

Design a Stock Trading Platform

like Robinhood

Robinhood:-

Robinhood is a brokerage and interfaces with external entities that actually manage order filling/cancellation.

Capability: 1) Order Processing = Synchronously places orders & tracks orders via request/response API.

2) Trade Feed : Offer subscribing to a trade feed for symbols. "pushy" data to the client every time a trade occurs, including the symbol, price per share, number of shares & the orderID.

Background : Financial Markets :-

- Symbol :- META, AAPL (Uniquely identify a stock)
- Order :- Buy or sell a stock. Can be market order or limit order.
- Market Order :- Immediate purchase or sale of stock at ~~at~~ current market prices
- Limit Order :- Purchase or sell a stock ~~at~~ stock at a specified price.

Functional Req :-

- ① See live price of stocks.
- ② Manage orders for stock.
(market/limit order, create/cancel order)

out of scope -

- User can trade outside of Market hours
- User can see orderbook in real time.
- Trade ETFs, option, crypt.

Non functional Req:-

- ① High consistency for orders.
- ② Scalable (lots of trades, lots of user)
- ③ Low latency (Price Updates, orders)
- ④ minimize exchange connections

out of scope -

- System connects to multiple exchanges for stock trading.
- System manages trading fees / calculations.
- System enforces daily limit on trading behaviour.
- System protects against bot usage.

Core Entity :-

- ① User
- ② Symbol
- ③ Order

The APIs :-

- ① GET /symbol/:name
Response : Symbol

- ② Post /order

Request : {
position: "buy",
symbol: "META",
PriceInCents: 52210,
numShares: 10}

Response : Order

- ④ GET /orders

Response : Order[]
(paginated)

- ⑤ DELETE /order/:id
Response : {
ok: true

}

#HLD

(1) User can see live price of stocks.

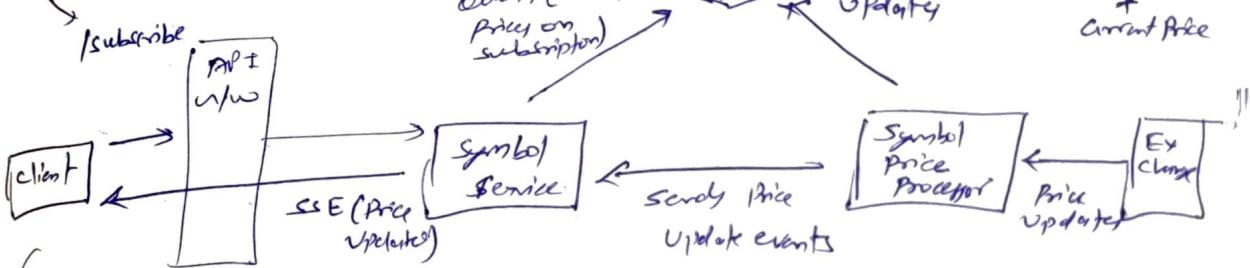
⇒ Bad solⁿ - Polling Exchange Directly

⇒ Good solⁿ - Polling Internal Cache

Server Sent Events (SSE) Price Updates.

(Similar to our Facebook live comment write up).

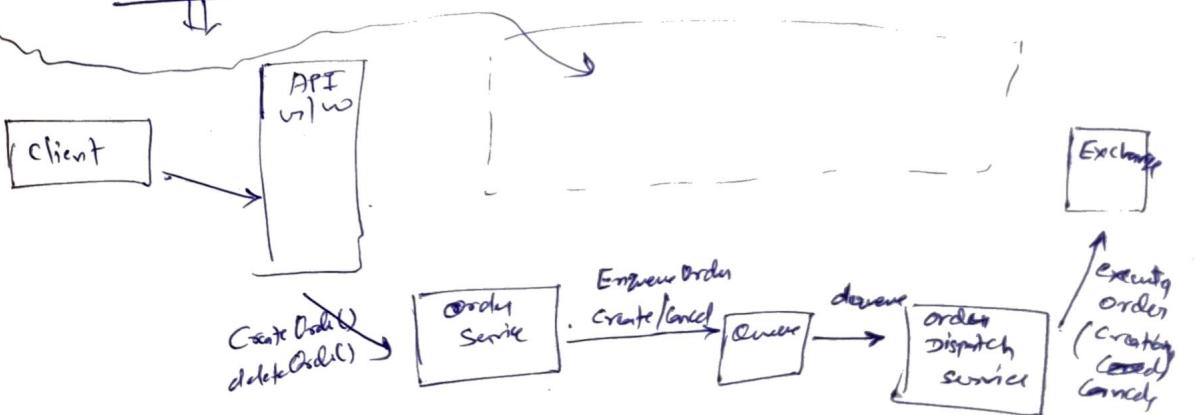
Initial price &
subscribes to feed



(2) User can manage orders for stocks.

⇒ Bad solⁿ - send order directly to exchange.

⇒ Good solⁿ - send order to Dispatch Service via Queue.



⇒ Create * We want SLA (n/msec)

* during very high trading traffic the queue before the order dispatch service has scaled up, order might take a while to be dispatched to exchange, which could violate our goal SLA.

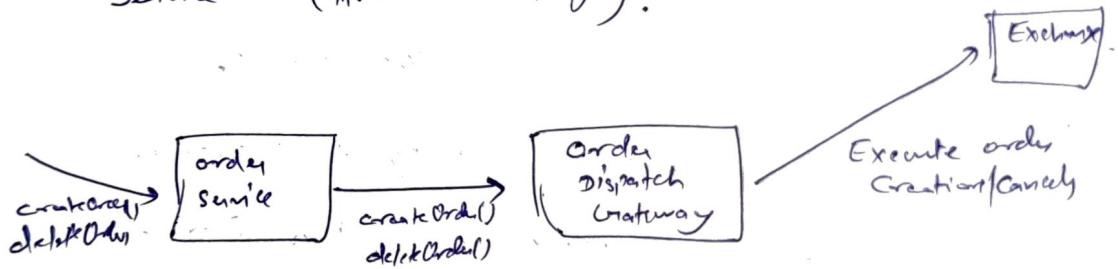
example ↴

+ Imagine a user want to quickly order stock or quickly cancel ~~order~~ an outstanding order. It would be unacceptable for them to be left waiting for our dispatcher to eventually handle their order, or for our service to start more machine up to scale up given increased queue load.

what other option do we have?

⇒ Create soM - Order Gateway

↳ This approach involve sending our orders directly from the order service to an order dispatch gateway. The gateway would enable external internet communication with the exchange via the order service. ("AWS NAT gateway").



⇒ In order to track orders, we can stand up an order database that is updated when orders are created or cancelled.

