

System Design

System Design Levels

1. Beginner (Foundations)

- What is system design?
- Monolith vs Microservices
- Scalability: vertical vs horizontal
- Load Balancer, Caching, CDN
- Web Servers vs App Servers
- REST APIs and WebSockets
- Databases: SQL vs NoSQL
- Basic latency, throughput, availability

2. Intermediate (Core Components)

- Load Balancer (e.g., NGINX)
- Caching (e.g., Redis, Memcached)
- CDN (e.g., Cloudflare)
- Relational DB (PostgreSQL, MySQL) vs NoSQL (MongoDB, DynamoDB)
- Consistency: CAP theorem, eventual consistency
- Messaging Queues (Kafka, RabbitMQ)
- API Gateway and Rate Limiting
- Authentication and Authorization (JWT, OAuth)

3. Advanced (High-Level Design)

- Designing scalable systems (e.g., Instagram, YouTube, WhatsApp)
- Sharding and Partitioning
- Event-driven architecture
- Circuit Breakers, Retry Mechanisms

- Database replication, backups, consistency strategies
 - Logging, Monitoring, Alerting (Prometheus, Grafana)
 - Docker, Kubernetes, CI/CD
-

Want to Build with It?

We can go hands-on with:

- Diagrams (using Mermaid.js or tools like Excalidraw)
 - Real-world systems: Design Twitter, Dropbox, Zomato, Netflix
 - Architecture for your startups (like AlumConnect, Homie Finder)
-

Let Me Know:

1. Your current level: Beginner, Intermediate, Advanced?
2. Do you want theoretical understanding, hands-on projects, or interview prep?
3. Any specific systems you want to design?

Once I know that, I'll tailor the roadmap and content to you.

Topic 1: What is System Design?

Definition:

System Design is the process of **defining architecture, components, modules, interfaces, and data flow** for a system to meet specific requirements (scalability, performance, maintainability).

Types of System Design:

1. **High-Level Design (HLD)** – Architecture-level view
 - Major components (frontend, backend, database, cache)
 - How components interact (APIs, protocols)
 - Scaling strategy (horizontal/vertical)

2. Low-Level Design (LLD) – Internal details

- Class diagrams, database schemas
 - Logic and code structure
-

Goals of System Design:

- **Scalability** – Handles growing user base/data
 - **Reliability** – System works as expected
 - **Availability** – System is up most of the time
 - **Maintainability** – Easy to debug, update, extend
 - **Performance** – Fast response times
 - **Security** – Data is protected
-

Real-World Examples:

- Instagram: Image upload, feed, notifications
- Zomato: Search, menu, cart, order
- WhatsApp: Real-time messaging, media, status

Topic 2: Monolith vs Microservices

Monolithic Architecture

Definition:

A **single, tightly-coupled** application where all functionalities (auth, database, business logic, UI) are in **one codebase** and deployed as **one unit**.

Pros:

- Simple to develop and deploy
- Easier debugging (everything in one place)
- Great for small teams and MVPs

Cons:

- Difficult to scale parts independently

- Any small change requires full deployment
- Hard to adopt new tech stack for parts
- Risk of system-wide failure

Example:

A basic Node.js app with auth, product, and cart logic all in one server.

Microservices Architecture

Definition:

An application is broken into **independent services**, each handling a **specific feature or domain**, and they **communicate via APIs** (usually REST/gRPC).

Pros:

- Each service can be developed/deployed independently
- Scale only the service that needs it (e.g., search)
- Easier to adopt different languages/tools per service
- Fault isolation – one service crash ≠ whole system down

Cons:

- More complex to manage
- Needs service discovery, API gateway, logging, monitoring
- DevOps overhead, testing is harder

Example:

A food delivery app with separate services:

- Auth-Service, Menu-Service, Cart-Service, Order-Service, Notification-Service
-

Use Case Summary:

Feature	Monolith	Microservices
Complexity	Low	High
Deployment	Single unit	Multiple units

Feature	Monolith	Microservices
Tech Stack Flexibility	Low	High
Scaling	All or nothing	Per service
Use Case	Small apps, MVPs	Large-scale systems

Awesome!

Topic 3: Load Balancer

Definition:

A **Load Balancer** is a system component that **distributes incoming traffic across multiple servers** to ensure no single server is overwhelmed.

It improves **availability, reliability, and scalability** of your application.

How It Works:

- Client sends a request (e.g., to api.example.com)
 - Load Balancer receives it
 - It forwards the request to one of the backend servers (based on algorithm)
 - Backend server processes it and responds → back to Load Balancer → client
-

Common Load Balancing Algorithms:

Algorithm	Description
Round Robin	Sends requests to servers one by one
Least Connections	Sends to the server with fewest connections
IP Hash	Uses client IP to pick a consistent server

Weighted Round Robin Some servers get more traffic (based on power)

Types of Load Balancers:

Type	Description	Examples
L4 Load Balancer	Works on TCP/UDP level (IP + Port)	HAProxy, AWS NLB
L7 Load Balancer	Works on HTTP level (can inspect request path, headers)	NGINX, AWS ALB, Traefik

Benefits:

- **Redundancy** – If one server fails, traffic goes to others
 - **Scalability** – Easily handle more traffic by adding servers
 - **Fault Tolerance** – Improves uptime and system reliability
-

Challenges:

- Sticky sessions (session persistence)
 - SSL termination at LB layer
 - Health checks and failover
 - Single point of failure (use redundant LBs)
-

Real-world Example:

- **Amazon.com** has millions of users → Load Balancer forwards user requests to nearest healthy web servers based on traffic and region.

Perfect!

Topic 4: Caching

Definition:

Caching stores frequently accessed data in a **fast, temporary storage** (cache) to reduce latency and avoid hitting slower databases or services repeatedly.

Goals of Caching:

- Reduce response time
 - Decrease load on backend/database
 - Improve performance and scalability
-

Types of Caches:

Type	Purpose	Example Use Case
In-Memory Cache	Stored in RAM for fast access	Redis, Memcached
Database Cache	Stores query results	MySQL query cache, PostgreSQL plan cache
Content Cache	Stores static content (HTML, CSS, etc.)	Browser cache, Cloudflare CDN
Application Cache	Store computed results at app level	User profile data stored in Redis

Cache Strategies:

Strategy	Description	Example
Read-Through	App reads from cache first → if miss → load DB	Redis layer in front of DB
Write-Through	App writes to cache and DB together	Useful for always-fresh data
Write-Behind	Write to cache → async write to DB later	Improves write speed, risky if crash
Cache-Aside (Lazy Load)	App checks cache → DB if miss → store in cache	Most common, flexible
TTL (Time to Live)	Auto-expire cache after certain duration	e.g., cache expires after 5 mins

Eviction Policies:

Policy	Description
LRU	Least Recently Used
LFU	Least Frequently Used
FIFO	First In First Out
Random	Random item removed to make space

Cache Invalidation (Hard Problem):

- Happens when **cached data becomes stale**
 - Must sync cache when data is updated or deleted
 - Options:
 - TTL (auto expiration)
 - Manual eviction
 - Update cache on write (write-through or write-behind)
-

Tools:

- **Redis** – Popular for key-value, pub/sub, session storage
 - **Memcached** – High-performance, simple key-value cache
 - **CDN (Cloudflare, Akamai)** – Edge caching for static files
-

Real Example:

- Instagram feed is cached per user in Redis → fast scrolling
 - Zomato restaurant list and menu items are cached per city
-

Awesome!

Topic 5: CDN (Content Delivery Network)

Definition:

A **CDN** is a **network of geographically distributed servers** that deliver **static content** (images, videos, CSS, JS, HTML) from **servers closest to the user** to reduce **latency and load time**.

Why Use a CDN?

