.com

Used heavily for:

- user accounts
- orders in e-commerce
- transactions
- product catalogs

They use **SQL** (Structured Query Language) to *query, insert, update, join* data.

4. When should you use a relational database?

Use RDBMS when your system needs:

Strong consistency / correctness

Banking, payments, orders, inventory updates — all must remain correct.

Structured queries

Complex filters, sorts, aggregations.

Joins (but carefully)

When you need to combine information from multiple tables.

✓ Transactions

Atomic multi-step operations (explained below).

That's why relational DBs are recommended for **product-focused interviews** (Instagram-like, e-commerce-like, etc.).

✓ 5. Key Features of Relational Databases You Must Know for Interviews

1 SQL JOINs

A JOIN lets you combine data from two tables.

Example:

You have a **users** table and a **posts** table.

To fetch all posts written by *Raj*, you join them on `user_id`.

JOINS are *powerful*, but they:

- can become slow on large tables
- can require indices to be fast
- sometimes force you to denormalize or cache data externally

In system design interviews:

“I will try to avoid heavy joins on high-traffic paths. I’ll use caching or pre-aggregated tables instead.”

2 Indexes

Indexes make lookups extremely fast.

Without index → Database must scan entire table
With index → Database jumps directly to relevant rows

Analogy:

A book without an index → you read every page
A book with index → jump straight to page number

Types:

- **B-Tree index** → range queries, sorting
- **Hash index** → exact match (key = value)
- **Multi-column index** → (city, age)
- **Special indexes** → geospatial, full-text search

Interview tip:

"I'll add an index on `email` to speed up login queries."

3 Transactions (ACID)

A transaction groups multiple operations into *one atomic step*.

Example:

When creating an order:

1. reduce inventory
2. create transaction
3. create order record

Either **all succeed** or **all fail**.
No partial data.

ACID stands for:

- **Atomicity** – all-or-nothing

- **Consistency** – DB moves from valid state to valid state
- **Isolation** – concurrent transactions don't interfere
- **Durability** – once committed, data is safe

This is extremely important for:

- payments
 - banking
 - inventory
 - user account creation
-



6. Common RDBMS Choices

- **Postgres** → Best default choice (rich features, strong transactions, JSONB support, indexing variety)
- **MySQL** → Lightweight, extremely popular, simple architecture

For interviews, Postgres is highly recommended because:

- supports multi-indexes
 - strong ACID
 - great JSON support
 - powerful query engine
-



7. NoSQL vs SQL – What NOT to Do

Avoid generic statements like:

- ✗ “SQL is for complex relationships.”
- ✗ “NoSQL is for scalability.”
- ✗ “NoSQL cannot do joins.”
- ✗ “SQL cannot scale horizontally.”

These are half-true and over-simplified.

Instead:

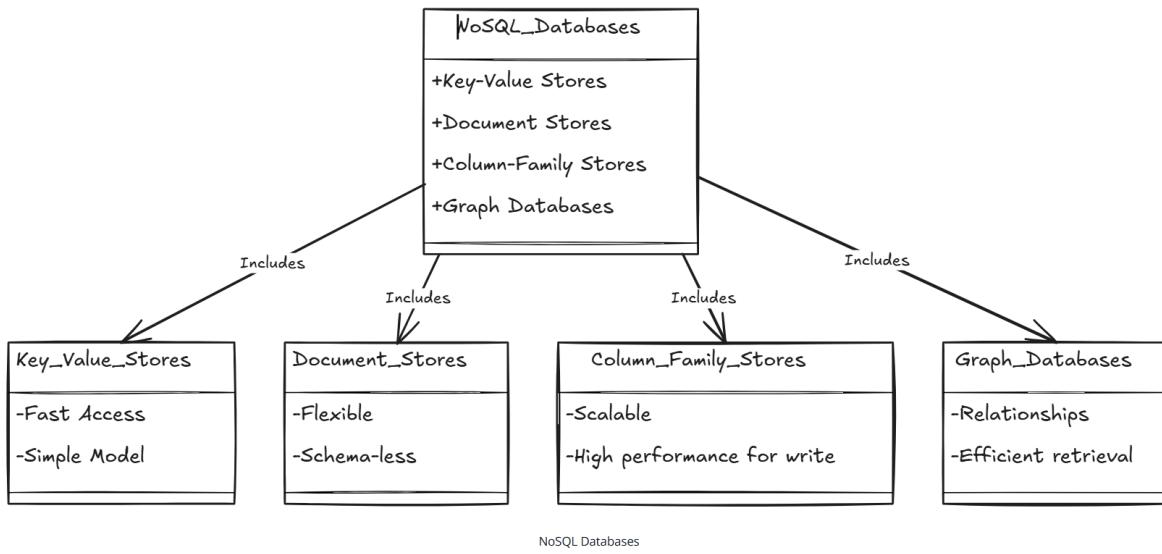
- ✓ Explain why your chosen DB fits your system
 - ✓ Talk about operational needs (consistency, latency, throughput)
 - ✓ Mention known trade-offs
-

8. Example of a **PERFECT** Interview Answer

“For storing user accounts, posts, and relationships, I’ll choose Postgres. Postgres gives me strong ACID guarantees, flexible indexing, and the ability to model normalized data cleanly. For parts of the system requiring very high-speed lookups, like session storage or feed caching, I’ll complement this with Redis.”

This shows clarity, confidence, and real-world understanding.

2.) NoSQL Databases



Below is a **clear, complete, and deeply structured explanation** of the diagram *and* the entire NoSQL explanation you pasted. I'll break it into:

1. **What NoSQL really is**
2. **When to use it**
3. **Why you shouldn't make wrong generalizations in interviews**
4. **Each of the 4 major NoSQL types explained (with examples + use cases)**
5. **Key internal concepts (indexes, consistency, scalability)**
6. **Common NoSQL databases and when to choose them**

I'll also describe the diagram visually in simple language.

1. What is a NoSQL Database? (Simple Definition)

NoSQL = Not Only SQL

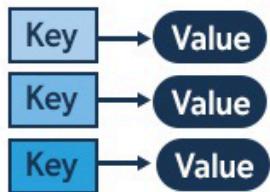
It represents a family of databases that:

- Do **NOT** follow fixed table-row-column structure
- Are usually **schema-less** (schema can evolve easily)
- Support **different data models** (key-value, document, column-family, graph)
- Designed for **horizontal scalability** (spread data across many servers)
- Handle **huge amounts of data & high traffic**

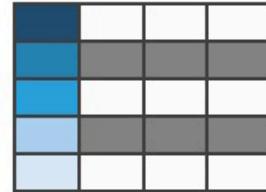
Think of NoSQL as a toolbox with **4 different types of tools**.

NoSQL

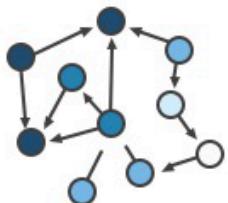
Key-Value



Column-Family



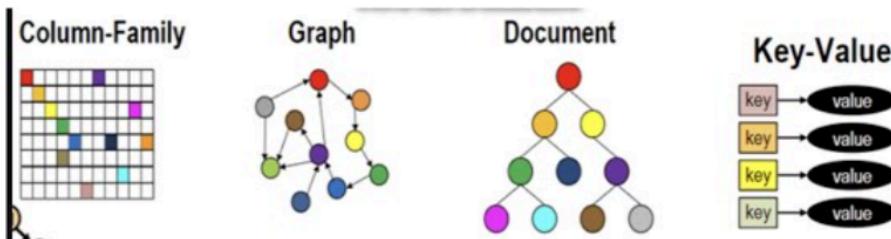
Graph



Document



NoSQL



✓ 2. Why and When Should You Use a NoSQL Database?

Use NoSQL when:

1 Flexible or Evolving Data Models

If data structure is frequently changing or varies per record.
Example: user profiles, logs, items with different attributes.

2 You Need Horizontal Scalability

You want to scale by adding more servers, not upgrading a single server.

3 Large-Scale, Real-Time, High-Speed Systems

Examples:

- Real-time analytics
- Logging pipelines
- IoT sensor data
- High-write systems (e.g., time-series data)

Important Note (interview warning):

The guide stresses an important point:

“Don’t say NoSQL is for scale and SQL cannot scale.”

“Don’t say NoSQL is for flexible schema and SQL cannot do that.”

Why?

Because **Postgres can scale** and **SQL can store JSON** too.

Instead say:

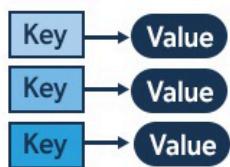
“For this system, I choose DynamoDB because it provides predictable scalability and fits our access patterns.”

✓ 3. Understanding the Diagram (Step-by-Step Explanation)

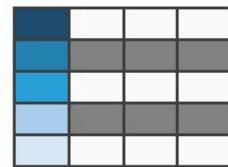
Here is the diagram you shared, explained clearly:

NoSQL

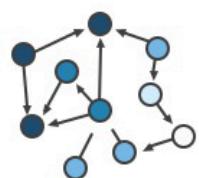
Key-Value



Column-Family



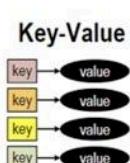
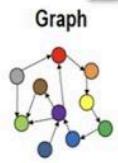
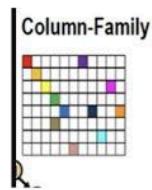
Graph



Document



NoSQL



The top box:

NoSQL Databases

Includes:

- Key-Value Stores

- Document Stores
- Column-Family Stores
- Graph Databases

This means “NoSQL” is not one database; it is an umbrella category.

4 boxes below explain each NoSQL type:

4. The 4 Types of NoSQL Databases (with examples)

1 Key–Value Stores

- ✓ Fastest access
- ✓ Simplest model (key → value)

The value can be anything: string, JSON, blob, object.

Examples: Redis, DynamoDB (internally a key-value store)

Use cases:

- Caching
- User session storage
- Feature flags
- Shopping cart (by user ID)

Why so fast?

O(1) lookup — just lookup by key.

2 Document Stores

✓ Flexible

✓ Schema-less

Stores objects as JSON-like documents

Each item can have different fields.

Examples: MongoDB, CouchDB, Firestore

Use cases:

- User profiles
- Product catalogs
- Content (blogs, comments)

Why popular?

Front-end JSON ↔ Backend JSON → very natural development model.

3 Column-Family Stores

✓ Highly scalable

✓ Great for high write throughput

✓ Inspired by Google Bigtable

Data is stored in **wide rows**, grouped by column families.

Examples: Cassandra, HBase

Use cases:

- Time-series data

- Sensor logs
- Event streams
- Write-heavy analytics systems

Why?

Append-only storage + distributed architecture → huge write performance.

4 Graph Databases

- ✓ Best for relationships
- ✓ Fast traversal of connections

Instead of tables or documents, it stores:

- Nodes (users)
- Edges (follows, friends, connections)

Examples: Neo4j, Amazon Neptune

Use cases:

- Social networks
- Recommendation graphs
- Fraud detection (link analysis)

Why?

Relational DBs struggle with deep relationship queries; graph DBs excel.

5. Things You MUST Know About NoSQL (For Interviews)

1 Consistency Models

SQL databases = strong consistency by default

NoSQL databases = variety:

- **Strong Consistency**

Reads always return the latest data (ex: DynamoDB strongly consistent reads)

- **Eventual Consistency**

After some time, all nodes have the same data

Faster, better for distributed systems

Explain like this:

“User timelines can use eventual consistency because slight delays are acceptable.”

