

System Design (/courses/system-design)

☐ Show Notes

/ System Design Interview Questions (/courses/system-design/topics/interview-questions/)

/ Design URL Shortener

**Design URL Shortener**

Bookmark

**Design a URL shortening service, like bit.ly**

Blog (<https://blog.interviewbit.com>) | About Us ([/pages/about\\_us/](/pages/about_us/)) | FAQ (</pages/faq/>) |

Contact Us ([/pages/contact\\_us/](/pages/contact_us/)) | Terms (</pages/terms/>) | Privacy Policy (</pages/privacy/>)

Programming — InterviewBit

<https://www.interviewbit.com/courses/programming/>

System Design Interview Questions (</courses/system-design/>) |

<http://bit.ly/2bRm86>

Google Interview Questions (</google-interview-questions/>) |

Facebook Interview Questions (</facebook-interview-questions/>) |

Amazon Interview Questions (</amazon-interview-questions/>) |

Microsoft Interview Questions (</microsoft-interview-questions/>) | Puzzles Questions (</puzzles/>)

**● Features.**Like Us (<https://www.facebook.com/interviewbit>)Follow Us ([https://twitter.com/interview\\_bit](https://twitter.com/interview_bit))

*This is the first part of any system design interview coming up with the features which the system should support. As an interviewee, you should try to list down all the features you can think of which our system should support. Try to spend around 2 minutes for this section in the interview. You can use the notes section alongside to remember what you wrote. ”*



**Shortening:** Take a url and return a much shorter url.  
 Got suggestions? We would love to hear your feedback.  Loved InterviewBit? Write us a testimonial.

Ex: <http://www.interviewbit.com/courses/programming/topics/time-complexity/>  
<http://go.gl/GU8w/> (<http://www.quora.com/What-is-your-review-of-InterviewBit>)

**Redirection:** Take a short url and redirect to the original url.

Ex: <http://goo.gl/GUKA8w> => <http://www.interviewbit.com/courses/programming/topics/time-complexity/>

**Custom url:** Allow the users to pick custom shortened url.

Ex: <http://www.interviewbit.com/courses/programming/topics/time-complexity/> => <http://goo.gl/ib-time>

**Analytics:** Usage statistics for site owner.

Ex: How many people clicked the shortened url in the last day?

**Gotcha:** What if two people try to shorten the same URL?

## ● Estimation:

“ This is usually the second part of a design interview, coming up with the estimated numbers of how scalable our system should be. Important parameters to remember for this section is the number of queries per second and the data which the system will be required to handle. Try to spend around 5 minutes for this section in the interview. ”

❓ ◀ **Q:** How many queries per second should the system handle?(Assuming 100 Million new URLs added each month)

**Hint:** Assuming average lifetime of a shortened URL is 2 weeks and 20% of websites creates 80% of the traffic, we see that we'll receive around 1 Billion queries in a month.

**A:** 400 queries per second in total. 360 read queries and 40 write queries per second.



Got suggestions ? We would love to hear your feedback.



Loved InterviewBit ? Write your testimonial.  
(<http://www.quora.com/What-is-your-review-of-InterviewBit>)

**Q:** How much data will we need to store so that we don't have to restructure our architecture for the next 5 years considering constant growth rate?

**Q:** How many URLs will we need to handle in the next 5 years?



**Hint:** Earlier we saw, we would see 100 Million new URLs each month. Assuming same growth rate for next 5 years, total URLs we will need to store will be  $100 \text{ Million} * 12 * 5 = 6 \text{ Billion}$ .

**A:** 6 Billion.

**Q:** What is the minimum length of shortened URL to represent 6 Billion URLs?



**Hint:** We will use (a-z, A-Z, 0-9) to encode our URLs. If we call  $x$  as minimum number of characters to represent 6 Billion total URLs, then will be the smallest integer such that  $x^{62} > 6 * 10^9$ .

**A:**  $\log(6 * 10^9)$  to the base 62 = 6

**A:** 3 Terabytes for URLs and 36 Gigabytes for shortened URLs

Note that for the shortened URL, we will only store the slug(6 characters) and compute the actual URL on the fly.



**Q:** Data read/written each second?

**Hint:** Data flow for each request would contain a shortened URL and the original URL. Shortened URL as we saw earlier, will take 6 bytes whereas the original URL can be assumed to take at most 500 bytes.

**A:** Written :  $40 * (500 + 6)$  bytes, Read :  $360 * (500 + 6)$  bytes



Got suggestions ? We would love to hear your feedback.



Loved InterviewBit ? Write us a testimonial.  
(<http://www.quora.com/What-is-your-review-of-InterviewBit>)

## ● Design Goals:

“

**Latency** - Is this problem very latency sensitive (Or in other words, Are requests with high latency and a failing request, equally bad?). For example, search typeahead suggestions are useless if they take more than a second.

**Consistency** - Does this problem require tight consistency? Or is it okay if things are eventually consistent?

**Availability** - Does this problem require 100% availability?

*There could be more goals depending on the problem. It's possible that all parameters might be important, and some of them might conflict. In that case, you'd need to prioritize one over the other.*”

❓ ◀ **Q:** Is Latency a very important metric for us?

**A:** Yes. Our system is similar to DNS resolution, higher latency on URL shortener is as good as a failure to resolve.



❓ ◀ **Q:** Should we choose Consistency or Availability for our service?

**A:** This is a tricky one. Both are extremely important. However, CAP theorem dictates that we choose one. Do we want a system that always answers correctly but is not available sometimes? Or else, do we want a system which is always available but can sometime say that a URL does not exist even if it does?

This tradeoff is a product decision around what we are trying to optimize. Let's say, we go with consistency here.



❓ ◀ **Q:** Not suggestions? We would love to hear your feedback. Try to list down other design goals?



Loved InterviewBit? Write us a testimonial.  
(<http://www.quora.com/What-is-your-review-of-InterviewBit>)

**A:** URL shortener by definition needs to be as short as possible. Shorter the shortened URL, better it compares to competition.



## ● Skeleton of the design:

“ The next step in most cases is to come up with the barebone design of your system, both in terms of API and the overall workflow of a read and write request. Workflow of read/write request here refers to specifying the important components and how they interact. Try to spend around 5 minutes for this section in the interview.

**Important :** Try to gather feedback from the interviewer here to indicate if you are headed in the right direction. ”

? ◀ **Q:** How should we define our APIs?

**A:**

ShorteningAPI(url) {store the url\_mapping and return hash(url)}

RedirectionAPI(hash) {redirect\_to url\_mapping[hash]}

Both APIs are very lightweight, their computation will not be the bottleneck plus we don't have to store any session information about users.

Basically, we are trying to build a service which serves as a huge HashMap



? ◀ **Q:** How would a typical write query look like?

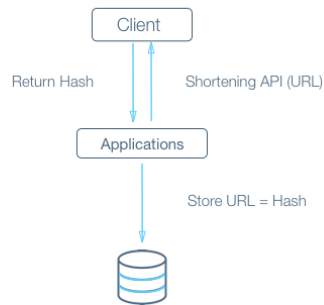
**A:** Components:



Client (Mobile app/Browser etc.) which calls ShorteningAPI(url) → Application server which interprets the API call and generates the shortened hash for the url

Database server which stores the hash => url mapping

#### HIGH LEVEL DESIGN (WRITE)



**Q:** How would a typical read query look like?

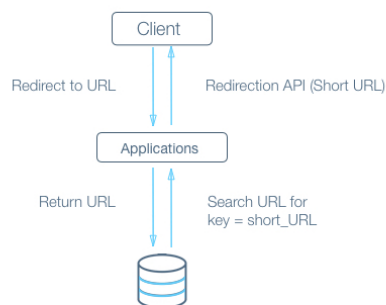
**A:** Components:

Client ( Mobile app / Browser , etc ) which calls RedirectionAPI(short\_url)

Application server which interprets the API call, extracts the hash from short\_url, asks database server for url corresponding to hash and returns the url.

Database server which stores the hash => url mapping

#### HIGH LEVEL DESIGN



Got suggestions ? We would love to hear your feedback.



Loved InterviewBit ? Write us a testimonial.  
(<http://www.quora.com/What-is-your-review-of-InterviewBit>)

## ● Deep Dive:

“ Lets dig deeper into every component one by one. Discussion for this section will take majority of the interview time(20-30 minutes). ”

❓ ◀ Q: How should we compute the hash of a URL?

**Gotcha:** How should we handle the case where two separate URLs get shortened to the same URL?



**A:** We can use a list of salts in case of collision.

For each read request, we can compute all possible shortened URLs using our list of salts and query for them in parallel to save time.

**A: `convert_to_base_62(md5(original_url + salt))[:6](first six characters)`**

Links: MD5-Wiki (<https://en.wikipedia.org/wiki/MD5>) Base 62 Conversion-Stackoverflow (<http://stackoverflow.com/questions/1119722/base-62-conversion-in-python>)

**Gotchas:**

Directly encoding URL to base\_62 will allow a user to check if a URL has been shortened already or not, reverse engineering can lead the user to the exact hash function used, this should not be allowed. Therefore randomization has to be introduced. Also, if two users shorten the same URL, it should result in two separate shortened URLs (for analytics)

Database ID encoded to base\_62 also won't be suitable for a production environment because it will leak information about the database. For example, patterns can be learnt to compute the growth rate (new URLs added per day) and in the worst case, copy the whole database.