Potential Deep Dive

- i) How can the system scale up live price updates?
- If many users are subscribing to price updates from several stocks, the system will need to ensure that price updates are successfully propagated to user via SSE.

The main problem we need to solve is:- "How do we route symbol price updates to the symbol service servers connected to users who care about these symbol updates."

→ To enable this functionality in a scalable way, we can leverage Redis Pub/Sub to publish/subscribe to price update. User can subscribe to price update via the symbol service & the symbol service can ensure it is subscribed to symbol price update via Redis for all the symbols the user cares about.

(2) How does the system track order updates?

When digging into the order dispatch flow, it's worth clarifying how we'll ensure our order DB is updated as orders are updated on exchange.

* problem - The trade processor cannot efficiently update the Order DB based on a trade, since it can't look up orders using just "externalOrderID", due to the order table being partitioned by "UUID".

* ② - Introduces a key-value store (e.g. Redis) mapping "externalOrderId" → (orderId, userId).

* who populates it & when -

↳ The Order Service will populate this store synchronously after submitting an order to the exchange.

* why this helps :-

↳ It allows the trade Processor to quickly determine if a trade involves an internal order and efficiently update the correct row in the sharded order table.

③ How does the system manage order consistency?

orders are stored in a horizontally partitioned relational database (e.g. PostgreSQL) by userId, ensuring that all reads & writes happen on the same shard for a given user.

A) Order workflow :-

when an order is created, the system -

- store the order with status Pending to ensure the client intent is recorded before external submission.
- submit the order synchronously to the exchange & receive an externalOrderId.
- insert this externalOrderId into a key-value store & update the order's status to submitted.
- Report to the client confirming the order.

This ensures that if a system failure occurs at any step, we always have a recoverable trace. If submission fails, the order is marked as failed.

- (ii) if ~~subsequent~~ failure occurs after submission but before status update, or background clean-up job
can use the clientOrderId to query the exchange, determine the true state of the order, & reconcile it (record the externalOrderId or mark it as failed)

(B) Cancel workflow -

when an order is canceled :-

- status is first updated to Pending Cancel.
- cancellation is submitted to the exchange.
- The cancellation is recorded in the database.
- The client is notified of success.

Failure at any of these steps is also mitigated through the clean-up job. For example, if cancellation fails or cancellation via the exchange API, updating local state accordingly.

Key Takeaways -

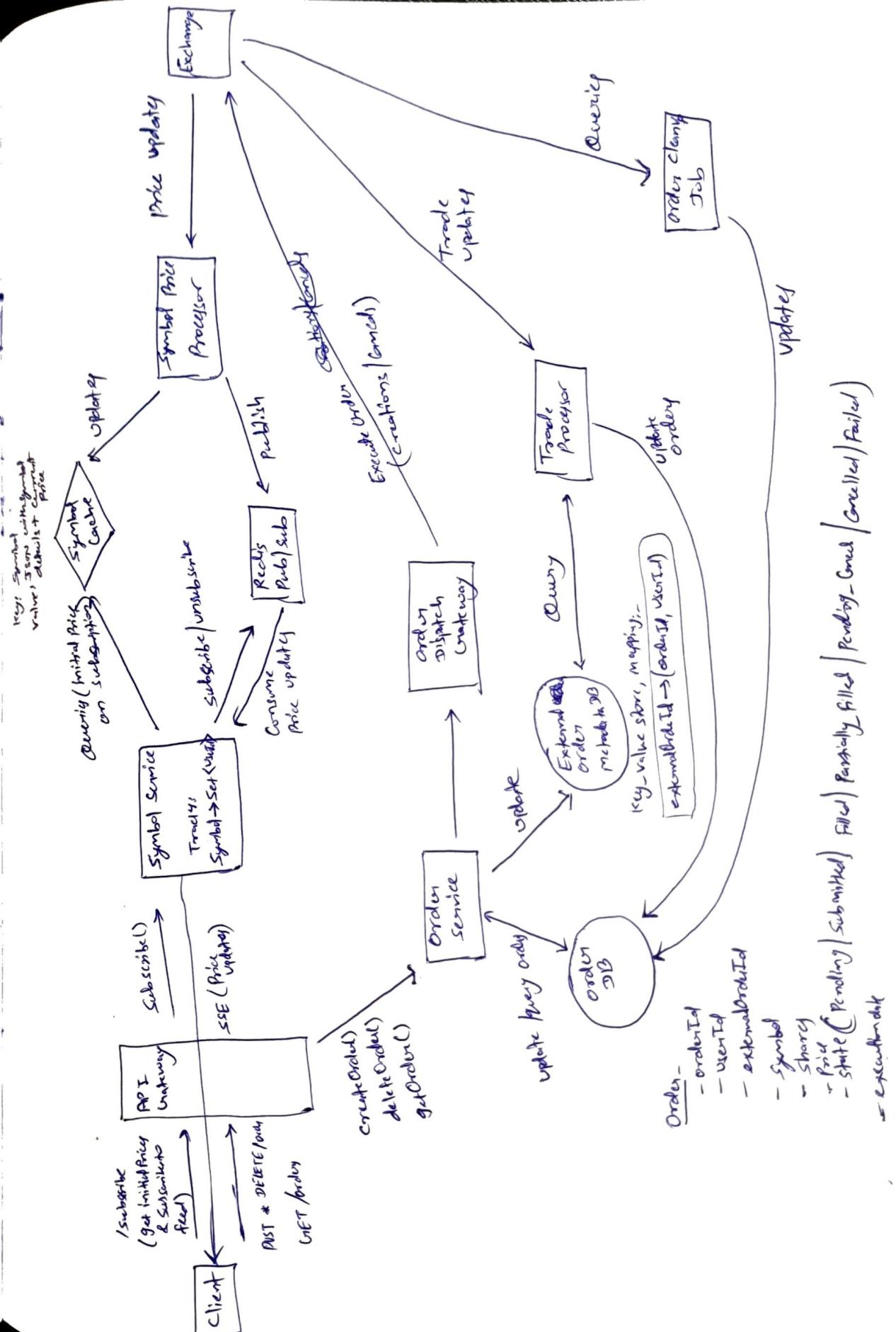
- * Issue - Partial failure during order/cancel flows may cause inconsistency.
- * Solution - work-ahead storage of order intent, use of a key-value store for mapping externalOrderId, and a clean-up job for reconciliation.
- * who populates :- The order service populates the key-value store & updates the database during normal execution.