“Payment systems must use strong consistency.”

2 Indexing in NoSQL

Like SQL, NoSQL also uses:

- B-tree indexes
- Hash indexes
- Secondary indexes (MongoDB, DynamoDB GSI)

But indexing differs between databases.

3 Horizontal Scalability

NoSQL databases typically use:

- **Sharding / Partitioning**

- **Consistent hashing**
- **Replication (for fault tolerance)**

Example answer:

“DynamoDB partitions data automatically based on partition key. Cassandra uses consistent hashing ring architecture.”

6. Most Common NoSQL Databases & When to Use Them

DynamoDB (AWS) — Recommended

When to use:

- Serverless apps
- Predictable high scale
- Key-value access patterns
- Low operational overhead

Why it's preferred:

- Auto-scaling
 - Global tables
 - Built-in caching (DAX)
-

Cassandra

When to use:

- Extremely high write throughput
- Large clusters across multiple datacenters
- Time-series data

Tradeoffs:

- No strong consistency by default
 - Limited query flexibility
-

MongoDB

When to use:

- Flexible document structures
- Analytics + operational workloads
- Developers want JSON-like documents

Tradeoffs:

- Requires careful schema design for large-scale workloads
-



FINAL SUMMARY (Interview-Ready)

NoSQL = a family of database models.

It includes:

- **Key-value (fastest)**
- **Document (flexible)**
- **Column-family (high writes)**
- **Graph (relationships)**

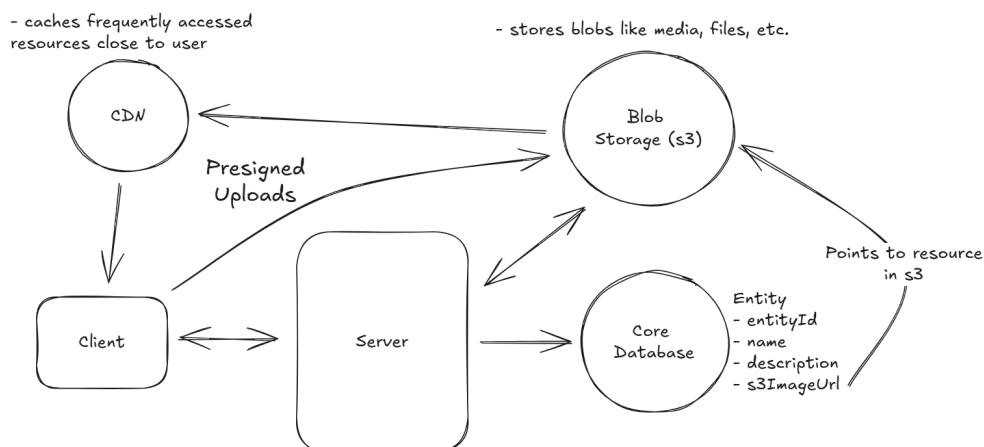
Use NoSQL when:

- Data model changes frequently
- You need horizontal scalability
- You're handling huge traffic or unstructured data
- Real-time applications require low latency

But:

Do NOT say “SQL cannot scale” or “SQL cannot handle JSON.”
 Explain **why the chosen NoSQL database fits your design.**

3.) Blob Storage



Basic Blob Storage Example

1. What Is Blob Storage? (Simple Definition)

Blob Storage = Storage for **Binary Large Objects**, i.e.:

- images
- videos
- PDFs
- audio
- raw files
- user uploads

Think of it as a giant **infinite hard disk in the cloud**.

Typical services:

- **Amazon S3**
- Google Cloud Storage
- Azure Blob Storage

Storing large files inside a **database** is:

-  Expensive
-  Slow
-  Inefficient
-  Hard to scale

So we store files in **blob storage** and only store **metadata + URL** in the database.

2. Why Use Blob Storage Instead of Database?

Because blob storage is:

- ✓ **Extremely cheap**

S3 = ~\$0.02 per GB/month
DB = ~\$1 per GB or more

✓ Extremely scalable

S3 is practically *infinite* storage.

✓ Highly durable

Your file is replicated across multiple data centers.
Durability is usually **11 9s (99.999999999%)**.

✓ Optimized for serving media

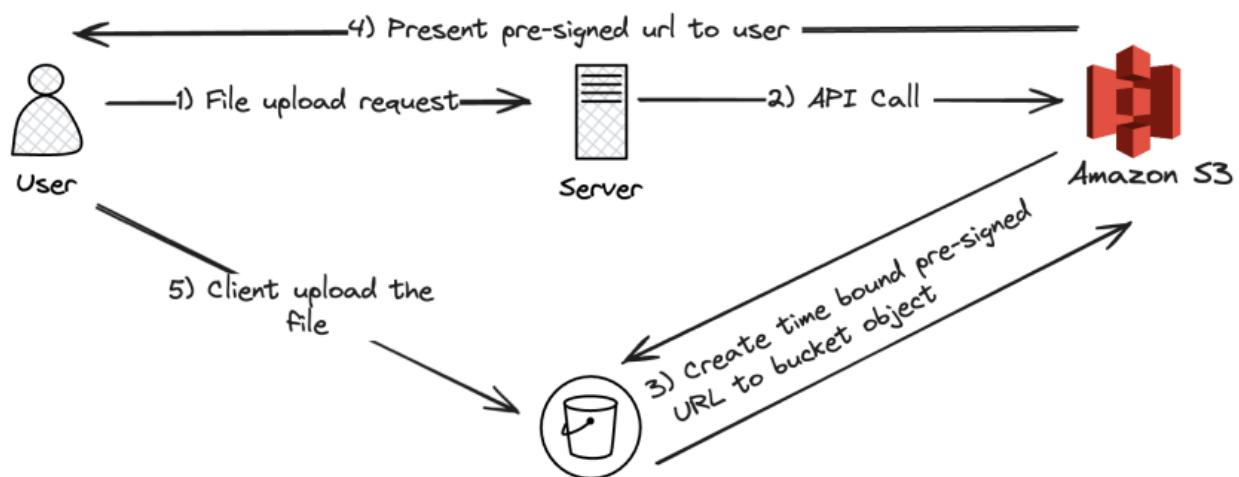
Especially when combined with a CDN.

✓ Perfect for applications with user-uploaded files

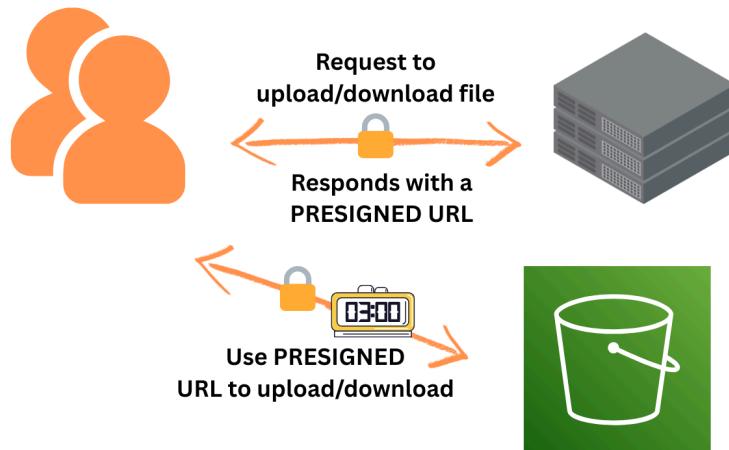
Instagram images, YouTube videos, Dropbox files, WhatsApp profile pictures.

🌟 3. Diagram Explanation (Step-by-Step)

I will describe your diagram in clean English:



Signed, Sealed, Secured: Amazon S3 Style



The architecture includes:

- Client (mobile/web app)
- Server / Backend API
- Blob storage (S3)
- Core Database (Postgres/DynamoDB)
- CDN (CloudFront / Cloudflare)

Let's walk through what each part does.

★ A. Upload Flow (Presigned URL Process)

1 Client wants to upload a file (image/video).

It sends a request to the **server**:

“Give me permission to upload.”

2 Server generates a *presigned URL*

A temporary, secure URL that allows the client to upload directly to S3.

The server also:

- creates a database record with a placeholder
- stores metadata (name, description, entityId)

3 Client uploads file *directly* to S3 using the presigned URL

Important:

- ✓ The file does NOT go through your backend
- ✓ Less load on server
- ✓ Faster for the user
- ✓ Cheaper for you

4 S3 notifies the server (optional)

“File upload complete.”

Server updates DB status → uploaded = true

★ B. Download Flow (Client Requests a File)

1 Client asks server for the file

Example:

“Get profile picture of user 123”

2 Server responds with URL

Either:

- A public CDN URL
OR
- A new **presigned URL** for secure download

3 Client downloads from CDN → not from your server

CDN = much faster, served from nearest edge location.

This improves:

- ✓ Performance
 - ✓ Reduces bandwidth cost
 - ✓ Reduces server load
 - ✓ Reduces S3 load
-

★ 4. Relationship Between Core Database and Blob Storage

Your database **does not store the file**.

It only stores:

`entityId`

`name`

`description`

`s3ImageUrl`

Meaning:

- Database stores metadata
- Database stores pointer (URL)
- Blob storage stores the actual large file

This gives the best performance and cost efficiency.

★ 5. When Should You Use Blob Storage?

✓ Instagram / Facebook

Store images & videos in S3, store metadata in DB.

✓ YouTube

Videos → S3

Metadata (title, duration, owner) → DB

✓ Dropbox / Google Drive

Files → S3

File metadata → DB

✓ E-commerce

Product images → S3

Product data → DB

Whenever the content is large and binary → store in blob storage.

★ 6. Key Concepts You MUST Know (Interview Points)

★ 1. Durability

Blob storage (like S3) uses:

- replication
- erasure coding
- multiple availability zones

This ensures your data is *nearly impossible to lose*.

★ 2. Scalability

Treat S3 as:

“Infinite storage that just works.”

You never scale it manually.

Interview Tip:

You never design sharding for S3. AWS does it for you.

★ 3. Cost

Huge difference:

Storage Type	Cost per GB
S3	~\$0.023 / GB
DynamoDB	~\$1.25 / GB
SQL DB	Even more expensive

That's why storing images/videos in database is a BIG mistake.

★ 4. Security

Blob storage supports:

- Encryption at rest
- Encryption in transit (HTTPS)
- Access control policies
- IAM roles
- Bucket policies

- Presigned URLs (temporary access)

Presigned URLs = *very important in interviews.*

★ 5. Chunking / Multipart Upload

Used for uploading large files by splitting into parts.

Benefits:

- Upload continues even if a chunk fails
- Upload parts in parallel
- Faster
- More reliable

S3 supports this natively.

★ 7. Examples of Blob Storage Services

- Amazon S3 ★ (most common)
- Google Cloud Storage
- Azure Blob Storage

Features include:

- Versioning
- Lifecycle policies (e.g., auto delete old files)
- Replication

- Access control
- Logging

Interview Tip:

If unsure, ALWAYS pick **Amazon S3**.
Every interviewer understands it.

⭐ 8. FULL Interview Summary You Can Speak

“Blob storage is ideal for storing large binaries like images, videos, and documents. Instead of storing files inside the database—which is expensive and slow—we store them in S3 and keep only metadata + a URL in the core database.”

