

Design a photo sharing app like

Instagram

Functional Requirements -

- ① Create Posts with media & caption.
- ② Follow users
- ③ view a feed of friends posts in chronological order.

out of scope

- like comment on posts.
- search for user, hashtags or locations
- User should be able to go live

Non functional Req -

- ① Highly available (Availability >> Consistency)
- ② Low latency feed loading (< 500ms)
- ③ Low latency media rendering (< 200ms)
- ④ Scalable to 500M DAU.

out of scope

- secure & protect user data (authentication, authorization, encryption)
- fault-tolerant & highly reliable (no database)
- analytics on user behaviour & engagement

Core Entity -

- ① User
- ② Post
- ③ Media
- ④ Follows

APIs -

- ① POST /posts → Post Id
 - { "media": { "photo or video type"},
 - "Caption": "my cool photo"

② POST /follows

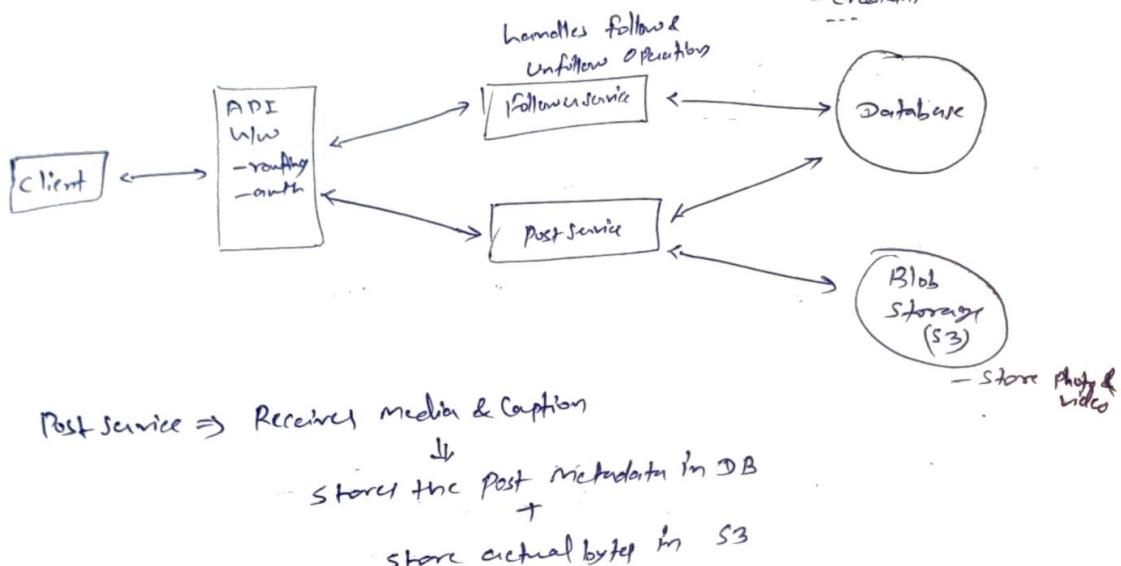
- { "followedId": "123"

- ③ GET /feed ?cursor={cursor} & limit={pageSize}

→ Post[]

HLD

- (1) Create post & follow user -
- (2)
- (3) view feed



For Feed -

- * query 'follow' table to get a list of users that the current user follows.
- * for each of those followed user, query post table.
- * merge & sort all posts.
- * return sorted posts to the client.

Back to indexing, given we're opting for DynamoDB, we'll need to add a few indexes to our tables to avoid full table scans.

(a) For 'follower' table -

$\underbrace{\text{Partition Key} \Rightarrow \text{followerId}}$ $\underbrace{\text{Sort Key} \Rightarrow \text{followerId}}$] \Rightarrow Allows us to efficiently query for all users that a given user follows.