Got suggestions ? We would love to hear your feedback.



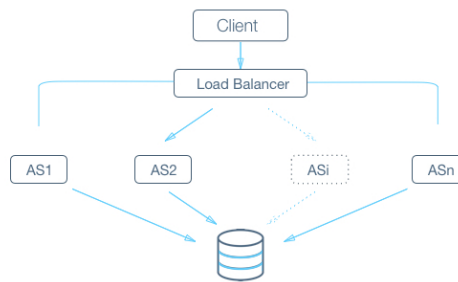
Loved InterviewBit ? Write your testimonial.  
(<http://www.quora.com/What-is-your-review-of-InterviewBit>)

**Q:** How would you take care of application layer fault tolerance?

**Q:** How do we handle the case where our application server dies?

**A:** The simplest thing that could be done here is to have multiple application servers. They do not store any data (stateless) and all of them behave the exact same way when up. So, if one of them goes down, we still have other application servers who would keep the site running.

HIGH LEVEL DESIGN



**Q:** How does our client know which application servers to talk to. How does it know which application servers have gone down and which ones are still working?

**A:** We introduce load balancers. Load balancers are a set of machines (an order of magnitude lower in number) which track the set of application servers which are active (not gone down). Client can send request to any of the load balancers who then forwarded the request to one of the working application servers randomly.



**A:** If we have only one application server machine, our whole service would become unavailable. Machines will fail and so will network. So, we need to plan for those events. Multiple application server machines along with load balancer.

**Got suggestions ? We would love to hear your feedback.**

**Loved InterviewBit ? Write us a testimonial.**

**(<http://www.quora.com/What-is-your-review-of-InterviewBit>)**



is the way to go.



**Q:** What all data should we store?

**A:** Data storage layer : Hash => URL mapping.

Billions of small sized(1kb) object. There is no relationship between objects.

We would also need to store data for analytics, for example, how many times was the url opened in the last hour?



**Q:** Should we choose RDBMS or NOSQL?

**Hint :** Things to consider :

Are joins required :

NoSQL databases are inefficient for joins.

In this case, assuming we don't need analytics, we only need to answer the query

"Given a hash, give me the corresponding URL" - a standard key to value lookup. As such, we don't need any joins here.

**Size of the DB :**

If the size of the data is so small that it fits on a single machine's main memory, SQL is a clear winner. SQL on a single machine has next to zero maintenance overhead and has great performance with right index built. If your index can fit into RAM, its best to go with a SQL solution. Let's analyze our current case :

*Assumptions :*

# of writes per month: 100 Million (Refer to estimations section)

Avg size of URL : 500 bytes

Provisioning for : 5 years

Space required :  $500 \text{ bytes} * 100M * 12 * 5 = 3TB$

3TB of data can fit on a single machine's hard disk. However, the index might not. If we store all data on a single machine, our write and read operations would be very slow (Page swaps for indices). Given read latency is critical for us, we can't store the data on a single machine.



**Got suggestions ? We would love to hear your**

feedback. So, a SQL solution will have a sharding overhead. Most NoSQL solutions however

are built with the assumption that the data does not fit on a single machine and (<http://www.quora.com/What-is-your-review-of-InterviewBit>)



**Loved InterviewBit ? Write us a testimonial.**

hence have sharding built in.

**A:** As discussed in hints, NoSQL would be a better fit for our case. Since we are optimizing for consistency over availability, we will choose a NoSQL DB which is highly consistent like HBase.



**Q:** What would the database schema look like?

**A:** We want to store the mapping from hash -> URL which is ideal for a key value store. In NoSQL domain, we need to be careful when designing the key as entries are indexed by the key. In our case, since we will only query by hash which will never update, our key will be hash with the value being the URL.



**Q:** How would we do sharding?

**A:** HBase inherently use consistent hashing for sharding. We explain it in detail at <https://www.interviewbit.com/problems/sharding-a-database/> (<https://www.interviewbit.com/problems/sharding-a-database/>). Since our total data size is a few TBs, we don't need to think about sharding across datacenters either.



**Q:** How would we handle a DB machine going down?

**A:** NoSQL again has single machine failure handling built in. A NoSQL DB like HBase would have an availability issue for a few seconds for the affected keys, as it tries to maintain tight consistency.



Got suggestions ? We would love to hear your feedback.



Loved InterviewBit ? Write us a testimonial. (<http://www.quora.com/What-is-your-review-of-InterviewBit>)

## ● Bonus Exercise :

🔍 ◀ **Q:** How can we optimize read queries?

**A:** Caching can be used to reduce average read time. Since most URL's are accessed for only a small amount of time after its creation, LRU cache fits our use case. Also, if a shortened URL goes viral, serving it through a cache will reduce the chances over loading the database. Redis and Memcached are examples of well known and widely used caches.

Redis-Tutorial (<http://try.redis.io/>) is an interactive tutorial for Redis.

Redis vs Memcached (<http://stackoverflow.com/questions/10558465/memcached-vs-redis>) lists down comparison between them.



🔍 ◀ **Q:** How would you share the data if you were working with SQL DB?

**Hint:** Consistent Hashing - Stackoverflow (<http://stackoverflow.com/questions/144360/simple-basic-explanation-of-a-distributed-hash-table-dht>)

**A:** We can use integer encoding of the shortened URL to distribute data among our DB shards.

Assuming we assign values from 0 to 61 to characters a to z, A to Z, and 0 to 9, we can compute the integer encoding of the shortened URL.

We can see that the maximum value of integer encoding will be less than  $10^{13}$ , which we can divide among our DB shards.

We will use consistent hashing to ensure we don't have to rehash all the data again if we add new DB shards later.



🔍 ◀ **Q:** How would you handle machine dying in the case of SQL DB?

**A:** This is a bit tricky. Obviously, for every shard, we need to have more than one machine

We can have a scheme better known as master-slave scheme, wherein there is one machine (master) which processes all writes and there is a slave machine

which just subscribes to all of the writes (<http://keepsquaringits.com/2014/01/01/what-is-your-review-of-interviewbit/>) so that if the master goes down, the slave can take over and start responding to the



Got suggestions? We would love to hear your



Loved InterviewBit? Write us a testimonial.


[System Design \(/courses/system-design/\)](/courses/system-design/)[Show Notes](#)[/ System Design Interview Questions \(/courses/system-design/topics/interview-questions/\)](/courses/system-design/topics/interview-questions/)[/ Design Search Typeahead](#)

## Design Search Typeahead

[Bookmark](#)

### Design a search typeahead ( Search autocomplete ) system at Google's scale.

Blog (<https://blog.interviewbit.com>) | [About Us \(/pages/about\\_us/\)](/pages/about_us/) | [FAQ \(/pages/faq/\)](/pages/faq/) | [Contact Us \(/pages/contact\\_us/\)](/pages/contact_us/) | [Terms \(/pages/terms/\)](/pages/terms/) | [Privacy Policy \(/pages/privacy/\)](/pages/privacy/)



**michael phelps**  
**michael system**  
**michael kors**

[System Design Interview Questions \(/courses/system-design/\)](/courses/system-design/) | [Google Interview Questions \(/google-interview-questions/\)](/google-interview-questions/) | [Facebook Interview Questions \(/facebook-interview-questions/\)](/facebook-interview-questions/) | [Amazon Interview Questions \(/amazon-interview-questions/\)](/amazon-interview-questions/) | [Microsoft Interview Questions \(/microsoft-interview-questions/\)](/microsoft-interview-questions/) | [Puzzles Questions \(/puzzles/\)](/puzzles/)

Press Enter to search.

#### Features:

[Like Us \(https://www.facebook.com/interviewbit\)](https://www.facebook.com/interviewbit)[Follow Us \(https://twitter.com/interview\\_bit\)](https://twitter.com/interview_bit)[Email \(mailto:hello@interviewbit.com\)](mailto:hello@interviewbit.com)

“ This is the first part of any system design interview, coming up with the features which the system should support. As an interviewee, you should try to list down all the features you can think of which our system should support. Try to spend around 2 minutes for this section in the interview. You can use the notes section alongside to remember what you wrote. ”



Got suggestions ? We would love to hear your

Q: How many typeahead suggestions are to be provided?

A: Let's assume 5 for this case.



Loved InterviewBit ? Write us a testimonial.

(<http://www.quora.com/What-is-your-review-of-InterviewBit>)

**Q:** Do we need to account for spelling mistakes ?

**A:** Example : Should typing *mik* give michael as a suggestion because michael is really popular as a query?

**Q:** What is the criteria for choosing the 5 suggestions ?

**A:** As the question suggests, all suggestions should have the typed phrase/query as the strict prefix. Now amongst those, the most relevant would be the most popular 5. Here, popularity of a query can be determined by the frequency of the query being searched in the past.

**Q:** Does the system need to be realtime ( For example, recent popular events like "Germany wins the FIFA worldcup" starts showing up in results within minutes ).

**A:** Let's assume that it needs to be realtime.

**Q:** Do we need to support personalization with the suggestions? ( My interests / queries affect the search suggestions shown to me).