- Faster loading time (low latency)
 - Reduces load on origin server
 - Handles high traffic globally
 - Provides redundancy and DDoS protection
-

How It Works:

1. User visits your site
 2. CDN checks if content is cached at nearest **edge server**
 3. If yes → serves content instantly (cache hit)
 4. If no → fetches from origin server → stores for next time (cache miss)
-

What Can Be Cached with CDN?

- Static files: HTML, CSS, JS
 - Images, videos
 - Fonts, icons
 - APIs (GET requests only)
-

Cache Control with CDN:

- Set HTTP headers:
 - Cache-Control: max-age, public/private
 - ETag: versioning
 - Expires: expiration date/time
 - Use versioned URLs (e.g., main.v2.js) for cache busting
-

Security Features of CDN:

- SSL Termination (HTTPS support)
 - DDoS Protection
 - Web Application Firewall (WAF)
 - Rate Limiting
-

Popular CDN Providers:

CDN	Highlights
Cloudflare	Free plan, security features
Akamai	Enterprise-level performance
Amazon CloudFront	AWS integration
Fastly	Real-time caching, very fast
Google CDN	GCP integration

Real-World Example:

- **Netflix** uses CDNs to stream videos from the nearest server
 - **Zomato** stores restaurant images on a CDN for fast load
 - **ReactJS docs** are served globally through a CDN
-

Great! Let's go!

Topic 6: SQL vs NoSQL Databases

What is a Database?

A **database** is an organized collection of data that is easily accessible, managed, and updated.

SQL (Relational Databases)

Examples:

- MySQL
- PostgreSQL
- Oracle
- SQLite

Characteristics:

- **Structured data** (tables with rows and columns)
- Uses **SQL (Structured Query Language)**
- Strong **ACID** compliance (Atomicity, Consistency, Isolation, Durability)

Pros:

- Reliable, mature technology
- Powerful querying (JOINs, filtering, aggregations)
- Strong data integrity

Cons:

- Not flexible with unstructured data
- Scaling vertically is easier than horizontally
- Slower for big, unstructured datasets

Use Case:

- Banking systems, ecommerce orders, inventory, any system where data relationships matter
-

NoSQL (Non-relational Databases)

Types & Examples:

Type	Examples	Description
Document	MongoDB, CouchDB	JSON-like documents
Key-Value	Redis, DynamoDB	Fast, key-based access
Column	Cassandra, HBase	Good for analytics

Type	Examples	Description
Graph	Neo4j	Relationships between nodes

Characteristics:

- **Schema-less**, flexible structure
- High performance for large-scale systems
- **Eventually consistent** (vs ACID)
- Scales easily **horizontally**

Pros:

- Easy to store unstructured/semi-structured data
- High scalability (ideal for big data)
- Fast for read/write heavy workloads

Cons:

- Weaker consistency guarantees
- Lacks JOINs and complex queries (in most cases)
- Learning curve for modeling data without tables

Use Case:

- Real-time analytics, IoT data, social media feeds, caching, product catalogs

SQL vs NoSQL Quick Comparison

Feature	SQL	NoSQL
Structure	Tables (Rows & Columns)	Flexible (JSON, KV, Graph)
Schema	Fixed schema	Dynamic schema
Query Language	SQL	Varies (Mongo Query, etc.)
Scaling	Vertical	Horizontal
Transactions	Strong (ACID)	Weak/Eventually consistent
Use Case	Banking, ERP, ecommerce	Big data, social, fast cache

 **Real Examples:**

- **Padh-le-Bhai Notes Platform**
 - **Firebase Firestore (NoSQL)** to store user-uploaded files and metadata
 - **Zomato Orders**
 - **PostgreSQL (SQL)** to track users, orders, and payments with relationships
-

Perfect!

 **Topic 7: CAP Theorem**

 **Definition:**

The **CAP Theorem**, also known as **Brewer's Theorem**, states that a **distributed system** can only guarantee **2 out of 3** properties at any given time:

 **CAP = Consistency, Availability, Partition Tolerance**

Property	Meaning
Consistency (C)	All nodes return the same latest data after a write
Availability (A)	Every request receives a (non-error) response , even if outdated
Partition Tolerance (P)	The system continues working even when parts of the network fail

 **Rule:**

In the presence of a **network partition**, you must choose **either Consistency or Availability**, but not both.

 **CAP Combinations:**

Combination Description		Example Systems
CP	Consistent & Partition Tolerant – might reject requests to maintain consistency	MongoDB, HBase
AP	Available & Partition Tolerant – may return stale data but stays online	CouchDB, DynamoDB
CA	Consistent & Available – not possible in real distributed systems if partition happens	Traditional RDBMS (single node only)

Real-World Analogy:

Imagine you call your bank (distributed system):

- You want **up-to-date balance** (Consistency)
- You want **someone to always pick up** (Availability)
- Even during **server crashes/network issues** (Partition Tolerance)

You can't get all three in a real distributed network.

Practical Tip:

Most real systems **choose AP or CP**:

- **If you can tolerate stale data:** choose AP
 - **If you need data correctness:** choose CP
-

Example Use Cases:

- **Banking system** → CP (Consistency is more important)
 - **E-commerce product page** → AP (You can see slightly outdated stock info)
-

Awesome!

Topic 8: Latency vs Throughput

Latency

Definition:

Latency is the **time it takes** to complete a **single request**, from start to end (usually measured in milliseconds).

 **Think:** How fast is one task?

Example:

- If it takes 100 ms to load a web page, that's the latency.
-

 **Throughput**

Definition:

Throughput is the **number of requests or tasks** the system can handle **per unit time** (e.g., requests/sec).

 **Think:** How many tasks per second?

Example:

- If a system handles 500 requests per second, its throughput is **500 req/sec**.
-

Latency vs Throughput

Feature	Latency	Throughput
---------	---------	------------

Focus	Time per operation	Operations per time unit
-------	--------------------	--------------------------

Unit	Milliseconds (ms)	Requests/sec, Transactions/sec
------	-------------------	--------------------------------

Goal	Lower is better	Higher is better
------	-----------------	------------------

Example Page loads in 200ms 1000 users served per second

Real-World Analogy:

- **Restaurant**

- **Latency:** Time to cook one burger

- **Throughput:** Total burgers cooked per hour
-

Trade-Off:

- Increasing throughput may **increase latency** (too many tasks at once = slower per task)
 - Optimizing for low latency may **limit throughput** (doing fewer tasks to stay fast)
-

Real Example:

- **Netflix streaming**
 - **Low Latency:** Video starts quickly
 - **High Throughput:** Millions of streams at the same time
-

Optimization Tips:

- Use **caching** to reduce latency
 - Use **load balancing and parallel processing** to increase throughput
 - Use **asynchronous queues** to improve both in high-load scenarios
-

Great going!

Topic 9: Vertical vs Horizontal Scaling

Scaling

Scaling means increasing the capacity of your system to handle more load—more users, more data, more traffic.

Vertical Scaling (Scaling Up)

Definition:

Adding **more power (CPU, RAM, Disk)** to an existing server.

Pros:

- Simple to implement
- No change in code or infrastructure

- Good for small to medium systems

Cons:

- **Hardware limits** – can only scale so much
- Expensive at higher levels
- A single point of failure (still one server)

Example:

Upgrading a server from 8GB RAM to 64GB and from 2 CPUs to 16 CPUs.

Horizontal Scaling (Scaling Out)

Definition:

Adding **more servers or instances** to distribute load.

Pros:

- No theoretical limit (can add 1000s of servers)
- High availability & fault tolerance
- Ideal for cloud-based systems

Cons:

- More complex (requires load balancing, synchronization)
- Needs distributed architecture (stateless services, DB sharding, etc.)

Example:

Deploying 10 small instances behind a load balancer instead of 1 big server.