(b) For 'Post' table -

$\underbrace{\text{Partition Key} \Rightarrow \text{userId}}$ (Since most of our queries will be to get the posts of a given user)

$\underbrace{\text{Sort Key} \Rightarrow \text{Composite of (createAt & postId)}}$ \Rightarrow (to ensure chronological ordering while maintaining uniqueness)

II Potential Deep Dive

① The system should deliver feed content with less latency ($< 500\text{ms}$):

→ Problem with Fan-out on Read:-

* Model:- Fetch all post from followed account when user request the feed.

* Issues:- → Thousands of DB reads per request.

→ Some posts fetched multiple times for similar followers.
→ Unpredictable Latency:- Varies by follow count & activity.

* Scale Problem:- → $500\text{M DAV} * 5 \text{ Feed refreshes} = 2.5 \text{B reads/day}$

→ Peak: 150k feed requests/sec

→ Cannot meet 500ms latency target reliably.

* Conclusion:- Doesn't scale; Computational work at read time is inefficient

⇒ Bad solution - Feed Cancellation with simple Cache

* Add a Redis cache layer before querying db.

* Key format: $\text{feedId:userId:cursor} \rightarrow \text{list of Post IDs}$

* Pros: → Faster reads for cache hits.

* Cons: → Cache size must be enormous to be effective.

→ Still performs expensive fan-out reads on cache misses.

* Conclusion: Helps but doesn't solve root problem;
still scales poorly.

⇒ Good solution: Precompute Feeds (Fan-out on write)

* Idea:- Generate feed when post is created, not when it's read.

* Pros:- (i) New post → Get all followers (via LSI on follower table)

(ii) for each follower → Add postID to their Redis feed (ZSET with time stamp)

* Feed Retrieval -

→ Get top N PostIDs from Redis.

→ Hydration options:-

1) Query each PostID from DB.

2) Cache full Post metadata in Redis.

3) Hybrid: Metadata Cache + Latch fallback to DB.

* Conclusion - Fast reads; write-heavy system now. Enables consistent low-latency reads.

⇒ challenges with Fan-out on write -

→ Write Amplification (Specially for celebrities).

→ Millions of writes for one celebrity post.

→ Redis's storage bloat + system overload risk.

⇒ Great Solution : Hybrid Approaches (Precompute + Realtime) :-

* Approach → Fan out on write for regular user (< 100K followers)

→ Fan out on Read for celebrities ($\geq 100K$ followers)

* Feed Retrieval → Get Precomputed parts from Redis.

→ Merge with recent posts from celebrities
(queried in real-time from DB)

* Benefit → Reduces write amplification.

→ Keeps read latency low for most users

* Challenges → Need to tune celebrity threshold (eg. 100K followers)

→ Added system complexity

→ Inconsistency experience for celeb-heavy users.

* Conclusion → Practical trade-off: Prove effective at Instagram scale.

NOTE - Redis Reliability Considerations

- * Redis is in-memory & not inherently durable.
- * Mitigation strategies:
 - 1) High availability → Redis Sentinel for failover.
 - 2) Durability → AOF (Append-only file) or RDB (snapshot) persistence.
 - 3) sharding → Redis cluster to distribute load.
- * Interview tip → Always mention durability & failover plans. (don't treat Redis as magical cache).

(2) Media upload & rendering for large photos/videos :-

we need to → instant rendering of media
[max size → photos up to 8 MB & video up to 4mB]

- * two major challenges
 - Efficient uploads
 - Low latency downloads/rendering

(A) Uploading Large Media :-

Problem :- HTTP requests are limited (< 2mB). need chunked upload.
Solution:- Multipart upload to S3.

* Use AWS S3 Multipart Upload API -

- (i) POST /posts create metadata & return an

Pre-signed S3 URL

- (ii) Client uploads file in 'chunks' to S3 directly.
- (iii) S3 assembles chunks into final file.

* Post Metadata Update -

- (i) media initially marked as pending.