Clients upload files directly to S3 using presigned URLs, which reduces backend load. We then serve files via the CDN for low-latency global access. Blob storage is extremely durable, highly scalable, cost-effective, and secure. It’s the standard solution for platforms like Instagram, YouTube, Dropbox, and many more.”

4.) Search Optimized Database

🔍 1. What Is a Search-Optimized Database? (Simple Definition)

A **search-optimized database** is a specialized system designed for **full-text search**—finding relevant documents based on words, phrases, partial matches, misspellings, or relevance scoring.

Examples:

- Elasticsearch (most popular)
- OpenSearch

- Apache Solr
- Meilisearch

Typical use-case questions:

- "Find all events containing the phrase 'coldplay concert'"
- "Search for users whose bio mentions 'machine learning'"
- "Search tweets containing 'openai'"

A **traditional SQL query** using:

```
SELECT * FROM documents WHERE document_text LIKE '%search_term%' ;
```

is extremely slow because:

- It scans the **entire table (full table scan)**
- It cannot use normal SQL indexes
- It does not scale beyond small datasets

So, we need a **special database designed for searching text**.



2. Why Do We Need a Search Database?

Search is fundamentally different from:

- exact matches
- numeric filtering
- primary key lookup

Search requires:

- matching incomplete text
- ranking by relevance
- handling typos
- tokenization of text
- fast lookups across millions/billions of documents

A normal SQL database **cannot efficiently perform**:

- full-text search
- relevance scoring
- fuzzy matching
- stemming (“run”, “running”, “runs”)
- multi-field ranking

This is where **Elasticsearch** and similar tools excel.



3. How Search Databases Work (Key Concepts)

Search engines rely on a set of core techniques:



1. Inverted Indexes (Most Important Concept)

This is the heart of search engines.

An **inverted index** maps:

`word → list of documents that contain the word`

Example:

"coldplay" → [doc1, doc7, doc9]

"concert" → [doc1, doc2, doc9]

"tickets" → [doc2, doc4, doc5]

Compare this to a normal database:

- Traditional DB:
"Document → words inside it"
- Search engine:
"Word → documents containing it"

This makes searching **instant**:

- No scanning
 - Immediate lookup
 - Scales well for millions of docs
-

★ 2. Tokenization

Tokenization breaks text into individual words (tokens).

Example text:

"Coldplay live concert in Mumbai"

Tokens:

- coldplay
- live

- concert
- mumbai

Search engines rely on tokenization to build inverted indexes.

★ 3. Stemming / Lemmatization

Words like:

- running
- runner
- runs

are all reduced to their root form:

- **run**

This ensures search:

- “run” finds “running”
- “runner” finds “runs”

Users get better results.

★ 4. Fuzzy Search (Handling Typos)

Users make mistakes:

- “iphon” → wants “iphone”

- “restaurent” → wants “restaurant”

Search engines use algorithms like:

- Edit distance
- Levenshtein automata

to find documents **close enough** to the query.

★ 5. Ranking & Relevance Scoring

Search engines score results based on:

- frequency
- proximity
- field boosting
- recency
- importance

This ensures that more relevant results appear first.

Traditional databases cannot do this.

★ 6. Distributed Scaling

Search databases scale horizontally:

- Data is **sharded** across multiple nodes
- Queries are executed in parallel

- Results are merged and ranked

This enables fast search across billions of documents.



4. When Should You Use a Search Database?

Use it whenever your product has **search** as a primary feature.

Examples:

✓ Ticketmaster

Search through thousands of concerts/events by:

- title
- description
- venue
- artist
- tags

✓ Twitter

Searching tweets by content, hashtags, username.

✓ Ecommerce (Amazon, Flipkart)

Search by product name, brand, category.

✓ Social Media (LinkedIn)

Searching user profiles or posts.

✓ Logging systems (ELK stack)

Elasticsearch powers search over logs.



5. Things You MUST Know for Interviews

These are the bullet points interviewers expect you to mention:

Inverted Indexes

- Core data structure
 - Maps words → documents
 - Fast lookups
-

Tokenization

- Break text into searchable tokens
 - Language-aware tokenizers (English, Hindi, etc.)
-

Stemming

- Reduces word variations to the same root
 - More relevant search
-

Fuzzy Search

- Supports typos
- Uses edit-distance algorithms

- Great UX improvement
-

Scaling

- Add more nodes
 - Elasticsearch automatically shards the index
 - Replication provides fault tolerance
-



6. Popular Search Databases (What to Say in Interviews)



Elasticsearch

Most widely used.

- Distributed
- REST API
- Uses Lucene internally
- High performance
- Supports analytics, aggregations
- Used by Netflix, Uber, Yelp, Slack

Interview Tip:

If unsure → always choose **Elasticsearch**.

⭐ Other Options

1. Postgres Full-Text Search (GIN Indexes)

- Good for simple search
- Prevents spinning up Elasticsearch cluster unnecessarily
- Not as powerful but very useful

2. Redis Search

Has modules like RediSearch
But weaker compared to Elastic.



7. What You Should Say in an Interview

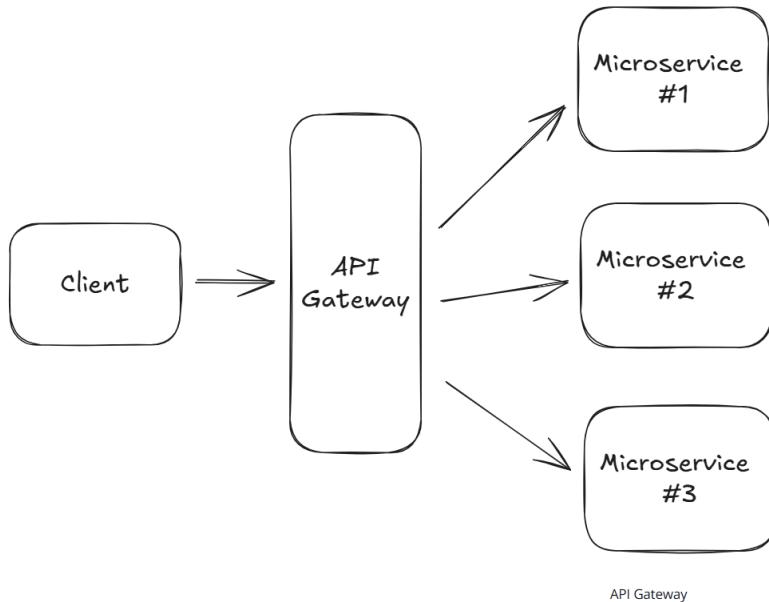
“For full-text search, I will use Elasticsearch.
It builds inverted indexes, supports tokenization and stemming, allows fuzzy search, and can scale horizontally.
This makes search fast and relevant even across millions of documents.
The primary database (Postgres) stores structured data, while Elasticsearch stores denormalized searchable documents.”

This immediately signals you understand:

- search architecture
- correct separation of responsibilities
- real-world patterns

5.) API Gateway

What is an API gateway and when should you use it?



Here is a **clear, deep, practical, interview-ready explanation** of the **API Gateway**, the diagram, how it works, why it exists, and when you should use it in system design. This is one of the **MOST** important components in microservice architecture.

1. What Is an API Gateway? (Simple Definition)

An **API Gateway** is the **front door** for all client requests in a microservice architecture.

It acts as a **single entry point** between:

- Clients (mobile app, web browser, backend services)
- Your internal microservices

Think of it as:

"A smart traffic controller that decides which microservice should handle which request.

🔍 2. Why Do We Need an API Gateway?

Without a gateway:

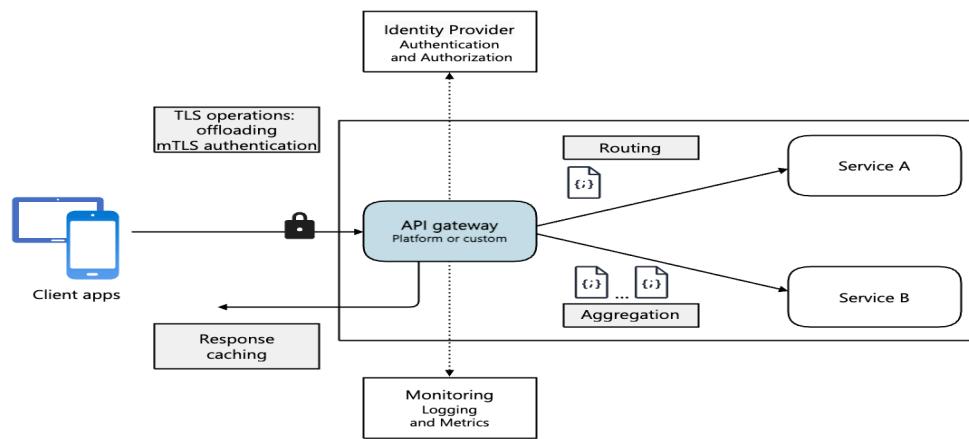
- Client must talk to **10, 20, or 50 different microservices**
- Client must handle **service discovery**
- Client must understand **internal URLs**
- Client must manage **auth, rate limiting, retries**, etc.

This is extremely messy.

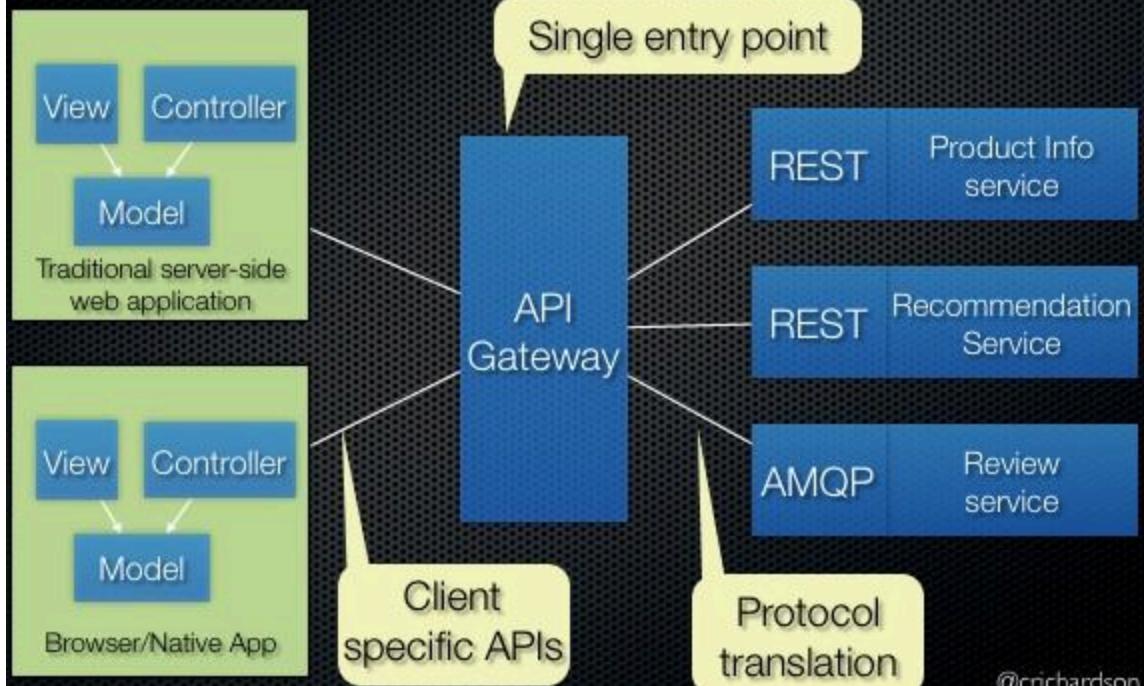
With a gateway:

- ✓ Client has **one unified endpoint**
- ✓ Gateway forwards the request to the right microservice
- ✓ Gateway adds **security, logging, monitoring, rate limiting**, etc.
- ✓ Microservices stay clean and focused on business logic

📘 3. Diagram Explanation (Step-by-Step)



Use an API gateway



In your diagram:

Client → API Gateway → Microservice #1 / #2 / #3

Meaning:

1. **Client sends a request**

Example:

`GET /users/123`

2. **API Gateway receives it as the single entry point**

3. **Gateway forwards it to the correct microservice**

Example:

`/users/123` goes to **User Service**

`/orders/555` goes to **Order Service**

4. **Gateway aggregates responses (if needed)**

Example:

`/user-profile` may require:

- User Service
- Order Service
- Payment Service

5. **Gateway returns final response to client**



4. What Responsibilities Does an API Gateway Handle?

This is the most important part for interviews.