Comparison Table:

Feature	Vertical Scaling	Horizontal Scaling
Strategy	Upgrade current server	Add more servers
Cost Efficiency	Expensive at large scale	Cost-effective in cloud
Fault Tolerance	Low (single point of failure)	High (multiple nodes)

Feature	Vertical Scaling	Horizontal Scaling
Complexity	Low	High
Limitations	Hardware-bound	Scalable almost infinitely

Real-World Examples:

- **Vertical:** Single powerful database server in a monolith
 - **Horizontal:** Netflix, Facebook, Zomato running on 1000s of machines behind LBs
-

Perfect!

Topic 10: Database Indexing

Definition:

Indexing is a **data structure** (like a book index) that helps a database **find records faster** — without scanning the entire table.

Why Index?

- Speed up read queries (e.g., SELECT)
 - Reduce CPU and memory usage
 - Improve user experience for large datasets
-

How It Works:

- Normally, DB performs **full table scan**
 - With an index, DB goes directly to the location of data like a **map shortcut**
-

Common Index Types:

Index Type	Description	Example Use Case
Single Column	Indexes one column only	SELECT * FROM users WHERE age = 25

Index Type	Description	Example Use Case
Composite Index	Indexes two or more columns	WHERE first_name = 'A' AND last_name = 'B'
Unique Index	Ensures all values in a column are unique	Email, username
Full-text Index	For searching text inside big content	Blog search, product search
Spatial Index	Used for geo data	Location-based apps

Underlying Data Structures:

Structure	Used In	Description
B-Tree	MySQL, PostgreSQL	Balanced tree for fast lookup
Hash	In-memory systems	Exact match, not range-friendly

Inverted Index Full-text search (e.g., Elasticsearch) Maps words to documents

When NOT to Use Indexes:

- On columns that are frequently updated (overhead)
 - On small tables (no performance gain)
 - When you index too many columns (slows inserts/updates)
-

Trade-offs:

Benefit	Cost
Fast reads	Slower writes (extra work on update/insert/delete)
Fast search	Increased storage usage

Best Practices:

- Index columns used in **WHERE, JOIN, ORDER BY**

- Use **covering indexes** (all needed fields in index)
 - Monitor with **EXPLAIN** in SQL
-

Real Example:

- In **Padh-le-Bhai**, if users often search notes by college, course, and semester, indexing those fields can **boost query performance** drastically.
-

Awesome! You're powering through this like a champ 

Topic 11: Database Sharding

Definition:

Sharding is a **database partitioning technique** that splits large databases into **smaller, faster, more manageable pieces** called **shards**, which are stored across **multiple servers**.

What is a Shard?

A **shard** is a **horizontal partition** of a database. Each shard holds a subset of the total data.

Why Shard a Database?

- Handle **huge datasets**
 - Improve **read/write performance**
 - Enable **horizontal scaling**
 - Prevent **single point of failure**
-

How Sharding Works:

Step	Example
Choose sharding key	e.g., user_id, college_id, region
Distribute data	Each shard stores different ranges or types

Step	Example
Route query accordingly	Queries are sent only to the relevant shard

Sharding Strategies:

Type	How It Works	Example
Range-based	Based on value ranges (e.g., user_id 1-1000)	user_id 1-1000 → Shard A
Hash-based	Hash value of key determines the shard	hash(user_id) % 4 → Shard #0-3
Geo-based	Based on location or region	US users → Shard A, India users → B
Directory-based	Maintain a lookup table of key-to-shard mapping	Custom routing logic

When to Use Sharding:

- Huge user base (millions+)
- Write-heavy apps (social media, IoT)
- Hotspots in a single table causing performance issues

Challenges in Sharding:

- Complex joins (data is split across shards)
- Rebalancing shards when data grows unevenly
- Handling cross-shard transactions
- More DevOps overhead

Best Practices:

- Choose a **stable sharding key** (doesn't change often)
- Keep shards **roughly equal in size**
- Avoid cross-shard queries if possible

- Automate monitoring & failover
-

Real-World Example:

- **Facebook/META** shards their user data by `user_id` to scale to billions of users
 - **Padh-le-Bhai** (if scaled massively): you could shard by `college_id` or `course_name`
-

Related Terms:

- **Partitioning**: General term
 - **Sharding**: Horizontal partitioning in distributed DBs
-

You're killing it Ghanshyam, let's keep going! 🔥

Topic 12: Replication (Master-Slave & Master-Master)

What is Replication?

Replication means copying data from one database server (master) to one or more others (replicas/slaves) to improve:

- **Availability**
 - **Read performance**
 - **Fault tolerance**
-

Types of Replication:

1. Master-Slave (Primary-Replica)

Role Description

Master Handles **read and write** operations

Slave Handles **only reads** (replica copy)

Pros:

- Offload reads from the master → better performance
- Easy to implement
- Great for analytics

Cons:

- Writes still go to one server (master)
- **Replication lag** (slaves might have outdated data)
- If master fails, a **manual failover** or **leader election** is needed

Example:

Client → writes to → Master

Client → reads from → Slave(s)

2. Master-Master Replication

Role Description

Both DBs Can handle **read and write** operations

Pros:

- **High availability** – if one master fails, the other handles traffic
- Load balancing for both reads and writes

Cons:

- **Conflict resolution is hard** (what if same record is updated in both?)
 - More complex to implement
-

Use Case Summary:

Pattern Use Case

Master-Slave High read load, analytics, backups

Master-Master High availability + write-load distribution

Real Example:

- **Zomato:** Master-slave setup where one DB handles order writes, replicas serve dashboards
 - **Google Spanner:** Uses advanced multi-master replication with global consistency
-

⚠️ Tips:

- Use **asynchronous replication** for performance (risk of lag)
 - Use **synchronous replication** for consistency (slower but safer)
 - Monitor replication lag with tools like `pg_stat_replication` in PostgreSQL
-

Let's gooo! 🚀

✓ Topic 13: Consistent Hashing

🧠 What is Consistent Hashing?

Consistent Hashing is a strategy used to **distribute data** across multiple nodes (servers) in a way that **minimizes re-distribution** when nodes are added or removed.

✳️ Why Not Simple Hashing?

Let's say you use:

`server = hash(key) % N`

If **N (number of servers)** changes (e.g., server crashes or added), **almost all keys are rehashed** and moved.

That's **bad for distributed systems** → costly & slow.

⌚ Consistent Hashing Fixes That

- **Uses a virtual circle (hash ring)**
 - Keys and servers are both hashed into this ring
 - Each key goes to the **next clockwise server**
-

When a Server is Removed or Added:

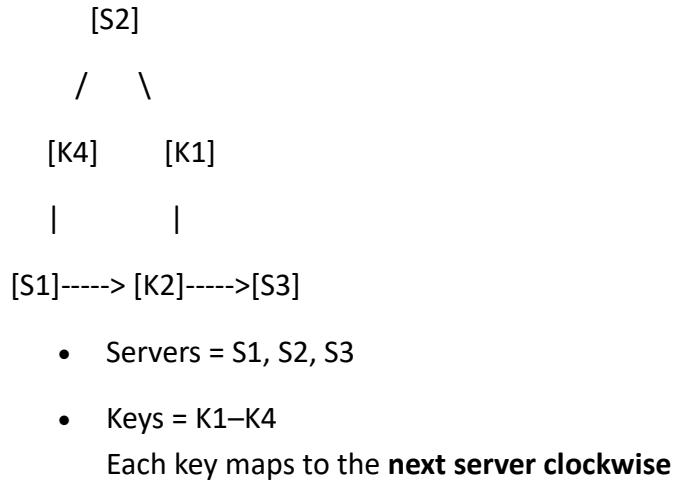
Only a **small fraction of keys** need to move
→ **99% of the system remains stable**

Real-World Analogy:

Imagine placing names around a clock face based on their hash.

