

DynamoDB

- Fully managed, highly scalable, key-value service provided by AWS.
- Not an open-source.

⇒ The Data Model :-

In dynamodb, data is organized into tables, where each table has multiple items that represent individual records.

Ex:- Using table in dynamodb:-

```
{ "personid": 101,  
  "name": "Fred"  
},  
{ "personid": 102,  
  "name": "Mary",  
  "address": {  
    "city": "Town",  
  }  
}
```

- dynamodb uses json for data transmission, it's merely a transport format.

Partition Key & Sort Key

DynamoDB tables are defined by a primary key, which can consist of one or two attributes:-

① Partition Keys:-

A single attribute that, along with the sort key (if present), uniquely identifies each item in table.

- DynamoDB uses the partition key's value to determine the physical location of item within the database.
- This value is hashed to determine the partition where the item is stored.

② Sort Key (optional) - An additional attribute that, when combined with the Partition Key, forms a Composite Primary Key. The Sort Key is used to order items with the same Partition Key value, enabling efficient range queries & sorting within a partition.

Primary key = {Partition Key} : {Sort Key}

↑
optional

→ But what is actually happening under the hood?

DynamoDB uses a combination of consistent hashing and B-tree, to efficiently manage data distribution & retrieval.

for sort keys

for Partition Key

~~Composite~~ * Composite Key operation - when querying with both keys, DynamoDB first uses the partition key's hash to find the right node, then uses the sort key to traverse the B-tree & find the specific items.

This two tier approach allows DynamoDB to achieve both horizontally scalability & efficient querying within partition.

Secondary Indexes

But what if you need to query your data by an attribute that isn't the Partition Key? This is where secondary index comes in. DynamoDB supports two types of secondary indexes:

(i) Global Secondary Index (GSI) :-

→ An index with a partition key & optional sort key that differs from the table's partition key. GSIs allow you to query items based on attributes stored on entirely different physical partitions from the base table & is replicated separately.

(ii) Local Secondary Index (~~GSI~~ (LSI)) :-

An index with the same partition key as the table's partition key but a different sort key.

→ LSIs enable range queries & sorting within a partition. Since LSIs use the same partition key as the base table, they are stored on the same physical partitions as the items they are indexing.

Accessing data in dynamodb

① Scan operation

② Query operation,

~~# Redis #~~

→ It is single threaded making it very fast & easy to reason about.

⇒ some of the most fundamental data structures supported by Redis:-

- (1) strings
- (2) Hashes (objects)
- (3) Lists
- (4) sets
- (5) sorted sets (Priority queues)
- (6) Bloom filters
- (7) Geospatial Indexing
- (8) Time series

In addition to simple data structures, Redis also supports different communication patterns like Pub/Sub & streams, partially standing in for more complex setups like Apache Kafka etc.

Capabilities

(1) Redis as a cache

(2) Redis as a distributed lock (Used in TicketMaster, Uber)

(3) Redis for leaderboards

→ Redis's sorted sets maintains ordered data which can be queried in log time which make them appropriate for leaderboard applications.

The high write throughput & low read latency make this especially useful for scaled applications where something like a SQL DB will start to struggle.

- In "Post-search Facebook" system design problem, we have a need to find the posts which contain a given keyword (e.g. "Tiger") which have the most likes.
- we can use Redis's sorted set to maintain a list of the top liked posts for a given keyword. Periodically, we can remove low-ranked posts to save space.

④ Redis for rate limiting

⑤ Redis for Proximity Search - (MEADD, MESEARCH command)

⑥ Redis for Event Sourcing

Facebook's Post Search

For this problem, we're zooming in on the search experience for Facebook. This question is prime for infrastructure-style interviews where your interviewer want to understand how deeply you understand data layout, indexing & scaling.

Functional Requirements -

core req:-

- 1) User should be able to create & like post
- 2) User should be able to search posts by keyword.
- 3) User should be able to get search results by recency or like count.

out of scope

- support fuzzy matching on terms (e.g. search for "bird" matches "sparrow")
- personalization in search result
- privacy rules & filter
- Images & media
- Realtime update to the search page as new post come in.

Non-functional Req:-

- 1) Fast queries < 500ms
- 2) High volume of requests
- 3) ~~Each~~ New post must be searchable in short amount of time < 1 minute.

- All parts must be discoverable, including old or unpopular parts. (we can take more time for these)
- The system should be highly ~~scalable~~ available.

Hot vs Cold Data If you have a lot of data in your system, there's a good chance you'll have "hot" and "cold" data. Hot data is readily accessed and needs to be served quickly, often from memory.

Cold data is infrequently accessed & can be served from disk, a remote database or even tape.

Scale Estimation -

volume of work?

$$\text{Post Created} = 1B * 1 \text{ post/day} / (100k \text{ sec/day}) = 10^4 = 10K \text{ post per sec.}$$

$$\text{Like Created} = 1B * 10 likes/day / (100k sec/day) = 100K likes per second.$$

Let's look at reads

$$(\text{Searches}) \quad 1B * 1 \text{ search/day} / (100k \text{ sec/day}) = 10K \text{ searches/sec.}$$

- while this is a lot (& may burst the 10x this value or more), our system is write-heavy or Read-heavy, another thing to note.

Storage Req

$$\text{posts searchable} : 1B \text{ posts/day} * 365 \text{ day/year} * 10 \text{ years} = 3.6 \text{ T posts}$$

$$\text{Raw size} : 3.6 \text{ posts} * 1 \text{ kb/post} = 3.6 \text{ PB.}$$

Core Endpoints

- ① User
- ② Post
- ③ Like

APIs:-

- ① /Search

GET /search ?query=QUERY & sort-order=Likes/Timestamp

- ② /Create a post

POST /Post

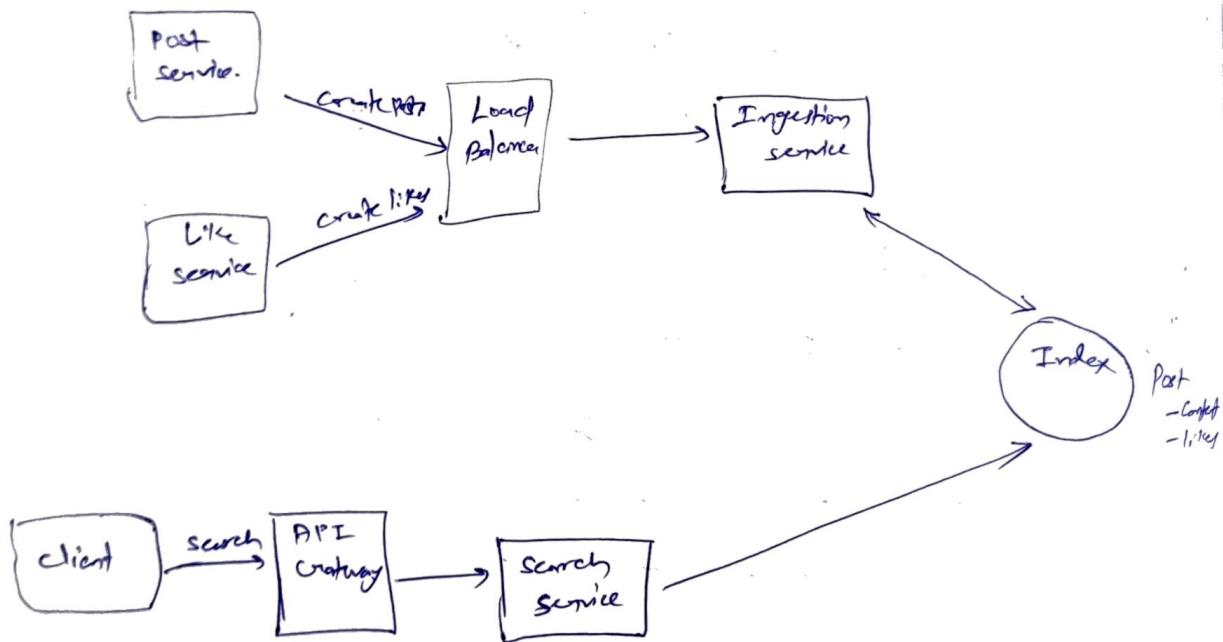
- ③ /Like a post

POST /post/{postId}/Like

②

High level design

- ① user should be able to create & like posts.
- ② user should be able to search posts by keyword.



Search post by Keyword - ways

- → Bud self — Scale an unindexed database

(Go through each content in DB & search)

Post
- Content
- Likes

~~→ cannot col~~ - Create an Inverted Index

Create dictionary, map keyword to document ids

key	value (post id)
Taylor	1, 5, 10
Board	2, 5, 13, 15

- we can use Redis for this (in-memory)

↓
There are durability concerns but can be handled
them with a durable alternative like
Memory DB or in a deep-drive.

- Challenges -
- * these post ids list are going to get very large, especially of common keywords. we will have to address this later
 - * we also need to ~~be~~ write to many keys for every post since a post might have 10-1,000 keywords. we'll need to handle some of scaling challenges associated with ~~search this~~.

- ③ User should be able to get search results sorted by recency or like count.

⇒ Bad boy - Request time sorting

↳ grab all post ids for a given keyword from index & then ~~loop~~ ~~repeatedly~~ look up all the posts for timestamp or like count, & then sort those in memory.

↳ challenges - there could be millions of posts for a common keyword - add latency to system.

⇒ Insert solution - Multiple Indexes

↳ different approach would be to have two separate indexes :-

- one ~~for~~ sorted by creation (creation index)
- one sorted by likes count (likes index)

* For 'creation index' - we can use standard Redis list. We are always going to append to this list & our query will only be taking from last element.

* For the 'Likes index' - each key can use a Redis sorted list ? (keep a list of items ordered by score)

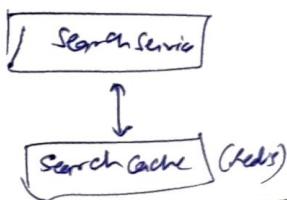
Challenge - ① doubled the storage for indexes.

② introduced a new problem - likes are happening quite frequently. Each like event requires us to update many scores so that the like indexes are up-to-date. This puts a lot of stress on our system, which we'll both note for our interview & plan to address later.

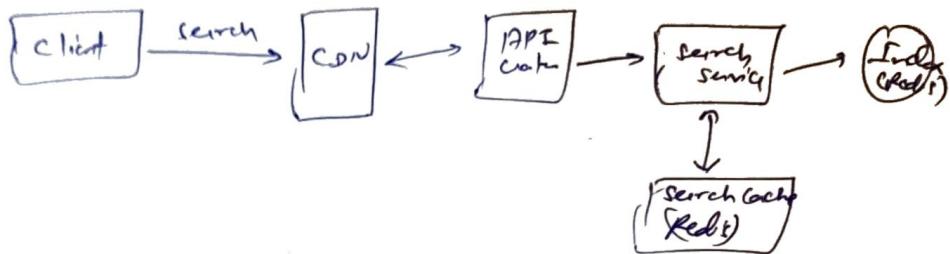
Potential deep dive

① How can we handle large volume of requests from user?

⇒ Brick wall - Use a distributed cache alongside our search service.



→ current solⁿ - Use a CDN to cache at the edge.



② How can we handle multi-keyword, phrase queries?

⇒ current solⁿ - Intersection & filter

↳ challenge ⇒ millions of post for each keyword

↳ challenge ⇒ millions of post for each keyword
↳ challenge ⇒ transferring megabytes of data
which need to put into a hash table
& intersected. This can be hard to
do while meeting our 500ms SLA.
- quite computationally extensive
(in worst case)

⇒ current solⁿ - Bigrams or Shingles

↳ create token for each pair of word & index.

Ex:- I saw taylor swift at concert.



↳ created this pairs to indexer.

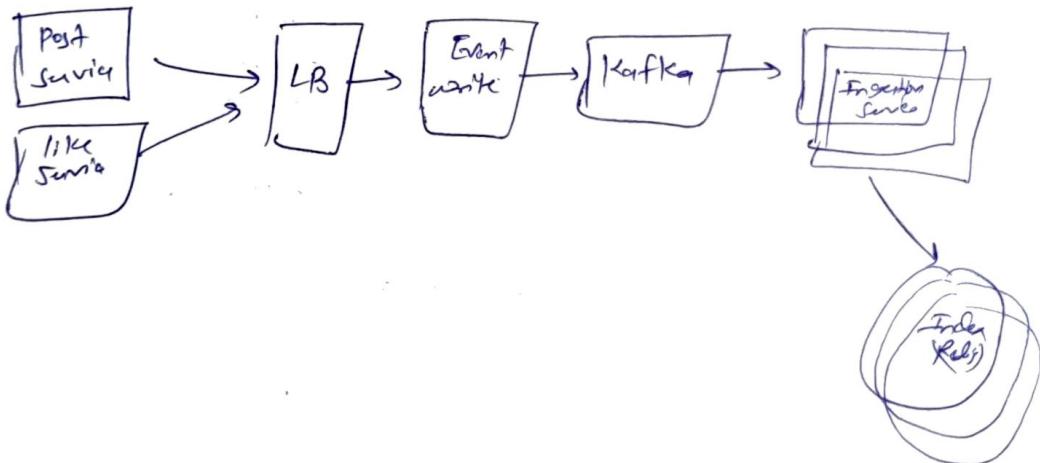
challenge ⇒ * dramatically increases the size of index,
so instead of 10M single-words keyword
we might end up 100M+ indexed keys

Potential remediation - only extract bigrams that are
likely to be searched & fall back
to intersection approach when there
is no entry in our index.