1. Routing

Routes each request to the correct backend microservice.



2. Authentication & Authorization

- Validate JWT tokens
- Check API keys
- Enforce user permissions
- Block unauthorized requests

This removes security logic from each microservice.



3. Rate Limiting

Prevents abuse:

- Max 100 requests/min per user
 - Max 1000/min per IP
 - Block bots
-

★ 4. Logging & Monitoring

Tracks:

- request latency
- failure rate
- performance metrics
- audit logs

Useful for debugging and observability.

★ 5. Request/Response Transformation

Examples:

- Convert HTTP → gRPC
 - Compress responses
 - Add/remove headers
 - Change API versions
-

★ 6. Caching

Can store frequent responses to speed up requests.

★ 7. Load Balancing

Distributes traffic among multiple instances of each microservice.

★ 8. Aggregation

A single API call might need data from multiple microservices.

API Gateway can:

Fetch from multiple services → merge → send one response

⚡ 5. When Should You Use an API Gateway?

✓ Always include it in system design interviews

Especially when designing:

- Instagram
- Twitter
- Netflix
- Uber
- E-commerce
- Social networking
- Microservice-based systems

It's a standard part of modern architectures.

You **MUST** include an API Gateway when:

- ✓ Using microservices
 - ✓ Serving multiple clients (web, Android, iOS, micro-frontend)
 - ✓ Want unified security and rate limiting
 - ✓ Want consistent logging
 - ✓ Want versioning of APIs
-

🚫 When Should You **NOT** Use an API Gateway?

Rare cases:

- Very small systems
- Simple monolithic architecture
- Internal-only services

But for interviews → **ALWAYS include it.**

6. Common API Gateway Technologies

⭐ Fully managed gateways

- AWS API Gateway
- AWS ALB + Lambda Authorizers
- Google Apigee
- Azure API Management

Open-source gateways

- Kong
- NGINX
- Traefik
- Ambassador
- Envoy

NGINX is extremely common (even used by Amazon in early days).

7. What to Say in a System Design Interview

Here is a perfect interview-ready statement:

“I will place an API Gateway at the entry point of my system.

It will handle routing, authentication, rate limiting, logging, request validation, and aggregation.

This keeps clients simple and keeps microservices focused on business logic.

Using a gateway also allows us to easily scale, version APIs, secure endpoints, and maintain observability.”

Interviewers LOVE this explanation.

8. Why Interviewers Don't Go Too Deep on API Gateways

Because:

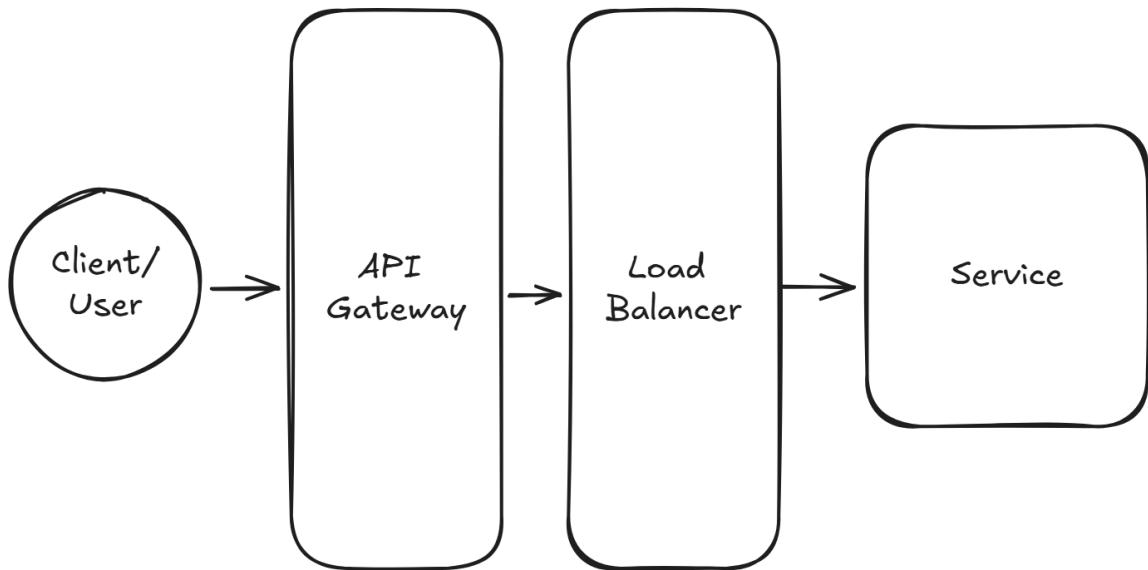
- Gateways are standard
- They're not the part of the design that changes
- Every modern microservice system uses a gateway

- They want to focus more on other components (cache, DB, queues)

So you just show awareness—no need to oversell it.

6.) Load Balancer

What is a load balancer and when should you use it?



A common setup for a horizontally scaled, authenticated service.

🧭 1. What Is a Load Balancer? (Simple Definition)

A **load balancer** is a component that distributes incoming traffic across multiple instances of the same service.

Think of it as:

“A traffic cop that prevents any single server from becoming overloaded.”

Whenever you scale a service **horizontally** (run multiple copies), you **must** place a load balancer in front of them.

2. Why Do We Use Load Balancers?

Systems must handle:

- High traffic
- Sudden bursts of users
- Server failures
- Uneven workload

If all traffic went to a single server:

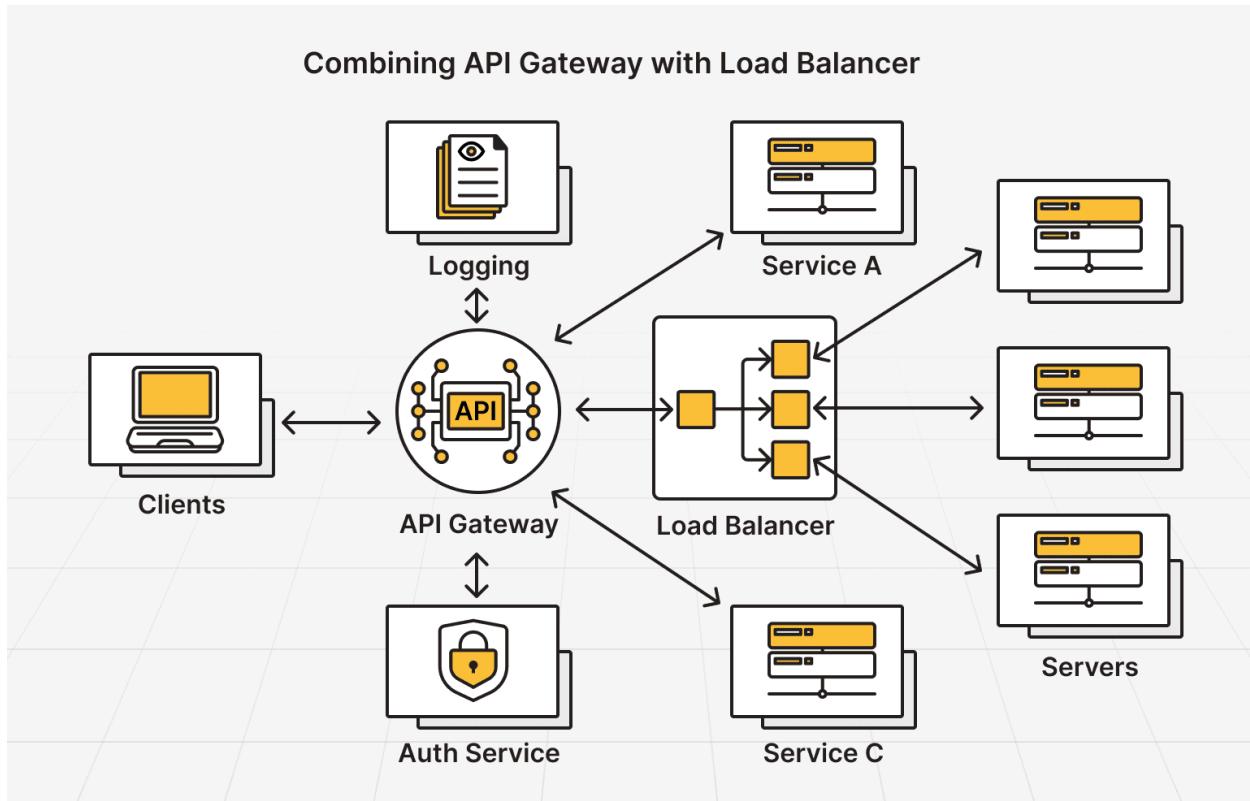
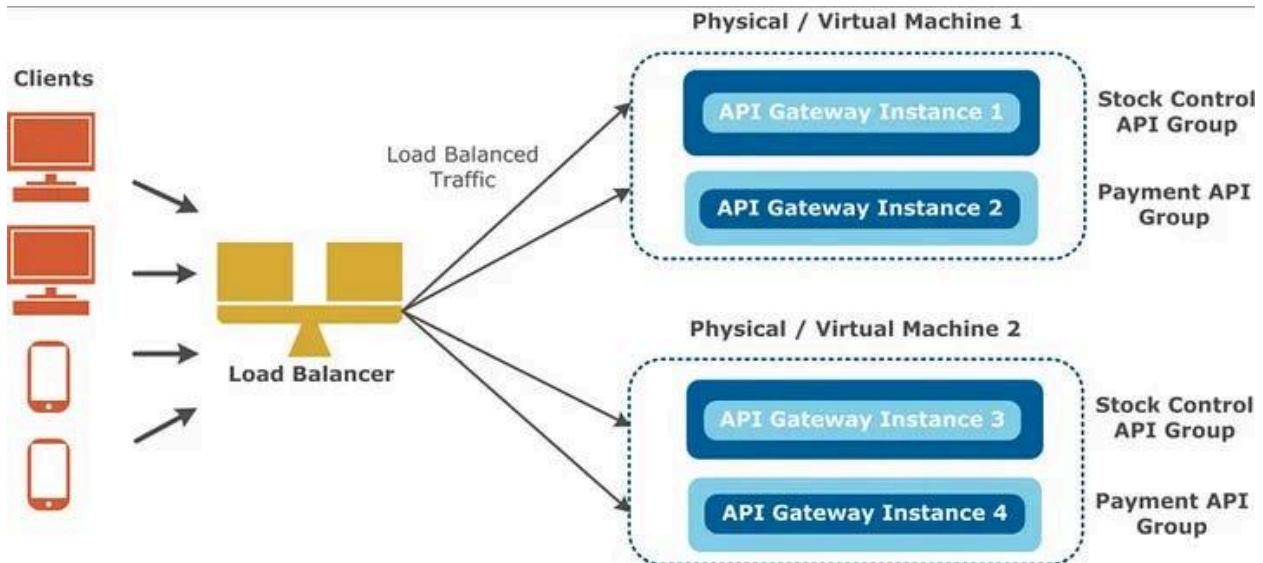
- ✗ That server would crash
- ✗ Users would see downtime
- ✗ You would have a single point of failure

A load balancer ensures:

- ✓ Even distribution
 - ✓ High availability
 - ✓ Scalability
 - ✓ Fault tolerance (if one server fails, LB removes it)
-



3. Diagram Explanation (Step by Step)



Your diagram shows:

Client → API Gateway → Load Balancer → Service

This architecture means:

1 The Client (User)

Makes HTTP requests.