**A:** Let's assume that we don't need to support personalization

## ● Estimation:

“ This is usually the second part of a design interview, coming up with the estimated numbers of how scalable our system should be. Important parameters to remember for this section is the number of queries per second and the data which the system will be required to handle. Try to spend around 5 minutes for this section in the interview. ”



Got suggestions ? We would love to hear your feedback.



Loved InterviewBit ? Write us a testimonial.

Clients can query my system for top 5 suggestions given a query prefix.  
(<http://www.quora.com/What-is-your-review-of-InterviewBit>)

Every search query done should feed into the system for an update.

Lets estimate the volume of each.

**Q:** How many search queries are done per day?

**A:** Assuming the scale of Google, we can expect around 2-4 Billion queries per day.



**Q:** How many queries per second should the system handle?

**A:** We can use the estimation from the last question here.

Total Number of queries : 4 Billion

Average length of query : 5 words = 25 letters ( Since avg length of english word is 5 letters ).

Assuming, every single keystroke results in a typeahead query, we are looking at an upper bound of  $4 \times 25 = 100$  Billion queries per day.



**Q:** How much data would we need to store?

**A:** Lets first look at the amount of new data we generate every day. 15% of the search queries are new for Google ( ~500 Million new queries ). Assuming 25 letters on average per query, we will 12.5G new data per day.

Assuming, we have accumulated queries over the last 10 years, the size would be  $12.5 \times 365 \times 10$  G which is approximately 50TB.



## Design Goals:



Got suggestions ? We would love to hear your feedback.



Loved InterviewBit ? Write us a testimonial.  
(<http://www.quora.com/What-is-your-review-of-InterviewBit>)

“

**Latency** - Is this problem very latency sensitive (Or in other words, Are requests with high latency and a failing request, equally bad?). For example, search typeahead suggestions are useless if they take more than a second.

**Consistency** - Does this problem require tight consistency? Or is it okay if things are eventually consistent?

**Availability** - Does this problem require 100% availability?

*There could be more goals depending on the problem. It's possible that all parameters might be important, and some of them might conflict. In that case, you'd need to prioritize one over the other. ”*

❓ ◀ **Q:** Is Latency a very important metric for us?

**A:** A big Yes. Search typeahead almost competes with typing speed and hence needs to have a really low latency.



❓ ◀ **Q:** How important is Consistency for us?

**A:** Not really important. If 2 people see different top 5 suggestions which are on the same scale of popularity, its not the end of the world. I, as a product owner, am happy as long as the results become eventually consistent.



❓ ◀ **Q:** How important is Availability for us?

**A:** Very important. If search typeahead is not available, the site would still keep working. However, it will lead to a much degraded experience.



**Got suggestions ? We would love to hear your feedback.**



**Loved InterviewBit ? Write us a testimonial.**  
(<http://www.quora.com/What-is-your-review-of-InterviewBit>)

## ● Skeleton of the design:

“ The next step in most cases is to come up with the barebone design of your system, both in terms of API and the overall workflow of a read and write request. Workflow of read/write request here refers to specifying the important components and how they interact. Try to spend around 5 minutes for this section in the interview.

**Important** : Try to gather feedback from the interviewer here to indicate if you are headed in the right direction. ”

❓ ◀ As discussed before, there are essentially 2 parts to this system :

Given a query, give me 5 most frequent search terms with the query as strict prefix

Given a search term, update the frequencies.

**Q:** What would the API look like for the client?

**A:**

Read: `List(string) getTopSuggestions(string currentQuery)`

Write: `void updateSuggestions(string searchTerm)`



❓ ◀ **Q:** What is a good data structure to store my search queries so that I can quickly retrieve the top 5 most popular queries?

**A:** For this question, we need to figure out top queries with another string as strict prefix. If you have dealt with enough string questions, you would realize a prefix tree (or trie) would be a perfect fit here.

Devil however lies in the details. We will dig deeper into the nitty gritty of this in the next section.



Got suggestions ? We would love to hear your feedback.



Loved InterviewBit ? Write us a testimonial.  
(<http://www.quora.com/What-is-your-review-of-InterviewBit>)



**Q:** How would a typical read query look like?

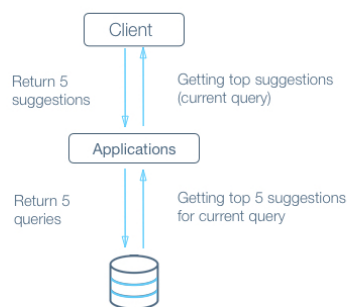
**A:** Components:

Client ( Mobile app / Browser, etc ) which calls  
getTopSuggestions(currentQuery)

Application server which interprets the API call and queries the database for the  
corresponding top 5 queries.

Database server which looks up the top queries in the trie.

HIGH LEVEL DESIGN (READ)



**Q:** How would a typical write query look like?

**A:** Components:

Client ( Mobile app / Browser, etc ) which calls updateSuggestions(searchTerm)

Application server which interprets the API call and forwards the searchTerm to  
database for update.

Database server which updates its trie using the searchTerm

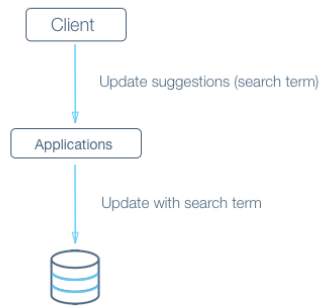


Got suggestions ? We would love to hear your  
feedback.



Loved InterviewBit ? Write us a testimonial.  
(<http://www.quora.com/What-is-your-review-of-InterviewBit>)

## HIGH LEVEL DESIGN



## ● Deep Dive:

“ Lets dig deeper into every component one by one. Discussion for this section will take majority of the interview time(20-30 minutes). ”

❓ ◀ Lets dig deeper into every component one by one.

### Application layer:

Think about all details/gotchas yourself before beginning.

**Q:** How would you take care of application layer fault tolerance?

**Q:** How do we handle the case where our application server dies?



**A:** The simplest thing that could be done here is to have multiple application server. They do not store any data (stateless) and all of them behave the



Got suggestions ? We would love to hear your feedback.

exact same way when up. So, if one of them goes down, we still have other application servers who would keep the site running.  
(<http://www.quora.com/What-is-your-review-of-InterviewBit>)



Loved InterviewBit ? Write us a testimonial.  
(<http://www.quora.com/What-is-your-review-of-InterviewBit>)

**Q:** How does our client know which application servers to talk to. How does it know which application servers have gone down and which ones are still working?



**A:** We introduce load balancers. Load balancers are a set of machines (an order of magnitude lower in number) which track the set of application servers which are active ( not gone down ). Client can send request to any of the load balancers who then forward the request to one of the working application servers randomly.

**A:** If we have only one application server machine, our whole service would become unavailable. Machines will fail and so will network. So, we need to plan for those events. Multiple application server machines along with load balancer is the way to go.



### Database layer:

Let's first dig deeper into the trie we talked about earlier.

**Q:** How would a read query on the trie work?

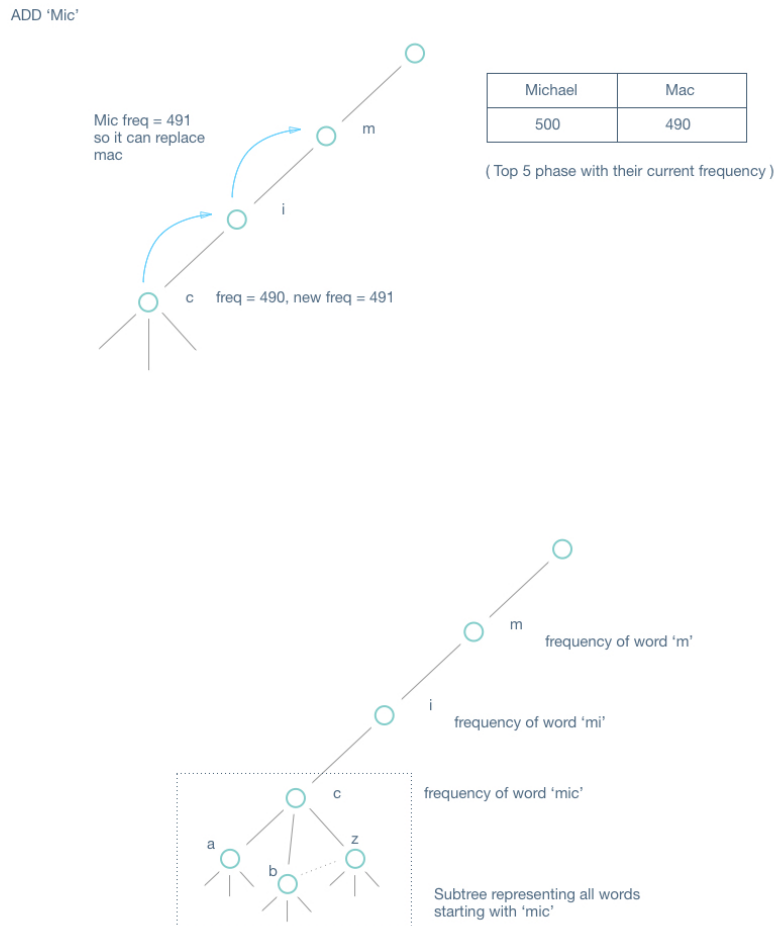
**A:** The read query would require us to fetch the top 5 results per query. A traditional trie would store the frequency of the search term ending on the node n1 at n1. In such a trie, how do we get the 5 most frequent queries which have the search term as strict prefix. Obviously, all such frequent queries would be the terms in the subtree under n1 ( as shown in diagram ).