③ How can we address large volume of writes?

(A) Post Creation -

- * Add more capacity to ingestion service & Partitioning the incoming requests.
- * use ~~log~~ log stream like Kafka → Partition load across ingestion service.
- * finally we can scale out our index by sharding the indexes by keyword. This way writes to the indexes are spread across many indexes.

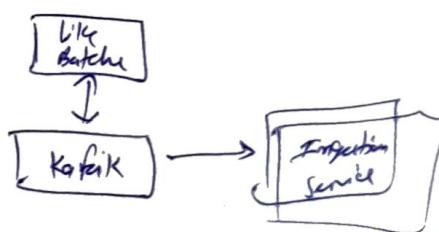


(B) Like Event -

→ Wood soi - Batch like before writing to the database.
we can batch likes for a given postid over a period (like 30 second). Then instead of needing to write 500 times for partially viral post, we can make 1 update with an increment.

Challenge -

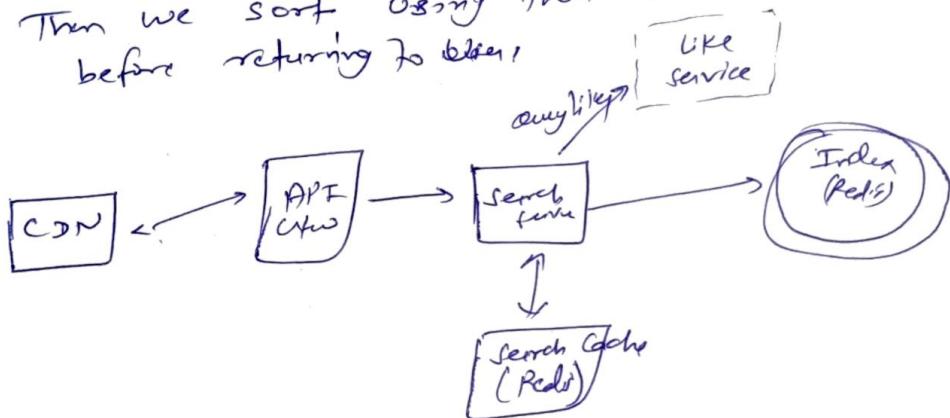
of 500.
It won't drastically reduce volume of writes since most post aren't viral
overhead added,



→ insert 50%: Two stage architecture

→ we can do even better by only updating the like count in our index when it passes specific milestones like power of 2 or 10. so we would write to the index when the post has 1 like, 2 likes, 4 likes etc. This reduces the number of writes exponentially.

→ but this means our index is inherently stale. If we want to retrieve N posts for a given query, we grab the top N^2 posts from like index. For each post in that list, we query the "LikeService" for up-to-date count. Then we sort using that fresh like count before returning to user.



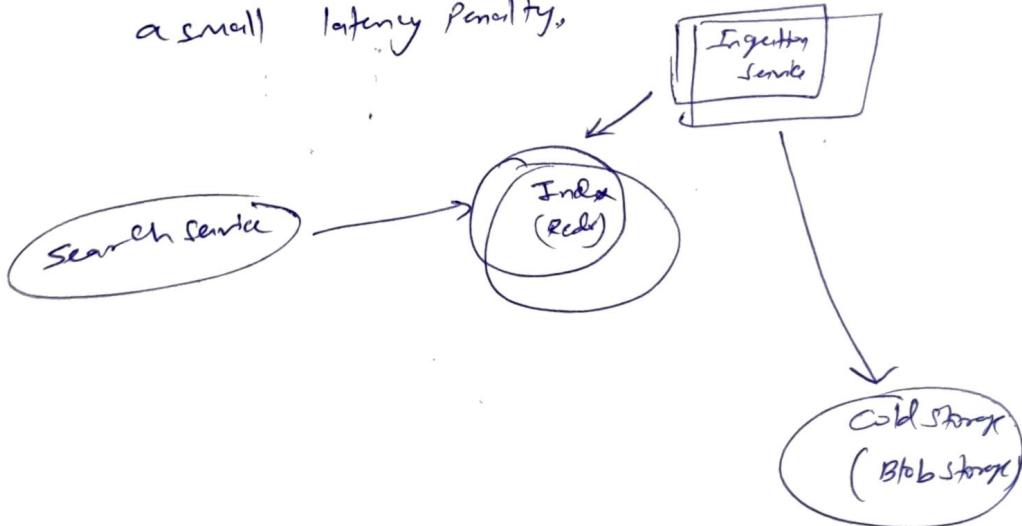
Challenger → more complex

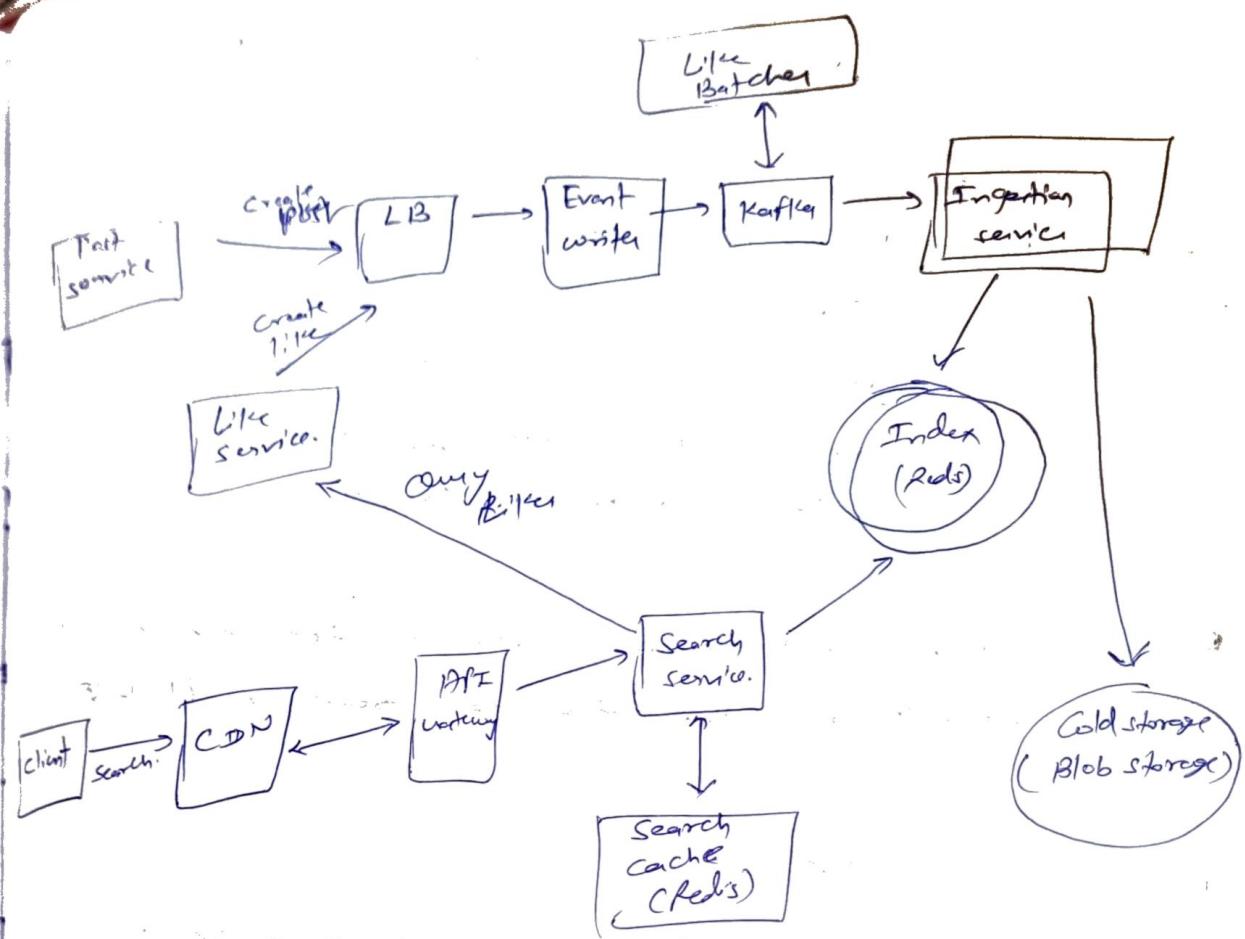
④ How can we optimize storage of our system?

→ first, we can put caps & limits on each of our inverted indexes. we ~~won't~~ probably won't need all 10M post ^{with} ~~content~~ "Mark" contained in their contexts. By keeping our index to 1K - 10K items, we can reduce the necessary storage by orders of magnitude.

Next, most keywords won't be searched often or even at all. Based on our search analytics, we can run a batch job to move rarely used keywords to a less frequently accessed but ultimately cheaper storage. One way to do this is to move those keyword indexes to cold, blob storage like S3 or R2.

- on a regular basis we'll determine which keywords were rarely used (or not at all) accessed in past month. we'll move these indexed from our in-memory Redis instance to a blob in our blob storage.
- when the index needs to be queried, we'll first try to query Redis. If we don't get our keyword there, we can query the index from our blob storage with a small latency penalty.





II Design a Web Crawler

Web crawler → Program which traverse web pages → download web pages & following link from one page to another.

- It is used to index the web for search engines, collect data for research or monitor websites for changes.

For our purpose, we'll design a web crawler whose goal is to extract from the web to train an LLM. ~~etc.~~

Functional Req.

- ① Crawl the web starting from a given set of seed urls.
- ② Extract text data from each web page & store the ~~rest~~ text for later processing.

out of scope

- Actual Processing of text data (LLM etc.)
- Handling non-text data (Audio, video)
- Handling of dynamic content
- Handling of Authentication (e.g. login required pages)

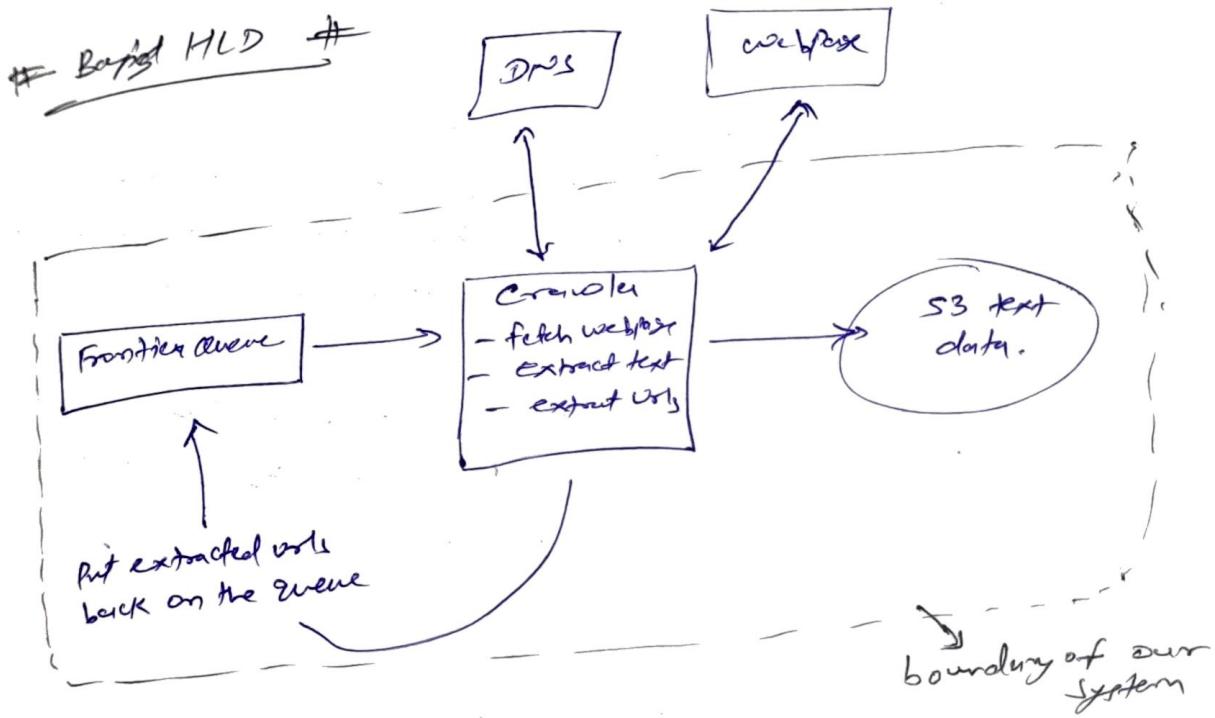
Non-functional

- ① Fault tolerance to handle failures gracefully & resume crawling without losing progress
- ② Politeness to adhere to robots.txt & not overloaded website with heavy traffic

- (3) Efficiency to crawl the web in under 5 day.
- (4) Scalability to handle 10B pages.

out of scope

- security to protect system from malicious actor
- cost to operate the system within budget constraint
- compliance to adhere to legal requirement & privacy regulation.