2 API Gateway

Handles:

- Authentication
- Rate limiting
- Logging
- Routing

3 Load Balancer

Distributes requests across **multiple instances** of the same service.

4 Service

Your microservice, usually running:

- many replicas
- on different servers
- behind the load balancer

Thus, the flow is:

User → Gateway → Load Balancer → Service Replica #1/#2/#3 ...



4. When Should You Use a Load Balancer?

You should use a load balancer whenever:

✓ You have multiple instances of a service

Example: 10 instances of User Service, 20 instances of Feed Service.

✓ You want fault tolerance

If one instance dies, the load balancer stops sending traffic to it.

✓ You need horizontal scaling

More replicas = more throughput.

✓ You need zero downtime deployment

Load balancer can switch between old and new versions.

✓ You need high availability

Even if one machine crashes, system stays alive.

Load balancers are essential for:

- Web servers
- Microservices
- Authentication services
- Real-time services
- APIs
- Any system expecting large traffic



⚠ 5. Load Balancers in Interviews — Important Advice

Interviewers don't want you to draw a load balancer everywhere.

Reason:

- It clutters the diagram
- It's implied when you say "the service is horizontally scaled"

Recommended approach:

- Mention that services are horizontally scaled
 - Add a load balancer **only once**, usually near the API gateway
 - Explain that all services behind it are scaled as needed
-



6. L4 vs L7 Load Balancers (Very Important)

Load balancers work at different layers of the OSI model.



L4 Load Balancer (Transport Layer)

Works at **TCP/UDP** level.

Features:

- Very fast
- Uses IP + port
- No knowledge of HTTP paths
- Best for **WebSockets, gRPC, raw TCP connections**

When to use:

- ✓ Persistent connections
- ✓ WebSocket servers
- ✓ Multiplayer gaming servers
- ✓ Real-time communication

★ L7 Load Balancer (Application Layer)

Understands **HTTP, URLs, cookies, headers**.

Features:

- Smart routing
- Can route based on REST path
- Can do A/B testing
- Can add/remove headers
- Can serve cached content

When to use:

✓ Normal web traffic (HTTP/HTTPS)

✓ Routing requests like:

`/users` → User Service

`/videos` → Video Service

✓ Content-aware routing

✓ API versioning

✓ Authentication gates

★ Rule of Thumb (from your text):

If you use WebSockets → L4 Load Balancer

If you use normal HTTP APIs → L7 Load Balancer



7. Features Load Balancers Often Provide

Beyond simple routing, load balancers also offer:

✓ **Health checks**

Detects if service instance is unhealthy.

✓ **Autoscaling integration**

Adds/removes instances based on traffic.

✓ **Sticky sessions (session affinity)**

Pin same user to same instance if needed.

✓ **SSL termination**

Offloads HTTPS encryption to LB.

✓ **Rate limiting**

Basic traffic control.

✓ **Connection pooling**

Reuse fewer TCP connections downstream.



8. Most Common Load Balancers

★ **AWS Elastic Load Balancer (ELB / ALB / NLB)**

- ALB (Application Load Balancer) = L7
- NLB (Network Load Balancer) = L4

★ **NGINX**

Extremely popular open-source load balancer.

★ **HAProxy**

Highly optimized, used in high-performance systems.

★ Hardware load balancers (for *extreme scale*)

Used by telecom companies, financial institutions.



9. What You Should Say in an Interview

Here's a complete interview-ready answer:

"Because this system needs to scale horizontally, I'll place a load balancer in front of my service replicas.

The load balancer will distribute requests evenly, perform health checks, and ensure fault tolerance. For standard HTTP traffic, an L7 load balancer works best because it allows path-based routing and request inspection.

If the system uses WebSockets or persistent connections, I'll switch to an L4 load balancer. We can use AWS ALB or NGINX depending on infrastructure."

This shows awareness of scale, routing, and performance.



10. Final Summary

A Load Balancer is essential when:

- You have **multiple service instances**
- You need **fault tolerance**
- You expect **high traffic**
- You want **high availability**

Load Balancers come in two types:

- **L4** → For persistent connections (WebSocket/gRPC)

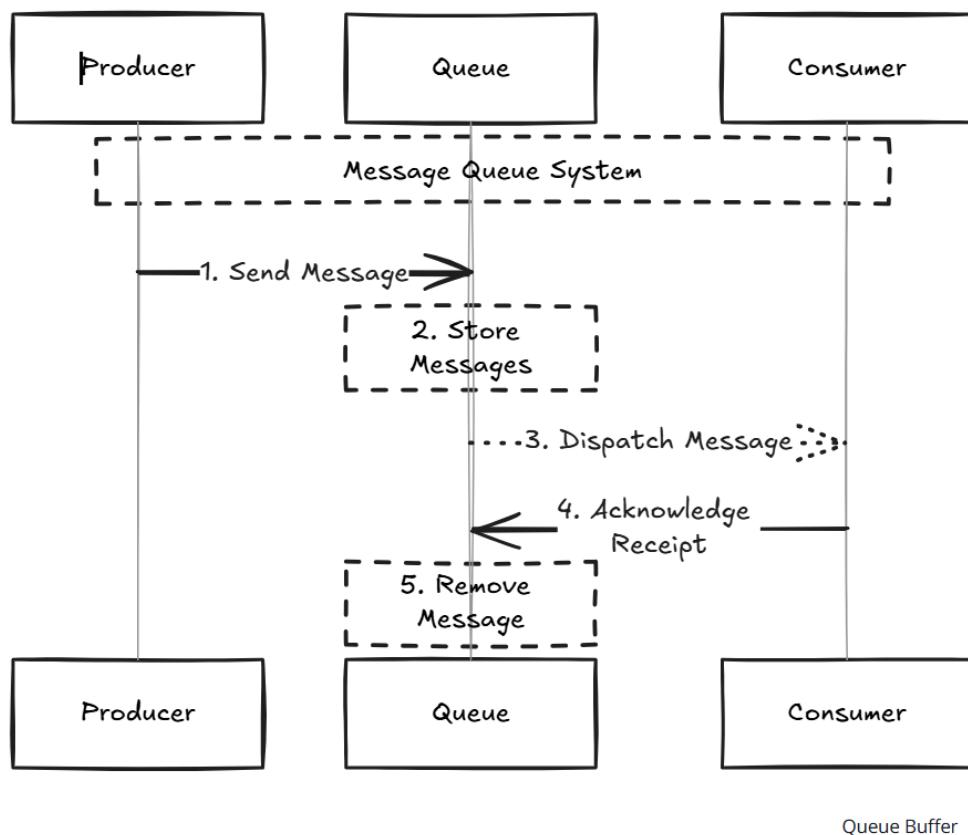
- L7 → For normal HTTP routing and smart routing

Most common choices:

- AWS ELB / ALB / NLB
 - NGINX
 - HAProxy

7.) Queue

What are queues and when should you use them?



✉ 1. What Is a Queue in System Design? (Simple Definition)

A **queue** is a buffer that sits between:

- a **producer** (who sends work/messages)
- a **consumer** (who processes work/messages)

A queue temporarily **stores messages** so that producers and consumers can work *independently* at their own pace.

Think of it like:

“A waiting line where messages wait until a worker is ready to process them.”

⭐ 2. Why Do We Use Queues?

Queues solve **three big problems** in distributed systems:

⭐ 1 Handling Bursty Traffic

If you get 1,000 requests at once but can only process 200 per second:

- Without a queue → **system crashes or drops requests**
- With a queue → extra 800 requests wait safely and get processed later

Nothing is lost, just delayed.

⭐ 2 Decoupling Producer & Consumer

Without a queue:

- Both systems must be up

- Both must be fast
- Both must be scaled equally

With a queue:

- ✓ Producer can send messages and immediately return
- ✓ Consumer can process at its own speed
- ✓ Both systems scale independently

This dramatically improves reliability.

★ 3 Distributing Work Across Multiple Workers

Queues allow multiple consumers to share the workload.

Example:

- Upload photos → photo processing service
- Each photo is a job
- 10 worker machines pull from the queue
- Workload is evenly distributed