So, a brute force way is to scan all the nodes in the subtree and find the 5 most frequent. Lets estimate the number of nodes we will have to scan this way.



Got suggestions? We would love to hear your feedback. (http://www.quora.com/What-is-your-review-of-InterviewBit) Writeups a testimonial. InterviewBit

Also, the high latency does not align with our design goals.



**Q:** How can we modify the trie so that reads become super efficient?

**Hint :** Store more data on every node of the trie.

**A:** Storage is cheap. Lets say we were allowed to store more stuff on each node. How would we use the extra storage to reduce the latency of answering the query.

A good choice would be storing the top 5 queries for the prefix ending on node n1 at n1 itself. So, every node has the top 5 search terms from the subtree below it. The read operation becomes fairly simple now. Given a search prefix,



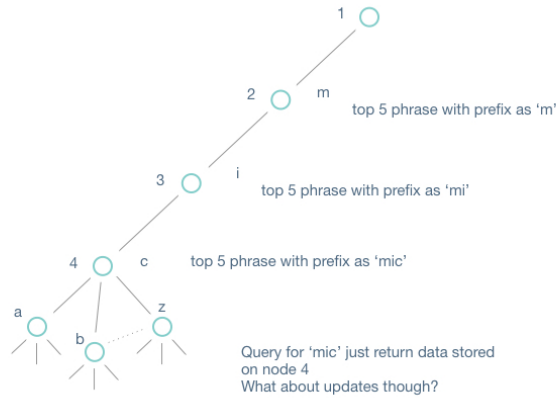
Got suggestions ? We would love to hear your

feedback.



Loved InterviewBit ? Write us a testimonial.

(<http://www.quora.com/What-is-your-review-of-InterviewBit>)



**Q:** How would a typical write work in this trie?

**A:** So, now whenever we get an actual search term, we will traverse down to the node corresponding to it and increase its frequency. But wait, we are not done yet. We store the top 5 queries in each node. Its possible that this particular search query jumped into the top 5 queries of a few other nodes. We need to update the top 5 queries of those nodes then. How do we do it then? Truthfully, we need to know the frequencies of the top 5 queries ( of every node in the path from root to the node ) to decide if this query becomes a part of the top 5.

There are 2 ways we could achieve this.

Along with the top 5 on every node, we also store their frequency. Anytime, a node's frequency gets updated, we traverse back from the node to its parent till we reach the root. For every parent, we check if the current query is part of the top 5. If so, we replace the corresponding frequency with the updated frequency. If not, we check if the current query's frequency is high enough to be a part of the top 5. If so, we update the top 5 with frequency. On every node, we store the top pointer to the end node of the 5 most frequent queries ( pointers instead of the text ). The update process would involve comparing the current query's frequency with the 5th lowest node's frequency and update the node pointer with the current query pointer if the new frequency is greater.



**Got suggestions ? We would love to hear your feedback.**





**Loved InterviewBit ? Write us a testimonial.**  
(<http://www.quora.com/What-is-your-review-of-InterviewBit>)

  **Q:** Can frequent writes affect read efficiency?

**A:** Yes, potentially. If we are updating the top 5 queries or the frequencies very frequently, we will need to take a lock on the node to make sure the reader thread does not get an inconsistent value. As such, writes start to compete with reads.



  **Q:** What optimizations can we do to improve read efficiency?

**Q:** Can we use sampling?



**A:** Yes. If we assume Google's scale, most frequent queries would appear 100s of times in an hour. As such instead of using every query to update, we can sample 1 in 100 or 1 in 1000 query and update the trie using that.

**Q:** Offline update?



**A:** Again if we assume that most queries appearing in the search typeahead would appear 100s of times in an hour, we can have an offline hashmap which keeps maintaining a map from query to frequency. Its only when the frequency becomes a multiple of a threshold that we go and update the query in the trie with the new frequency. The hashmap being a separate datastore would not collide with the actual trie for reads.

**A:** As mentioned earlier, writes compete with read. Sampling writes and Offline updates can be used to improve read efficiency.



**Got suggestions ? We would love to hear your**

  **Q:** What if I use a separate trie for updates and copy it over to the active one  
feedback.



**Loved InterviewBit ? Write us a testimonial.**

(<http://www.quora.com/What-is-your-review-of-InterviewBit>)

periodically?

**A:** Not really, there are 2 major problems with this approach.

You are not realtime anymore. Lets say you copy over the trie every hour. Its possible a search term became very popular and it wasn't reflected for an hour because it was present in the offline trie and did not appear till it was copied to the original trie

The trie is humungous. Copying over the trie can't be an atomic operation. As such, how would you make sure that reads are still consistent while still processing incoming writes?



**Q:** Would all data fit on a single machine?

**A:** Refer to estimations section. We would need to store more than 50TB of data.

Ideally, we would want most of it in memory to help with the latency. Thats a lot to ask from a single machine. We will go with a "No" here.



**Q:** Alright, how do we shard the data then?

**Q:** Would we only shard on the first level?



**A:** The number of shards could very well be more than the number of branches on first level(26). We will need to be more intelligent than just sharding on first level.

**Q:** What is the downside of assigning one branch to a different shard?



Got suggestions ? We would love to hear your



Loved InterviewBit ? Write us a testimonial.

**A:** Load imbalance. Storage imbalance. Some letters are more frequent than the others. For example, letters starting with 'a' are more likely than

feedback

(<http://www.quora.com/What-is-your-review-of-InterviewBit>)

letters starting with 'x'. As such, we can run into cases of certain shards running hot on load. Also, certain shards will have to store more data because there are more queries starting with a certain letter. Another fact in favor of sharding a little more intelligently.

**A:** Lets say we were sharding till the second or third level and we optimize for load here. Lets also say that we have the data around the expected load for every prefix.

We keep traversing the 2 letter prefixes in order ('a', 'aa', 'ab', 'ac',...) and break when the total load exceeds an threshold load and assign that range to a shard. We will need to have a master which has this mapping with it, so that it can route a prefix query to the correct shard.



**Q:** How would we handle a DB machine going down?

**A:** As we discussed earlier, availability is more important to us than consistency. If thats the case, we can maintain multiple replica of each shard and an update goes to all replicas. The read can go to multiple replicas (not necessarily all) and uses the first response it gets. If a replica goes down, reads and writes continue to work fine as there are other replicas to serve the queries.

The issue occurs when this replica comes back up. There are 2 options here : If the frequency of the replica going down is lower or we have much higher number of replicas, the replica which comes back up can read the whole data from one of the older working replica while keeping the new incoming writes in a queue.

There is a queue with every server which contains the changelog or the exact write query being sent to them. The replica can request any of the other replicas in its shard for all changelog since a particular timestamp and use that to update its trie.



Got suggestions ? We would love to hear your feedback.



Loved InterviewBit ? Write us a testimonial.  
(<http://www.quora.com/What-is-your-review-of-InterviewBit>)





You have now mastered this problem!

## Discussion

Post a comment ([https://discuss.interviewbit.com/session/sso?return\\_path=https://discus](https://discuss.interviewbit.com/session/sso?return_path=https://discus)

T

thecoder2016

about 1 year ago

**What would be the schema for storing the trie ? Would we use SQL or no-SQL databa 2**

[Reply](https://discuss.interviewbit.com/session/sso?return_path=https://discuss.interviewbit.com/t/what-would-be-the-schema-for-storing-the-trie-would-we-use-sql-or-no-sql-databa) ([https://discuss.interviewbit.com/session/sso?return\\_path=https://discuss.interviewbit.com/t/what-would-be-the-schema-for-storing-the-trie-would-we-use-sql-or-no-sql-databa](https://discuss.interviewbit.com/session/sso?return_path=https://discuss.interviewbit.com/t/what-would-be-the-schema-for-storing-the-trie-would-we-use-sql-or-no-sql-databa))

H

himanshu\_setia

over 1 year ago

**But How would we store Trie in DB? 2**

[Reply](https://discuss.interviewbit.com/session/sso?return_path=https://discuss.interviewbit.com/t/but-how-would-we-store-trie-in-db) ([https://discuss.interviewbit.com/session/sso?return\\_path=https://discuss.interviewbit.com/t/but-how-would-we-store-trie-in-db](https://discuss.interviewbit.com/session/sso?return_path=https://discuss.interviewbit.com/t/but-how-would-we-store-trie-in-db))



Got suggestions ? We would love to hear your

D

feedback.



Loved InterviewBit ? Write us a testimonial.

(<http://www.quora.com/What-is-your-review-of-InterviewBit>)

System Design (/courses/system-design/)

Show Notes

/ Storage Scalability (/courses/system-design/topics/storage-scalability/)

/ Highly Consistent Database

## Highly Consistent Database

Bookmark

Design a distributed key value store which is highly consistent and is network partition tolerant.