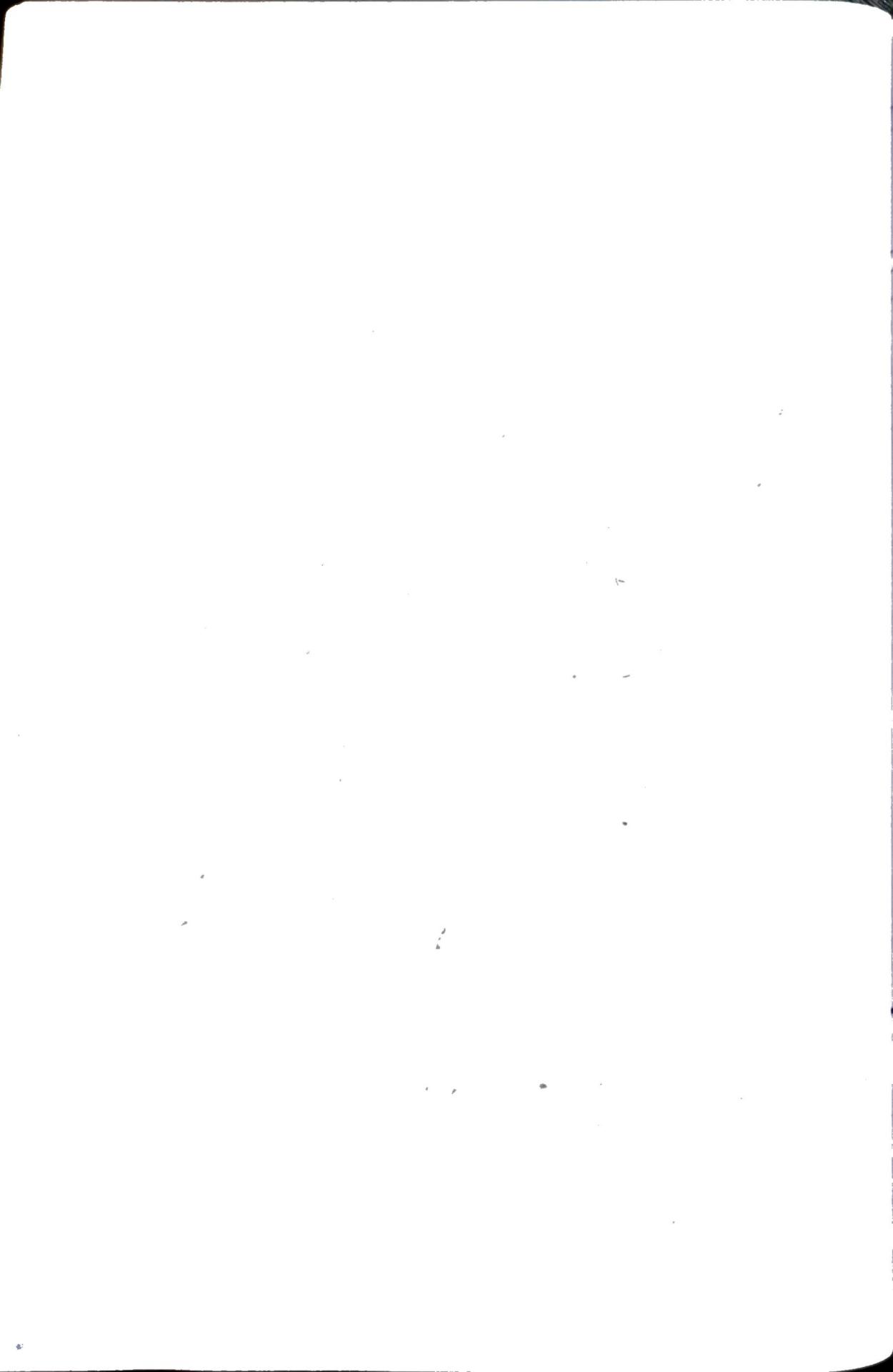
* Benefit - Ensures strong fault-tolerance & eventual consistency, allowing mid-flight failures without losing track of client activity or exchange state.

H2D design [on next page]

⇒ some additional sleep dives you might consider

- ① Excess price updates
- ② Limiting exchange correspondence
- ③ Live order updates
- ④ Historical price / portfolio value data





Design an Ad click Aggregator

Functional Req -

- ① User clicks on ad & get redirected to the advertiser's website.
- ② Advertisers can query click metrics over time w/ 1 minute min granularity.

out of scope

- Ad targeting & serving
- Cross device tracking
- Integration with offline channels

Non functional Req -

(clicks per second)

- ① Scalable to support peak 10K CPS.
- ② Low latency analytics queries < 1 sec.
- ③ Fault tolerance
- ④ Data integrity
- ⑤ As real-time as possible
- ⑥ Idempotency of ad click.

out of scope

- Spam Detection
- Demographic Profiling
- Conversion tracking

Scale

- ① 10M ads at any given time
- ② 10 K ad clicks per second.

HLD :-

- ① User can click on ads & be redirected to the target.

⇒ Browser - Client Side redirect

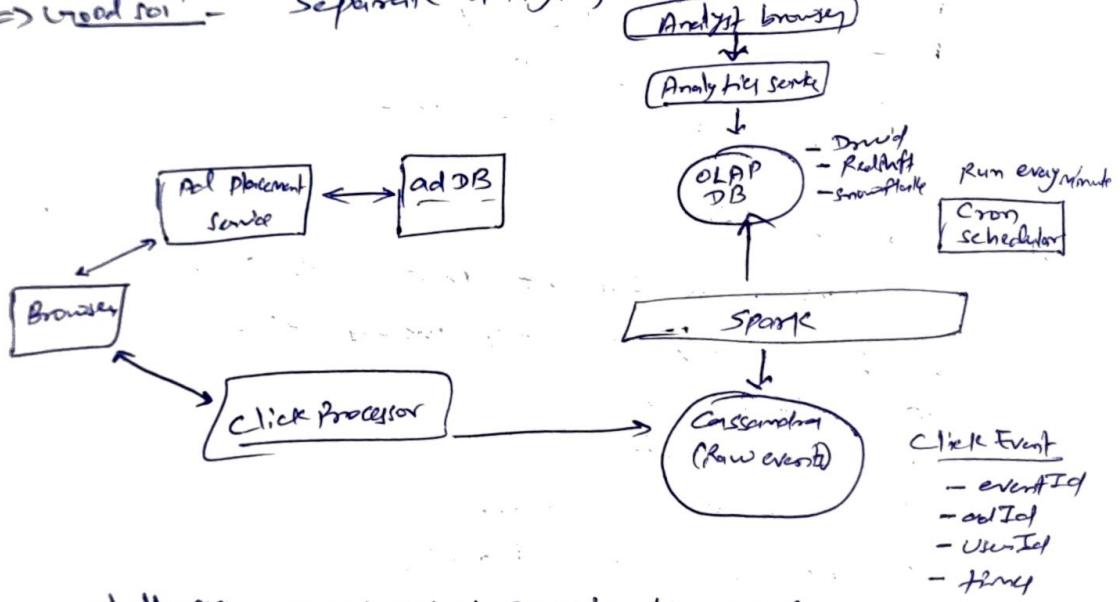
⇒ Server - Server Side redirect

② Advertising on query and click metrics over time at 1 minute intervals.

