

Grokking the System Design Interview

(/collection/5668639101419520/56490502)

87% completed



Q Search Course

- Designing a URL Shortening service like TinyURL (/courses/grokking-the-system-design-interview/m2ygV4E81AR)
- Designing Pastebin (/courses/grokking-the-system-design-interview/3jyvQ3pg6KO)
- Designing Instagram (/courses/grokking-the-system-design-interview/m2yDVZnQ8IG)
- Designing Dropbox (/courses/grokking-the-system-design-interview/m22Gymjp4mG)
- Designing Facebook Messenger (/courses/grokking-the-system-design-interview/B8R22v0wqJo)
- Designing Twitter (/courses/grokking-the-system-design-interview/m2G48X18NDO)
- Designing Youtube or Netflix (/courses/grokking-the-system-design-interview/xV26VjZ7yMI)
- Designing Typeahead Suggestion (/courses/grokking-the-system-design-interview/mE2XkgGRnmp)
- Designing an API Rate Limiter (/courses/grokking-the-system-design-interview/3jYKmrVAPGQ)
- Designing Twitter Search (/courses/grokking-the-system-design-interview/xV9mMjj74gE)
- Designing a Web Crawler (/courses/grokking-the-system-design-interview/NE5LpPrWrKv)
- Designing Facebook's Newsfeed (/courses/grokking-the-system-design-interview/gxpWJ3ZKYwl)
- Designing Yelp or Nearby Friends (/courses/grokking-the-system-design-interview/B8rpM8E16LQ)
- Designing Uber backend (/courses/grokking-the-system-design-interview/YQVkj548NM)
- Design Ticketmaster (*New*) (/courses/grokking-the-system-design-interview/YQyq6mBKq4n)
- Additional Resources (/courses/grokking-the-system-design-interview/...

Storage Capacity: Since we have 400 million new tweets every day and each tweet on average is 300 bytes, the total storage we need, will be:

$$400M * 300 \Rightarrow 120GB/day$$

Total storage per second:

$$120GB / 24hours / 3600sec \approx 1.38MB/second$$

4. System APIs

#

We can have SOAP or REST APIs to expose the functionality of our service; following could be the definition of the search API:

```
search(api_dev_key, search_terms, maximum_results_to_return, sort, page_token)
```

Parameters:

api_dev_key (string): The API developer key of a registered account. This will be used to, among other things, throttle users based on their allocated quota.

search_terms (string): A string containing the search terms.

maximum_results_to_return (number): Number of tweets to return.

sort (number): Optional sort mode: Latest first (0 - default), Best matched (1), Most liked (2).

page_token (string): This token will specify a page in the result set that should be returned.

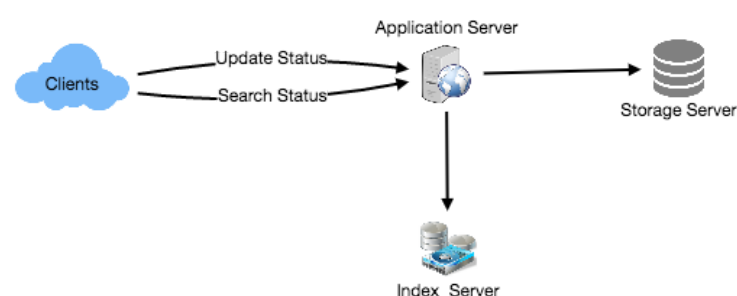
Returns: (JSON)

A JSON containing information about a list of tweets matching the search query. Each result entry can have the user ID & name, tweet text, tweet ID, creation time, number of likes, etc.

5. High Level Design

#

At the high level, we need to store all the statuses in a database and also build an index that can keep track of which word appears in which tweet. This index will help us quickly find tweets that users are trying to search.



High level design for Twitter search

🌐

Explore

(/explore)

📁

Tracks

(/tracks)

📖

My Courses

(/mycourses)

☕

Edpresso

(/edpresso)

👤+

Refer a Friend

(/refer-a-friend)