Blog (https://blog.interviewbit.com/) | About Us (/pages/about\_us/) | FAQ (/pages/faq/)

Contact Us (/pages/contact\_us/) | Terms (/pages/terms/) | Privacy Policy (/pages/privacy/)

### Features:

“ System Design Interview Questions (/courses/system-design/) |  
This is the first part of any system design interview, coming up with the  
features which the system should support. As an interviewee, you should try  
to list down all the features you can think of which our system should support.  
Try to spend around 2 minutes for this section in the interview. You can use  
the notes section alongside to remember what you wrote. ”

Microsoft Interview Questions (/microsoft-interview-questions/) | Puzzles Questions (/puzzles/)

Q: What is the amount of data that we need to store?



Like Us (https://www.facebook.com/interviewbit)



Follow Us (https://twitter.com/interview\_bit)



Email (mailto:hello@interviewbit.com)

Q: Do we need to support updates?

A: Yes.

Q: Can the size of the value for a key increase with updates?

A: Yes. In other words, it's possible a sequence of keys could co-exist on one



Got suggestions? We would love to hear your

feedback. server previously, but with time, they grew to a size where all of them don't fit on

a single machine.



Loved InterviewBit? Write us a testimonial.

(http://www.quora.com/What-is-your-review-of-InterviewBit)

**Q:** Can a value be so big that it does not fit on a single machine?

**A:** No. Let's assume that there is an upper cap of 1GB to the size of the value.

**Q:** What would the estimated QPS be for this DB?

**A:** Let's assume around 100k

## ● Estimation:

“ This is usually the second part of a design interview, coming up with the estimated numbers of how scalable our system should be. Important parameters to remember for this section is the number of queries per second and the data which the system will be required to handle. Try to spend around 5 minutes for this section in the interview. ”

- ❓ ◀ Total storage size : 100 TB as estimated earlier  
Total estimated QPS : Around 10M

**Q:** What is the minimum number of machines required to store the data?

**A:** Assuming a machine has 10TB of hard disk, we would need minimum of  $100\text{TB} / 10\text{ TB} = 10$  machines to store the said data. Do note that this is bare minimum. The actual number might be higher if we decide to have replication or more machines in case we need more shards to lower the QPS load on every shard.



Got suggestions ? We would love to hear your

## ● Design Goals:



Loved InterviewBit ? Write us a testimonial.

(<http://www.quora.com/What-is-your-review-of-InterviewBit>)



“

**Latency** - Is this problem very latency sensitive (Or in other words, Are requests with high latency and a failing request, equally bad?). For example, search typeahead suggestions are useless if they take more than a second.

**Consistency** - Does this problem require tight consistency? Or is it okay if things are eventually consistent?

**Availability** - Does this problem require 100% availability?

*There could be more goals depending on the problem. It's possible that all parameters might be important, and some of them might conflict. In that case, you'd need to prioritize one over the other.*”

  **Q:** Is Latency a very important metric for us?

**A:** No, but it would be good to have a lower latency.



  **Q:** Consistency vs Availability?

**A:** As the question states, we need tight consistency and partitioning. Going by the CAP theorem (Nicely explained at <http://robertgrain.com/2014/08/cap-theorem-revisited/> (<http://robertgrain.com/2014/08/cap-theorem-revisited/>)), we would need to compromise with availability if we have tight consistency and partitioning. As is the case with any storage system, data loss is not acceptable.



## ● Deep Dive:

“

*Lets dig deeper into every component one by one. Discussion for this*



Got suggestions? We would love to hear your feedback. (Loved InterviewBit? Write us a testimonial. section will take majority of the interview time(20-30 minutes).)

feedback.

(<http://www.quora.com/What-is-your-review-of-InterviewBit>)

“ Note : In questions like these, the interviewer is looking at how you approach designing a solution. So, saying that I'll use a NoSQL DB like HBase is not an ideal answer. It is okay to discuss the architecture of HBase for example with rationale around why some components were designed the way they were.”

**Q:** Is sharding required?

**A:** Let's look at our earlier estimate about the data to be stored. 100TB of data can't be stored on a single machine.

Let's say that we somehow have a really beefy machine which can store that amount of data, that machine would have to handle all of the queries ( All of the load ) which could lead to a significant performance hit.

“ Tip: You could argue that there can be multiple copies of the same machine, but this would not scale in the future. As my data grows, its possible that I might not find a big beefy enough machine to fit my data. ”

So, the best course of action would be to shard the data and distribute the load amongst multiple machines.



❓ ◀ **Q:** Should the data stored be normalized?

<http://www.studytonight.com/dbms/database-normalization.php>

(<http://www.studytonight.com/dbms/database-normalization.php>)

**Q:** Can I shard the data so that all the data required for answering my most frequent queries live on a single machine?



**A:** Most applications are built to store data for a user ( consider messaging for example. Every user has his / her own mailbox ). As such, if you shard

based on every user as a row, its okay to store data in a denormalized fashion so that you won't have to query information across users. In this

case, lets say we go with storing data in denormalized fashion.



Got suggestions ? We would love to hear your

feedback

Loved InterviewBit ? Write us a testimonial.

(<https://www.quora.com/What-is-your-review-of-InterviewBit>)

**A:** If the data is not normalized, then we need to join across tables and across rows to fetch data. If the data is already shared across machine, any join across machines is highly undesirable ( High latency, Less indexing support ).

With storing denormalized information however, we would be storing the same fields at more than one place. However, all information related to a row ( or a key ) would be on the same machine. This would lead to lower latency.

However, if the sharing criteria is not chosen properly, it could lead to consistency concerns ( After all, we are storing the same data at multiple places ).



**Q:** How many machines per shard? How does a read / write look in every shard?

**Q:** Can we keep just one copy of data?



**A:** Since there is only one copy of the data, reading it should be consistent. As long as there are enough shards to ensure a reasonable load on each shard, latency should be acceptable as well. Reads and writes would work exactly how they work with a single DB just that there would be a row -> shard -> machine IP ( Given a row, tell me the shard it belongs to and then given the shard, give me the machine I should be querying / writing to ) resolution layer in between.

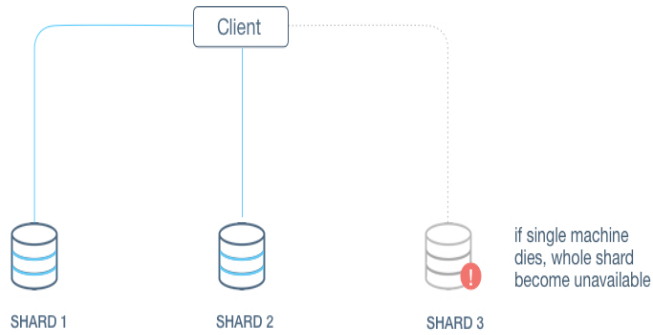
There is just one tiny problem with this model. What if the machine in the shard goes down? Our shard will be unavailable ( which is fine as governed by the CAP theorem ). However, what if the machine dies and its hard disk becomes corrupt. We suddenly run into the risk of losing the data which is not acceptable. Imagine losing all your messages because your shard went down and the hard disk got corrupted. That means we definitely need more than one copy of data being written with us.



Got suggestions ? We would love to hear your feedback.



Loved InterviewBit ? Write us a testimonial.  
(<http://www.quora.com/What-is-your-review-of-InterviewBit>)



**Q:** What problem may arise if we keep multiple copies of data?



**A:** Let's say we keep 3 copies of data ( The probability of all 3 machines dying and having a corrupted disk is negligible ). Now, the issue is how do we maintain all of the copies in absolute sync ( Consistency, remember ? ). One naive way would be that a write would not succeed unless its written to all 3 copies / machines. That would make our write latency go up significantly apart from making writes very unreliable ( My write fails if it fails on any of the machines ). Let's see if we can make this a bit better .

If we have to allow writes succeeding even if the write has been written on a majority of machines (2 out of 3, let's say), to maintain consistency, its important that there is a master machine which keeps track of this information. This master machine can track which machines have a particular block in each shard. This means that every read will go through this master machine, figure out the machines with the block and query for the required block. The machines which do not have the block can check with this master machine to see which block are not present on it, and catch up to replicate the block on it.

However , now if this master machine dies, our whole shard is unavailable till this machine comes back up. If this machine has a corrupted hard disk, then the unavailability becomes indefinite ( Note that we do not lose data in this case, as total data is the union of data present on 3 nodes ). This is not an



Got suggestions? We would love to hear your: improvement ideas in the questions

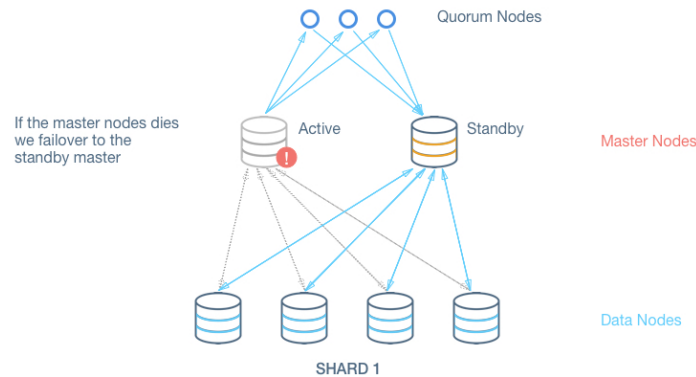
feedback.