We'll expand on the click processor path that we introduced above by breaking some options for how a click is processed & stored.

③ Bad soln - store & query from the same database.

④ \Rightarrow Good soln - separate analytics db with batch processing



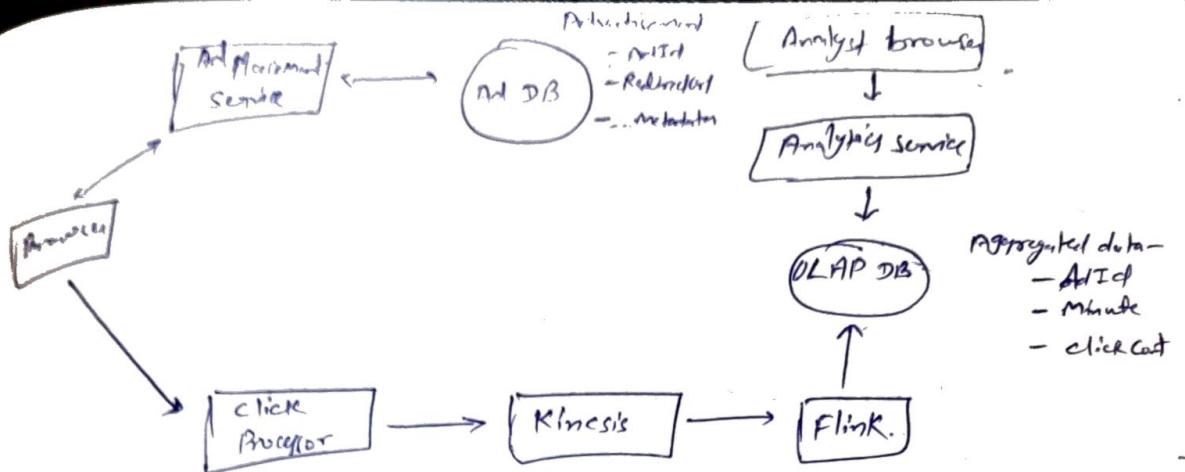
Challenges \rightarrow * batched Processing job introduce a significant delay in the data pipeline.

* not able to handle load, if we have a sudden spike in clicks.

↓
Introduce queue b/w Click Processor & event store could help with scalability issue, it would not solve the latency problem.

⑤ Current soln - Real-time Analytics with Stream Processing





Now when a click comes in:-

- The click processor service write events to a stream like Kafka or Kinesis.
 - A stream processor like Flink or spark streaming reads the events from the stream & aggregates them in real-time.
 - The aggregated data is stored in our OLAP database for querying.
 - Advertisers can query the OLAP database to get metrics on their ads in near real-time.
- ⇒ It's much more realistic for us to decrease the Flink aggregation window than it is for us to increase the spark job frequency given the overhead of running spark jobs.

This means we could aggregate clicks every couple of seconds with this streaming approach to decrease latency, while we'd ~~not~~ be hard-pressed to run spark jobs every couple of seconds.

Even better, Flink has flush intervals that can be configured, so we can aggregate on a minute boundary but flush the results every couple of seconds. This way, we get the best of both worlds.

Potential Deep Dive

(1) How can we scale to support 10K clicks per second?

(a) Click Processor Service -

* Bottleneck → Processing incoming clicks.

* Solution → Scale horizontally with multiple instances behind a load balancer, use cloud auto-scaling.

(b) Stream (e.g. Kafka / Kinesis) -

* Bottleneck : Limited throughput per shard (e.g. Kinesis: 1MB/s or 1000 records/s per shard).

* Solution - Add more shards & shard by AdId to allow parallel processing.

(c) Stream Processor (e.g. Apache Flink) -

* Bottleneck : Limited parallel processing capability.

* Solution - Scale Flink horizontally with more tasks; assign one job per shard for aggregation. We'll have a separate Flink job reading from each shard doing the aggregation for the AdIds in that shard.

(d) OLAP database -

* Bottleneck : Query performance at scale.

* Solution : Scale horizontally, shard by Advertisers for efficient querying across all their ads.

(e) Hot Shards (skewed traffic problem) -

* Problem :- Popular ads (e.g. Nike x LeBron) overload specific shards.

* Solution - Add randomness to Partition Keys for hot ads, e.g. use Advertiser to distribute load across multiple shards.

② How can we ensure that we don't lose any click data?

* we are already using stream like Kafka/Pmem's
↳ distributed, fault-tolerant,
highly available

* we can configure a retention period of 7 days, for example, so that if, for some reason, our stream processor goes down, it will come back up & can read the data that it lost from the stream again.

* skip checkpointing for small windows - for small aggregation windows (e.g., 1min), Flink can reprocess from stream without needing checkpoint recovery.

* Periodic Reconciliation -

- * store raw events in S3.
- * run batch jobs to re-aggregate & compare with stream results for accuracy

⇒ combines real-time + batch processing for reliability & correctness

③ How can we prevent a user from clicking on ads multiple times?

⇒ Buckets - add userId to click event payload. (Copy All very need to login.)

⇒ Impression - generate a unique impression id.

④ How can we ensure that advertisers can query metrics at low latency?

In Real time

- In real time stream processing, data is already aggregated & stored in the OLAP database making the query fast.

- where this query can still be slow is when we are aggregating over larger time windows, like a days, weeks or even years. In this case, we can pre-aggregate the data in the OLAP database. This can be done by creating a new table that stores the aggregated data at a higher level of granularity, like daily or weekly.

↓
This can be done via a nightly cronjob that runs a query to aggregate the data & store it in the new table.