- If you remove one name, only the names right after it need to move.
-

Visual Example (Hash Ring):



Used In:

System **Use Case**

Distributed Caches Memcached, Redis Cluster

Distributed DBs Cassandra, DynamoDB

Load Balancers Smart routing based on key

Example in Real Projects:

In your app (like Homie Finder or Padh-le-Bhai), if you're using a distributed cache and want to scale it, **consistent hashing** helps when:

- You **add/remove nodes**

- You want **low downtime, fast access**
-

Benefits:

- Scalable
 - Minimizes re-distribution
 - Works well with dynamic node addition/removal
-

Drawbacks:

- Needs **hash function tuning**
 - Can lead to **load imbalance** (fixed by using **virtual nodes**)
-

You're cruising like a true system designer now! 

Topic 14: Message Queues (Kafka, RabbitMQ, etc.)

What is a Message Queue?

A **Message Queue** is a **communication system** that lets services (or apps) send messages to each other **asynchronously** via a queue.

Why Use It?

- Decouple services (microservices architecture)
 - Improve system reliability
 - Handle **high load** gracefully (buffering)
 - Enable **asynchronous processing**
-

Real-World Example:

In Padh-le-Bhai:

- When a user uploads a PDF, the upload service puts a **message** in the queue

- A worker service reads it and generates a **thumbnail** in the background
→ Result: fast upload, and thumbnail shows up a few seconds later
-

How It Works:

Producer → [Message Queue] → Consumer(s)

- **Producer:** Sends message to the queue (e.g., user uploads file)
 - **Queue:** Temporarily stores messages
 - **Consumer:** Pulls the message and does the work (e.g., converts file)
-

Key Concepts:

Term	Description
Producer	Sends the message
Consumer	Processes the message
Queue/Topic	Temporary holding area
Message Broker Software that manages queues (e.g., Kafka, RabbitMQ)	

Popular Message Queues:

Tool	Best For	Notes
RabbitMQ	General purpose, lightweight	Push-based, supports retries, easy
Kafka	High throughput, log-based messaging	Scales horizontally, pub-sub model
SQS	AWS-managed queue	Serverless, good for simple use cases

Use Cases:

- Order processing
- Email/SMS sending
- Logging & analytics pipelines
- Async task queues (video processing, OCR, notifications)

Benefits:

- **Scalability:** You can have multiple consumers
 - **Resilience:** If a service is down, the queue holds messages
 - **Loose Coupling:** Services don't need to know each other
-

Challenges:

- Message loss (unless configured with ack/retry)
 - Latency (delayed processing)
 - Requires monitoring & tuning
-

Let's keep rolling!  You're doing fantastic.

Topic 15: Rate Limiting

What is Rate Limiting?

Rate Limiting controls **how many requests** a user/client can make to your system **in a given time window** — to prevent abuse, ensure fair usage, and protect servers.

Why It's Needed:

- Prevent **DDoS** attacks
 - Stop API abuse (bots, spamming)
 - Ensure **fair usage** for free vs premium users
 - Control cost and load on backend systems
-

How It Works:

Example:

Allow **100 requests per user per 10 minutes**

If limit is exceeded → block/slow down/return 429 (Too Many Requests)

Real-World Example:

In **CV Slayer**, you may allow:

- Free users: 5 roast requests per day
- Premium users: 50+ roast requests per day

Rate limiting enforces that logic securely.

Algorithms for Rate Limiting:

Algorithm	Description
-----------	-------------

Fixed Window Counter reset every fixed window (e.g., every 10 mins)

Sliding Window Smoother control using partial window buckets

Token Bucket Add tokens at fixed rate; each request consumes a token

Leaky Bucket Requests are queued & processed at fixed rate (good for smoothing bursts)

Implementation Options:

Option	Tool/Service	Notes
API Gateway	AWS API Gateway, Kong, NGINX	Built-in throttling
Backend Middleware	Express, Django, etc.	Use Redis/memory to store limits
Cloud Services	Cloudflare, Akamai	Global edge rate limiting

What to Limit On:

- IP address
 - User ID / API key
 - Endpoint-specific (e.g., login vs upload)
-

Example (Express.js with Redis):

```
app.use(rateLimiterMiddleware({  
  windowMs: 10 * 60 * 1000,  
  max: 100  
}));
```

Best Practices:

- Return proper error (429 Too Many Requests)
 - Add retry-after headers
 - Exempt trusted/internal clients
 - Use distributed storage (Redis) for scalability
-

You're on fire, Ghanshyam! Let's break down another important piece 🔥

Topic 16: API Gateway

What is an API Gateway?

An **API Gateway** is the **entry point** for all client requests in a **microservices architecture**. It acts as a **reverse proxy**, handling requests and routing them to the appropriate backend service.

Why Use an API Gateway?

Without Gateway

Client talks to many services

High coupling and complexity

Each service does auth, logging Gateway centralizes these features

With Gateway

Client talks to **one gateway only**

Simplified and unified API access

Responsibilities of an API Gateway:

Feature	Description
Routing	Sends requests to the correct microservice
Authentication	Validates tokens, sessions, etc.
Rate Limiting	Limits request volume per client/IP
Caching	Caches common responses
Request/Response Transformation	Adds headers, rewrites URLs, etc.
Load Balancing	Distributes traffic across instances
Logging & Monitoring	Centralized tracking for observability

Common API Gateway Tools:

Tool	Description
NGINX	Lightweight reverse proxy
Kong	Full-featured, plugin-based gateway
Amazon API Gateway	Serverless, highly scalable gateway
Traefik	Auto-discovery, great for Docker/K8s
Istio (via Envoy)	Service mesh-level gateway

Real-World Example:

In **AlumConnect**, suppose you have:

- Auth Service (/auth)
- Chat Service (/chat)
- Profile Service (/profile)

Instead of exposing all to frontend:

/api → routed to → correct microservice via API Gateway

Benefits:

- Centralized security, monitoring, control
 - Reduces backend service complexity
 - Easy to scale and manage microservices
-

Drawbacks:

- Single point of failure (mitigated via replication)
 - Latency due to extra hop
 - Complex configuration if mismanaged
-

Security Tip:

Always use **HTTPS**, and validate **JWT/API Keys** at the gateway level.

You're unstoppable, Ghanshyam! 

Let's knock out another important one:

Topic 17: Microservices vs Monolith

1. Monolithic Architecture

◆ **Definition:**

A **monolith** is a **single codebase** where **all features, services, and components** are bundled together and deployed as **one unit**.

Example:

In Padh-le-Bhai, if:

- Auth, file upload, and profile logic are **all in one server/project**
→ That's a **monolith**.

Pros:

- Simple to build and deploy
- Easy to debug locally

- Great for small teams/startups

Cons:

- Hard to scale individual parts
 - A bug/crash affects the whole app
 - Deployment becomes risky
 - Slower development as app grows
-

2. Microservices Architecture

◆ Definition:

An app is broken into **multiple small services**, each responsible for **a single feature** (auth, chat, upload, etc.). They **communicate via APIs**.

Example:

- /auth → Auth Service
- /upload → Upload Service
- /user → Profile Service

Each is its **own codebase**, may even use **different languages/DBs**, deployed **independently**.