- (ii) Two options to update metadata after upload -



- ↓
- (a) client-driven \Rightarrow Client calls 'PATCH' with 'Object-Key' & 'complete' status.
 (Can't rely much on client)
- (b) server-driven \Rightarrow S3 triggered a 'Lambda' to update metadata.
- (Conclusion \rightarrow server-driven offers better consistency).

B) Serving Media (Rendering) :-

(i) Naive approach - Direct S3 serving -

- * clients download media directly from S3.
- * Problems:-
 - High latency for users far from S3 region
 - No Caching \rightarrow High Cost
 - one-size-fits-all delivery (wasteful on slow/mobile devices)

(ii) Good solⁿ - Using CDN (e.g. CloudFront) -

- * CDN caches media at global edge locations
- * Benefits
 - lower latency
 - Reduced load on S3
 - supports file-type-based Caching & Compression
- * Limitations
 - still serves the same file to all users, regardless of device/network.
 - Cache evictions can cause orphan fetches.

(iii) Current solⁿ \rightarrow CDN + Dynamic Media Optimization -

- * Generate 'multiple media variants' (renditions, formats) on upload
 - use tools like Cloudinary, Imgix.
 - Adaptive Streaming for videos

* Device-aware delivery :-

- client sends device/network info
- CDN serves optimized version (e.g. webP on mobile)

* Advanced Caching:-

- Popular media cached aggressively.
- Proactive cache warming for trending content

* Challenges:-

- more storage (multiple versions)
- media processing - Pipeline adds complexity.

* Conclusion: Best performance & scalability; used by large platforms.

③ System should be scalable to support 500M DAU.

① Hybrid feed generation :- Precompute feeds for regular users, merge celebrity posts at read time. Balances read speed & write load.

② CDN for media delivery :- Use a global CDN (e.g. CloudFront) to cache & serve media close to user. Supports fast, low-latency access worldwide.

③ chunked media uploads - Upload large files (e.g. 4GB) in parts using S3 multipart upload. Improves readability & allows parallel uploads.

④ Scalable Metadata storage - Store metadata in DynamoDB with proper indexing for fast access & horizontal scaling.