Potential Deep Dive

- ① How can we ensure we are fault tolerant & don't lose progress?

Currently Crawler Service
is doing many things

- hitting DNS
- Fetching web pages \Rightarrow (most likely to fail)
- extracting text data
- extracting new urls to add to the frontier queue.

4

This is not fault tolerance. If there is a failure in any single task, all progress will be lost.

- we should break crawler service into smaller, pipelined stages:-

(a) URL fetcher \rightarrow fetch HTML of web pages from external server. If there is failure, we can retry without losing progress on the rest of data. (storing data in blob storage)

(b) Text & URL Extraction \rightarrow Extract text data from HTML, any linked url to add to frontier queue.

\Rightarrow what about if we fail to fetch a URL?

(i) Bad sol \Rightarrow In memory timer $\{$ This isn't robust. if crawler goes down, we will lose timer.
Also it will not likely to work in few more seconds need some sort of exponential backoff

(ii) Bad (good sol) \Rightarrow Kafka with manual exponential backoff

$\{\text{Complex to implement}$ & maintain $\}$ \hookrightarrow separate topic for failed urls & separate service that reads from this topic & retries the fetch with exponential backoff.

- (iii) what col? - SES with exponential backoff
- ↳ SES support retry with configurable exponential backoff out of the box. (No need to implement own logic)
 - ↳ Retry attempt - 30sec, 2 min, 5 min, up to 15min.

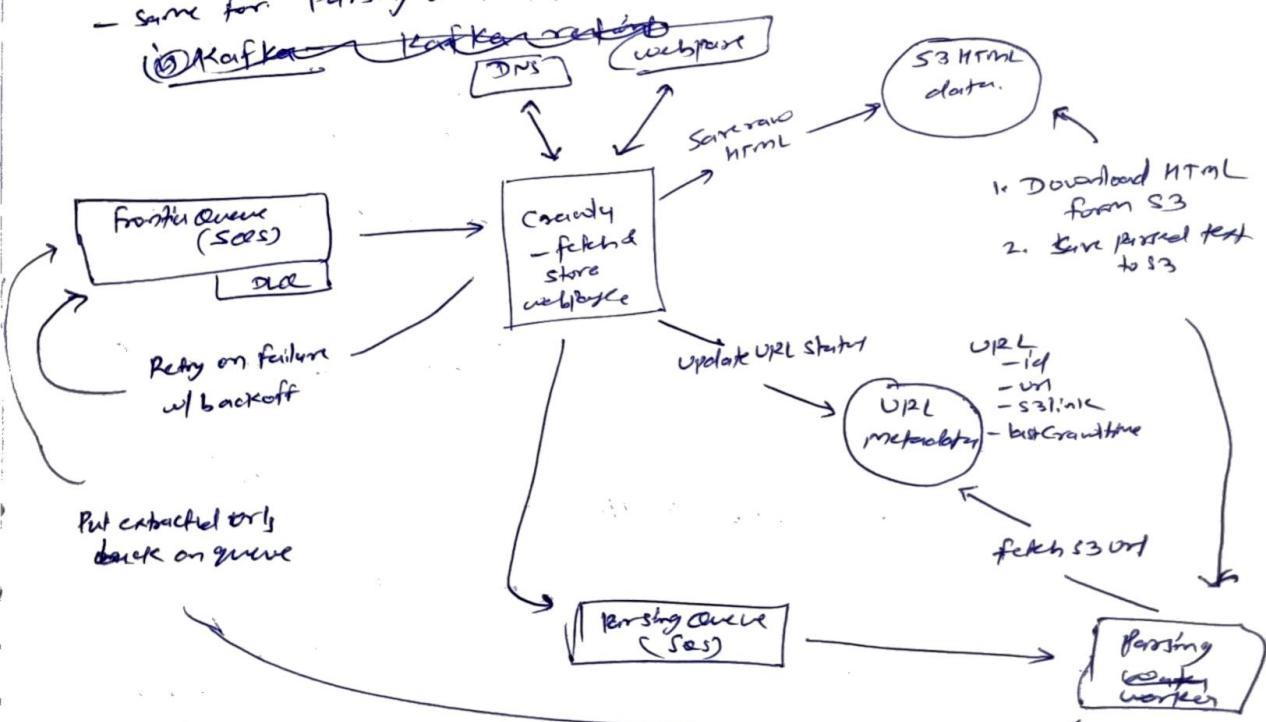
what if crawler goes down?

Answer is simple! - we spin up a new one. We'll just have to make sure that the half-finished URL is not lost.

Good news is the URL will stay in queue until it is confirmed to have been fetch by a crawler & the HTML is stored in blob storage. This way if a crawler goes down, the URL will be picked up by another crawler & process will continue.

We can accomplish this in Apache Kafka & SES both.

- same for Parsing worker goes down.



(2) How can we ensure politeness & adhere to robots.txt

~~robots.txt~~

robots.txt → file that website use to communicate with web crawler.

↳ similar to

user-agents: *
Disallow: /private
CrawlDelay: 10

→ specifies how many seconds the crawler should wait b/w requests.

To ensure politeness & adhere to robots.txt, we will need to do two things:

(i) Respect robots.txt - Before crawling a page, we will need to check the robots.txt file to see if we are allowed to crawl the page.

(ii) Rate limiting - (implement a global, domain-specific rate limiting mechanism using a centralized data store (like Redis) to track requests b/wt per domain per second).

(3) How to scale to 10B pages & efficiently crawl them in under 5 days.

In AWS ecosystem, a network optimized instance can handle about 400 Mbps.

$$* \frac{400 \text{ Mbps}}{8 \text{ bits/byte}} / 2 \text{ MB/page} = 25000 \text{ pages/second}$$

$$* \text{utilizing } 30\% \Rightarrow 7500 \text{ pages/second}$$

$$* \frac{10 \text{ B pages}}{7500 \text{ pages/second}} = 1,333,333 \text{ second} = 15.4 \text{ days}$$

for single machine

⇒ 4 machines needed

DNS → 1) DNS Caching

2) Multiple DNS providers

⇒ Last thing, watch out for crawler traps -

* Crawler traps are pages that are designed to keep crawlers on the site indefinitely. They can be created by having a page that links to itself many times or by having a page that links to many other pages on site. If we're not careful, we could end up crawling the same site over & over again & never finish.

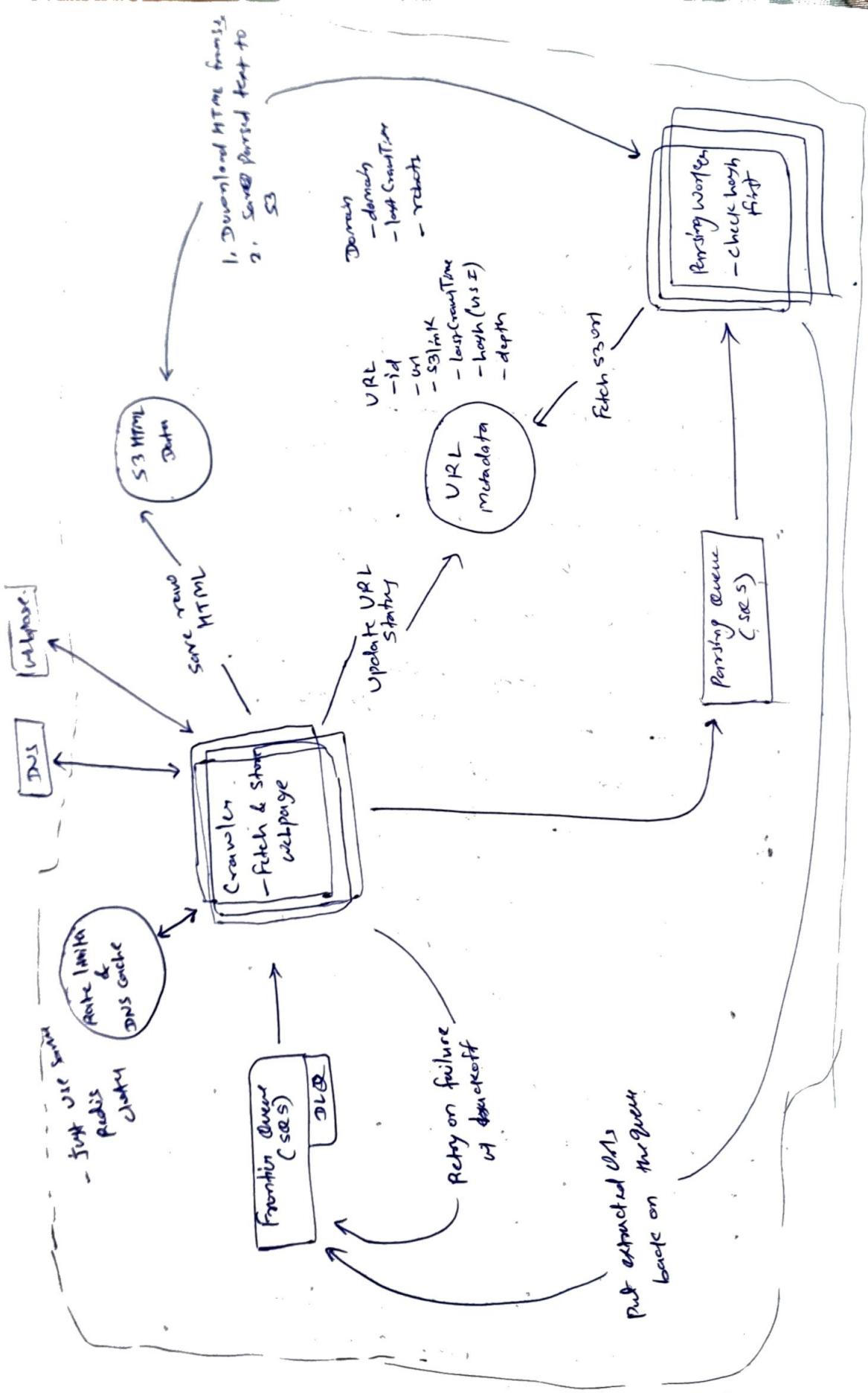
* Solution is pretty straightforward, we can implement a maximum depth for our crawler. We can add a 'depth' field to our DB table in Metadata DB & increment this field each time we follow a link. If the depth exceeds a certain threshold, we can stop crawling the page. This will prevent us from getting stuck in a crawler trap.

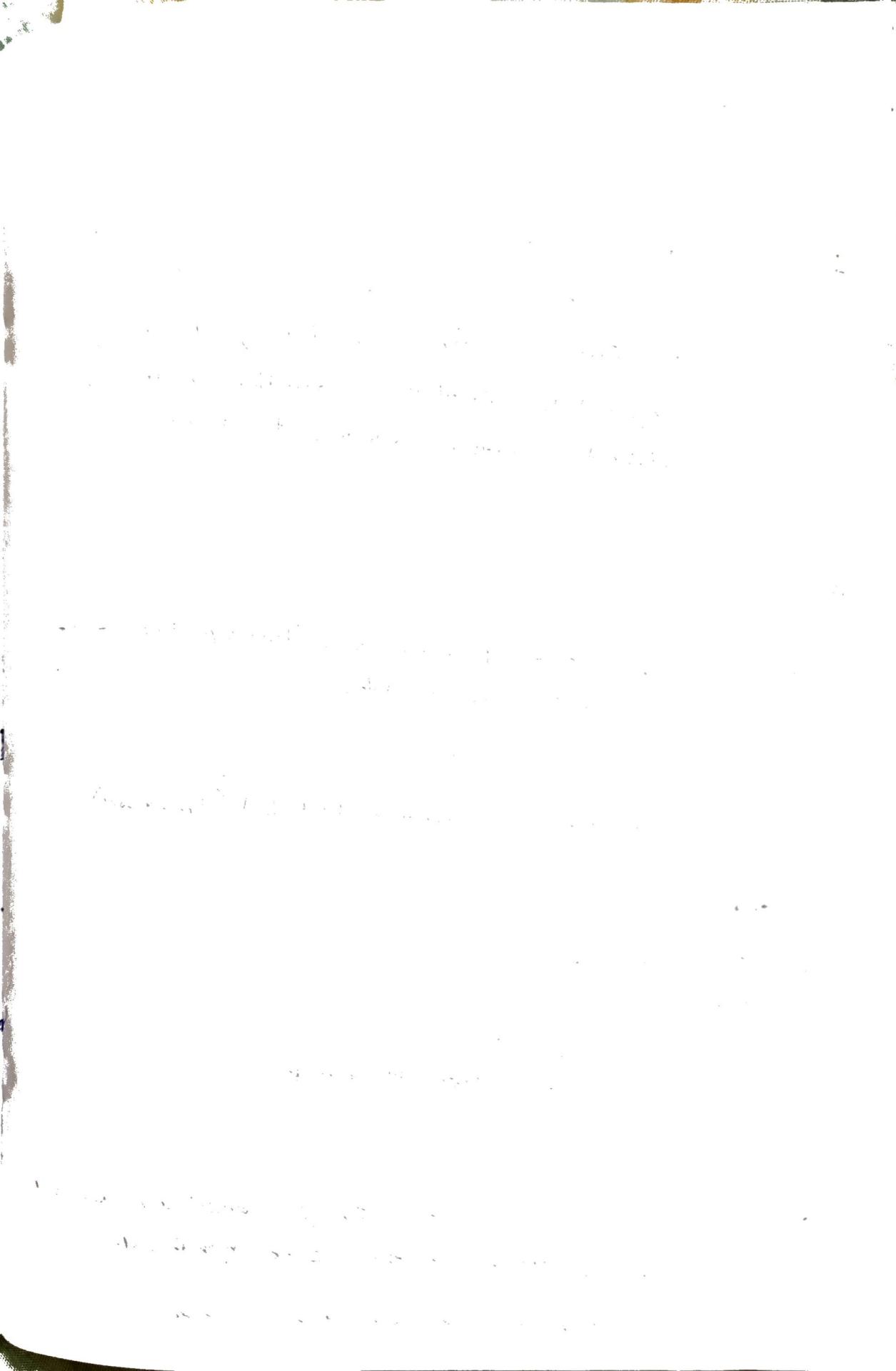
⇒ How to handle dynamic content?