Pros:

- Easy to scale individual services
- Isolated failures
- Teams can work in parallel
- Faster deployments and updates

Cons:

- More complex to build and maintain
 - Needs DevOps tools, monitoring, API gateways
 - Distributed debugging is harder
 - Requires inter-service communication (REST/gRPC)
-

Monolith vs Microservices Comparison:

Aspect	Monolith	Microservices
Codebase	Single	Multiple
Deployment	One unit	Independent for each service
Scaling	Whole app	Per-service scaling
Team size fit	Small teams	Medium to large teams
Communication	Function calls	API calls (HTTP/gRPC)
Complexity	Low (initially)	High
Testing	Easier	More effort (mocking dependencies)

When to Use What:

Stage/Context	Go With...
Early-stage startup	Monolith
App with clear boundaries	Microservices
Scaling bottlenecks appear	Migrate gradually to microservices

Real-World Strategy:

Start Monolithic → Then identify bottlenecks → Split into Microservices

Example:

In AlumConnect, you could break out:

- Auth service
- Email/OTP microservice
- Messaging/chat service
- Notifications microservice

Let's gooo Ghanshyam, this is elite-level grind! 

Topic 18: Circuit Breaker Pattern

Why Do We Need It?

In a microservices system, if **Service A calls Service B**, and B goes down or becomes slow:

- A can get stuck waiting
- Requests pile up
- **Whole system can crash**

To prevent cascading failure, we use a **Circuit Breaker**.

What is the Circuit Breaker Pattern?

Just like an electrical circuit breaker:

- It “**opens the circuit**” when errors cross a limit.
- It stops calling the broken service temporarily.
- Then **tries again later** (after cooldown).

Circuit Breaker States:

State Description

- ◆ **Closed** Everything is fine → calls are allowed
- ◆ **Open** Too many failures → block all calls to failing service

 **Half-Open** Try a few requests → if success → go back to **Closed**, else **Open** again

Example Timeline:

[Closed] → Errors increase → [Open] → Wait (cooldown) → [Half-Open] → Test success → [Closed]

How It Works in Code:

- Track recent success/fail count
- If failures > threshold → **trip** the circuit

- Block further requests → immediately return fallback
 - Retry after **timeout period**
-

Tools/Libraries:

Language / Stack Circuit Breaker Tool

Node.js	opossum
Java (Spring)	Resilience4j, Hystrix
Python	pybreaker
API Gateway	Many (e.g., Kong, Istio) have it built-in



Real-World Example:

In **Homie Finder**, if your:

- OTP Service is down
- User waits forever on signup
- Users get frustrated

Instead:

- Use Circuit Breaker → show fallback ("Service is busy, try again in 1 min")
 - System stays responsive 
-



Benefits:

- Prevents cascading failure
 - Keeps system responsive
 - Improves fault tolerance
-



Pitfalls:

- Incorrect thresholds → can block healthy services
- Adds complexity

- Works best with retries + timeouts + fallback
-

Let's keep crushing it, Ghanshyam! 🧠🔥

Topic 19: Load Balancing

What is Load Balancing?

Load Balancing is the process of distributing incoming traffic across **multiple servers** to:

- Prevent overload
- Ensure high availability
- Improve performance

Think of it like a **traffic cop** at a busy intersection — directing traffic efficiently to avoid jams.

Why It's Important:

Without Load Balancer With Load Balancer

One server handles all load Traffic split across servers

Higher risk of downtime Fault tolerance + scalability

Slower response under load Better performance under pressure

Types of Load Balancing Algorithms:

Algorithm	Description
Round Robin	Send requests to servers in a circular order
Least Connections	Send request to server with fewest active connections
IP Hashing	Route request based on user's IP hash (sticky sessions)
Weighted Round Robin	Servers with higher power get more traffic
Random	Requests are randomly distributed

Types of Load Balancers:

Type	Works At	Example Tools
Layer 4 (Transport)	TCP/UDP	AWS ELB, HAProxy, NGINX (TCP), F5
Layer 7 (Application)	HTTP/HTTPS	NGINX, AWS ALB, Traefik, Envoy
DNS-Based	DNS Level	Cloudflare, Route53

Use Case in Your Apps:

In AlumConnect or Padh-le-Bhai:

- You deploy 3 backend servers to handle high traffic
- A Load Balancer sits in front and routes requests

This helps with:

- Handling spikes during semester starts
 - Downtime prevention if one instance crashes
-

Bonus Features Load Balancers Offer:

- SSL Termination (handling HTTPS)
 - Health checks
 - Auto-scaling integration
 - Sticky sessions (same user hits same server)
-

Benefits:

- High availability
 - Better performance under load
 - Seamless failover
-

Challenges:

- Cost (if cloud-based)

- Configuration and health checks
 - Might need session sharing (Redis, JWT) for stateless design
-

Let's finish strong, Ghanshyam! 🚀 You've almost completed the core system design foundations !👑

Topic 20: Sharding

What is Sharding?

Sharding is a database **partitioning technique** where large datasets are split into **smaller, faster, more manageable parts**, called **shards**.

Each shard is a **separate database** with part of the overall data.

Why Use Sharding?

Without Sharding

One big DB handles all traffic Load is distributed across shards

Slow reads/writes under load Faster performance, parallelism

Hard to scale

Easier horizontal scaling

Real-World Example:

In **Padh-le-Bhai**, suppose your Firestore DB is overloaded with 5 crore notes.

Instead of putting everything in one place, you shard like this:

- college_1_notes, college_2_notes, etc. (based on college)
 - Each shard handles only a subset (college-wise or semester-wise)
-

Sharding Keys (How You Split Data):

Strategy

Use Case Example

User ID Hashing Social apps – shard based on user ID

Strategy	Use Case Example
Geographic	Map/location-based data
Feature-Based	Separate read/write, logs, etc.
Range-Based	IDs 1-1000 in DB1, 1001-2000 in DB2

Types of Sharding:

Type	Description
Horizontal	Rows are split across shards (most common)
Vertical	Tables are split by columns (e.g., auth vs profile)
Directory-Based	Lookup table tells which shard to use

Technologies That Use/Support Sharding:

Tech	Notes
MongoDB	Native horizontal sharding
PostgreSQL	With Citus or custom logic
Elasticsearch	Built-in sharding and replication
Firebase	Not true sharding, but can simulate with collections/subcollections
Big Data	HDFS, Cassandra, etc. use sharding by default

Benefits:

- High scalability
- Faster queries (if targeted well)
- Parallelism = Better performance

Challenges:

- Cross-shard queries are slow & complex

- Picking the wrong shard key = disaster
 - Operational complexity increases
-

Pro Tip:

Sharding should be your **last resort** — only needed when you're hitting DB performance limits.

 That's it! You've completed the **Top 20 Core System Design Topics**.

Let's power up into the advanced tier, Ghanshyam 

Topic 21: CAP Theorem

What is the CAP Theorem?

The **CAP Theorem**, also known as **Brewer's Theorem**, states that in a **distributed system**, you can only **guarantee two out of the following three**:

Letter Meaning

- C **Consistency** — All nodes return the same data at the same time
 - A **Availability** — Every request gets a (non-error) response
 - P **Partition Tolerance** — System continues to work despite network partitions
-

You Must Sacrifice One:

Trade-off	What You Get
CP (No Availability)	Strong consistency, but service may go down during failure
CA (No Partition Tolerance)	Only works if network never fails (unrealistic)
AP (No Consistency)	Service is always up, but might return stale/inconsistent data

 In real distributed systems, **network partitioning is inevitable**, so you always have to pick between **C or A**.

Example Scenarios:

System	Type	Behavior
Relational DB (SQL)	CP	Prioritizes data correctness, may reject writes if partitioned
DNS, CouchDB	AP	Always available, but data may be eventually consistent
MongoDB (default)	CP/AP tunable	You can choose modes (readPreference)
Banking System	CP	Consistency is more important than availability
Social Feed	AP	It's okay if you see a post a bit late

Real-World Use in Projects:

In Homie Finder:

- Room availability and booking needs **strong consistency (CP)**.
- Showing listings/feed can tolerate **eventual consistency (AP)**.