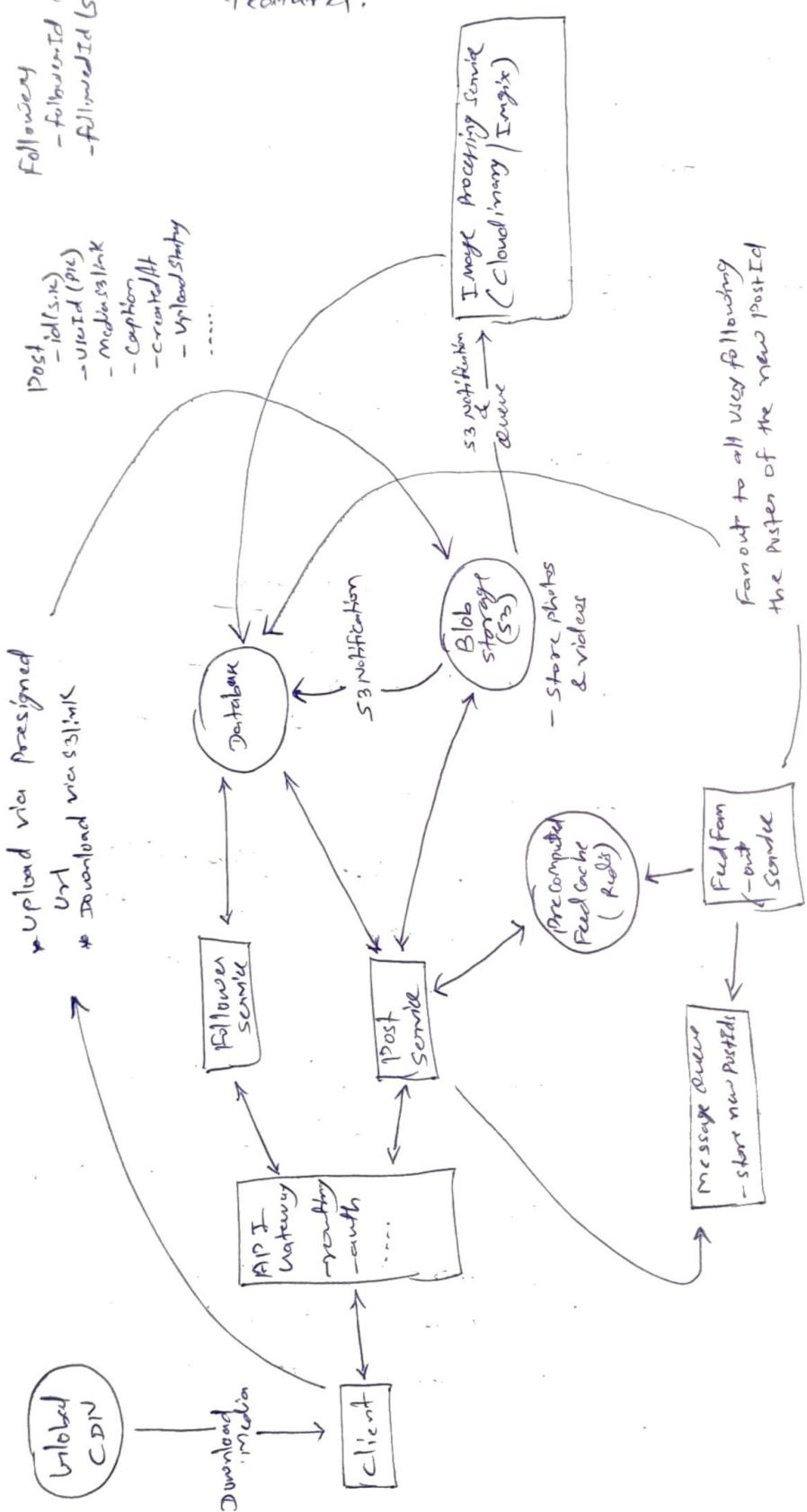
⑤ Smart storage tiering:-

↳ Move rarely accessed media/metadata to cheaper storage (e.g. S3 → Glacier) to reduce cost.

(vi) Auto-scaling Microservices:-

with traffic, using load features.

Microservices scale horizontally
balancers & cloud auto-scaling



II Apache Kafka

① Introduction :-

- * It is a distributed event streaming platform.
- * Can be used as a "message queue" or for "Stream Processing".
- * Provides high performance, durability, & scalability.
- * Handles real-time, high-volume data reliably.

② Kafka in the World Cup Use Case :-

① Producer : Sends events (goals/bookings) to Kafka

↓
Think it as: A reporter who captures every match update & submits it to central Kafka system

② Consumer : Reads events & update website

↓
Think it as: A reader who reads match updates from Kafka & displays them on the live site.

③ Scaling : Needed when games/events increase massively.

* What's the challenge?

- initially, one game = manageable no. of events
- But imagine a 1000 team tournament \Rightarrow with hundreds of matches running at once.
- Now Kafka is receiving a firehose of events - too much for a single broker or consumer to handle.

* Solution : Kafka supports horizontal scaling by:

- Adding more brokers
- splitting data across multiple partitions
- using consumer groups to parallelize processing

(D) Partitioning : Assign events by game → maintain order

- * Partition messages by game id.
- * All events for the same game go to the same partition → hence, they stay in correct order.
- * Think it as: - A private log for each match, so that nothing is jumbled.

(E) Consumer Groups ; Ensure each message is processed only once.

- * What's the issue?
 - ↳ You want more consumers to handle the load, but:-
 - don't want all consumers processing the same messages (unless it's pub/sub).
 - You want "each message to be processed once" (not duplicated).

* Kafka's solution :-

- Use "Consumer groups".
- Kafka assigns each partition to ~~"one consumer"~~ "only one consumer in a group".