→ we'll need to use a headless browser like puppeteer to render the page & extract the component.

⇒ How to handle large files?

→ we can use 'Content-Length' header to determine size of the file before downloading it & skip files that are too large.





Design an Online Auction Platform Like eBay

Intro:- An online auction service lets users list items for sale while others compete to purchase them by placing increasingly higher bids until the auction ends, with the highest bidder winning the item.

Functional Req:-

- ① User should be able to post/put an item up for auction with a starting price & end date.
- ② Bid on auctions by user
- ③ view an auction & current highest bid (1 max bid)

out of scope

- Search for items
- Filter items by category
- Sort items by price
- view the auction history of an item.

Non-Functional Req:-

- ① System should maintain strong consistency for bids to ensure all user see the same highest bid.
- ② Fault tolerant & durable. we can't drop any bids

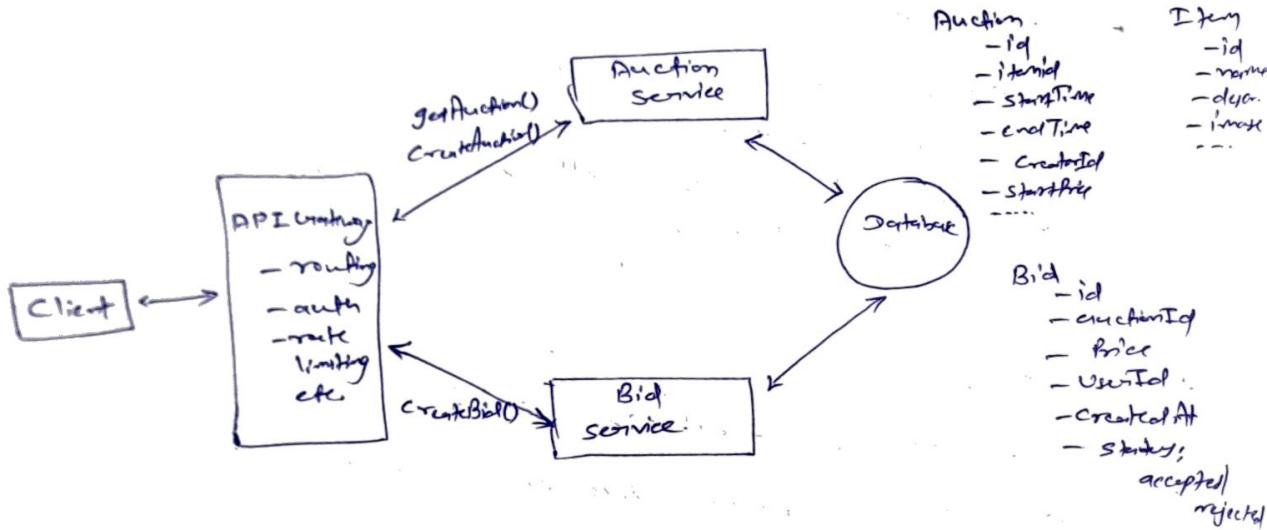
- ③ Real-time max bid display on client.
 ④ scale to 10M concurrent auction.
- out of scope
- system should have proper observability & monitoring.
 - secure & protect user data.
 - system should be well tested & easy to deploy
(CI/CD Pipelining)

- # Entity :-
- ① Auction :- ~~start date~~, ^{price} end date, item being auctioned
 - ② Item :- name, description, image
 - ③ Bid :- amount bid, user placing the bid, the auction being bid on
 - ④ User :- who starts auction or bids on an auction.

API | System Interface

- ① POST /auctions → Auction & Item
 - { item: Item,
start date: Date,
end date: Date,
starting price: number }
- ② POST /auctions/: auctionId / bid → Bid
 - { Bid }
- ③ GET /auctions/: auctionId → Auction & Item.

High Level Design



* why two services:-

→ Bidding Traffic is typically much higher volume than auction creation - we expect ~100x more bids than auction. Having a separate service allows us to scale the bidding infrastructure independently.

* Bidding Service will :-

- validate incoming bid (Ex:- check bid amount is higher than current max)
- Update the auction with new highest bid.
- Store bid history in the database
- Notify relevant parties of bid updates.

⇒ Let's assume user request for update for an auction with current max bid amount. what happens next is more interesting:-

- If we never reflect maximum bid price, then the user will bid ~~too~~ based on a stale amount & be confused (& frustrated) when they are

told their bid was not accepted. Especially in an auction with a lot of activity, this is a problem.

To solve this, we can simply poll for the latest maximum bid price every few seconds. While imperfect, this ensures at least some degree of consistency & reduces the likelihood of a user being told their bid was not accepted when it actually was.

Potentially deep dives

① How can we ensure strong consistency for bids?

② Back story - Row Level locking with Bids every n

- * [approach] -
 - Begin a transaction
 - Lock all bids for auction using select... for update
 - query the current maximum bid from the bids table
 - Compare the new bid against it
 - update the new bid if accepted
 - commit the transaction

- * [challenge] - Performance & scalability issue.
 - ↳ locking all bid rows for an auction serializes bid processing creating a major bottleneck

- Poor User Experience
 - ↳ delay introduced by lock contention deserialized processing result in slow response or timeout when placing bids.

① Gridsol1 - Cache max bid externally.

↳ Problem with above solⁿ (a) is lock on large no. of bid rows.

↳ Solution
we can cache the max bid in memory in something like Redis.

Flow :- * Read Cache (for max bid)
* Update Cache (with new bid #)
* write bid (write bid to db with status of accepted/rejected).

Issue → Cache-DB inconsistency (Cache update succeeds but db write fails, leading to inconsistent state.)

→ Rollback complexity - (Rolling back cache after DB failure adds operational complexity & potential race condition)

② Gridsol1 - Cache max bid in database

↳ storing max bid in the database.

New flow

- lock auction row for given auction (just one)
- read max bid for given auction from Auction table
- write bid to database with status of accepted or rejected
- update max bid in Auction table if new bid is higher
- Commit transaction.

↳ we now only lock a single row & for a short duration

② How can we ensure that system is fault-tolerant & durable?

→ we need to guarantee durability & ensure that all bids are recorded & processed, even in the face of system failure.

The best approach here is to introduce a durable message queue & get bids into it as soon as possible. This offers several benefits—

a) Durable storage

b) Buffering Against Load Spikes

→ Imagine popular auction

↑ thousands of bids per second

without queue → Drop bids → unacceptable,
Crash under the load
massively overprovision our server (expensive)

c) Guaranteed ordering

→ This is important for factored—if two users bid the same amount, we want to ensure the first bid wins. The queue gives us this ordering for free.

⇒ Kafka is suitable here because—

- High throughput

- Durability

- Partitioning (by auctioned)

③ How can we ensure that the system displays the current highest bid in real-time?

⇒ bad - Polling (two slow, inefficient)

⇒ good - long polling for max bid

⇒ great - Server-Sent Events (SSE)

→ when a user views an auction, their browser establishes a single SSE connection. The server can then push new ~~not~~ maximum bid values through this connection whenever they change. This creates a true real-time experience.

→ on server side, we maintain a set of active SSE connections for each auction. When a new bid is accepted, we push the updated maximum bid to all connected clients for that auction.

Challenge - Scaling issue when multiple servers come to picture.

- we will talk about solution in next chapter

④

How can we ensure that system scales to support 10M concurrent auctions?

- while it comes to discussing scale, you will typically want to follow a similar process for every system design question: working left to right, evaluate each component in your design asking the following questions:

a) what are the resource requirement at peak?
Consider storage, compute & bandwidth.

- b) Does the current design satisfy the requirement at peak?
- c) If not, how can we scale the component to meet the new requirement?

Start with some assumptions \Rightarrow

10M concurrent auctions, 100 bids.

$\frac{1}{4}$
1 Billion bids / day

\downarrow
10K bids per second

\Rightarrow Message Queue

\hookrightarrow decent hardware can handle 10K request per second.
 \hookrightarrow No issue here.

\Rightarrow Bid Service

\hookrightarrow In this case with almost all stateless servers, we can horizontally scale.

\Rightarrow Database :-

- Auction = 1KB
- Bid = 500 bytes

- Avg auction runs for a week

$$10M \times 52 = 520M \text{ auction per year}$$

$$520M \times (1KB + (0.5KB + 100)) \\ = 25 \text{ TB of storage per year}$$

$\frac{1}{4}$
modern SSDs can handle
~~10TB~~ 100 TB of storage, no issue here.

\Rightarrow SSE :-

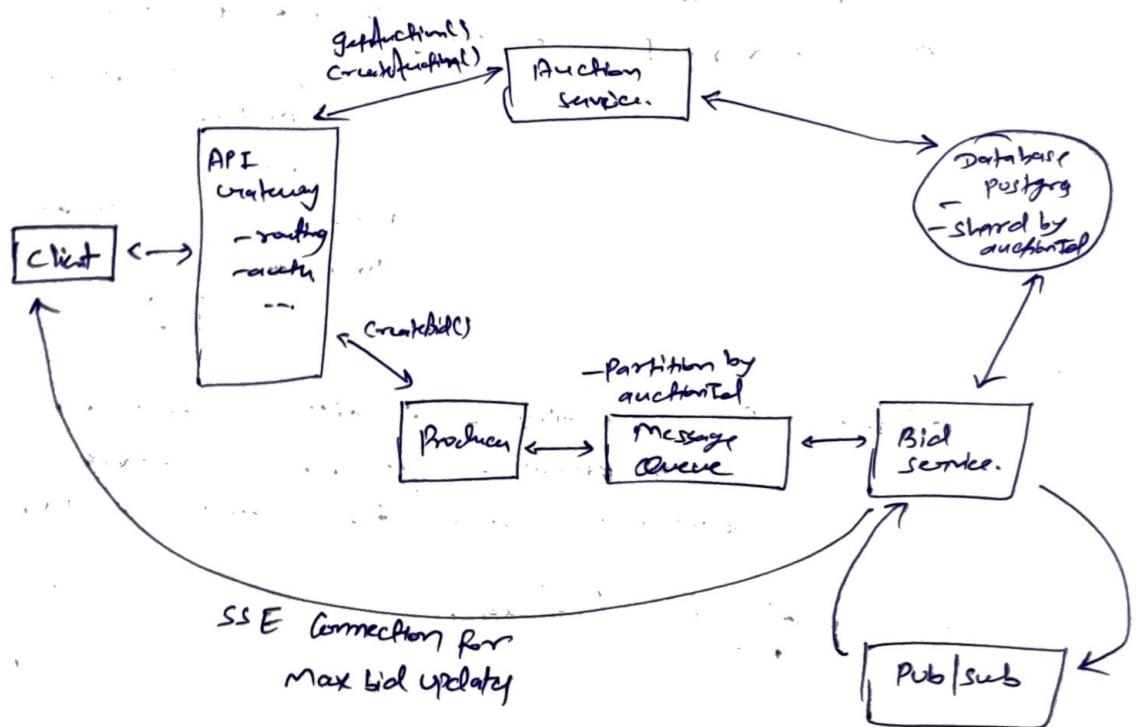
\hookrightarrow issue - All connections won't fit on same server.
~~so we need~~ a way for Bid service servers to coordinate with one another.