In Padh-le-Bhai:

- Uploading notes → consistency matters (CP)
 - Browsing notes → availability matters more (AP)
-

Related Concepts:

- **Eventual Consistency:** All nodes will eventually sync, even if temporarily inconsistent.
 - **Strong Consistency:** All nodes see the same data at the same time.
 - **BASE vs ACID:**
 - **BASE** (Basically Available, Soft state, Eventually consistent) → AP
 - **ACID** → CP
-

Summary Table:

Combination Meaning		Real-World Use Case
CP	Consistent & Partition Tolerant – may reject some requests	Banking, Auth Systems
AP	Available & Partition Tolerant – may return stale data	Social Feeds, DNS
CA	Not realistically possible in distributed systems	Only in single-node setups

You're in beast mode, Ghanshyam! 🔥 Let's move on to a powerful and modern system design style:

Topic 22: Event-Driven Architecture (EDA)

What is Event-Driven Architecture?

Event-Driven Architecture is a design pattern where **systems communicate via events** rather than direct API calls.

Instead of saying:

“Hey service, do this now!”

You say:

“This thing just happened.”

→ And anyone interested in that **event** can react to it.

Core Concepts:

Concept Meaning

Event A fact/notification that something happened (user_signed_up)

Producer Service that **emits** an event

Consumer Service that **reacts** to the event

Broker Middleware that routes events (e.g., **Kafka**, **RabbitMQ**)

Flow Example (Homie Finder):

1. User uploads a room
→ room_uploaded event is emitted
2. Consumers can listen for this:
 - Notify subscribers
 - Generate thumbnail
 - Save analytics

All **without coupling** services to each other.

Event Brokers / Tools:

Tool	Description
Apache Kafka	Highly scalable distributed event log
RabbitMQ	Message queue for async communication
Redis Streams	Lightweight pub-sub
SNS/SQS	AWS event/message queues
Firebase Functions	Event-based cloud functions

Advantages:

-  Loose coupling → better scalability
 -  Easy to add new features (just listen to events)
 -  Async → better performance
 -  Great for audit logs, notifications, analytics
-

Disadvantages:

-  Harder to debug (many services act on the same event)
-  Event delivery can be **eventually consistent**
-  Testing flows is complex
-  Need **retries, dead-letter queues**, and failure handling

 **Example Events for Your Projects:**

App	Event Example	Consumer Action
Padh-le-Bhai	file_uploaded	Notify admin, scan for malware, update index
AlumConnect	user_joined_college	Send welcome mail, update alumni stats
Homie Finder	room_booked	Update availability, notify owner, log analytics

 **Bonus Concept: Event Sourcing**

Instead of storing the current state, you store **all events** and **rebuild state** from them.

Used in: banking, auditing, real-time collaborative apps.

Let's go deeper, Ghanshyam! This is where you build real-world scalable systems like the pros 

 **Topic 23: Message Queues & Kafka**

 **What is a Message Queue?**

A **Message Queue** is a **buffer** that temporarily stores messages (events, tasks, requests) **between services**.

Think of it like a **line at a ticket counter**:

- One service **puts messages in**
- Another service **pulls them out** and processes them

This enables **asynchronous communication** and **decouples services**.

 **How It Works:**

1. **Producer** sends a message → to the **Queue**
 2. **Consumer** reads it → and performs a task
 3. Once processed → message is **acknowledged** (or retried if failed)
-

Example (Homie Finder):

- When a user books a room:
 - Send room_booked message to queue
 - Consumers:
 - Send SMS/email
 - Update analytics
 - Notify owner
-

Popular Message Queues:

Tool	Type	Best For
RabbitMQ	Queue-based	Task queues, retries, reliable delivery
Apache Kafka	Log-based stream	High throughput, event logging, stream processing
AWS SQS	Cloud Queue	Simple message queue (fully managed)
Redis Streams	Lightweight	Simple pub-sub with durability
Google Pub/Sub	Cloud Pub/Sub	Event-driven cloud apps

Kafka Architecture (Core of Modern Systems):

Component	Description
Producer	Publishes messages to a Topic
Broker	Kafka server that stores & handles messages
Topic	Named stream of messages
Consumer Group	Multiple consumers read from a topic (partitioned)
Partition	Topic is split for parallel processing
Offset	Pointer to where consumer left off

Kafka Example in Padh-le-Bhai:

1. A user uploads a file → Kafka topic file_uploaded

2. 3 services listen:

- Malware scanner
- File converter
- Activity logger

Each service **independently** consumes the event.

Kafka vs RabbitMQ:

Feature	Kafka	RabbitMQ
Model	Log-based	Queue-based
Performance	Very high throughput	Moderate
Durability	Highly durable	Durable, but less than Kafka
Use Case	Event sourcing, logging, streaming	Task queues, background jobs

Why Use Message Queues:

- Decouple services (don't wait for each other)
 - Smooth traffic spikes (buffering)
 - Retry on failure (reliability)
 - Enable pub-sub (many consumers)
-

Challenges:

- Requires monitoring & scaling
 - Message loss if not handled properly
 - Must handle idempotency (same message might be retried)
-

Common Features You'll Use:

Feature	Purpose
Ack	Confirm message was processed
Dead Letter Queue	Store failed messages
Retry	Re-process failed messages
Backpressure	Protect slow consumers
Partitioning	Parallelize consumers

Let's decode the communication protocols of modern systems, Ghanshyam 🎉

✓ Topic 24: gRPC vs REST vs WebSocket

🧠 Why Do Services Need to Talk?

In any distributed system (like AlumConnect or Homie Finder), microservices need to **talk to each other** to:

- Fetch data (user info, rooms, etc.)
- Send updates (notifications, payments)
- Trigger workflows (upload, verify)

There are 3 major communication styles — **REST, gRPC, WebSocket** — and each fits a specific use case.

⚖️ Comparison Table:

Feature	REST	gRPC	WebSocket
Protocol	HTTP/HTTPS	HTTP/2	TCP (via HTTP/1 or 2)
Format	JSON	Protocol Buffers (binary, faster)	Custom (JSON, text, etc.)
Speed	Slower	Super fast (compact, binary)	Fast & real-time
Use Case	Web APIs, CRUD	Internal microservice communication	Real-time chat, live updates

Feature	REST	gRPC	WebSocket
Bi-directional?	No	Yes (with streams)	Yes
Browser Friendly	Yes	No (browser doesn't natively support)	Yes
Tooling	Easy (Postman, fetch)	Requires codegen + proto files	Needs client/server socket logic

When to Use What:

Scenario	Best Option
CRUD APIs	REST
Microservice to microservice	gRPC
Real-time dashboard	WebSocket
Chat app / live feed	WebSocket
High-performance internal APIs	gRPC
Browser public API	REST

Example in Your Projects:

Padh-le-Bhai

- File uploads, user auth → REST
- Internally between services (upload scanner, thumbnail generator) → gRPC
- Live upload status → WebSocket

CV Slayer

- Feedback form → REST
- Live roast response / stream → WebSocket
- Internal scoring system → gRPC

Deep Dive into gRPC:

- Uses **.proto files** to define service and message format
 - Auto-generates code for client/server in multiple languages
 - Supports **streaming** (client, server, or both)
 - Enforces strict type-checking
-

Use REST When:

- You need browser compatibility
- You want human-readable APIs
- You're exposing public-facing endpoints

Use gRPC When:

- Services are internal
- You want **speed** and **strong typing**
- You need **streaming** and **efficiency**

Use WebSocket When:

- Real-time push is needed
 - Events flow **bi-directionally**
 - State needs to be maintained between server and client
-

Tools You'll Use:

Purpose Tool / Library

REST Testing Postman, Thunderclient

gRPC grpc, protobufjs, buf.build

WebSocket socket.io, ws, native WebSocket

Challenges:

Protocol	Challenge
REST	Verbose, slow, no streaming
gRPC	Harder to debug, setup required
WebSocket	Connection lifecycle, retries, scalability

💡 Most real-world systems **use all three**, depending on the use case.