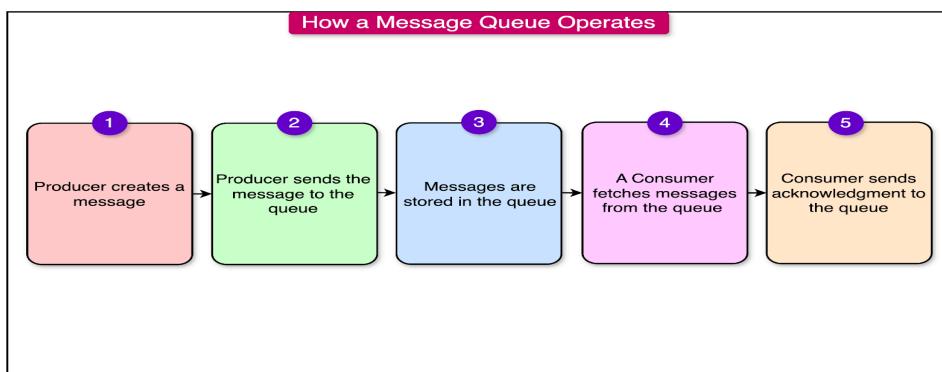
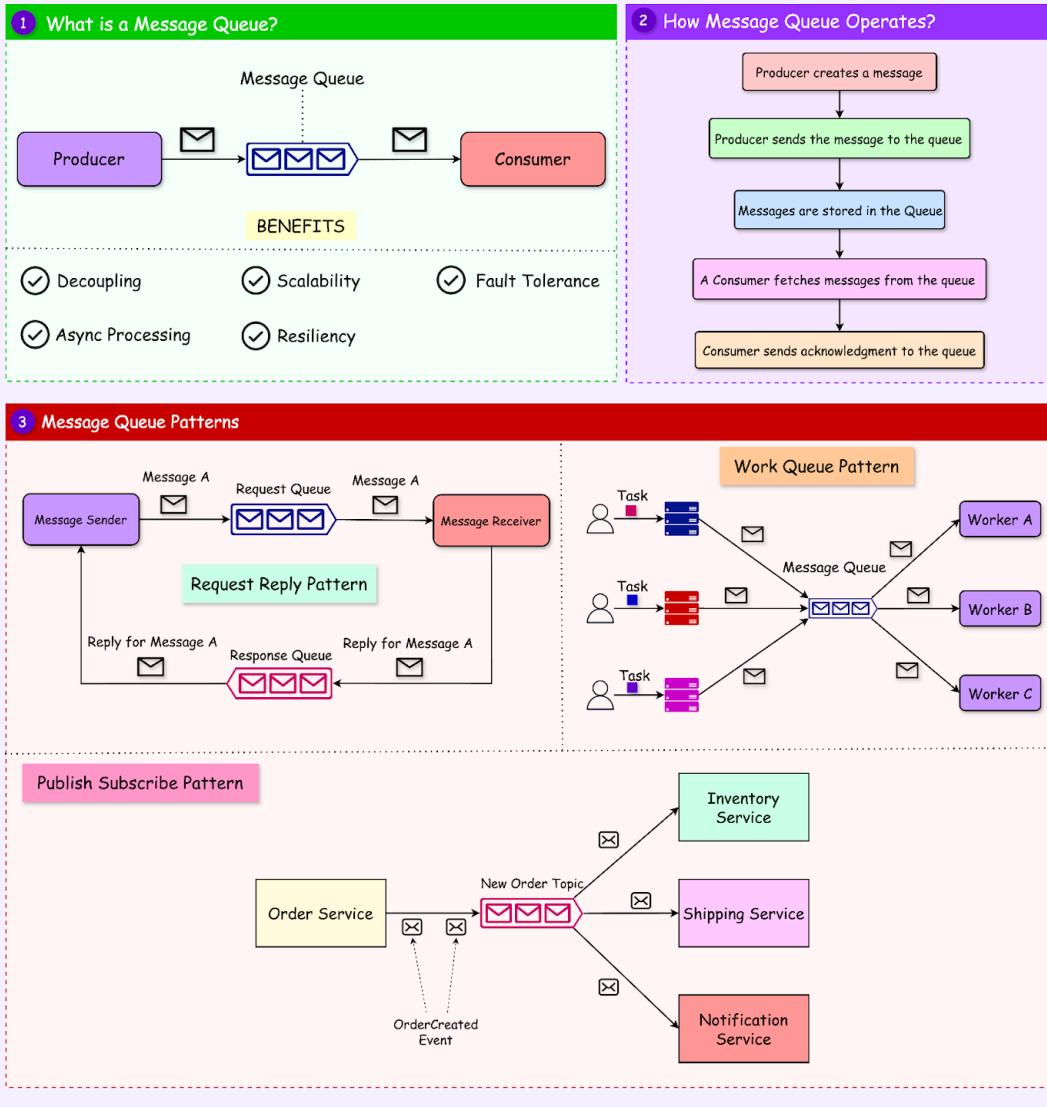
This is the backbone of scalable asynchronous systems.



3. Diagram Explanation (Step-by-Step)

Here's what your diagram represents:

A Cheatsheet On Message Queues



Flow:

1 Producer → Send Message

Producer sends a message to the queue.

2 Queue → Store Message

Queue stores the message safely.

3 Queue → Dispatch Message

Queue sends the message to a consumer (worker).

4 Consumer → Acknowledge

Consumer processes the message and confirms (ACK):

“Job completed.”

5 Queue → Remove Message

Queue deletes the message permanently.

This ensures:

- No duplication
 - No message loss
 - At-least-once processing semantics
-



4. When Should You Use a Queue?

Use a queue whenever:



A. You must handle unexpected bursts of traffic

Example: Uber surges, Black Friday sales, flash events.

★ B. Work is expensive and slow

Example: image processing, video encoding, ML model inference, email sending.

These tasks should NOT block user requests.

★ C. You want asynchronous behavior

User uploads a photo → processing happens later → user doesn't wait.

★ D. You want to scale workers independently

Example:

- 1 producer → 50 workers
- 20 producers → 100 workers

Queues allow dynamic scaling.

🎯 5. Real-World Use Cases

✓ Uber / Ola

Queue user ride-requests during peak hours.

✓ Instagram / Snapchat

Photo/video processing pipelines.

✓ Netflix

Transcoding videos into many formats.

✓ **Amazon**

Order fulfillment, event processing.

✓ **Email / Notification Services**

Send millions of emails asynchronously.

6. Key Concepts You MUST Know for Interviews

1. Message Ordering

Most queues are **FIFO** (first-in-first-out).

But some systems like **Kafka** support ordered processing *within partitions*.

2. Retry Mechanisms

If a consumer fails:

- The queue retries delivery
- Retries can be immediate or delayed
- Some systems use exponential backoff

This prevents occasional failures from losing data.

3. Dead-Letter Queues (DLQ)

If a message fails too many times, it's moved to a **DLQ**.

Purpose:

- Debug errors
 - Inspect malformed messages
 - Prevent poison messages from blocking the queue
-

★ 4. Scaling with Partitions

A single queue can be split into **partitions** across servers.

Each partition:

- Is processed by different workers
- Increases throughput
- Maintains ordering *within* that partition

Kafka does this extremely well.

★ 5. Backpressure (VERY IMPORTANT)

Queues are great, but dangerous if unbounded.

Example:

- System can process 200 req/s
- It receives 300 req/s
- 100 req/s accumulate in the queue
- Eventually, the queue becomes huge
- You will **NEVER catch up**

- Latency explodes → starts delaying hours of work

Backpressure solves this:

- Reject new messages
- Slow producers
- Apply rate limiting
- Drop low-priority tasks
- Auto-scale consumers

This prevents silent overload.



7. Most Common Queue Technologies

★ Kafka (distributed, high throughput, partitions, streams)

Best for:

- event streaming
- real-time analytics
- log aggregation
- high volume pipelines
- microservice communication

★ AWS SQS (simple, managed, reliable)

Best for:

- serverless architectures

- background tasks
- decoupled workloads

Others include:

- RabbitMQ
- Google Pub/Sub
- Redis Streams

But for interviews, the top choices are:

- Kafka
 - SQS
-



8. Perfect Interview-Ready Explanation

Here's a concise but powerful statement you can speak:

"I'll introduce a queue between the producer and consumer to buffer bursty traffic and decouple the two systems.

The producer sends messages to the queue and returns quickly, while a worker pool processes messages asynchronously at its own pace.

The queue provides durability, retries, dead-letter queues, and can be partitioned to scale horizontally.

This ensures stable performance even when incoming traffic spikes."

This immediately demonstrates strong understanding.



9. Final Summary

- ✓ Queues smooth out bursty traffic
- ✓ Decouple producer and consumer
- ✓ Allow asynchronous, scalable processing
- ✓ Support FIFO ordering, retries, DLQs
- ✓ Require backpressure to avoid overload
- ✓ Kafka & SQS are industry-standard

8.) Streams / Event Sourcing

What are streams and when should you use them?

Here is a deep, clear, beautifully structured, interview-ready explanation of Streams / Event Sourcing, including what streams are, how they differ from queues, why we use them, event sourcing, real-time analytics, chat messaging, partitioning, consumer groups, windowing, fault tolerance, and the main technologies (Kafka, Kinesis, Flink).

This is one of the MOST powerful concepts in modern system design, especially for real-time systems, auditability, financial systems, analytics, and distributed architectures.

1. What Are Streams? (Simple Definition)

A **stream** is a continuous, ordered, append-only log of events that can be **read multiple times** by consumers.

Key Points:

- Messages/events are **not removed** after consumption

- Consumers can **re-read** events from any point in time
- Streams retain data for **hours, days, or weeks**
- Many consumers can read **independently**

Streams are used when you need:

- ✓ Real-time processing
 - ✓ Replay capability
 - ✓ Multiple consumers
 - ✓ Auditing / event sourcing
 - ✓ High throughput event ingestion
-

2. How Streams Differ from Queues (Very Important)

Feature	Queue	Stream
Message deleted after consumption	✓ Yes	✗ No
Replay old messages	✗ No	✓ Yes
Multiple consumers	✗ Hard (competing workers)	✓ Easy (consumer groups)
Event ordering	Sometimes	Very strong (per partition)
Ideal for	Background jobs	Real-time pipelines, analytics, event sourcing

Queues = one-and-done processing

Streams = persistent event log for replay & analysis



3. What Is Event Sourcing?

Event Sourcing = store every change to application state as an event.

Instead of storing only the **current state**, you store the **sequence of events** that *led* to that state.

Example (Bank Account):

Instead of storing:

`balance = 2000`

You store:

`Deposit +1000`

`Deposit +800`

`Withdraw -200`

`Deposit +400`

Using the stream, you can:

- Rebuild account balance anytime
- Recompute history
- Undo events
- Run audits
- Ensure correctness in financial systems

This is extremely useful for:

- Banking
 - Trading
 - Ledgers
 - Accounting
 - Revision history
 - Distributed workflows
-



4. When Should You Use Streams?

Let's break down the three main use cases with examples.

★ A. When you need to process huge amounts of data in real time

Example: **Instagram analytics**

- Every like
- Every comment
- Every share
- Every view

Users generate millions of events per second.

Using streams:

1. Events flow into Kafka

2. Flink or Spark Streaming processes them
3. Real-time dashboards update instantly

Perfect for:

- Real-time leaderboards
 - Trending topics
 - Fraud detection
 - Monitoring & alerting
-

★ B. Event Sourcing (Financial & transactional systems)

Example: **Banking**

Each transaction is an event:

- deposit
- withdrawal
- transfer

These events are stored in a stream like Kafka.

Benefits:

- ✓ Replay to rebuild account state
- ✓ Audit ability — who did what when
- ✓ Undo operations
- ✓ Multiple services can consume the same event

This is core to **high-integrity systems**.

★ C. Publish-Subscribe (Chat applications / notifications)

Example: **WhatsApp group chat**

- User sends a message
- Message is published to a stream
- All group members receive the message simultaneously
- Each consumer (each user's device) reads the same event

Streams allow:

- Multiple subscribers
- Ordered messages
- Persistent history

This is the **Pub/Sub** pattern.

🔧 5. Important Concepts You MUST Know for Interviews

These are essential if you're using Kafka or any streaming system in your design.

★ 1. Partitioning (Scaling Streams)

Streams are divided into **partitions**.

Each partition:

- Stores events in order
- Can be processed by a different consumer
- Allows **horizontal scaling**

You must choose a **partition key**, like:

- userId
- chatRoomId
- accountId

All related events must go to the same partition to maintain ordering.

★ 2. Multiple Consumer Groups

A killer feature of streams:

Each consumer group gets its **own independent read position**.

Example:

- Group A → updates real-time dashboard
- Group B → stores events in analytics DB
- Group C → triggers notifications

All reading **the same stream**, independently.

This is impossible with queues.

★ 3. Replication (Fault Tolerance)

To avoid data loss, streams replicate data across nodes.

Kafka uses:

- **leader replica**
- **follower replicas**

If one fails → consumers automatically failover.

This ensures:

- high availability
 - durability
 - fault tolerance
-

★ 4. Windowing

Windowing groups events by:

- **time window** (5 min, 1 hour, 24 hours)
- **count window** (every 100 events)

Example:

“Calculate average ride time in last 24 hours updated every minute.”

Used in dashboards, analytics, metrics.



6. Most Common Stream Technologies

★ Kafka (Industry Standard)

- Distributed streaming platform
- High throughput
- Partitioned logs
- Retention-based storage
- Used by Netflix, Uber, LinkedIn
- Best for **event sourcing + analytics + Pub/Sub**

This is the #1 stream you should mention in interviews.

★ Flink (Stream Processing Framework)

Not a queue or stream storage —

Flink reads from Kafka and processes data in real-time.