~~one~~

(F) Topics : Logically separate streams like soccer / basketball.

* Why topics?

- Now your system grows: You're covering not just soccer, but also basketball, cricket etc
- You don't want soccer updates showing up on basketball dashboards.

* Kafka Solution :-

- Use "topics" to separate data streams:
 - soccer events, basketball events etc.

↳ Producers send to appropriate topic, & Consumers subscribe only to the topics they care about.

* Think of a topic as:- A dedicated channel for each sport.

⇒ Problem first: Why do we need Consumer groups?

Imagine you're running a live world cup site where millions of events (goals, bookings, substitutions) are coming in at the same time.

The challenge: - A single consumer can't keep up - it's too slow.
→ so you decide to add more consumers
→ But if each consumer reads all messages, every update is duplicated - that's not what we want!

* Kafka's answer: Consumer Groups :-

(a) What is a Consumer Group?

- A Consumer group is a set of consumers working together to process messages from a Kafka topic.
- Kafka guarantees:
 - * Each partition is read by only one consumer in the group.
 - * So each message is processed only once, by ~~one~~ one consumer.
 - * Consumers in the same group share the workload, ie load-balanced.

(b) How it works?

Example:- You have a 'soccer-events' topic with 3 partitions (P0, P1, P2).

- You create a consumer group "soccer-updates" with 3 consumers (C1, C2, C3)

Kafka will do this:-

P0 → C1

P1 → C2

P2 → C3

} Now each consumer gets a unique partition = No duplication, full coverage

if one consumer fails (e.g. C2), Kafka will rebalance:-

P0 → C1

P1 → C3 ← reassigned

P2 → C3

② what if we want multiple consumers to get the same message?

That's when you use multiple consumer groups.

- Each group gets a copy of the topic's message.
- This gives you pub-sub behavior.

Example:-

- * analytic group → consumes messages for dashboards
- * alert group → consumes some message & triggers alerts

③ ~~By~~

Kafka Core Terminology & Architecture

① Cluster → Group of brokers (servers).

② Broker → Stores partitions & serves clients

③ Partition → Append-only, ordered log of messages

④ Topic → Logical grouping of partitions.

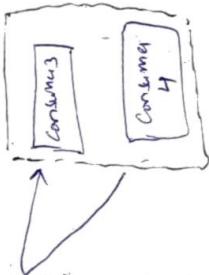
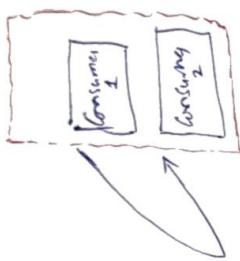
⑤ Producer → Sends data to topics.

⑥ Consumer → Reads from topics.

② Message Queue vs Stream →

- * queue → message is acknowledged Post Processing
- * stream → continuous consumption with acknowledgement.

Consumer Group 1
Subscribed to Topic A



Consumer Group 2
Subscribed to Topic B

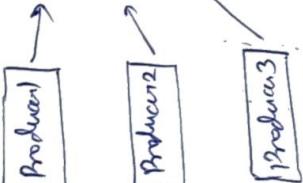
Consumer API



Broker 1
Consumer Cluster



Producer API



④ Kafka Message -

Each message (also called a record) has:-

- * value : the actual data (eg "user clicked on ad")
- * key : decide which partition it goes to.
- * timestamp : when the event happened.
- * headers : metadata (like HTTP headers)

- without a key, Kafka picks a ~~random~~ partition randomly.
- with a key, it hashes it and sends similar keys to the same partition (preserving order).

⑤ How Kafka decides where Data goes:-

- Kafka first "decides the partition" using the message key and a hashing algorithm.
- Then, it uses "cluster metadata" to know "which broker" holds that partition & sends the message there.

⑥ Consumers and offsets in

- Each message in a ~~message~~ Partition has an offset, like a unique message number.
- Consumers Pull (not pushed!) message & commit offsets after processing - this is how Kafka knows where to resume if something crashes.