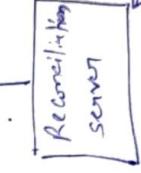
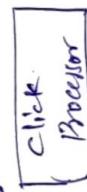
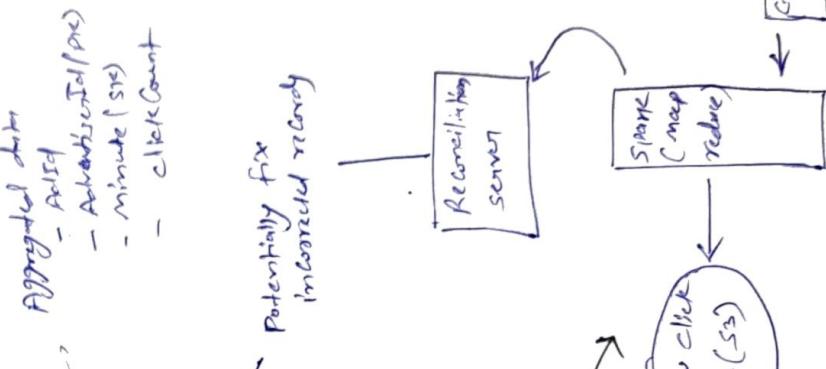
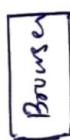
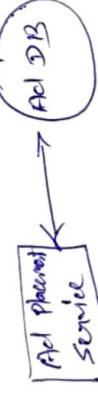
when an advertiser queries the data, they can query the pre-aggregated table for higher level of granularity & then drill down to the lower level of granularity if needed.

Advertisement :-

- Ad ID
- Random ID
- Impression
- 2. Store in Cache.

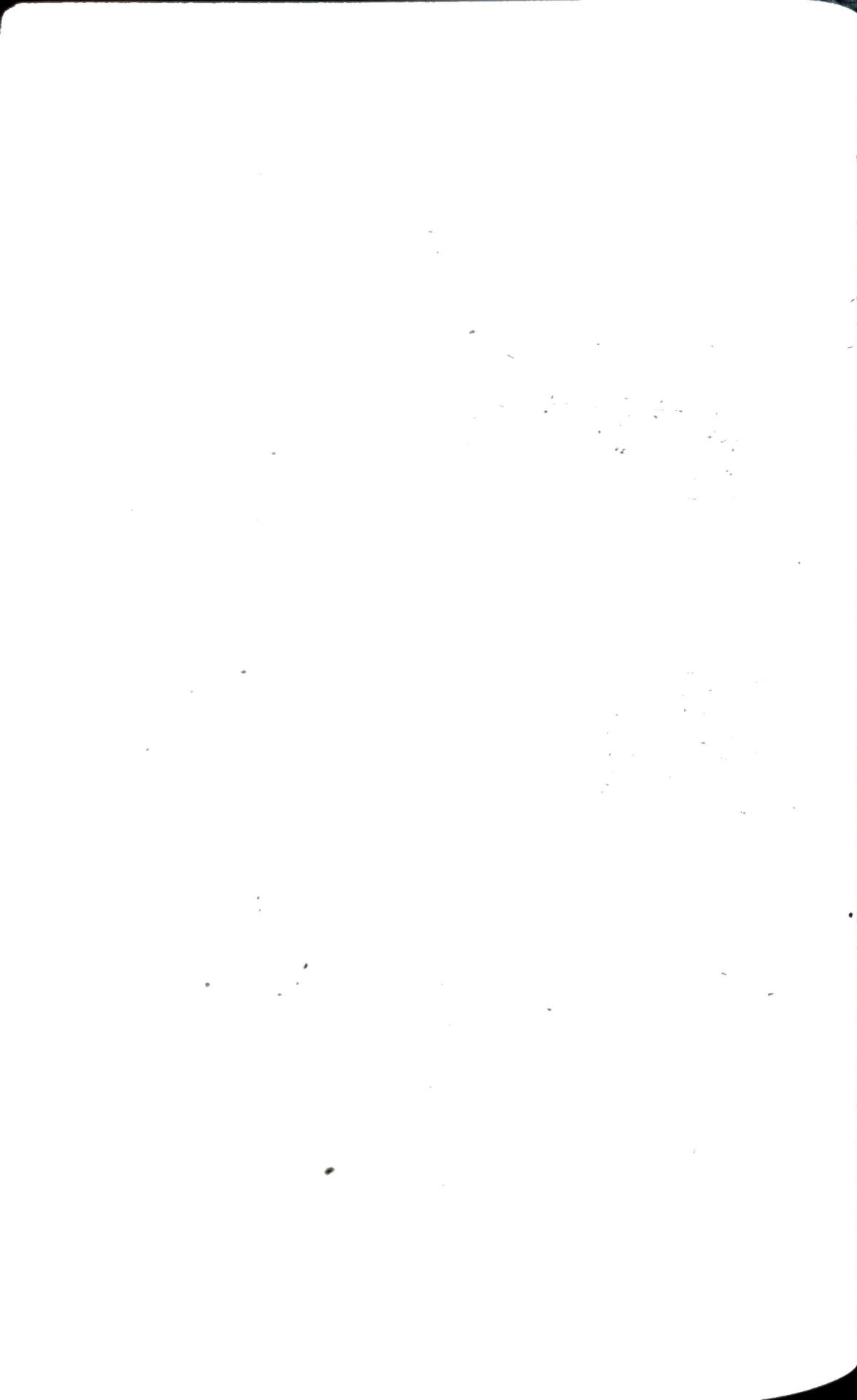
- generate ad impression key id as idempotency key

2. store in Cache.



Both Kinesis & Flink partition by AdID
with celebrity solution for hot
shards

* 7 day retention period on Kinesis



Design a Collaborative Document Editor like Google Docs

Functional Req

- (1) User can create new documents.
- (2) User can edit documents.
- (3) User can get real-time updates.
- (4) User can see each other's cursors.

Out of Scope

- Assuming simple text editor (no sophisticated document structure).
- Permissions & collaboration levels (e.g. who has access to document).
- Document history & versioning.

Non-functional Req

- (1) Documents are eventually consistent.
- (2) Updates are low latency (100ms).
- (3) Millions of users, billions of documents.
- (4) No more than 100 concurrent users.
- (5) Durable documents.

Core Entities -

- (1) Editor : A user editing a document.
- (2) Document : A collection of text managed by an editor.
- (3) Edit : A change made to the document by an editor.
- (4) Cursor : The position of editor's cursor in the document.

APIs :-

- (1) POST /docs
 - { title: String }
 - } → ↗ docId
 - 3

② ws /docs/docId?

```
SEND {  
    type: "insert"  
}
```

```
SEND {  
    type: "updateCursor",  
    position: ...  
}
```

```
SEND {  
    "type": "delete",  
}  
RECV {  
    type: "update"  
}
```

~~# Hold~~

① User should be able to Create new documents.

API → Create new document, take a title & docId.



② Multiple User should be able to edit the same document Concurrently.

we have Consistency & scaling issues regarding this
will see in deep-dive

Initial problem - sending snapshots (wrong approach)

↳ on every edit, entire document is sent to a central document service & stored (e.g. In33).

Issue ↳ Inefficient

Inconsistent → if two users update simultaneously, final document depends on which updates arrive last - causing loss of changes from the other user.