Used for:

- aggregations
 - windowed computations
 - analytics
 - continuous data pipelines
-

★ AWS Kinesis

AWS alternative to Kafka.

Built-in integrations with AWS services.



7. Perfect Interview-Ready Explanation

Here's what you can confidently say in an interview:

"Since this system requires real-time processing and replay capability, I'll use a streaming platform like Kafka.

Streams retain events for a configurable window and allow multiple consumer groups to independently read and process events. This makes them ideal for analytics pipelines, event sourcing, and scenarios with multiple subscribers.

Kafka partitions the stream so it can scale horizontally, ensures fault tolerance via replication, and supports advanced processing techniques like windowing for time-based aggregations."

This instantly signals you understand advanced distributed systems.

8. Final Summary (Super Clean)

- ✓ Streams retain data (unlike queues)
- ✓ Consumers can re-read events
- ✓ Ideal for real-time pipelines, analytics, event sourcing, pub/sub
- ✓ Partitions allow massive scaling
- ✓ Consumer groups allow multiple independent processing paths
- ✓ Windowing enables real-time aggregations
- ✓ Kafka + Flink + Kinesis are the most common tools

9.) Distributed Lock

What are distributed locks and when should you use them?

let's break **Distributed Locks** down in the clearest, most intuitive way with **real-world examples, why they're needed, how they work, and what interviewers expect.**

This is one of the most misunderstood but extremely important system-design topics — especially for e-commerce, ride matching, ticketing systems, auctions, and distributed cron jobs.



1. What Is a Distributed Lock? (Simple Definition)

A **distributed lock** ensures that **only one machine / service / process** can access a shared resource at a time *in a distributed system*.

Think of it as:

“A system-wide lock that works across multiple servers.”

This is necessary because in distributed systems:

- You have **multiple servers handling traffic**
- Many users may try to modify the **same resource** at the same time
- Traditional DB locks are **too short-lived** (milliseconds), not suitable for 5–10 minute “reservations”

Distributed locks solve this problem.



2. Why Do We Need Distributed Locks? (Real Example: Ticketmaster)

Imagine Coldplay tickets go live.

- 10,000 people try to buy **the same seat**
- Your system is running on **200 servers**
- All 200 servers receive booking requests for the same seat (Seat A12)

We need:

- ✓ Only **ONE** user to reserve A12
- ✓ Others must be blocked or redirected
- ✓ Lock must last **minutes**, not milliseconds
- ✓ Lock must automatically expire if user disappears

A normal SQL transaction lock lasts only until the transaction ends — usually seconds or milliseconds — NOT enough for:

- payment
- OTP
- user deciding
- entering card details

Distributed locks fill this gap.

3. How Distributed Locks Work (Clear Example)

Let's say we use **Redis** as our distributed lock manager.

You want to lock ticket → `ticket-123`.

Step 1 — Try to set the lock

```
SET ticket-123 locked NX EX 600
```

Meaning:

- `NX` → Only set if NOT exists
- `EX 600` → Expire after 600 seconds (10 minutes)

If successful → you own the lock.

If not → somebody else already locked it.

Step 2 — User proceeds with checkout

User sees “seat reserved for 10 minutes”.

Step 3 — Payment completes → release lock

`DEL ticket-123`

Now another person may buy the seat.

Step 4 — If user crashes or closes browser

After 10 minutes → lock automatically expires.

- ✓ No stuck reservations
 - ✓ No manual cleanup
 - ✓ No human involvement
-



4. Real-World Use Cases (With Explanation)



A. E-Commerce Checkout (High-Demand Items)

Example: Limited edition sneakers drop.

- User A adds sneakers to cart
- System locks the product ID
- Product stays reserved for 10 minutes
- Other users cannot grab it
- If user A pays → lock released
- If not → lock expires → product free again

This avoids:

-
- ✖ Two users checking out the same item
 - ✖ Overselling limited stock
 - ✖ Cart-race conditions
-

★ B. Ride-Sharing (Uber/Ola Driver Assignment)

When a rider requests a trip:

- System finds 5 nearby drivers
- It selects **one driver**, say Driver D5
- It applies a distributed lock on **driver-D5**
- During this period, **no other rider** can be assigned to D5
- If D5 accepts → lock is released after assignment
- If D5 ignores or rejects → lock expires → next rider uses him

Prevents double-booking the same driver.

★ C. Distributed Cron Jobs

Imagine your service is running on 10 servers.

Every night at 12:00 AM → run a **daily report job**.

Without a lock:

- All 10 servers run the job
- You get **10 duplicate reports**
- DB gets hammered

With distributed lock:

- Server #3 acquires lock `daily-report`
- #3 runs the job
- All others skip
- After job completes → lock is released

Only **one** server runs the scheduled job.

★ D. Online Auctions (Last-Second Bidding)

Near auction end, users rapidly place bids.

For item ID `item-777`:

- System locks the item for a short window (2–3 seconds)
- Updates highest bid
- Prevents two users from placing “highest bid” at same millisecond

Prevents:

- ✗ Tie bids
 - ✗ Incorrect highest bidder
 - ✗ Double writes
-

🏛️ 5. Where Are Distributed Locks Stored? (Redis, Zookeeper)

Most common systems:

★ Redis

- Simple
- Very fast

- Supports atomic **SETNX** and expiry
- Used by Shopify, Uber, Stripe

Algorithm: **Redlock**

ZooKeeper

- Very strong consistency
 - Used in Hadoop, Kafka, Kubernetes leader election
 - Ideal for long-lived locks
-

6. Things Interviewers Expect You to Know

1. Lock Expiry (TTL)

If a process crashes while holding a lock → you don't want the resource stuck forever.

Solution:

- Every lock must have an expiry (e.g., 10 minutes)

Perfect for checkout and booking systems.

2. Locking Granularity

You can lock:

- A **single item** (ticket-123)
- An entire **section** (section-A)

- A **driver's availability** (driver-42)
- A **user account** during sensitive updates

Granularity improves correctness and prevents race conditions.

★ 3. Deadlocks (What they are & how to avoid)

A deadlock occurs when:

- Process 1 locks A → waits for B
- Process 2 locks B → waits for A

Both wait forever.

To avoid:

- Always acquire locks in a fixed order
 - Keep lock duration short
 - Use TTL (expiry)
 - Prefer single-lock operations when possible
-

★ 4. Correctness Guarantee

Distributed locks must ensure:

1. **Only one owner** at a time
2. **Lock expiry** prevents permanent lockouts
3. **Lock renewal** (heartbeat) for long tasks

4. Atomic acquisition across multiple nodes



7. Full Example: Ticket Reservation Using Redis Lock

User wants seat **A12**.

Step 1 — Try to lock

```
SET seat-A12 user123 NX EX 600
```

If success → reserved

If failure → someone else already took it

Step 2 — Continue checkout

Show countdown timer (10 minutes).

Step 3 — Payment done

```
DEL seat-A12
```

Step 4 — If user disappears

Lock expires automatically after 10 minutes.



8. Perfect Interview-Ready Explanation

Here's how **YOU** should explain distributed locks in an interview:

"When multiple servers may access the same resource, I use a distributed lock so that only one request can modify the resource at a time."

For example, in a ticketing system, when a user tries to buy a specific seat, I acquire a Redis-based distributed lock on that seat for ~10 minutes. If another request tries to acquire the same lock, it fails, preventing double-booking.

The lock automatically expires to avoid stuck reservations if the user disconnects.

Distributed locks are also useful in ride-sharing assignments, preventing two riders from being matched to the same driver, or in distributed cron jobs to ensure only one machine runs the job.

Systems like Redis (Redlock) and ZooKeeper provide strong consistency and safe lock acquisition.”

This is a **10/10 answer** for any SDE-2/SDE-3 system design round.

10.) Distributed Cache

What is a distributed cache and when should you use it?

This is one of the MOST important system design building blocks, and understanding it deeply will take your level up significantly.

I'll explain **what a distributed cache is, why we use it, real-world examples, cache eviction, cache invalidation, write strategies, and how Redis data structures help you design better systems.**

Let's begin. 

1. What Is a Distributed Cache? (Super Simple Definition)

A **distributed cache** is:

A fast, in-memory storage layer (usually Redis or Memcached) spread across multiple servers, used to speed up reads and reduce load on your database.

Key points:

- Much faster than a database
- Stores frequently accessed or expensive-to-compute data
- Can scale horizontally

- Reduces repeated expensive work
- Helps meet low-latency requirements (like <10ms)

A typical database query might take **40–60 ms**.

A Redis lookup takes **0.3–1 ms**.

That's why caches transform performance dramatically.

2. Why Do We Use a Distributed Cache?

Caches solve **three main performance problems**:

A. Save Expensive Computations (Aggregated Metrics)

Example: An analytics dashboard shows:

- total sales today
- new signups
- active users
- revenue by category

These metrics take:

- heavy DB queries
- aggregation
- scanning large datasets

Instead of computing this **every time**, we:

1. Compute metrics **once per hour**
2. Store in Redis as a key-value map
3. Dashboard reads from Redis → super fast

- ✓ Reduces DB load
 - ✓ Gives instant results
 - ✓ Saves CPU time
-

★ B. Reduce Database Load (User Sessions)

When a user logs in:

- Instead of storing sessions in SQL DB
- Store them in Redis (very fast and cheap)

Why?

- Sessions read/write frequently
- Thousands/millions of concurrent users
- Queries must be instant

Example of session data:

```
key: session:12345
value: { userId: 99, name: "Raj", loginTime: "2025-12-09T10:01" }
```

Every incoming request checks Redis instead of DB → huge performance win.

★ C. Speed Up Expensive Queries

Example: Twitter feed

Query to get posts of everyone you follow requires:

- joining `users` + `follows` + `posts`
- sorting by timestamp
- applying visibility filters

Running this query on SQL **for every request** is too slow.

Instead:

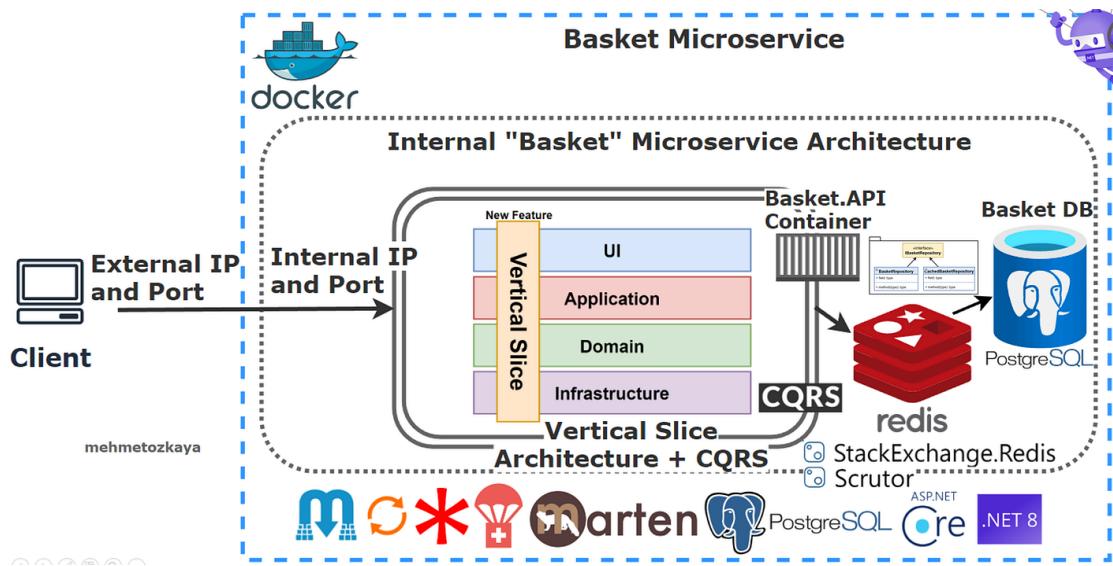
- Precompute timeline
- Store in Redis (usually a `list` or `sorted set`)
- Serve instantly

This turns a **100ms DB query** into a **1ms Redis call**.