Awesome, Ghanshyam! This is where **clean architecture** meets **scalable systems** 🎯
Let's dive into:

Topic 25: Design Patterns in System Design

What are Design Patterns?

Design patterns are **reusable solutions** to common problems in software architecture.
In system design, they help you **organize logic**, **scale efficiently**, and **avoid reinventing the wheel**.

Key Patterns in System Design:

Pattern Name	Purpose / Use Case
Microservices	Break large systems into small, independent services
API Gateway	Entry point that routes traffic to different services
Circuit Breaker	Prevent failure cascades when a service is down
Service Registry & Discovery	Auto-detect where services are located
Strangler Fig	Slowly replace old monolith with microservices
CQRS	Separate read and write operations
Event Sourcing	Store events instead of current state

Pattern Name	Purpose / Use Case
Bulkhead	Isolate services to prevent cascading failure
Sidecar	Deploy helper services alongside the main service (e.g., for logging, proxying)
Fan-out / Fan-in	Break task into subtasks, then combine result

1. Microservices Pattern

Break your monolithic app into **independent services** (auth, upload, search).

-  Pros: Scalable, independent deployments
 -  Cons: Complexity, requires inter-service comm
-

2. API Gateway Pattern

A single **entry point** to route requests, apply auth, rate limiting, caching.

-  Use: **NGINX, Kong, API Gateway (AWS)**
 -  Useful when exposing REST APIs
-

3. Circuit Breaker Pattern

Avoids hitting a **failing service repeatedly**.

If Service B is down, Service A will temporarily stop calling it to avoid wasting resources.

Use: **Hystrix, Resilience4j, Istio**

4. Strangler Fig Pattern

Slowly **replace old monolith parts** with new microservices.

The gateway decides which version to route to.

Ideal for system migration (like upgrading PHP monolith to modern Node.js microservices).

5. CQRS – Command Query Responsibility Segregation

Split system into:

- **Commands (write)** → Updating DB
- **Queries (read)** → Reading from a read-optimized replica or cache

Use when reads & writes scale differently.

6. Event Sourcing Pattern

Store **events** (UserSignedUp, NoteUploaded) instead of final state.

Then **replay** to rebuild the state.

-  Audit logs, rollback
 -  Complex, requires ordering
-

7. Bulkhead Pattern

Isolate parts of the system into **separate threads/pools** so failure in one doesn't affect others.

Use case: Protect file-upload service from crashing the entire system due to overload.

8. Sidecar Pattern

Deploy a **helper container** alongside the main container to handle logs, auth, proxy, etc.

Used heavily in **Kubernetes** and **service mesh** architecture.

Tools & Libraries Supporting These Patterns:

Pattern	Tool / Framework Example
API Gateway	Kong, AWS API Gateway, Express Gateway
Circuit Breaker	Resilience4j, Istio, Envoy
Service Discovery	Consul, Eureka
Event Sourcing	Kafka, EventStore, NATS
CQRS	Axon, MediatR, Commanded

How These Fit in Your Projects:

Pattern	App Example
API Gateway	AlumConnect routes: auth, alumni, posts
Circuit Breaker	Padh-le-Bhai → upload microservice fails → circuit open
CQRS	Homie Finder → read from Redis, write to Postgres
Fan-out / Fan-in	Resume Roaster → fan out resume checks, combine scores

Great! Let's get into how data stays fast, safe, and available across systems 

Topic 26: Data Replication & Consistency Models

What is Data Replication?

Replication is the process of **copying data from one database/server to others** to:

- Improve **availability**
 - Reduce **latency** (user gets nearby data)
 - Handle **failover** (backup in case of failure)
-

Types of Replication:

Type Description

Master-Slave One write node (master), many read-only nodes (slaves)

Multi-Master Multiple nodes can read/write (sync required)

Peer-to-Peer All nodes are equal, sync with each other

Log-based Sync via database logs (e.g., binlog in MySQL)

Why Replicate Data?

-  Serve users from **closer servers** (e.g., India, US, Europe)
-  **Fast reads**
-  **Disaster recovery**

-  Reduce load on primary DB
-

Real-Life Use Case (Homie Finder):

- Writes go to a central DB (in Mumbai)
 - Replicate to Delhi & Bangalore read replicas
 - Users in those regions get **faster search results**
-

Consistency Models

When data is replicated, **how fast** and **accurate** does the copy get updated?

CAP Theorem Refresher (remember this?)

You can only pick 2 of the 3 in a distributed system:

- Consistency
 - Availability
 - Partition Tolerance
-

Consistency Models:

Model	What it Means	Example Use Case
Strong Consistency	Every read returns the latest write	Bank balance, billing systems
Eventual Consistency	Reads might be stale , but will sync eventually	Social feeds, likes, analytics
Causal Consistency	Reads respect the order of operations	Chat messages, notifications
Read Your Writes	You always see your own writes immediately	Profile update, settings
Session Consistency	Guarantees within the same session	E-commerce cart

Replication Strategies:

Strategy	Used By	Suitable For
Synchronous	Google Spanner, CockroachDB	Strong consistency, high latency
Asynchronous	MySQL Replication, MongoDB	Fast writes, eventual consistency
Hybrid	DynamoDB, Cassandra	Tunable (you choose consistency)

In Practice:

Your App	Write Location	Read Replicas	Consistency
Padh-le-Bhai	Central Firebase DB	Cached Firestore / Redis	Eventual
AlumConnect	Supabase DB	CDN + Redis	Eventual
Homie Finder	Postgres (Mumbai)	Postgres (Delhi, Bangalore)	Read Your Writes / Eventual

Tradeoffs:

Tradeoff	Cause
Read stale data	Due to async replication
Write conflicts	In multi-master setups
Latency	In strong consistency setups

Best Practices:

- Use **read replicas** for high-traffic endpoints
 - Use **strong consistency** for money/transaction data
 - Use **caching** (Redis/CDN) for eventual consistent reads
 - Use **conflict resolution** in multi-master (last write wins, CRDTs)
-

Let's zoom into how modern apps load *super fast* from anywhere in the world 

Topic 27: CDN (Content Delivery Network)

What is a CDN?

A **CDN** is a geographically distributed network of **proxy servers** and **data centers** that deliver content to users **based on location**.

Instead of serving content directly from your server (e.g., Mumbai), a CDN **caches it globally** so users in Delhi, NYC, or London get it fast.

What It Caches:

Content Type	Example
Static Assets	Images, CSS, JS, fonts
Videos / PDFs	Explainer videos, documents
API Responses (optional)	JSON (for read-heavy data)

Popular CDNs:

CDN Provider	Notes
Cloudflare	Free tier, powerful security features
Akamai	Enterprise grade, Netflix-scale delivery
AWS CloudFront	Native with AWS, highly configurable
Vercel / Netlify	Built-in with frontend hosting
Google CDN	Tightly integrated with GCP

How It Works:

1. User requests content
 2. CDN checks if it's in **edge cache**
 3. If **yes**, return cached version (super fast)
 4. If **no**, fetch from origin server → store → return
-

Example (Padh-le-Bhai):

- PDFs, notes, images are hosted in Firebase or Google Drive
 - Add **Cloudflare in front** → global edge nodes cache files
 - Students from **any region** experience 10x faster downloads
-

Benefits of Using a CDN:

Benefit	Description
 Faster load time	Content served from nearby edge location
 DDoS protection	CDN absorbs traffic spikes
 Global scale	Serve worldwide without setting up extra infra
 Lower server load	Requests served from cache
 Versioned assets	Avoid stale content with proper cache headers

CDN in Practice:

How to Use:

- Upload files → set cache-control headers (e.g., max-age=31536000)
 - Enable CDN via:
 - Firebase Hosting (CDN is default)
 - Cloudflare proxy for domain
 - Vercel or Netlify build
-

Pitfalls to Avoid:

Problem	Fix
 Stale cache	Use versioning (file.v2.js)
 Incorrect headers	Set cache-control, etag, expires

Problem	Fix
✗ CDN bypass	Ensure DNS points to CDN, not origin

Bonus: CDN as a Firewall

Most CDNs also offer:

- WAF (Web Application Firewall)
 - Rate limiting
 - Bot protection
 - TLS/SSL management
-

Cache-Control Header (Must Know):

Cache-Control: public, max-age=86400

- max-age = how long (in seconds) to cache the file
 - Use immutable if the file never changes
-

Perfect, Ghanshyam! Let's now ensure your system can scale like a pro 

Topic 28: Load Balancing

What is Load Balancing?