+

Create

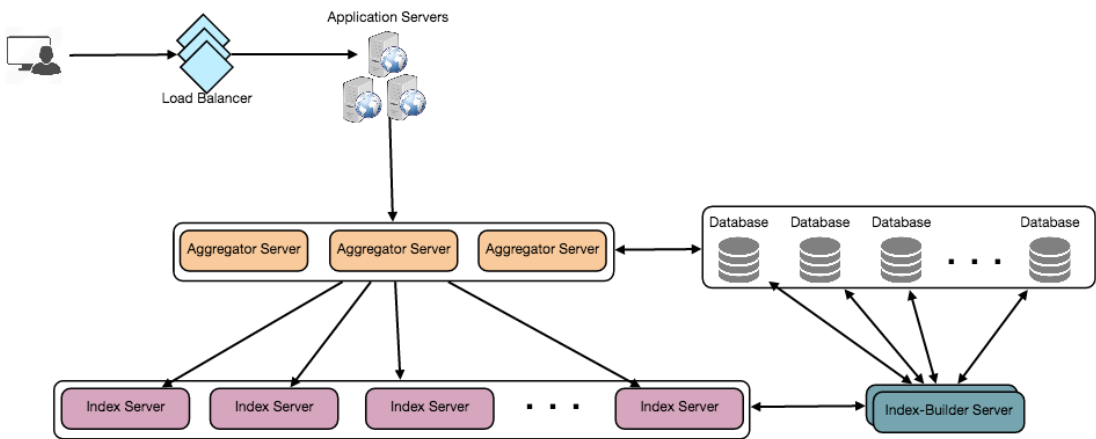
courses/grokking-the-system-design-interview/xV9mMjj7

Grokking the System Design Interview

(/collection/5668639101419520/56490502)

87% completed

- Designing a URL Shortening service like TinyURL (/courses/grokking-the-system-design-interview/m2ygV4E81AR)
- Designing Pastebin (/courses/grokking-the-system-design-interview/3jyvQ3pg6KO)
- Designing Instagram (/courses/grokking-the-system-design-interview/m2yDVZnQ8IG)
- Designing Dropbox (/courses/grokking-the-system-design-interview/m22Gymjp4mG)
- Designing Facebook Messenger (/courses/grokking-the-system-design-interview/B8R22v0wqJo)
- Designing Twitter (/courses/grokking-the-system-design-interview/m2G48X18NDO)
- Designing Youtube or Netflix (/courses/grokking-the-system-design-interview/xV26VjZ7yMI)
- Designing Typeahead Suggestion (/courses/grokking-the-system-design-interview/mE2XkgGRnmp)
- Designing an API Rate Limiter (/courses/grokking-the-system-design-interview/3jYKmrVAPGQ)
- Designing Twitter Search (/courses/grokking-the-system-design-interview/xV9mMjj74gE)
- Designing a Web Crawler (/courses/grokking-the-system-design-interview/NE5LpPrWrKv)
- Designing Facebook's Newsfeed (/courses/grokking-the-system-design-interview/gxpWJ3ZKYwl)
- Designing Yelp or Nearby Friends (/courses/grokking-the-system-design-interview/B8rpM8E16LQ)
- Designing Uber backend (/courses/grokking-the-system-design-interview/YQVkj548NM)
- Design Ticketmaster (*New*) (/courses/grokking-the-system-design-interview/YQyq6mBKq4n)
- Additional Resources (/courses/grokking-the-system-design-interview/NE5LpPrWrKv)



Detailed component design

7. Fault Tolerance

#

What will happen when an index server dies? We can have a secondary replica of each server and if the primary server dies it can take control after the failover. Both primary and secondary servers will have the same copy of the index.

What if both primary and secondary servers die at the same time? We have to allocate a new server and rebuild the same index on it. How can we do that? We don't know what words/tweets were kept on this server. If we were using 'Sharding based on the tweet object', the brute-force solution would be to iterate through the whole database and filter TweetIDs using our hash function to figure out all the required tweets that would be stored on this server. This would be inefficient and also during the time when the server was being rebuilt we would not be able to serve any query from it, thus missing some tweets that should have been seen by the user.

How can we efficiently retrieve a mapping between tweets and the index server? We have to build a reverse index that will map all the TweetID to their index server. Our Index-Builder server can hold this information. We will need to build a Hashtable where the 'key' will be the index server number and the 'value' will be a HashSet containing all the TweetIDs being kept at that index server. Notice that we are keeping all the TweetIDs in a HashSet; this will enable us to add/remove tweets from our index quickly. So now, whenever an index server has to rebuild itself, it can simply ask the Index-Builder server for all the tweets it needs to store and then fetch those tweets to build the index. This approach will surely be fast. We should also have a replica of the Index-Builder server for fault tolerance.

8. Cache

#

To deal with hot tweets we can introduce a cache in front of our database. We can use Memcached (<https://en.wikipedia.org/wiki/Memcached>), which can store all