This is why Twitter uses Redis heavily for timelines.

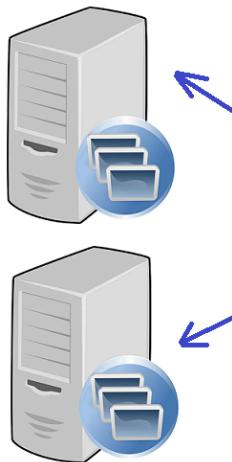
3. Understanding the Diagram (Distributed Cache Architecture)

Here's a conceptual view similar to your diagram:



Distributed Caching

App Server 1



Distributed Cache



App Server 2

Flow:

1. **Read request**

- Check cache
- If found → return instantly
- If not → fetch from DB, store in cache, then return

2. Write request

- Strategies vary (explained below: write-through, write-back etc.)

This pattern is known as **cache-aside** (lazy loading) and it's the most common in system design interviews.

⚡ 4. Things You MUST Know About Caches (Interview Focus)

These are the details interviewers LOVE to hear.

⭐ 1. Eviction Policies (When Cache Gets Full)

Caches have limited memory (RAM), so old items must be removed.

Common policies:

✓ LRU (Least Recently Used)

Evicts items that haven't been accessed in the longest time.

Most popular and intuitive.

✓ LFU (Least Frequently Used)

Evicts items rarely accessed.

Useful for stable workloads.

✓ FIFO (First In, First Out)

Evicts oldest items.

Simplest but often suboptimal.

Mention **LRU** in interviews unless asked otherwise.

★ 2. Cache Invalidation (VERY Important)

This is about keeping cache **consistent** with the database.

Example:

- You cache event details in Ticketmaster
- Venue of event changes
- Cache must be invalidated (removed)
- Next request repopulates it from DB

Cache invalidation is hard because:

"There are only two hard things in computer science: cache invalidation and naming things."
— Phil Karlton

But in interviews, you can mention techniques like:

- delete on write
- TTL (time-based expiry)
- versioning

★ 3. Cache Write Strategies

✓ Write-Through

Write to DB + cache at same time.

Pros: always consistent

Cons: slower writes

✓ Write-Around

Write only to DB; populate cache later.

Pros: avoids cache pollution

Cons: next read is slower (miss)

✓ Write-Back (Write-Behind)

Write only to cache first → asynchronously to DB.

Pros: fastest writes

Cons: risk of data loss if cache fails

Used for analytics or non-critical writes, NOT for payments.

★ 4. Use the Right Redis Data Structure (Key Interview Tip)

Most candidates say:

“I will store events in Redis.”

This is incomplete.

Instead say:

“I will store the top 100 events in a **Redis Sorted Set** sorted by popularity score.”

Examples of Redis data structures:

✓ **Strings → simple values**

✓ **Hashes → objects (`user:123 → {name, email, age}`)**

✓ **Lists → timelines, queues (`LPUSH, LRANGE`)**

✓ **Sets → unique tags, followers**

✓ **Sorted Sets → leaderboards, trending hashtags**

✓ **Bitmaps → daily activity flags**

✓ **HyperLogLog → unique counts (fast approximate counts)**

Mentioning the right structure shows **deep understanding**.



5. When Should You Use a Distributed Cache?

Let's use real-world examples.



A. Instagram Feed Generation

- Store feed in Redis list
 - Query is expensive, but read is super fast
 - Reduces DB load massively
-



B. E-Commerce Product Page

Thousands of users check same product.

Cache product record:

`productId → product details`

Avoids DB overload.



C. Leaderboards (e.g., gaming platforms)

Use Redis sorted sets:

```
ZADD leaderboard score userId  
ZRANGE leaderboard 0 100 WITHSCORES
```

Massive performance benefits.

★ D. Session Storage

Distributed cache stores:

- auth tokens
- user session
- preference data

Used by: Netflix, Amazon, Facebook.

★ E. Popular Search Results

If many users search “iPhone 15 Price”:

- First query hits DB
 - Next 10,000 hit Redis
-

🔧 6. Most Common Distributed Cache Technologies

★ Redis (MOST IMPORTANT)

Supports rich data structures:

- strings
- lists
- hashes

- sets
- sorted sets
- bitmaps
- hyperloglogs
- geospatial

Use Redis in **90% of interview problems.**

Memcached

Simpler key-value RAM store.

Used for basic caching.

7. Interview-Ready Explanation

Say something like:

“To reduce latency and scale the system, I’ll introduce a distributed cache such as Redis.

I’ll store precomputed feed results in a Redis sorted set so I can return user timelines in under 5 ms instead of running expensive DB joins.

I’ll use LRU eviction and a cache-aside strategy with TTL-based invalidation to ensure freshness.”

This shows **depth + correctness + practical thinking.**

8. Final Summary

✓ **Distributed cache = fast in-memory data store (Redis)**

- ✓ Reduces latency and DB load
- ✓ Used for sessions, feeds, expensive metrics
- ✓ Know eviction, invalidation, write strategies
- ✓ Use correct Redis data structures
- ✓ Redis & Memcached are standard

11.) CDN

What is a CDN and when should you use it?

Let's make **CDNs (Content Delivery Networks)** extremely simple, intuitive, and deeply practical with **real-world examples**.

This is one of the easiest yet most impactful system design concepts, and interviewers love when you explain it clearly.



1. What Is a CDN? (Crystal-Clear Definition)

A **CDN (Content Delivery Network)** is:

A globally distributed network of servers that cache and deliver content to users from the location closest to them.

Instead of your server in **Mumbai, Berlin, or US-East** serving all global users:

- A user in Japan gets content from Tokyo CDN Node
- A user in Germany gets content from Frankfurt CDN Node
- A user in USA gets content from California CDN Node

This **reduces latency, improves speed, and reduces load on your servers**.

Why Do CDNs Exist?

Because *distance* → *latency* → *slow websites* → *unhappy users*.

Example:

- If someone in London tries to load a file from a server in San Francisco → latency ~200ms
- But if the file is cached in a CDN node in London → latency ~10ms

That's **20x faster**.



2. How a CDN Works (Step-by-Step)

First request (cache miss)

1. User requests image: `/profile-pics/raj.png`
2. CDN checks local edge server → *not found*
3. CDN fetches from the **origin server** (your backend)
4. CDN caches the image in that region
5. Returns image to user

Subsequent requests (cache hit)

1. Next user in same region requests the same image
2. CDN serves the cached version **instantly**
3. No hit to origin server

4. Zero load on your backend

This massively reduces:

- latency
 - bandwidth usage
 - server cost
 - server load
-



3. When Should You Use a CDN? (Real-World Examples)

CDNs shine in **global systems** or **content-heavy platforms**.



A. Instagram / Facebook / Snapchat (Images & Videos)

If you open Instagram from India and the image is only stored in the US:

- It would load slowly
- Every request hits the origin server
- Latency and cost explode

So Instagram uses CDN to cache:

- profile pictures
- stories
- reels
- thumbnails

- post images

Users get instant loading, regardless of location.

★ B. YouTube (Video Streaming)

YouTube stores videos in multiple CDN edges.

When a user plays a video:

- They do NOT stream from YouTube's main servers
- They stream from the nearest CDN node (Akamai, Google Edge)

This keeps latency low even during:

- World Cup
 - Olympics
 - Music video premieres
 - Live streams
-

★ C. E-Commerce Websites (Flipkart, Amazon)

CDNs cache:

- product images
- category images
- CSS & JS files
- static HTML templates

On Black Friday or Big Billion Days, without a CDN:

- your servers would crash under load
- customers see slow pages
- revenue drops

With CDN:

- most traffic never reaches your backend
 - pages load fast worldwide
-

★ D. News Websites & Blogs

A blog article that updates once per day can be cached via CDN for:

- 1 hour
- 3 hours
- 1 day

This leads to:

- lightning-fast loads
 - massive savings in backend traffic
-

★ E. API Responses

This is less obvious, but very powerful.

If your API returns responses that:

- don't change constantly
- are read-heavy
- need global performance

...you can **cache them in CDN**.

Examples:

- Product catalog
- Location of ATMs
- Trending topics (refreshed every 30 seconds)
- Configuration files for mobile apps

Even Twitter and TikTok use CDN for **precomputed feeds** or metadata.



4. What You **MUST** Know About CDNs (Interview Essentials)



1. CDNs are not just for static content

They can cache:

- API responses
- dynamic pages
- personalized content with edge logic
- GraphQL responses

Using:

- **Edge compute** (Cloudflare Workers, Lambda@Edge)
 - **Edge caching**
 - **Route-specific rules**
-

★ 2. TTL (Time to Live)

CDN content can expire after:

- 5 minutes
- 1 hour
- 1 day

You control this with headers:

`Cache-Control: max-age=3600`

Interview tip:

"I'll set a 5-minute TTL for trending results so they remain fresh but reduce backend traffic."

★ 3. Cache Invalidation

You must be able to remove or refresh CDN content when underlying data changes.

Example:

- You update a blog post
- The CDN still serves old version

- You invalidate the cached version

Methods:

- Purge specific URL
 - Purge entire directory
 - Cache-busting using versioned URLs (`image_v5.png`)
-

★ 4. CDN is a distributed cache, not a full backend

It only caches and delivers content; it doesn't store your original data long-term.

★ 5. CDN + Blob Storage is a standard combination

The most common architecture:

`Client → CDN → S3 (Blob Storage) → Server (if needed)`

CDN caches files from S3, reducing costs & latency.

★ 6. CDNs provide extra features

- DDoS protection

- Web Application Firewall (WAF)
- Bot mitigation
- Rate limiting
- Edge compute functions

Cloudflare and Akamai especially excel here.



5. Popular CDNs

★ Cloudflare (fastest, most modern, global)

Used by: Discord, Shopify, Medium, Coinbase

★ Akamai (oldest, used by enterprise)

Used by: Netflix, Microsoft, Adobe

★ Amazon CloudFront

Best for AWS users.

Other examples: Fastly, Google CDN.



6. Interview-Ready Explanation

Here is a perfect answer to use:

“To reduce latency and serve global users efficiently, I’ll add a CDN like CloudFront or Cloudflare.

The CDN caches static assets such as images, videos, CSS, and even API responses at edge locations near the user.

This ensures users get sub-10ms response times, reduces load on the origin server, and helps the system scale during spikes.

I will set appropriate TTL values and implement cache invalidation when content changes.”

This shows mastery of CDN usage.



7. Final Summary

- ✓ **CDN = globally distributed cache**
- ✓ **Serves content from nearest location → very fast**
- ✓ **Reduces backend load**
- ✓ **Used for images, videos, HTML, CSS, JS, APIs**
- ✓ **Key concepts: TTL, invalidation, edge caching**
- ✓ **Cloudflare, Akamai, CloudFront are top choices**