Improvement - Sending Edits (Better, but not enough)

- How it works ⇒ send only specific edit (INSERT, DELETE)
- Benefit ⇒ solve inefficiency
- New problem ⇒ order of operations matter

Collaborative Editing Approach

(a) Operational Transformation (OT)

↳ An approach where each user's edit is transformed based on prior edits to maintain consistency in collaborative editing. For example, if one user inserts text before another's deletion, the delete position is adjusted accordingly. ~~OT is~~

- OT is fast & memory-efficient but requires a central server for ordering operations & is complex to implement correctly, limiting scalability.

(b) Conflict-Free Replicated Data Types (~~CRDT~~ CRDT)

↳ CRDTs allow edits to be applied in any order & still reach the same final result.

- They use real numbers to represent text positions & keep deleted characters as hidden "tombstones". This

this avoids the need for a central server & works well offline or in peer-to-peer setups.

- However, CRDTs can use more memory & may be slower due to storing all edits & handling conflicts

~~get editing~~ ~~the design~~

- * we'll use OT for our design, like Google Docs, leveraging a central server for low-memory, fast friendly collaboration.
- * For large scale or peer-to-peer scenarios, CRDTs could be an alternative, though this design focuses on OT.
- * Edits are sent as operators to the servers, transformed in order, & then served.
- * A "Cassandra-backed, append only document operations DB" ensures fast writes & durability, partitioned by documentId & ordered by timestamp.

(3) User should be able to view each other's changes in real-time

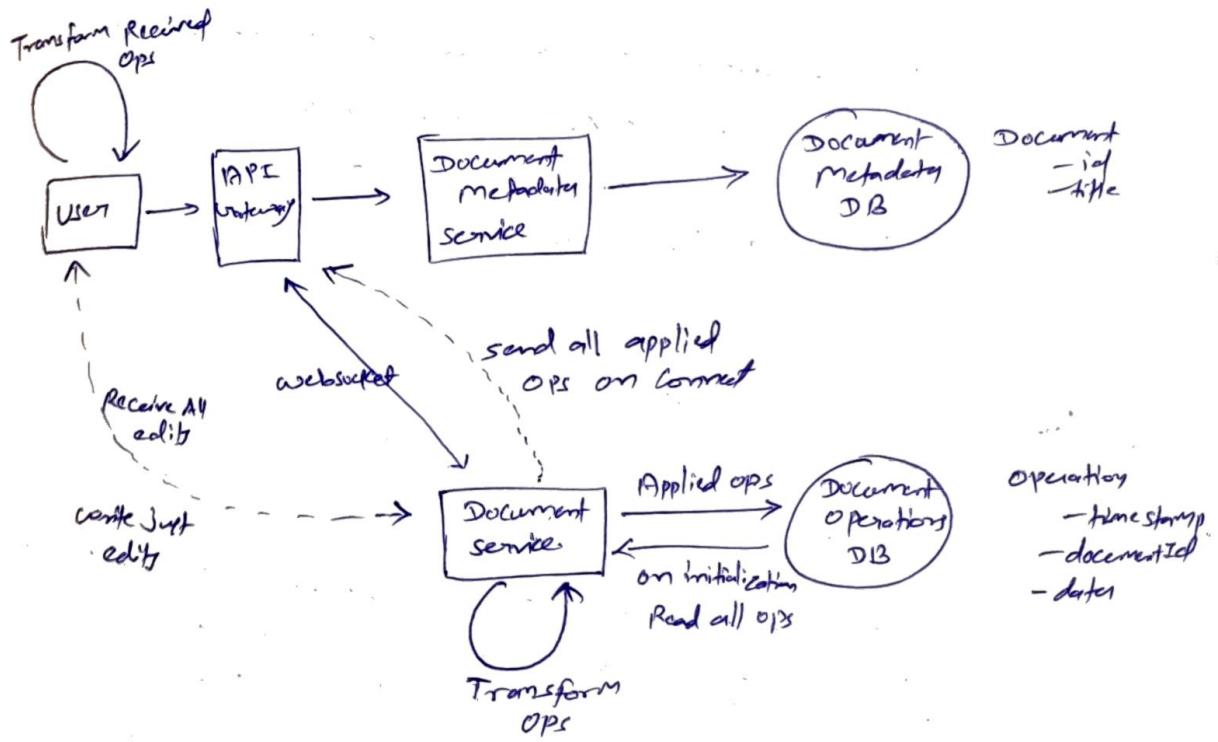
To allow users to see each other's changes in real-time you need to manage two key aspects of reading data:

(a) Initial document loading

↳ when a user first connects or reconnects to a document, the system pushes all prior operations from the database to their "websocket connection". This ensures everyone starts with the latest document version!

⑥ Real Time Updates & Consistency -

- when a collaborator makes an edit, the server records the operation in the database & then immediately sends it to all other connected clients for that document.
- client also performs operational transform (OT) on incoming updates. This is crucial because users expect their own changes to appear instantly.
- OT handles out-of-order edits, ensuring that no matter the sequence in which operations are received by different clients or the server, everyone's document ultimately becomes consistent.
- This means even if User A's edit arrives before User B's on the server, & User B's local edit was applied first on their end, OT guarantees that both users will see the exact same final document state.



④ User should be able to see the cursor position & presence of other users.

↳ This feature is crucial for avoiding simultaneous edits to the same area & prevented duplicated work.

Given that cursor position & user presence are ephemeral (only relevant in ~~editing~~ the moment & while connected), this data isn't stored within document itself. Instead—

- user report changes to their cursor position to document service.
- Document service stores this info 'in-memory' & immediately broadcasts it to all other connected users via the same websocket connection.
- when a new user connects, document service retrieves current cursor position & presence of other user from memory & send them to the new client.
- The document service also monitors disconnections: when a user leaves, their info is removed from memory & a broadcast is sent to remaining user to reflect their departure.

Potential Deep Dive

① How do we scale to millions of websocket connections?

— we need to scale number of 'Document Service' serving to the number of concurrent connections: when I connect to 'DocumentService', I either (a) need to be collocated with all other user connected to some document or (b) know where all the other users are

connected. (see whatsapp for a similar problem).

To ensure millions of concurrent users & to ensure high availability, the core strategy is to "horizontally scale the Document service", which handles real-time websocket connections.

Key points includes:-

- * websocket Edge Termination - websockets are managed at the Document service "edge" to maintain statelessness in other backend services.
- * consistent hashing for distribution :- A consistent hashing ring, coordinated by Apache Zookeeper, is used to distribute document connections. This ensures all collaborators on a single document are routed to and co-located on the "same Document Service server".
- * client connection flow - Clients connect to an initial server, are ~~re~~ redirected if necessary to the correct server based on the hash ring & then establish a websocket connection. This co-location enables efficient "broadcasting" of real-time updates directly to all connected users for that document.
- * scalability challenges - ~~with~~ significant "state migration" (connection, document operations) is required when servers are added/removed.
 - ↳ This necessitates robust handling of server failure, client reconnection logic & careful capacity planning.