(<http://www.quora.com/What-is-your-review-of-InterviewBit>)

**Q:** What if the master keeping track of where the blocks are stored dies?



**Answer:** To overcome this problem we keep a standby master which in the failover process becomes the acting master and keeps unavailability to minimum. Now, to keep the standby master up to date we can have a shared network file system. When any namespace modification is performed by the Active master, it durably logs a record of the modification to an edit log file stored in the shared directory. The Standby node constantly watches this directory for edits, and when edits occur, the Standby node applies them to its own namespace. In the event of a failover, the Standby will ensure that it has read all of the edits from the shared storage before promoting itself to the Active state. This ensures that the namespace state is fully synchronized before a failover occurs.



**A:** Going back to our design goals, latency and consistency are our design goals.

A simple way to resolve this is to make sure we only have one machine per shard. Reads and writes would work exactly how they work with a single DB. However, if the machine holding the only copy dies and its hard disk becomes corrupt, we suddenly run into the risk of losing the data which is not acceptable. That means we definitely need more than one copy of data being written with us. Let's say that number is 3. Now, the issue is how do we maintain all of the copies



in absolute sync (Consistency, remember?).

Got suggestions? We would love to hear your feedback.

Loved InterviewBit? Write us a testimonial.

One naive way would be that a write would not succeed unless it's written to all 3 copies / machines. That would make our write latency go up significantly apart

(<http://www.quora.com/What-is-your-review-of-InterviewBit>)



from making writes very unreliable ( My write fails if it fails on any of the machines ).

If we have to allow writes succeeding when the write has been written on a majority of machines ( 2 out of 3, lets say ), to maintain consistency, its important that there is a master machine which keeps track of this information. This master machine can track which machines have a particular block in each shard. However, now if this master machine dies, our whole shard is unavailable till this machine comes back up. If this machine has a corrupted hard disk, then the unavailability becomes indefinite.

There are couple of ways to keep unavailability to minimum using a standby master keeping track of master node data through a shared file system(Explained in detail in the last hint).



**You have now mastered this problem!**

## Discussion

Post a comment ([https://discuss.interviewbit.com/session/sso?return\\_path=https://discuss.interviewbit.com/t/what-are-the-cons-if-we-use-last-modified-time-solution-given-in-highly-available](https://discuss.interviewbit.com/session/sso?return_path=https://discuss.interviewbit.com/t/what-are-the-cons-if-we-use-last-modified-time-solution-given-in-highly-available))

**P**

piscado

🕒 over 1 year ago

**What are the cons if we use last modified time(solution given in highly available 3**

[Reply](https://discuss.interviewbit.com/session/sso?return_path=https://discuss.interviewbit.com/t/what-are-the-cons-if-we-use-last-modified-time-solution-given-in-highly-available) ([https://discuss.interviewbit.com/session/sso?return\\_path=https://discuss.interviewbit.com/t/what-are-the-cons-if-we-use-last-modified-time-solution-given-in-highly-available](https://discuss.interviewbit.com/session/sso?return_path=https://discuss.interviewbit.com/t/what-are-the-cons-if-we-use-last-modified-time-solution-given-in-highly-available))



**Got suggestions ? We would love to hear your**

**K**

feedback.



**Loved InterviewBit ? Write us a testimonial.**

(<http://www.quora.com/What-is-your-review-of-InterviewBit>)

System Design (/courses/system-design/)

 Show Notes

/ Storage Scalability (/courses/system-design/topics/storage-scalability/)

/ Highly Available Database

**Highly Available Database**

Bookmark

## Design a distributed key value store which is highly available and is network partition tolerant

[Blog \(https://blog.interviewbit.com/\)](https://blog.interviewbit.com/) | [About Us \(/pages/about\\_us/\)](/pages/about_us/) | [FAQ \(/pages/faq/\)](/pages/faq/) |

[Contact Us \(/pages/contact\\_us/\)](/pages/contact_us/) | [Terms \(/pages/terms/\)](/pages/terms/) | [Privacy Policy \(/pages/privacy/\)](/pages/privacy/)
**Features:**

“ System Design Interview Questions (/courses/system-design/) |  
 This is the first part of any system design interview, coming up with the  
 features which the system should support. As an interviewee, you should try  
 to list down all the features you can think of which our system should support.  
 Try to spend around 2 minutes for this section in the interview. You can use  
 the notes section alongside to remember what you wrote. ”

[Microsoft Interview Questions \(/microsoft-interview-questions/\)](/microsoft-interview-questions/) | [Puzzles Questions \(/puzzles/\)](/puzzles/)
**Q:** What is the amount of data that we need to store?Like Us (<https://www.facebook.com/interviewbit>)Follow Us ([https://twitter.com/interview\\_bit](https://twitter.com/interview_bit))Email (<mailto:hello@interviewbit.com>)**Q:** Do we need to support updates?**A:** Yes.**Q:** Can the size of the value for a key increase with updates?

**A:** Yes. In other words, it's possible a sequence of keys could co-exist on one server previously, but with time, they grew to a size where all of them don't fit on a single machine.

**Got suggestions ? We would love to hear your feedback.**

**Loved InterviewBit ? Write us a testimonial.**  
 (<http://www.quora.com/What-is-your-review-of-InterviewBit>)

**Q:** Can a value be so big that it does not fit on a single machine?

**A:** No. Let's assume that there is an upper cap of 1GB to the size of the value.

**Q:** What would the estimated QPS be for this DB?

**A:** Let's assume around 100k.

## ● Estimation:

“ This is usually the second part of a design interview, coming up with the estimated numbers of how scalable our system should be. Important parameters to remember for this section is the number of queries per second and the data which the system will be required to handle. Try to spend around 5 minutes for this section in the interview. ”

- ❓ ◀ Total storage size : 100 TB as estimated earlier  
Total estimated QPS : Around 100k

**Q:** What is the minimum number of machines required to store the data?

**A:** Assuming a machine has 10TB of hard disk, we would need minimum of  $100\text{TB} / 10\text{ TB} = 10$  machines to store the said data. Do note that this is bare minimum. The actual number might be higher if we decide to have replication or more machines in case we need more shards to lower the QPS load on every shard.



Got suggestions ? We would love to hear your

## ● Design Goals:



Loved InterviewBit ? Write us a testimonial.

(<http://www.quora.com/What-is-your-review-of-InterviewBit>)



“

**Latency** - Is this problem very latency sensitive (Or in other words, Are requests with high latency and a failing request, equally bad?). For example, search typeahead suggestions are useless if they take more than a second.

**Consistency** - Does this problem require tight consistency? Or is it okay if things are eventually consistent?

**Availability** - Does this problem require 100% availability?

*There could be more goals depending on the problem. It's possible that all parameters might be important, and some of them might conflict. In that case, you'd need to prioritize one over the other.*”

  **Q:** Is Latency a very important metric for us?

**A:** Since we want to be available all the time, we should try to have lower latency.




  **Q:** Consistency vs Availability?

**A:** As the question states, we need good availability and partition tolerance. Going by the CAP theorem ( Nicely explained at <http://robertgrainier.com/2014/08/cap-theorem-revisited/> ( <http://robertgrainier.com/2014/08/cap-theorem-revisited/> ) ), we would need to compromise with consistency if we have availability and partition tolerance.

We can however aim at having eventual consistency. As is the case with any storage system, data loss is not acceptable.



 **Deep Dive:** Got suggestions? We would love to hear your feedback.



Loved InterviewBit ? Write us a testimonial.  
(<http://www.quora.com/What-is-your-review-of-InterviewBit>)

“ Lets dig deeper into every component one by one. Discussion for this section will take majority of the interview time(20-30 minutes). ”



“ Note : In questions like these, the interviewer is looking at how you approach designing a solution. So, saying that I'll use a NoSQL DB like Cassandra is not an ideal answer. It is okay to discuss the architecture of Cassandra for example with rationale around why some components were designed the way they were.. ”

**Q:** Is sharing required?

**A:** Lets look at our earlier estimate about the data to be stored. 100TB of data can't be stored on a single machine.

Lets say that we somehow have a really beefy machine which can store that amount of data, that machine would have to handle all of the queries ( All of the load ) which could lead to a significant performance hit.

“ Tip: You could argue that there can be multiple copies of the same machine, but this would not scale in the future. As my data grows, its possible that I might not find a big beefy enough machine to fit my data. ”

So, the best course of action would be to share the data and distribute the load amongst multiple machines.



**Q:** Should the data stored be normalized?

(<http://www.studytonight.com/dbms/database-normalization.php>)

(<http://www.studytonight.com/dbms/database-normalization.php>)

**A:** If the data is normalized, then we need to join across tables and across rows to fetch data. If the data is already shared across machine, any join across machines is highly undesirable ( High latency, Less indexing support ).



Got suggestions? We would love to hear your feedback. (http://www.quora.com/What-is-your-review-of-InterviewBit)

However, if the sharding criteria is not chosen properly, it could lead to consistency concerns (After all, we are storing the same data at multiple places). However, for this case, we are more concerned with availability and ready to compromise on consistency as long as things become eventually consistent. In this case, it seems like having denormalized rows makes sharding easier for us and suits our use case better.