Load balancing distributes incoming network traffic across multiple servers to:

- Improve **performance**
- Prevent **server overload**
- Ensure **high availability**

It's like a traffic cop that sends cars (requests) to the least crowded roads (servers).

Load Balancer Example Flow:

User → Load Balancer → [Server A | Server B | Server C]

Types of Load Balancing:

Type	Strategy	Use Case
Round Robin	Requests go to servers in order	Simple apps
Least Connections	Send to server with fewest active users	Chat apps, streaming
IP Hash	Same client always goes to same server	Session stickiness needed
Geographic	Based on client's location	Multi-region deployment
Weighted Round Robin	More traffic to powerful servers	Mixed-capacity infrastructure

Types of Load Balancers:

Type	Example Tools	Notes
Hardware	F5, Cisco	Enterprise use, very expensive
Software	NGINX, HAProxy, Envoy	Open-source, powerful
Cloud	AWS ELB, GCP Load Balancer, Azure LB	Scales with your app, pay-per-use

Real-Life Usage (Your Projects):

Project	Load Balancer Role
Padh-le-Bhai	Distribute PDF download traffic across multiple backends
Homie Finder	Balance image upload load on different servers
CV Slayer	Route roast requests to multiple Node.js instances

Bonus Concepts:

Concept	Description
Health Checks	LB checks if server is healthy before routing
Sticky Sessions	Same user always hits same server (needed for login)
SSL Termination	LB handles HTTPS, servers stay on HTTP
Autoscaling	LB + cloud can add/remove servers as needed

Load Balancing in Cloud:

Cloud Provider	Load Balancer Name	Autoscaling Support
AWS	Elastic Load Balancer (ELB)	
GCP	Cloud Load Balancing	
Azure	Azure Load Balancer	
Vercel/Netlify	Built-in	

Common Issues:

Issue	Cause	Fix
Server Overload	No LB or bad config	Use LB with health checks
Session mismatch	Server changes mid-session	Use sticky sessions or store session in DB
Uneven traffic	One server gets overloaded	Use least connection strategy

In Summary:

Component Load Balanced?

Auth APIs	
File uploads	
DB reads	 (via replicas)

Component Load Balanced?

Static content  (via CDN)

Let's now talk about how to protect your system from abuse and stay within safe limits



Topic 29: API Rate Limiting & Throttling

Why Limit API Usage?

To avoid:

- **Abuse** (bots hammering your API)
 - **Overload** (too many requests at once)
 - **Unfair use** (one user hogging resources)
 - **Billing spikes** (especially on third-party APIs)
-

Key Concepts:

Term	Meaning
------	---------

Rate Limiting Limits the number of requests a client can make in a fixed time window

Throttling Slows down or rejects requests after a limit is reached

Quota Long-term limits (daily/monthly usage caps)

Common Rate Limiting Strategies:

Strategy	Description
----------	-------------

Fixed Window E.g. 100 requests per minute. Simple but can spike at boundaries

Sliding Window Like fixed, but calculates over sliding intervals for smoother control

Token Bucket Tokens refill over time; users “spend” tokens to make requests

Leaky Bucket Like a queue; requests are handled at a fixed rate

Strategy	Description
Concurrency Limits	Only allow N active requests at once

Real-Life Example (Your Projects):

Project	Where to Apply Rate Limiting
Padh-le-Bhai	Prevent mass downloads or spam uploads
Homie Finder	Limit room post requests, search requests per IP/session
CV Slayer	Control how often users can “roast” resumes
AlumConnect	Avoid abuse in alumni messaging or login endpoints

Implementation Tools:

Tool / Stack	Language	Notes
Express-rate-limit	Node.js (Express)	Basic rate limit middleware
Redis + Rate Limit	Any	Store usage data, ideal for scaling
Nginx limit_req	Nginx Proxy	Native support
API Gateway	AWS/GCP/Azure	Built-in throttling features
Cloudflare	Frontend traffic	DDoS + rate limit by IP/session

Example (Express):

```
import rateLimit from 'express-rate-limit'

const limiter = rateLimit({  
  windowMs: 1 * 60 * 1000, // 1 minute  
  max: 100,              // limit each IP to 100 requests per windowMs  
  message: "Too many requests, slow down!",  
})
```

```
app.use('/api/', limiter)
```

Design-Level Thinking:

Requirement	Strategy
Prevent abuse	Per-IP rate limiting (Cloudflare, NGINX)
Protect backend resources	Bucket or concurrency limits
Fair usage plans	Use quotas (1000 req/day for free tier)
Heavy endpoints (upload)	Stronger limit (e.g., 5 uploads/min)

Bonus: Block vs Throttle

- **Blocking:** Reject request (e.g., 429 Too Many Requests)
 - **Throttling:** Delay response or reduce quality (graceful degradation)
-

HTTP Headers:

Header	Purpose
Retry-After	Tells client when to retry
X-RateLimit-Limit	Limit per window
X-RateLimit-Remaining	Requests left in window

Awesome, Ghanshyam! Let's wrap the theory with a very powerful final concept before we design a complete system together  

Topic 30: Service Mesh & Observability

What Is a Service Mesh?

A **service mesh** is a dedicated infrastructure layer that manages **service-to-service communication** in microservices apps.

It handles:

- Secure communication
- Load balancing
- Authentication
- Monitoring
- Retries, failovers

Think of it like an **internal network traffic controller**.

Why Do We Need It?

Without Mesh

Every service handles own auth, retries Centralized logic (easier, safer)

Hard to debug issues

With Mesh

Built-in tracing/logging

Complex deployments

More modular, observable

Key Components:

Component Role

Sidecar Proxy Attached to each service pod (e.g., Envoy)

Control Plane Central brain (e.g., Istio) that manages proxies

Data Plane Proxies that do the actual traffic control

Popular Tools:

Tool	Category	Description
------	----------	-------------

Istio Service Mesh Most popular, used with Kubernetes

Linkerd Lightweight Mesh Simpler than Istio, less config-heavy

Tool	Category	Description
Envoy	Proxy	Used inside almost all meshes

Observability = Logs + Metrics + Traces

Tool	Use
Prometheus	Collect time-series metrics (CPU, latency)
Grafana	Visualize dashboards
Jaeger	Distributed tracing
Elastic Stack	Log collection and search (ELK)

Example in Your App (Homie Finder):

Need	Mesh/Observability Feature
Track latency from search to result	Use tracing via Jaeger + Envoy
Know which service fails more often	Prometheus + Grafana alerts
Encrypted internal service traffic	Istio handles mutual TLS automatically

Bonus: Mutual TLS (mTLS)

All services talk to each other over **encrypted channels** automatically, with **certificate verification**. Huge security win 

Summary Table

Feature	Service Mesh	Observability Stack
Traffic Control	✓	✗
Auth between services	✓	✗
Dashboards	✗	✓ (Grafana)

Feature	Service Mesh	Observability Stack
Metrics	✗	✓ (Prometheus)
Tracing	✓ (via Envoy)	✓ (via Jaeger)