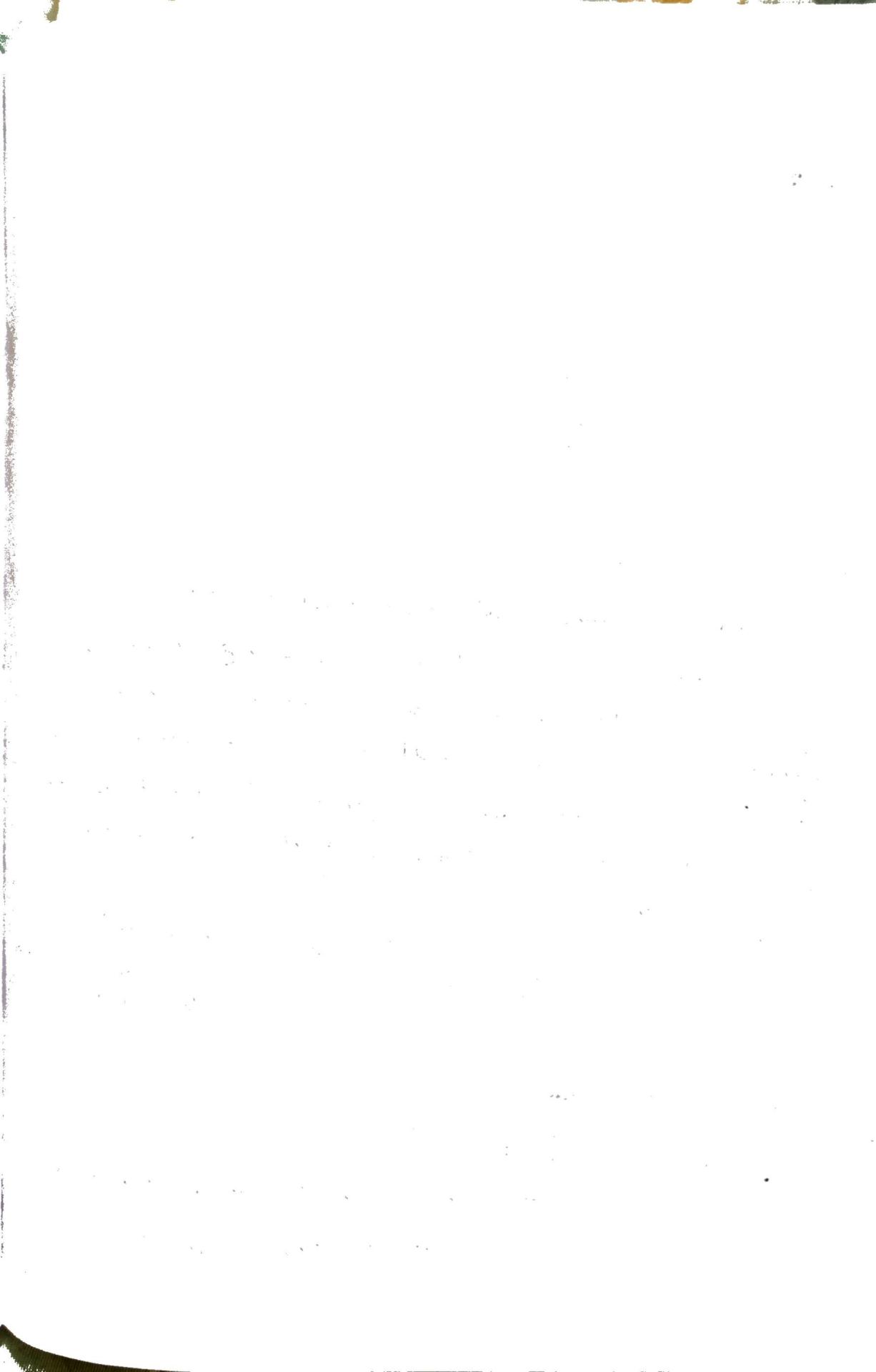
The solution is Pub/Sub.



↓
 whether using a technology like Redis's pub/sub, Kafka, or even a custom solution, we need to broadcast the new bid to all relevant servers so they can push the update to their connected client.

The way this works is simple: when server A receives a new bid, it publishes that bid to the pub/sub system, essentially putting it on a queue. All other instances of Bid Service are subscribed to some pub/sub channel & receive the message. If the bid for an auction that one of their current client connections is subscribed to, they send that connection the new bid data.





AutoComplete | Word Completion

- AutoComplete is involved in many search features like search features on various platforms like Facebook, Instagram, Google search, etc.
- For Facebook - one use case which involves keeping user's friends, mutual connections in mind while giving suggestions in the typeahead.

Functional Req. :-

- ① ~~suggestions~~ suggestions should be displayed on type of user.
*if search query
is in trending
add it to
do
so that include
it in
autocomplete
suggestions*
- ② It is supposed to be real-time. Let's say if currently there is an India vs Pakistan match going on, then ~~should~~ our autocomplete should include this as a suggestion in the typeahead if searching "India vs Australia" is a trending query.
- ③ Rank our suggestions (based on frequency)
*no of times searched
by our user*

or without loss

- \approx Billion ~~second~~ queries per day.
- we need 7 autocomplete suggestions for each character we enter in the search box.

Non-functional Req :-

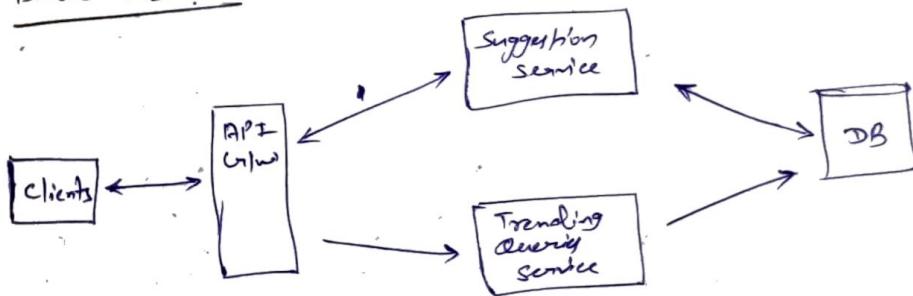
- ① Latency of suggestions must be very less.
- ② Handle large amount of data. (Partition tolerance)
- ③ Highly available \rightarrow
- ④ Eventual consistency is fine.
- ⑤ Search results must be ordered by some ranking score
- ⑥ The results must be relevant to whatever prefix the user has entered.

System APIs :-

① get_suggestions (Prefix)

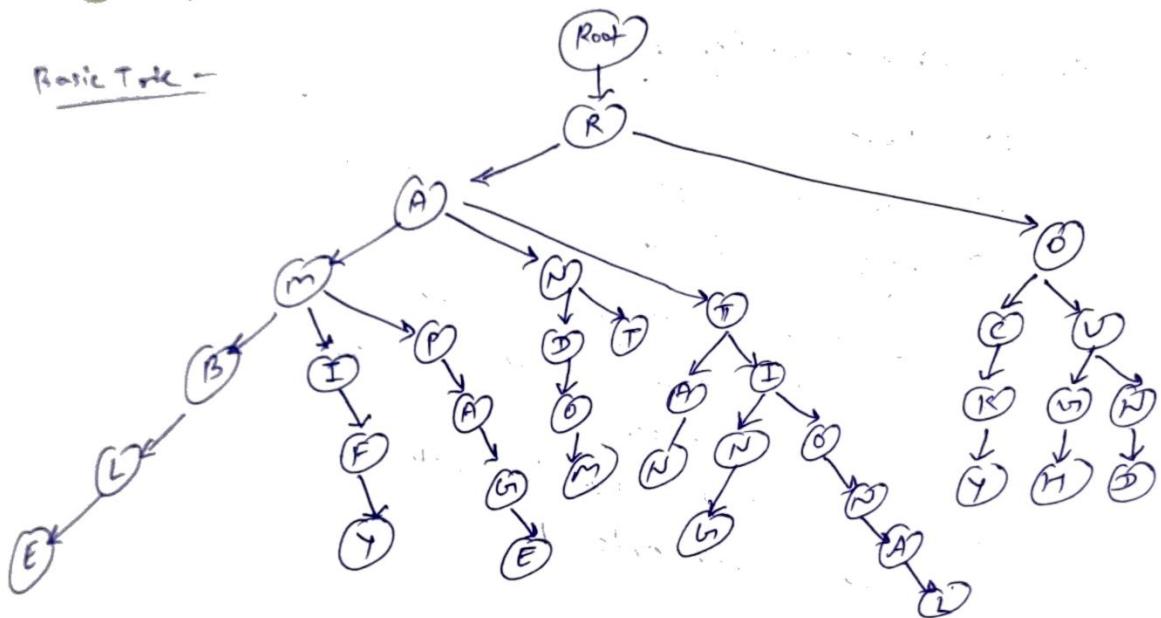
② add_trend (query) - a new unique trending query which has been searched above a certain threshold will be stored in our DB

Basic HLD



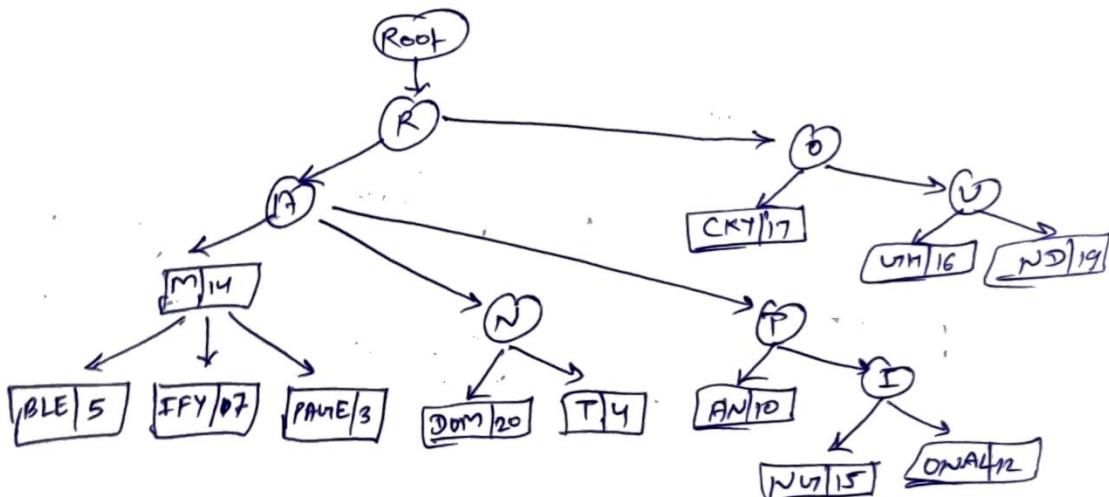
\Rightarrow Data structure to store search queries :-

Basic Trie -



In order to optimize space, we can merge nodes that have only one branch. Now tree will be stored like:-

This is called a suffix tree. (Compressed trie of all the suffixes).
Further, we can store the frequency of each search result in the tree itself as follows:-



Time complexity to go to the prefix node - $O(L)$

$L \rightarrow$ length of prefix

Let's say we entered 'RAT' - all possible suggestions would be found by going through the entire subtree under RAT & all nodes which have frequency 1, i.e. which are complete words will be fetched.

RAT we get -- RATAN-10, RATIONAL-12, RATINU-15

Time complexity to find all possible nodes under the prefix node i.e. scanning the entire subtree - $O(n)$
number of nodes under the prefix node.

- Then we have to sort it in descending order so that the word with the highest frequency is shown at the top.

- then we get suggestion in the following order

for RAT - RATINU-15, RATIONAL-12, RATAN-10

- Time complexity to sort output - $O(K \log K)$
no. of top suggestions.

$$\Rightarrow \text{total time complexity} = O(L) + O(N) + O(K \log K)$$

\Rightarrow with this time complexity \rightarrow latency would be big.

↓
we need to do better

↓

we can precompute all the possible top suggestions for each & every prefix node in tree data structure.

↓

Now the time complexity to find suggestion for a given prefix would be - $O(L)$, where L is length of the prefix ~~entered~~.

How can we do this?