⑦ Kafka uses pull model (why?)

- Is Kafka Consumers Poll brokers for data, why?
- so they don't get overwhelmed (good for slow consumers)
 - They can control their pace
 - Allows efficient batching & better fault handling.

⑧ When to use Kafka?

① Message Queue :-

- * For asynchronous jobs: Like sending a video link to be processed later.
- * If order matters.
- * If Producers & Consumers need to scale independently.

② Streams :-

- * For real-time processing: Like showing live ad click or FB live comments.
- * When multiple consumers need to ~~read~~ read the same stream.

③ Hot Partitions Problem (or Fixy) :-

Sometimes, one partition gets too much traffic (e.g. popular ad like Nike). That's bad - it can crash that broker.

Fixy:-

- a) Don't use a key: → Random P distribution (but lose ordering)
- b) Salt the key: add random suffix to it (e.g. -ad123-r1, ad123-r2) → Sporadic load
- c) Compound Key: Combine adId + region, or adId + guid group.
- d) Use back pressure:- make the producers wait if the partition is lagging.

⑩ Durability in Kafka :-

ⓐ Replication Mechanism:-

↳ Each Kafka Partition is replicated across multiple brokers.

↳ one broker is the leader, others are followers.

↳ messages are written to leader & ~~then~~ then replicated to followers.

⇒ This ensures that even if a broker fails, the data remains available.

ⓑ Producer acknowledgement :-

↳ [ack = all] : message is acknowledged only after all replicas receive it.

↳ Ensures "maximum durability".

ⓒ Replication factor :-

↳ Defining how many replicas a partition has.

↳ Common value: 3 (1 leader + 2 followers)

↳ If one broker fails, a follower can take over.

⑪ Consumer Fault Tolerance :-

ⓐ Offset management :-

→ Consumers commit offsets after processing messages.

→ on restart, Consumer resumes from the "last committed offset", avoiding data loss or duplication.

ⓑ Rebalancing :-

→ If a consumer in a group fails, Kafka reassigns partitions to remaining consumers.

- ⑥ Offset Commit Timing:-
- ↳ Commit after critical processing (e.g. starting timer)
 - ↳ Prevents data loss in the event of consumer failure.
 - ↳ Minimize consumer's responsibility to reduce retry overhead.

⑫ Handling Failure & Retries:-

ⓐ Producer Retries :-

- * Causes of failure: network issues, breakdown etc.
- * Kafka supports automatic producer retries:-

```
const Producer = Kafka.Producer({
    retry: { maxRetries: 5, initialRetryTime: 100 },
    idempotent: true, // avoid duplication
});
```

ⓑ Consumer Retries :-

- * Kafka does not support consumer retries natively.
- * Common solution:-
 - Move failed messages to a "retry topic".
 - Process them with a separate consumer.
 - After max retry, move them to a "Dead Letter Queue (DLQ)" for manual inspection.

(13) Kafka vs SQS (for retry use cases):-

In systems like Web Crawlers, AWS SQS may be preferred:-

- * Built-in retry & DLQ support.
- * Reduces complexity by eliminating need to manage retry manually.

(14) Performance Boostery :-

- * Batching : Send many messages together.
- * Compression : Use GZIP, Snappy etc. smaller payload = faster processing.
- * Partition wisely : Good partition key = even load = faster processing.

(15) Retention Policy :-

Kafka does not

delete messages immediately after

Consumption :-

- messages are retained based on policy, not consumption start.
- This allows "replaying" messages if needed (e.g., for recovery or new consumers).

(i) Types of Retention :-

(i) Time-based retention (retention.ms) :-

- | |
|---|
| → messages are kept for a fixed time duration (e.g. 7 days). |
| → After this period, Kafka deletes the message, regardless of consumption status. |

(ii) Size-based Retention (retention.bytes) :-

- | |
|--|
| → messages are retained until the total size of data in a topic exceeds a set threshold. |
| → Oldest messages are deleted to stay within the size limit. |

(b) why Customize Retention Policy?