(2) How do we keep storage under control?

To manage storage effectively with potentially billions of documents (assume each document is $50KB \Rightarrow$ total $50TB$) & Counter operations, & to keep active document memory usage low, periodic snapshotting/compaction of operations is crucial. This reduces stored operations & data transferred to clients.

two approaches -

(a) Good approach - offline snapshot / compact (with compaction service)

* **Approach** → A dedicated compaction service periodically reads operations from the document operations DB, compacts them (e.g. into a single insert) & writes them to a new documentVersioned in the DB.

* **Atomicity** → To ensure consistency for live documents, the compaction service tells the Document Service to "flip" the documentVersioned. If the document is currently active/loaded in the Document Service, this flip command is deferred to avoid corrupting in-flight operations.

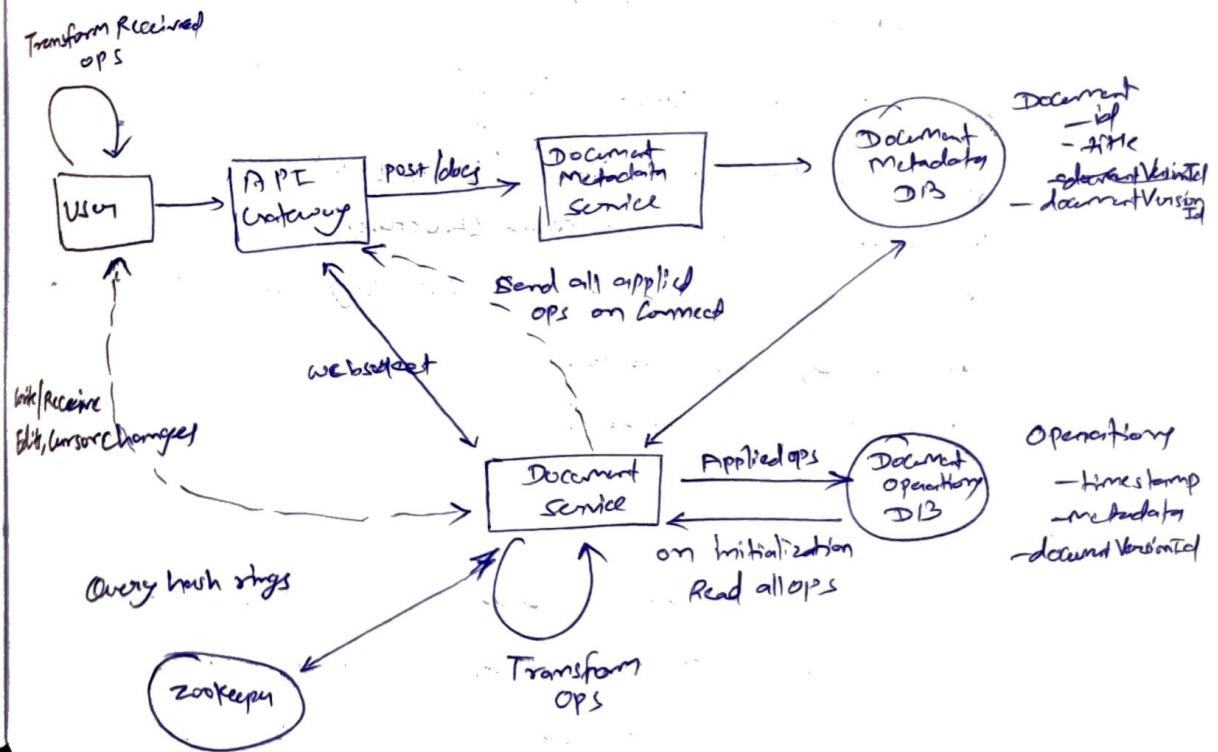
* **Challenges** → This involves complex distributed transaction, especially with databases like Cassandra that have limited transaction scope (row-level only). Coordinating b/w services to ensure atomicity is tricky.

(b) current approach — online snapshot/Compact (within document service)

* [Approach] → The Document service itself periodically performs compaction. This is ideal when a document becomes idle (no active connections), as the Document service already has all operations in memory & has exclusive ownership.

* [Process] → when the last client disconnects from a document, the Document service offloads compaction to a separate, low-priority CPU process. The results are written to the DB under new documentVersionId, which is then flipped in the Document Metadata DB.

* [Challenge] → The main risk is increasing latency for live document operations if compaction tasks hog CPU & resources. Running compaction in a separate, lower-priority process helps mitigate this.



Additional deep dive -

① Read-only mode

② Versioning

③ memory (primary use of Document Service can be bottleneck)
How can we further optimize it?

④ offline mode.

① Read only mode -

To scale read only access significantly, it bypasses the complex OT system. This involves: "separate access path", "serving static snapshots via CDNs", "using read replicas", and for "live" read-only, employing a dedicated read-only stream service that receives pre-computed updates from the main document service & broadcasts them without client-side OT.

(SSE over websockets)

② memory optimization (document service)

- * Smart Document Eviction (LRU, LRU)
- * Efficient in-memory data structures
- * offload & asynchronous processing (for compaction process)
- * Language/Runtime tuning (C/C++)
- * Sharding (inherent)

③ Versioning :-

To implement versioning, each compacted document snapshot is preserved as a distinct, timestamped historical version. The version metadata (e.g. ID, Timestamp, author, parent) is stored

for a dedicated Version History DB.
past document states are reconstructed by loading the
relevant snapshots & then applying any subsequent
granular operation up to the desired point.

The user interface allows ~~Browsing~~ Browse, Previewing &
restoring these historical versions of the current
live document.

Apache Cassandra

(1) Introduction -

Type :- open source, distributed NoSQL DB
Origin :- Facebook (inspired by Dynamo + Bigtable)

Strengths :- Scalable, fault-tolerant, highly available, schema-flexible
Use cases :- Netflix, Discord, Apple, Bloomberg

(2) Data Model -

- * Keyspace - Like RDBMS database; contains tables, defining replication.
- * Table - store rows.
- * Row - Identified by a "Primary Key."
- * Column - Flexible schema; supports VDTs, timestamps & sparse data.

- (Think of Cassandra like a JSON of tables with flexible columns per row.)