Store in Hash table
↓

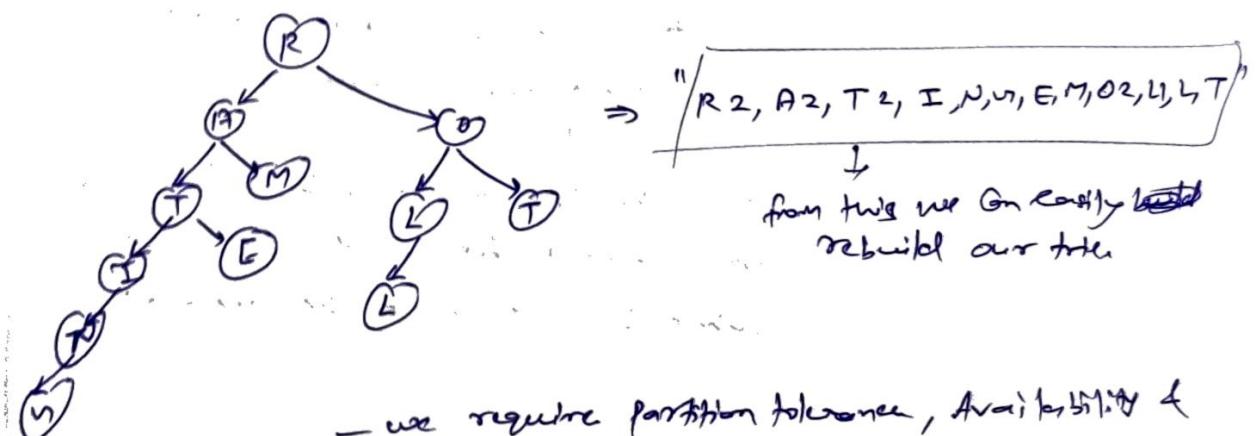
(Redis)

Keys Prefix
Value \rightarrow Suggestions

Prefix	Suggestion
R	{ RANDOM, ROUND, ROCKY } 20 19 17
RA	{ RANDOM, RATINN, RAM } 20 15 14
RAM	{ RAM, RAMIFY, RAMBLE } 14 17 5
RAN	{ RANDOM, RANT } 20 4
RAT	{ RATINN, RATIONAL, RATT } 15 12 10
RO	{ ROCKY } { ROUND, ROCKY, ROUND } 19 17 16
ROC	{ ROCKY }
ROUT	{ ROUNDY }

Storage

→ Storing entire tree in database can be done by taking snapshot of tree.



→ we require partition tolerance, Availability & scalability. Hence, most databases like Cassandra can be used.

\Rightarrow Traffic & storage estimation -

(1) Traffic estimation -

* 4B query/day

* each query consist 3 words (word \rightarrow 5 letters)

total read requests for one query = 3×5 request/query

total read request for one day = $3 \times 5 \times 4B = 60B$ queries

read req per second = $60 \times 10^9 / 10^5 (\text{sec}) = 60K \text{ req/sec.}$

(2) Storage estimation -

* 2 byte to store a character, total 30 bytes for 15char

* 20% of 4B queries to be unique.

$\therefore 800M \times 30 \text{ bytes} = 24 \text{ GB need to store.}$

- Assume 10y. new search query generated every day.

- 10y. of 24 GB = 240 GB

- For a year = $24 \times 365 = 8760 \text{ GB}$

- 10 years $\Rightarrow 8.76 \text{ TB of data in 10 years.}$

\Rightarrow Scaling trie database -

(2) will the trie data fit in one server?

\therefore initially 240 GB can fit in a server. but later may not because of increasing new search queries

\therefore need more servers with trie data stored in them.

Let's call these db servers "trie servers".

\Rightarrow Sharding - sharding the data, but how?

(approach) - dedicate tree for every alphabet

A trie serves for queries starting with A,
another " " " " " " " " " " " " B.

Problem - May lead to an imbalance in load. Many query may start with A but very less with

Start from Δ .

(approach 2) — we can add another server when one fills & then split the one file. whenever space on the server gets occupied, we start using the next free server.

for eg:- one server gets filled from A-abcd, then
next server or store from babc-cdef. E.g.
so on.

If user type A → it will hit two forney
aggregate results & return top 7 results to user.

Problems - Can still lead to an unbalanced load.

Approach 3) - Hash the words & $(hash\text{y}, \text{number of shard})$ will give seven numbers. This way query will be stored randomly. We can introduce Consistent hashing.

Problem - We have balanced tree servers, but we will need to aggregate the search results from each tree server & then aggregate root results & provide it to client.

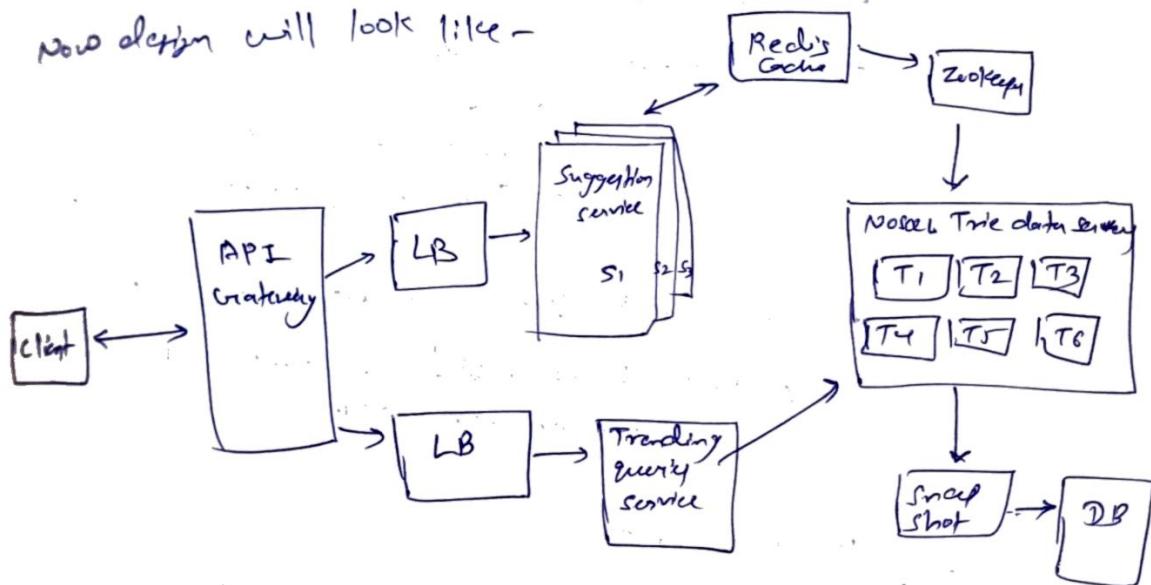
so here, solution 2 seems better. Aggregating from a few servers is better than aggregating results from each server. Hence we will use a mapping to see which server stores which range of prefix.

- we can use Zookeeper.

Zookeeper is very highly available & is very good with lots of reads & some writes.

⇒ Also we can cache popular searched prefixes.

Now system will look like -



② How will the Trending Query Service work?

It collects popular or recent queries & if the frequency of these new queries is more than threshold, it is stored in our trie server.

⇒ Log will look like

Query	Time	Weight
RAM	6th July 12pm to 1pm	173848
RAMIAY	6th July 1pm to 2PM	34869
PATNA	6th July 2 to 3PM	258488
ROUND	5th July	8234

Computation -

(Step 2) Collect Query Log

- * Queries are logged hourly with weights based on frequency & recency.
- * Older & low-frequency queries are discarded.
- * Daily aggregates are created for long-term efficiency.

(Step 2) MapReduce + Aggregator

- * "MapReduce jobs" run hourly to compute frequency ^{Score}
- * Aggregator service contains:
 - Old rank store
 - New frequency store
 - Time decay factor (recent query favored)
- * Queries with low score are discarded.
- * Output: "Rank table" containing trending queries with scores.

(Step 3) Apply to Trie

- * Specialized components called 'Appliers' update trie structures:-

- A1 → queries from qfr
- P2 → fs-01
- A2 → om-2

* Each Applier:-

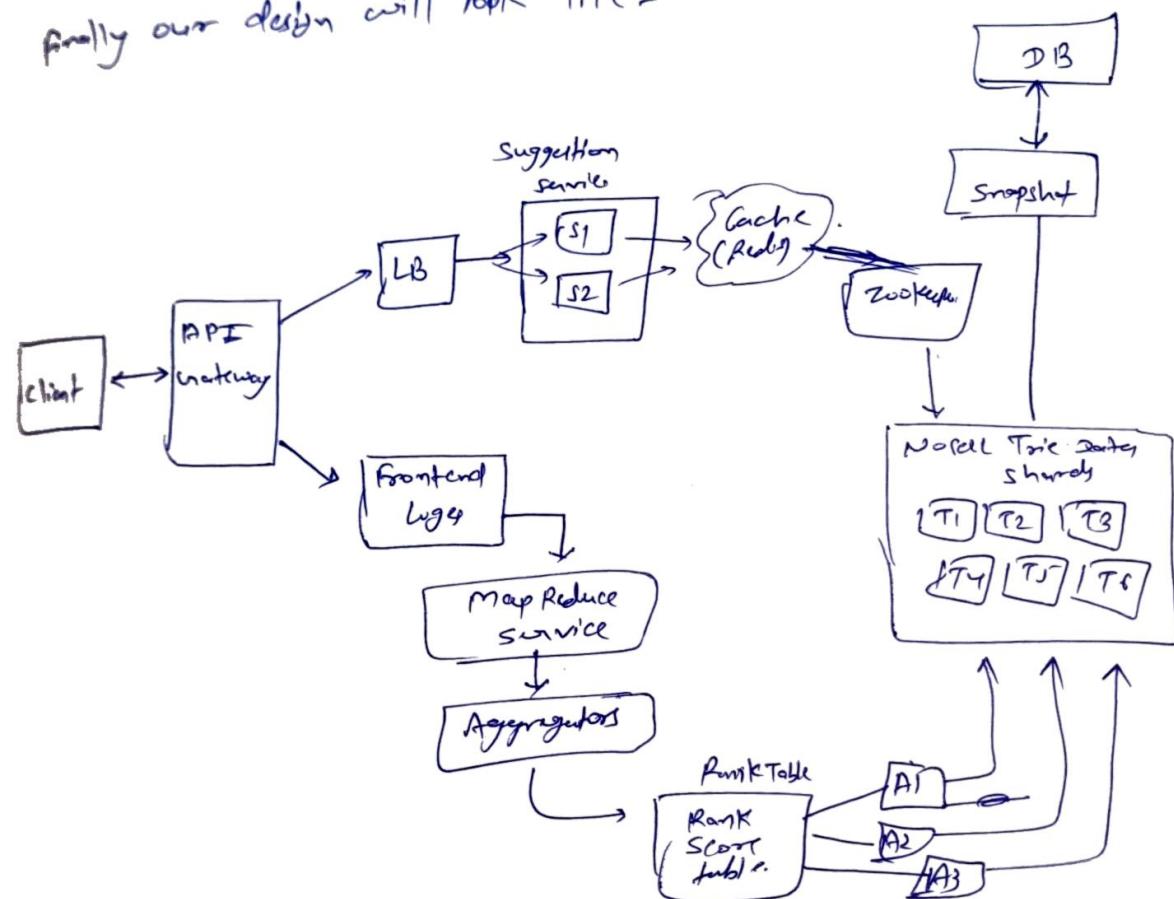
- Use rank table to assign weight
- Build internal trie with top-K suggestions per node.
- Dump into corresponding trie server.

Step 4) Update trie servers -

* done offline using master-slave configuration -

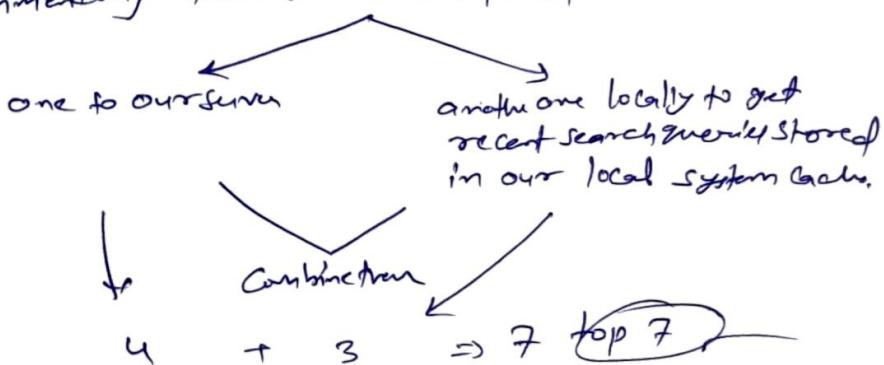
- Slave is updated while master serves traffic.
- After update, slave becomes the new master.

finally our design will look like -



③ Incorporating user's historical search queries:-

while recommending two queries are fired





Design a Distributed Job Scheduler like Airflow

What is a job scheduler:-

A job scheduler is a program that automatically schedules & executes jobs at specified times or intervals. It is used to automate repetitive tasks, run scheduled maintenance, or execute batch processes.

Key Terms:-

- (A) Task:- A task is the abstract concept of work to be done. For example, "send an email". Tasks are reusable & can be executed multiple times by different jobs.
- (B) Job:- A job is an instance of a task. It is made up of tasks to be executed, the schedule for when the task should be executed, & parameters needed to execute the task.

The main responsibility of a job scheduler is to take a set of jobs & execute them according to the schedule.

Functional Req:-

- ① User should be able to schedule jobs to be executed immediately, at a future date, or on a recurring schedule (i.e. every day at 10:00AM).