(i) Storage Limitations:-

↳ free up disk space by controlling how long/how much data is stored.

(ii) Business requirements :-

↳ Example: Audit logs may need to be retained for months or years for compliance.

(16) Scaling Kafka :-

Kafka is fast, but don't treat it like a database!

A single broker can handle:-

* ~1TB data

* ~1 million messages/sec. (Ballard)

To scale:-

* Add more brokers.

* Add more partitions.

* Choose partition key smartly to avoid hot spots

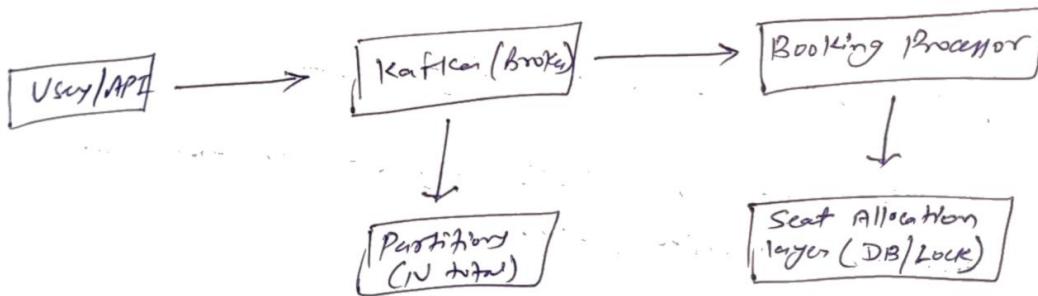
Example - Don't put entire video in Kafka. Put it in S3 and send the video URL via Kafka.

1) Scalable Ticket Booking System with First-Come, First-served Semantics

Goal :-

Design a highly scalable ticket booking system that preserves first-come, first-served (FCFS) semantics while handling "millions of concurrent requests", leveraging Kafka for message ~~partition~~ ingestions.

Key Components :-



How it works?

Step(1) : Request Ingestion at API Layer :-

- * API assigns each incoming request a 'unique', 'time-ordered token' (e.g. UUID or timestamp).
- * This ensures that even across partitions, requests can be ordered based on time.

Step(2) : Write to Kafka Topic with multiple partitions :-

- * Partitioning can be done by event ID (e.g. ConcertID) or user region.

- * Multiple partitions enable 'horizontal scalability'.
- * Kafka ensures ordering 'within each partition'.

Step (3) : Booking Processor layer :-

- * Multiple consumers ~~consume~~ consume from different partitions in parallel.
- * Each request carries its timestamp/token for global ordering.
- * Consumers place incoming requests into a priority queue based on token. (Can be Redis sorted set or in-memory structure).

Step (4) Seat Allocation (Critical Section) :-

- * Booking processor picks requests in timestamp order from the priority queue.
 - * For each request :-
 - Try atomic booking using DB constraint or distributed lock.
- UPDATE Seats
SET status = 'booked'
WHERE seatid=? AND status = 'available';
- only one will succeed, others get a failure response.

Step (5) Response Handling :-

- * successful bookings return confirmation.
- * others receive appropriate error or retry suggestion.

Guarantees & trade-offs :-

<u>Aspect</u>	<u>Technique</u>
(a) ordering :-	Timestamp / VLIID tokens + Priority Queue
(b) Scalability :-	Kafka multi-partition + parallel consumers
(c) Atomicity :-	DB constraint or distributed lock (e.g. Redis' RedLock)
(d) performance :-	In-memory queue for FCFS, no global kafka ordering requirement
(e) Fault-tolerance :-	Kafka retry + idempotent processors.

Summary

- * Don't rely on single Kafka partition - it limits throughput.
- * Use token-based ordering, partitioned ingestion, and atomic seat allocation.
- * Kafka provides durability & scalability; ordering logic is moved to processors.