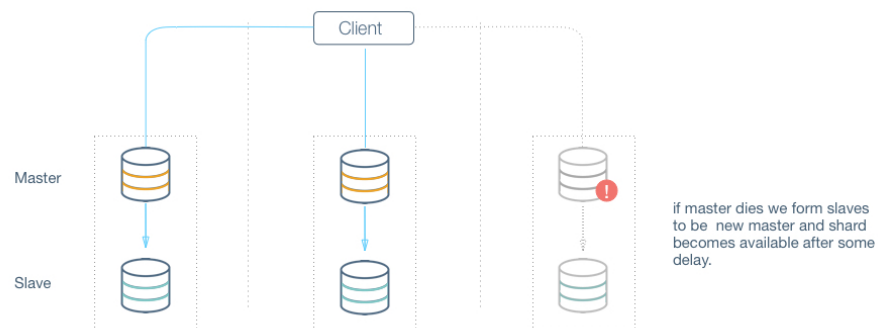


**Q:** How many machines per shard? How does a read/write look in every shard?

**A:** Going back to our design goals, low latency and high availability are our design goals.

Lets assume we have somehow sharded the rows into shards. Hence, lets first look at how the architecture might look at within a shard.

### Master Slave



One simple solution might be to have a master node in each shard which has a slave node which reads all new updates from master and keeps updating itself (The slave in this case might not have the same view as master and would lag a little bit). Clients can read from either the master or the slave depending on which responds earlier (or being slightly more intelligent with the reads to give more preference to the master, and fallback to slave after the read call to master). That could lead to inconsistent views on newer entries across master and client, but would ensure high read availability.



**Got suggestions? We would love to hear your**

feedback. However, what happens to writes when the master goes down. The writes start failing since only master was taking up the writes.

**Loved InterviewBit? Write us a testimonial.**

(<http://www.quora.com/What-is-your-review-of-InterviewBit>)

We can argue that we can have the slave become the new master in such a case. However, even that implies unavailability for the period of failover from master to the slave as new master.

Also, if the slave is lagging a bit, and then the master has a hardware failure, we run the risk of losing data.

This means that we definitely need more than one machine taking the write traffic if we are to be available all the time.

### Multi Master

Lets say we modify the previous design where both machines accept write AND read traffic. Lets name the machine m1 and m2.

If m1 accepts write without depending on m2, then it is possible m1 is doing write on a row state which is not the same as m2. That could lead to huge consistency concerns and the DB might become forever inconsistent. DBs might get operations out of order and it can cause eventual inconsistency if the order of operation matters (double the column value and then increment it vs increment the column value and then double it).

From above examples we see that any system with a *master* node will not be highly available, therefore we move to peer to peer systems.



**Q:** Can a peer to peer system be highly available in case of a DB machine dying?

**Hint:** Single point of failure!

**A:** We define a peer to peer system where every node is equally privileged and any two nodes can communicate. Yes, since we don't have a single point of failure anymore, therefore our system can theoretically be available even in presence of dying DB machines. Dynamo and Cassandra are examples of examples of such systems, both of them lack the *master* node and therefore have no single point of failure. Our highly available datastore will be highly based on Dynamo and Cassandra, as a reader you don't need to know about them.



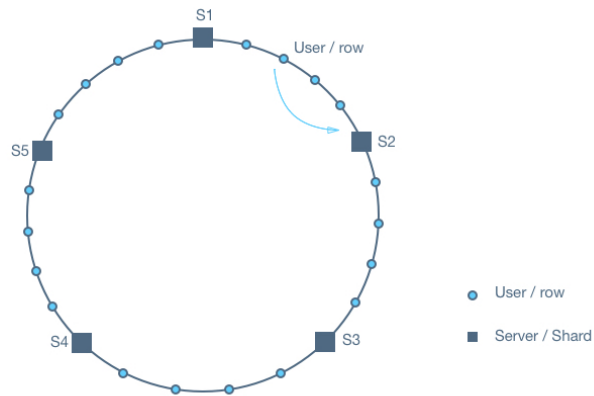
Got suggestions ? We would love to hear your feedback.



Loved InterviewBit ? Write us a testimonial.  
(<http://www.quora.com/What-is-your-review-of-InterviewBit>)

**Q:** How will we exactly share data for a peer-to-peer system?

**A:** Refer to <https://www.interviewbit.com/problems/sharing-a-database/> (<https://www.interviewbit.com/problems/sharing-a-database/>) for a detailed answer.



**Q:** How do we store redundant data?

**A:** We will need to store duplicate data to prevent data loss in case of some of our DB machines getting corrupted. To store the data we can use consistent hashing which assigns every data to a particular node on the ring. Let's call our replication factor (we will store  $P$  copies of data). Now for a data  $D$ , we have to choose  $P$  nodes where we will store copies of  $D$ .

Now, how do we choose these  $P$  nodes? We will choose the  $P$  clockwise consecutive nodes starting from the node where  $D$  is mapped by our hashing function.

An important point to discuss here is that even though any data might be mapped to a particular virtual node, the virtual node is not the master node for this data for either read or write. A client can request to read or write a data from any node they want. This is essential in creating a highly available system.

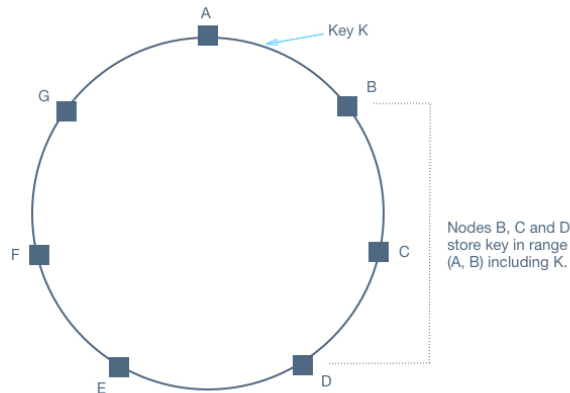


Got suggestions ? We would love to hear your feedback.



Loved InterviewBit ? Write us a testimonial.  
(<http://www.quora.com/What-is-your-review-of-InterviewBit>)





**Q:** How does a write/read happen in our system?

**A:**

#### Write request:

A client can contact any node in the ring with a `put()` request to write data, this node acts as a coordinating node for this request. Coordinating node then forwards the request to the mapping nodes for the data (hashing) and waits for acknowledgement from them. When it receives **W** (explained later) acknowledgements, it returns a write-accepted message to the client.

#### Read request:

To perform a `get()` request, client can connect to any node in the ring which becomes the coordinating node for this request. The coordinating node then asks all replica nodes for this data and returns consolidated data to the client when **R** of them replies back.

#### Read and Write consistency:

**W** and **R** are called write and read consistency number respectively. To recap, **W** is the minimum number of nodes from which the coordinating node should get an ack before making a write successful and **R** is the minimum number of nodes from which the coordinating node should get back read values to return them



Got suggestions? We would love to hear your

**R, W** together for ms quorum of the system. For a read to be consistent (return the latest write), we need to keep  $W + R > P$ .



Loved InterviewBit? Write us a testimonial.

(<http://www.quora.com/What-is-your-review-of-InterviewBit>)

Depending on the feature requirement  $W$  and  $R$  can be adjusted, for example to have very fast writes we can keep  $W = 1$  and  $R = P$ . If our system is read heavy we can keep  $R = 1$  and  $W = P$ . If read and write are equally distributed, we can keep both  $R$  and  $W$  as  $(P+1)/2$ .



**Q:** What if a machine goes down?

**A:** Since no node is only responsible for a piece of data, it's going down won't really affect our writes/reads. As long as  $W$  out of  $P$  nodes are available for some key, it can be updated (similarly for read).

Note that in case of less than  $W$  nodes available to write for some data, we can relax our write consistency (sacrificing data consistency for availability).



**Q:** What kind of consistency can we provide?

**A:** If we keep  $W = P$ , we can provide strong consistency but we won't be available for some writes even if one of our DB machine dies.

Earlier we saw in master-master configuration that in network partition cases, our masters may diverge to provide availability. In our current system, essentially all of our nodes are master and the point that they will diverge should be taken for granted and we should build our system considering it.

Therefore we should build for the case where  $W$  is less than  $P$ , hence our writes will be propagated i.e. some machines will have an updated view of data and some will not, therefore they are not consistent. The best we can guarantee here is eventual consistency, that is in due time, all the changes will be applied to every server and in the end they will all have a consistent view of the data.

To achieve eventual consistency, we need to be able to resolve differences between data on two servers. There are a couple of detect and resolve data conflicts that may arise.

First, if data(key, value) we store is such that value is just a single column, we



Got suggestions? We would love to hear your

can use a simple criteria of LWW (last write wins) to resolve conflicts. So if two servers have different view of a key, in the resolve step we can update the server with the stale with the new data and therefore become consistent.



Loved InterviewBit? Write us a testimonial.

(<http://www.quora.com/What-is-your-review-of-InterviewBit>)

The other way is to store augmented data for each row indicating all the coordinating nodes for the row till now. Now, to detect and understand conflict we can compare the augmented data. If one is a subset of the other (all the writes seen by one of the row has been seen by the other row) we can safely ignore the one with smaller augmented data. Otherwise, we have a conflict for our row and need application level logic to resolve it. This way is usually required when our value is composed of more than one independent column.



You have now mastered this problem!