(3) Primary Key -

- * Primary Key = (Partition Key) + [Clustering Key(s)]
- * Partition Key = Decides data distribution
- * Clustering Key = Decides row sorting within a partition.
 - Use clustering keys for ranges scans & sorting.

(4) Partitioning -

- * Uses Consistent Hashing on a ring.
- * Data stored on first vnode clockwise from the hash.
- * Vnodes (virtual nodes) :- Distribute load evenly across physical nodes.
 - Adds/removes nodes with minimal data reshuffling.

⑤ Replication -

- (a) Simple Strategy → For dev/test. clockwise replica placement.
- (b) Network Topology Strategy → Prod use. Data center + rack aware

|| simple
 $\text{REPLICATION} = \{ \text{'class': 'SimpleStrategy'}, \text{'replication-factor': 3} \}$

|| network aware
 $\text{REPLICATION} = \{ \text{'class': 'NetworkTopologyStrategy'}, \text{'dc1': 3, 'dc2': 2} \}$

⑥ Consistency -

- * Turnable consistency: ONE, TWO, QUORUM, ALL
- * QUORUM = $\lceil (N/2) + 1 \rceil$ → overlap in reads/writes
 = read-your-write-guarantee
- * Eventual consistency is default behaviour.
- * No full ACID: only atomic + isolated row-level writes.

⑦ Query Routing +

* Any node can be a coordinator.

- * Coordinator:-
 - knows all nodes (via gossip)
 - routes requests to appropriate replicas using hash rings & replication info.

⑧ Storage Model -

→ (log structured Merge Tree)

→ (write ahead log)

- * Uses LSM Tree (not B-tree):-
 - CommitLog (WAL) → memtable (RAM) → SSTable (disk, immutable)
 - "Tombstones" mark deletions
- * Read flow:-
 - memtable → Bloom filter → SSTable (new → old)

⇒ Compaction: Merges SSTables, clears tombstones

⇒ SSTable Indexing: Maps keys to byte offsets for fast disk access.

⑨ Gossip Protocol -

- * peer-to-peer info sharing (live nodes, schema etc)
- * seed nodes:- Anchor gossip to prevent split cluster.
- * vector clock:- Track node state freshness.

⑩ Fault Tolerance -

- * Failure Detection :- Phi Accrual Detector (Probabilistic)

* Hinted Handoff :-

→ stores temporary "hints", if a node is down,
→ hints sent when node is back.

- * Node is only "dead" after manual decommission,

⑪ Data Modelling Strategy -

Cassandra = Query-first design (not entity-relationship driven)

- * Partition key :- Drives distribution
- * Clustering key :- Drives sort/order
- * Denormalize data to support access patterns
- * Avoid large partitions (Performance bottleneck)

⑫ Example : Discord Messages -

Final schema:-

```
CREATE TABLE message(
    channel-id  bigint,
    bucket      int,
    message-id  bigint,
    ---
    PRIMARY KEY ((channel-id,bucket),message-id))
```

) WITH CLUSTERING ORDER BY (message-id DESC)

- Snowflake ids → Avoid timestamp collisions
- Bucket is 10-day window avoids unbounded partitions
- Partition = (channelid, bucket)

Postgre SQL

① why Postgre SQL?

- * ACID-compliant with strong consistency.
- * Ideal for systems needing complex relationships & structured + semi-structured (JSONB) data.
- * Rich built-in features: Full-text search, geospatial features, (PostGIS), JSONB querying.
- * Scales horizontally via replication, partitioning & sharding.
- * well-supported ecosystem & strong community.

② Read Performance & Indexing

- * Default: B-tree indexes (equality, range, sort)
- * Custom: Create indexes for foreign keys, filter conditions.
- * Covering Indexes: INCLUDE all needed columns to avoid table lookups.
- * Partial Indexes: Smaller + faster indexes for filtered data.

```
CREATE INDEX idx_user_email ON users(email);  
CREATE INDEX idx_post_user_include ON posts(user_id) INCLUDE (title, created_at);  
CREATE INDEX idx_active_user ON users(email) WHERE status='active';
```

③ Full-text & JSONB search -

- * Full-text (tsvector + tsin) - fast word search, relevance ranking, stemming, stopword removal
- * JSONB + tsin - store flexible metadata (mentions, tags etc)

Geospatial with PostGIS -

- * Add extension: `CREATE EXTENSION postgis;`
- * Use GIST index for fast proximity queries.

Query Optimization -

- * Keep working set in RAM: Huge boost in speed.
- * Avoid over-indexing: Slow writes, increase space usage.
- * Use Covering + Partial indexes for critical query paths.

Write architecture -

- * WAL (Write-Ahead Log): Durable disk write.
- * Buffer cache: In-memory data update.
- * Background writer: Flushes dirty pages to disk.
- * Indexes: Every write updates relevant indexes (Costly with many indexes)

Write optimization Techniques -

(a) Vertical scaling - (Upgrade hardware: NVMe SSD, more RAM, better CPU)

Batch Processing -

- * Instead of 1000 separate inserts:-

```
INSERT INTO likes (post_id, user_id) VALUES (1,101), (1,102), ...;
```

Write offloading -

Async write via Kafka or queues for logs/analytics/metrics

Table Partitioning

```
CREATE TABLE posts (...) PARTITION BY RANGE (created_at);
```

```
CREATE TABLE post_2024_01 PARTITION OF posts FOR VALUES FROM('2024-01-01') TO ('2024-02-01');
```

These two statements together define one partition of posts table.

(Create main table)

(Create one partition)

- (e) Sharding -
- ↳ shard by user-id, avoid cross-shard queries.
 - ↳ use tools like City for easier sharding

(f) Replication :- (Scaling Reads + High availability)

- * Asynchronous : Faster, risks of lag.
- * Synchronous : Slower, stronger consistency
- * Used read replicas to offload queries
- * Auto-failover in managed solutions (AWS RDS, MySQL)

(g) Data Consistency & Transactions -

ACID:

- * Atomicity - All or nothing writes
- * Consistency - DB remains valid state
- * Isolation - Each transaction as if alone
- * Durability - Committed = won't be lost

Isolation levels: Read Committed, Repeatable Read, Serializable

(h) When to Consider alternatives -

Need	Alternative
(1) $\geq 1M$ writes/sec.	Cassandra, DynamoDB
(2) Global writes (Multi-region)	CockroachDB
(3) Simple KV access	Redis, DynamoDB

Interview tips :-

- (a) Start with PostgreSQL, unless requirements demand otherwise.
- (b) Justify with ACID, indexing, write/read scaling, search, & partitioning.
- (c) Don't just say "PostgreSQL is ACID" - show how you use it for consistency.
- (d) Mention trade offs : eg- (write throughput vs number of indexes), (read replicas vs lag).