- ② User should be able to monitor the status of their jobs.

out of scope

- Cancel or reschedule jobs

Non Functional Req :-

Core Requirements:-

- ① Availability > Consistency
- ② System should execute jobs within 2s of their scheduled time.
- ③ Scalable to support up to 10K jobs per second.
- ④ The system ensure at-least-once execution of jobs.

out of scope

- enforce security policy
- CI/CD Pipeline

Core Entities -

① Task

② Job

③ Schedule

④ User

APIs -

① POST /jobs

```
{
  "jobId": "send-email",
  "schedule": "0 10 * * *",
  "parameters": {
    "to": "john@example.com",
    "subject": "Daily report"
  }
}
```

3 3

② GET /jobs ? userId=101 & status=[status] & startAndEnd
& endTimestamp

Data Flow

- ① A user schedules a job by providing the task to be executed, schedule for when the task should be executed, & the parameters needed to execute the task.
- ② The job is persisted in the system.
- ③ The job is picked up by a worker & executed at the scheduled time.
 - if the job fails, it is retried with exponential backoff
- ④ Update the job status in the status

~~#~~ High level design

- ① User should be able to schedule jobs to be executed immediately, at a future date, or on a recurring schedule.

:-

we will have 2 tables:-

a) Job table

```
{ "job_id": "123e4567-e89b", // Partition key for easy lookup by job_id  
  "user_id": "user-123",  
  "task_id": "sendemail",  
  "schedule": {  
    "type": "CRON" | "DATE"  
    "expression": "0 10 * * *", // Every day at 10:00 AM for Cron, specific date for DATE  
  },  
  "parameters": {  
    "to": "john@example.com",  
    "subject": "Daily Report"  
  }  
}
```

- ② Executions table tracks each individual time a job should run:-

```
{ "timebucket": 1715547600, // Partition key (Unix timestamp rounded down to hour)  
  "execution_time": "1715548800-123e4567-e89b",  
  "job_id": "user-123",  
  "status": "PENDING",  
  "attempt": 0  
}
```

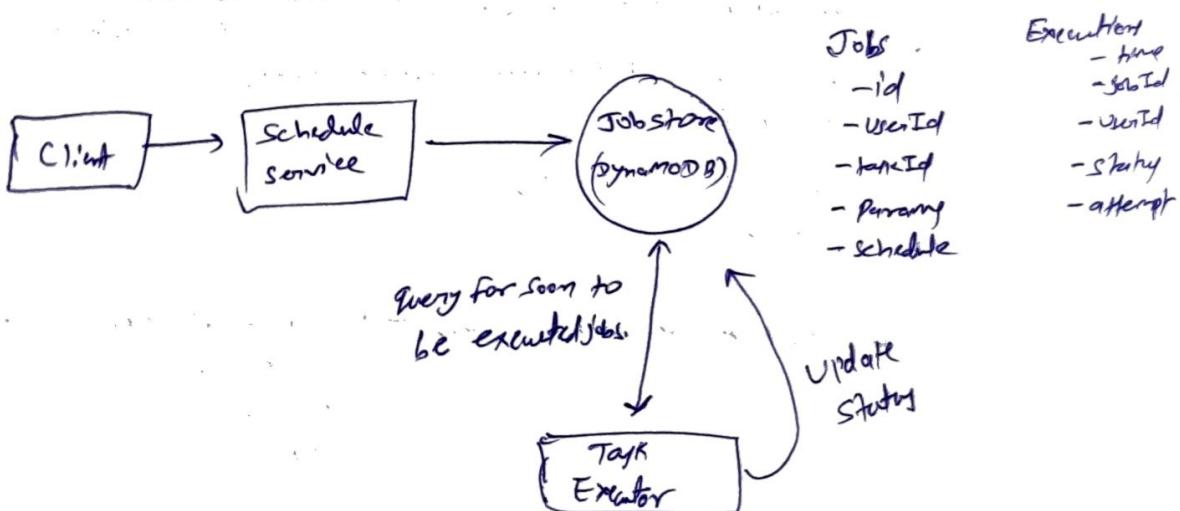
By using timebucket as partition key → achieved efficient querying while avoiding hot partition issues



for example, to find jobs that need to run in the next few minutes, we only need to query the current hour's bucket & possibly the next hour's bucket. The time bucket can be easily calculated:-

$$\text{time-bucket} = (\text{executionTime} / 3600) * 3600 \quad \begin{array}{l} \text{\# Round down} \\ \text{Return to nearest hour} \end{array}$$

- * First, by querying only 1-2 partitions to find upcoming jobs, we maintain efficient read operation. The hourly bucket ensure jobs are well-distributed across partitions, preventing hotspots.
 - * Second, when a recurring job completes, we can easily schedule its next occurrence by calculating the next execution time & creating a new entry in the "Execution" table. The job definition stays the same, but we keep creating new execution instances.
- ⇒ when a worker node is ready to execute jobs, it simply queries the "Execution" table for entries where:-
- * 'executionTime' is within the next few minutes,
 - * 'status' is PENDING.



② User should be able to monitor the status of their jobs, -

when a job is executed → need to update status on execution table.

COMPLETED,
FAILED,
IN-PROGRESS,
RETRYING

But how can we query for status of all jobs for a given user?

with current design ⇒ two-step process

- 1) find all job id for user from Jobs table,
- 2) query Execution table to find status of each job.

To solve this, we "add a Global Secondary Index (GSI) on the Execution table!"

- * Partition key + user_id
- * sort keys: execution_time + job_id

This GSI allows us to efficiently find all executions for a user, sort them by execution time, support pagination, & filter by status if needed.

Potential Deep Dive

① How can we ensure the system executes jobs within 2 sec. of their scheduled time?

⇒ To ensure jobs are executed with 2 seconds of their scheduled time, we must move away from our current model of querying the database every few minutes. Frequent querying (e.g. every 2 seconds) to

meet the timing precision would cause excessive database load & latency due to large payloads, making it impractical.

Instead, a more efficient solution is a two-layered scheduler architecture: (1) periodically query the database (e.g. every 5 minutes) for upcoming jobs, and (2) push these jobs into a message queue ordered by execution time. Workers then pull jobs from this queue for timely execution.

⇒ This design reduces database strain while leveraging message queue's ordering & throughput to achieve sub-second execution precision. However, handling newly created jobs scheduled within the next few minutes presents a challenge, as they could be missed by the periodic database query.

~~direct~~ Directly inserting these jobs into the queue doesn't help either due to FIFO constraints in systems like Kafka. To address this, the queue system must support dynamic prioritization based on execution time, allowing later-arriving jobs to be inserted in the correct position & executed on schedule.

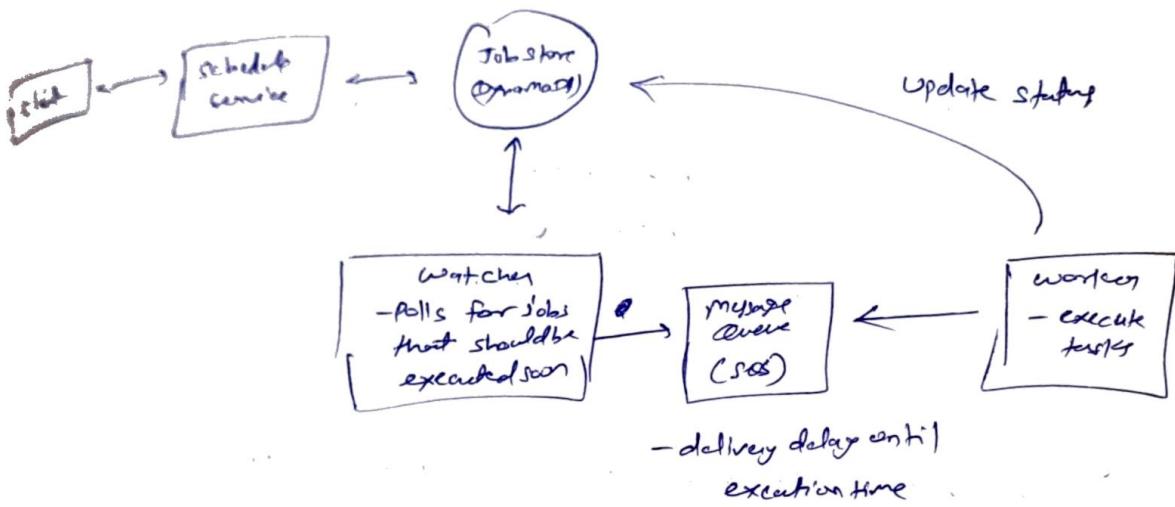
⇒ Cloud soln - Redis Sorted Set (issues → operations complexity in distributed system, implement own retry logic)

⇒ Cloud soln - RabbitMQ [issues → less reliable at scale]

⇒ Cloud soln - Amazon SQS (no native support for delayed message delivery.)

For e.g. schedule a job in 10 sec. ⇒ send message to SQS with a delay of 10sec ⇒ SQS will deliver msg to worker after delay. Easy as

Now we implement two layered schedule architecture.



② How can we ensure the system is ~~support~~ scalable to support up to 10K jobs per second?

- we can go left to right in design then look any bottlenecks & addressing them one-by-one.

③ Job creation - to handle high job create rate
↓
Introduce message queue. (allow to horizontally scale to schedule service)

④ Jobs DB -
↳ used DynamoDB with Partition key optimally
↳ All set here.

(10K job/sec 300 sec)

⑤ Message Queue Capacity -
for each 5 min window = process around 3 million jobs

- see ~~each~~ message byte upto 256 KB. (our message will be less ~ 200 bytes)

- total 600 MB data per 5 minute window

- default limit = 300 message/second (Amazon SQS)

↑ AWS can increase quota after specific request.

(d) workers -

two main choice of compute options -

- Container
- Lambda functions

→ Container with ECS & auto-scaling groups, ~~but~~
can be used as containers give us the best balance of
cost efficiency & operational simplicity while
meeting our performance needs.

(3) How can we ensure at-least-once execution of jobs?

Jobs can fail within a worker for one of two reasons:-

* Visible failure: - most likely bug in code or incorrect input parameters.

* Invisible failure: - most likely worker itself went down.

Bad job \Rightarrow Health check Policy

Crashing \Rightarrow Job leasing

creation of SQS visibility timeout (built-in mechanism)

when a worker receives a message from the queue, SQS automatically makes that message invisible to other workers for a configurable period. The worker processes the message & deletes it upon successful completion. If worker crashes or fails to process the message within the visibility timeout period, SQS automatically makes the message visible again for other workers to process.

- we can have relatively short visibility timeout (~30 sec).

After trying
3 times
again failed
then
� Shutdown
at FAILED

⇒ LastHy, one consequence of at-least-once execution is that we need to ensure our task code is idempotent. In other words, running the task multiple times should have the same outcome as running it just once.

* ways to handle this:-

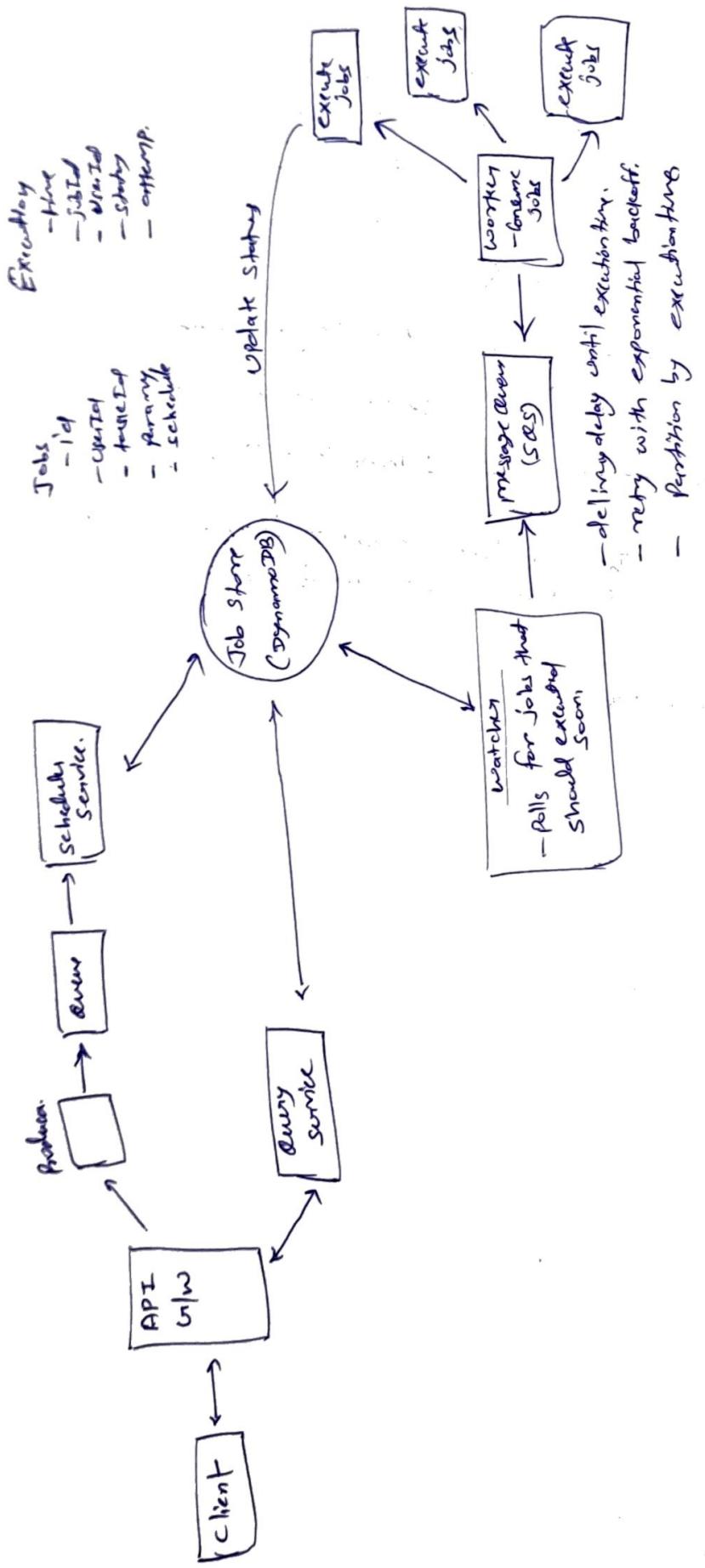
① ⇒ Bad Soln - No idempotency controls

⇒ Good Soln - Deduplication table

⇒ Correct Soln { Idempotent Job Design }

This involves incorporating idempotency keys & using conditional or state-aware operations.

⇒ HLD is on next page-----



Design a Distributed Cache # Like Redis

A distributed cache is a system that stores data as key-value pairs in-memory across multiple machines in a network.

Functional Req :-

- (1) Set, get, delete key-value pairs
- (2) User should be able to configure expiration time for key-value pairs.
- (3) Data should be evicted according to LRU policy.

Below the line

- User should be able to configure the cache size

Non-Functional Req :-

- (1) Highly available. Eventual consistency is acceptable.
- (2) Low latency operations (< 10ms for get/set operation)
- (3) Scalable to support the expected 1TB of data & 100K requests per second.

out of scope

- Durability
- Strong consistency guarantees
- Complex querying capability
- Transaction support

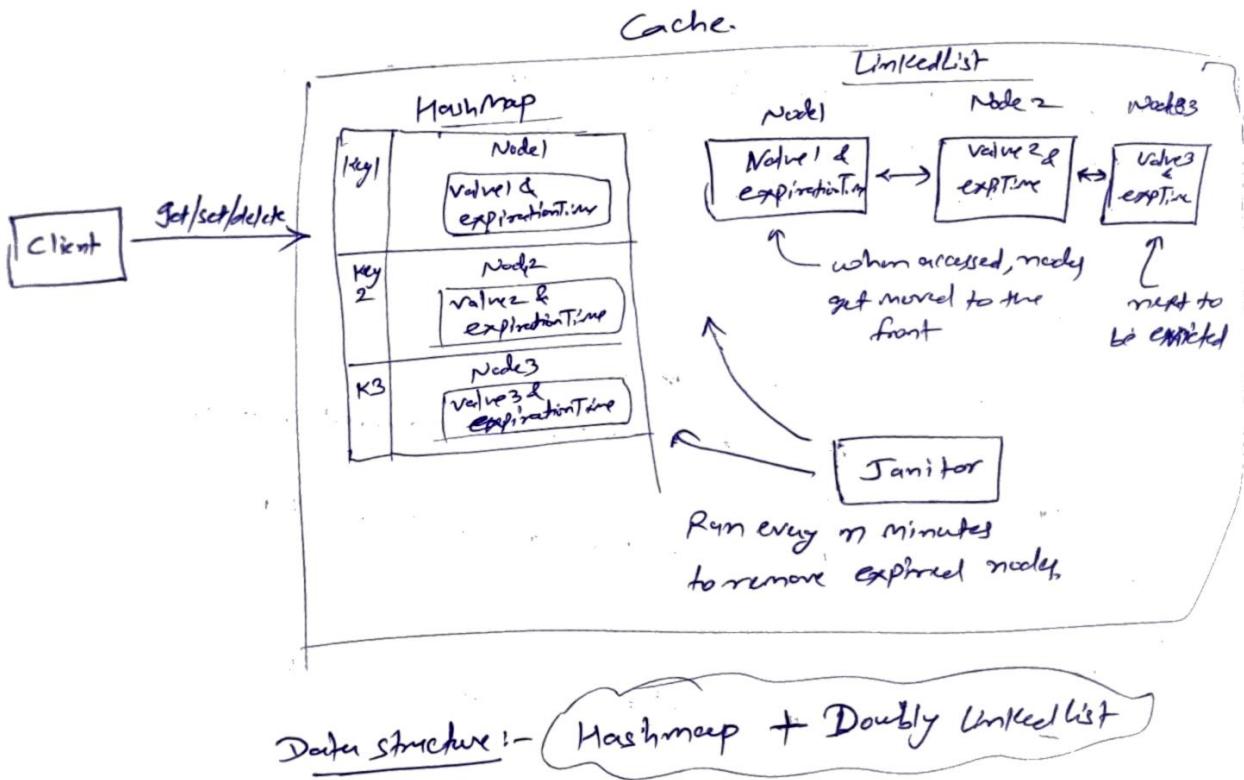
Entity

- Keys
- Values

APIs

- (1) POST /key
 1 "value": "..."
 2 "
- (2) PUT /key → {value: "..."}
- (3) DELETE /key

Basic HLD



Data structure:- Hashmap + Doubly Unlinked

Potential Deep Dive

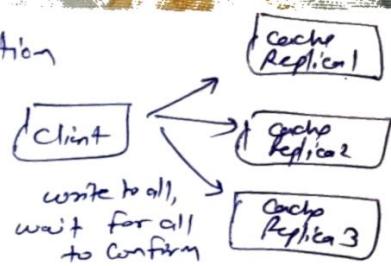
① How do we ensure our cache is highly available & fault-tolerant?

* Single node create issue when node fails.

We need multiple copies of our data spread across different nodes. But this opens up a whole new set of operations!

- How many copies should we keep?
- Which nodes should store the copies?
- How do we keep the copies in sync?
- What happens when node fail or quit communicating?

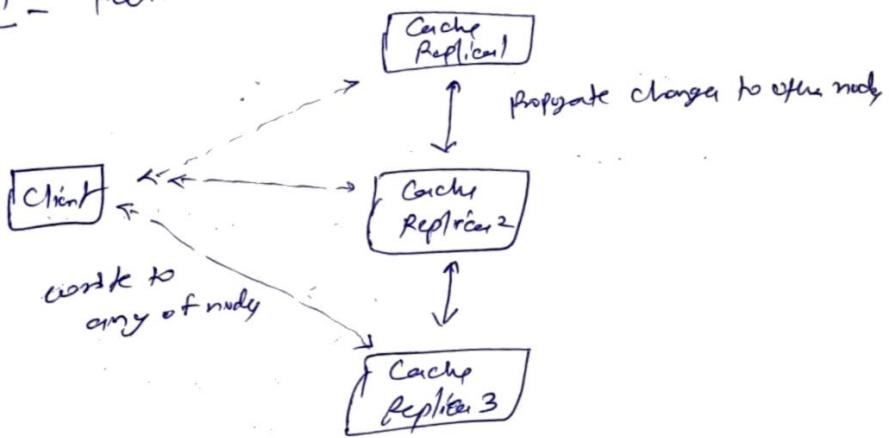
\Rightarrow Bad solⁿ - Synchronous replication



\Rightarrow Good solⁿ - Asynchronous replication

(New primary node)
(Replica may temporarily have stale data
need to select it until changes fully propagate through system)
(Primary fails)

\Rightarrow Great solⁿ - Peer-to-peer replication (using Gossip Protocol)



- maximum scalability / availability
- any node can read/write request serve

Cache is Scalable?

② How do we ensure our

\hookrightarrow load sharing

\downarrow
How we ensure even distribution of
keys across our nodes

\downarrow
Consistent Hashing

\downarrow
 $hash(key) \% n$

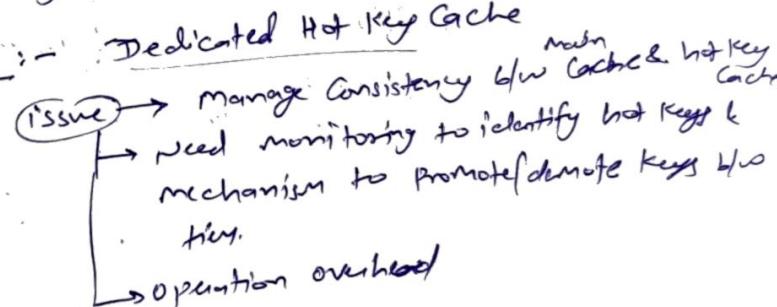
③ what happens if you have a hot key that is being read from a lot?

there are two distinct types of hot key problem we need to handle:-

- Hot Reads - key that receive an extremely high volume of read requests.
- Hot writes - keys that receive many concurrent write requests.

ways to handle Hot Reads :-

a) Vertical Scaling all nodes (Bad solⁿ)

b) Good solⁿ :- Dedicated Hot Key Cache


- Issue → Manage Consistency b/w main Cache & hot key Cache,
- need monitoring to identify hot keys & mechanism to promote/demote keys b/w tiers.
- operation overhead

c) Good solⁿ - Read Replicas ⇒ create multiple copies of the same data across diff nodes.

(Issue) is overhead in terms of storage & network bandwidth since entire nodes need to be replicated.

d) Great solⁿ - Copies of Hot Keys :-

↳ unlike read replicas which copy entire nodes, this approach selectively copies only the specific keys that are experiencing high read traffic.

How does it work? -



* First, system monitors key access pattern to detect hot keys that are frequently read.

* When a key becomes "hot", instead of having just one copy at User:123, the system creates multiple copies with different suffixes:-

User:123#1 → Node A Stores a Copy

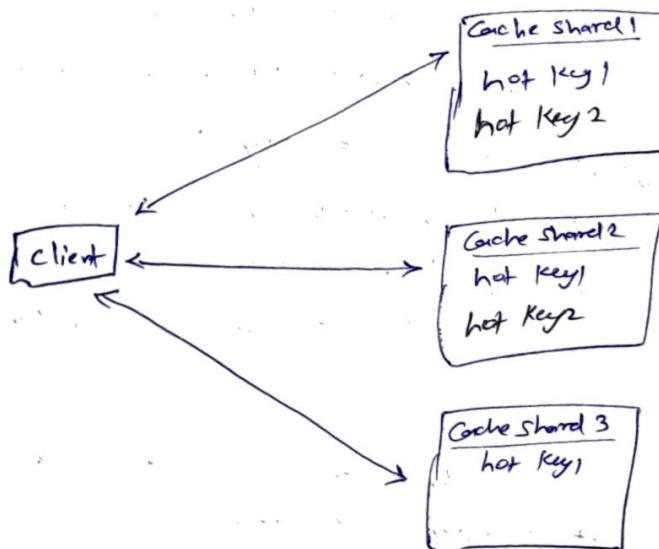
User:123#2 → Node B " " "

User:123#3 → Node C " "

* These copies get distributed to different nodes via Consistent hashing.

* For reads, clients randomly choose one of the suffix keys, spreading read load across multiple keys.

* For writes, the system must update all copies to maintain consistency.



- overhead in monitoring to detect hot keys & managing life cycle of copies - when to create them & when to remove them if a key is no longer hot. The approach works best when hot keys are primarily read-heavy with minimal writes.

(4) what happens if you have a hot key that is being written to a lot?

→ Great solⁿ - write Batching.

↳ Consider a viral video, receiving 10,000 views Update per second. Rather than execute 10,000 separate operation to set new value (eg view=1, view=2,...etc), write batching might collect these updates for 100ms, then execute a single operation to set the final value 1000 views.

↳ It reduces the work pressure on database.

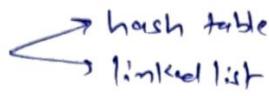
⇒ Great solⁿ - Sharding Hot Key with suffix.

- Instead of having a single counter or value that receives all writes, the system spreads write across multiple shards using a suffix strategy. For eg., a hot counter key "views:video123" might be split into 10 shards: "views:video123:1" through "views:video123:10". When a write arrives, the system randomly selects one of these shards to update.

- This approach distributes work load across multiple nodes in cluster.

- while read, need to sum values from all shards.

⇒ Trying it all together:-

Two data structures 

when it comes to scaling, it depends on which of the above approaches you choose, as you can end up with a slightly different design. Assuming you opted for:-

- 1) Asynchronous replication for high availability & handling hot key reads.
- 2) Consistent hashing for sharding & routing.
- 3) Random suffixes for distributing hot key writes across nodes.
- 4) write batching & connection pooling for decreasing network latency/overhead.

Final design:-

Client
- Batch writes
- Connection Pooling

- Consistent hashing to determine the right node
- random suffixes to spread hot key load across nodes