## Discussion

Post a comment ([https://discuss.interviewbit.com/session/sso?return\\_path=https://discus](https://discuss.interviewbit.com/session/sso?return_path=https://discus)

H

henry\_henry

almost 2 years ago

**How exactly do peer-to-peer systems work? Doesn't there still have to be a high-I** 2

[Reply](https://discuss.interviewbit.com/session/sso?return_path=https://discuss.interviewbit.com/t/how-exactly-do-peer-to-peer-systems-work-doesnt-there-still-have-to-be-a-high-l) ([https://discuss.interviewbit.com/session/sso?return\\_path=https://discuss.interviewbit.com/t/how-exactly-do-peer-to-peer-systems-work-doesnt-there-still-have-to-be-a-high-l](https://discuss.interviewbit.com/session/sso?return_path=https://discuss.interviewbit.com/t/how-exactly-do-peer-to-peer-systems-work-doesnt-there-still-have-to-be-a-high-l))

S



Got suggestions ? We would love to hear your feedback.



Loved InterviewBit ? Write us a testimonial.  
(<http://www.quora.com/What-is-your-review-of-InterviewBit>)

System Design (/courses/system-design)

Show Notes

/ System Design Interview Questions (/courses/system-design/topics/interview-questions/)

/ Design Twitter

## Design Twitter

Bookmark



Let's design a Twitter like system.

InterviewBit Retweeted

**Krishna @awesomkk** · 31 May 2015  
[Blog \(https://blog.interviewbit.com\)](https://blog.interviewbit.com) | [About Us \(/pages/about\\_us/\)](/pages/about_us/) | [FAQ \(/pages/faq/\)](/pages/faq/) |  
[Contact Us \(/pages/contact\\_us/\)](/pages/contact_us/) | [Terms \(/pages/terms/\)](/pages/terms/) | [Privacy Policy \(/pages/privacy/\)](/pages/privacy/)

InterviewBit Retweeted

**Sanjeev Tamang @tmz\_xanziv** · 30 May 2015  
[System Design Interview Questions \(/courses/system-design/\)](/courses/system-design/) |  
[Google Interview Questions \(/google-interview-questions/\)](/google-interview-questions/) |  
[Facebook Interview Questions \(/facebook-interview-questions/\)](/facebook-interview-questions/) |  
[Amazon Interview Questions \(/amazon-interview-questions/\)](/amazon-interview-questions/) |

InterviewBit Retweeted

**Prabhakar Undurthi @prabhakar\_u** · 30 May 2015  
 Your awesome resource to become the ninja programmer with everything you need. Check this out [.interviewbit.com](https://www.interviewbit.com) @interview\_bit

[Microsoft Interview Questions \(/microsoft-interview-questions/\)](/microsoft-interview-questions/) | [Puzzles Questions \(/puzzles/\)](/puzzles-questions/)

Like Us (<https://www.facebook.com/interviewbit>) | Follow Us ([https://twitter.com/interview\\_bit](https://twitter.com/interview_bit))  
 Email (<mailto:hello@interviewbit.com>)

## ● Features:

“ This is the first part of any system design interview, coming up with the features which the system should support. As an interviewee, you should try to list down all the features you can think of which our system should support.



Try to spend around 2 minutes for this section in the interview. You can use the notes section alongside to remember what you wrote. ”  
 Got suggestions ? We would love to hear your feedback. Loved InterviewBit ? Write us a testimonial. (<http://www.quora.com/What-is-your-review-of-InterviewBit>)

**Q:** What are some of the Twitter features we should support?

**A:** Let's assume that we are looking at posting tweets, following people and favoriting tweets. A user should also be able to see a feed of tweets of his/her followers.

**Q:** Do we need to support replies to tweets / grouping tweets by conversations?

**A:** Let's assume we don't need to for this case.

**Q:** How about privacy controls around each tweet?

**A:** Not required. Let's assume for this case that all tweets are public.

**Q:** Do we need to support trending tweets? If so, do we need to support localization and personalization?

**A:** For this case, let's just assume we are not focussing on building the trending tweets feature.

**Q:** How about Direct messaging?

**A:** No. Let's leave that out for this question. That could be another question by itself.

**Q:** How about mentions / tagging?

**A:** Let's assume we don't need to support mentions/tagging.

**Q:** Do we need to support a notification system?

**A:** For the purpose of this question, no.



Got suggestions? We would love to hear your  
● **Estimation:**  
feedback.



Loved InterviewBit? Write us a testimonial.  
(<http://www.quora.com/What-is-your-review-of-InterviewBit>)

“ This is usually the second part of a design interview, coming up with the estimated numbers of how scalable our system should be. Important parameters to remember for this section is the number of queries per second and the data which the system will be required to handle. Try to spend around 5 minutes for this section in the interview. ”

- ❓ ◀ Lets estimate the volume of tweets. Assume that our system would be the second most popular tweeting service after Twitter.

**Q:** What is the number of users and traffic that we expect the system to handle?

**A:** Twitter does around 500 million tweets per day with 100 million daily active users. Lets assume similar numbers.



- ❓ ◀ **Q:** How many followers does every user have?

**A:** The behavior should be similar to Twitter here. Each user has on average 200 followers, with certain hot users having a lot more followers. For example, users like Justin Bieber would have millions of followers.



- ❓ ◀ **Q:** How many times is a tweet favorited?

**A:** Assuming the same behavior as Twitter, we can assume that each tweet is favorited twice. However, in this case as well, there will be outliers. There are certain tweets which might be favorited by millions of people.



- ❓ ◀ **Q:** Assuming the network of users, how many user to follower edge would exist?



Got suggestions? We would love to hear your

**A:** We had 100 million active users with 200 followers on average. This means

100M x 200 edges = 20B edges



Loved InterviewBit? Write us a testimonial.

(<http://www.quora.com/What-is-your-review-of-InterviewBit>)



## ● Design Goals:

“

**Latency** - Is this problem very latency sensitive (Or in other words, Are requests with high latency and a failing request, equally bad?). For example, search typeahead suggestions are useless if they take more than a second.

**Consistency** - Does this problem require tight consistency? Or is it okay if things are eventually consistent?

**Availability** - Does this problem require 100% availability?

*There could be more goals depending on the problem. It's possible that all parameters might be important, and some of them might conflict. In that case, you'd need to prioritize one over the other. ”*

❓ ◀ Q: Is Latency a very important metric for us?

A: Yes. A twitter like system needs to be fast, especially when you are competing with Twitter .



❓ ◀ Q: How important is Consistency for us?

A: Not really. Assuming a lot of activity on this system, if I miss out on a tweet of a person I am following every now and then, its not the end of the world.

Compare this to direct messaging where consistency is extremely important.



❓ ◀ Got suggestions ? We would love to hear your feedback.



Loved InterviewBit ? Write us a testimonial.

(<http://www.quora.com/What-is-your-review-of-InterviewBit>)

**A:** Yes. If Twitter becomes unavailable, it becomes a news. As a product, it needs to be highly available.



## ● Skeleton of the design:

“ The next step in most cases is to come up with the barebone design of your system, both in terms of API and the overall workflow of a read and write request. Workflow of read/write request here refers to specifying the important components and how they interact. Try to spend around 5 minutes for this section in the interview.

**Important :** Try to gather feedback from the interviewer here to indicate if you are headed in the right direction. ”

❓ ◀ As discussed before, there are 4 major operations we need to support:

Posting new tweets

Following a user

Favoriting a tweet

Get the feed for a user

**Q:** What would the API look like for the client?

**Q:** What data would need with every Tweet we fetch?



**A:** We should have the content of the tweet, the person who posted the tweet, the timestamp when tweet was created and number of favorites for the tweet.



Got suggestions ? We would love to hear your feedback.



Loved InterviewBit ? Write us a testimonial.  
(<http://www.quora.com/What-is-your-review-of-InterviewBit>)



**Q:** Would we need all the user profiles of users who have favorited a tweet?



**A:** Given that's a lot of data to fetch, we can be more intelligent about it and just fetch top 2 people in the list. In this scheme, we would show every tweet as 200 favorites which on hover shows Favorited by X, Y and 198 others

**Q:** How many tweets should we fetch at a time?



**A:** At a time, only a certain number of tweets will be in the viewport (lets say 20). Lets call it a page of tweets. For a user, we would only want to fetch a page of tweets at a time.

*Gotcha:* Would the page size remain constant across different situations? Probably not. The page size would be different across clients based on screen size and resolution. For example, a mobile's page size might be lower than that of a web browser's.

**A:** The first 3 operations end up doing a write to the database. The last operation does a read. Following is an example of how the API might look like :

Posting new tweets : `addTweet(userId, tweetContent, timestamp)`

Following a user : `followUser(userId, toFollowUserId)`

Favorite a tweet : `favoriteTweet(userId, tweetId)`

`TweetResult getUserFeed(user, pageNumber, pageSize, lastUpdatedTimestamp)`

where `TweetResult` has the following fields :

```
TweetResult {
    List(Tweets) tweets,
    boolean isDeltaUpdate
}
```



Got suggestions ? We would love to hear your feedback.



Loved InterviewBit ? Write us a testimonial.  
(<http://www.quora.com/What-is-your-review-of-InterviewBit>)

```

        content,
        timestamp,
        numFavorites,
        sampleFavoriteNames
    }

```

There could be other APIs as well which would help us fetch the most recent tweets of a user, or fetch the followers for a tweet.



**Q:** How would a typical write query (addTweet) look like?

**A:** Components:

Client ( Mobile app / Browser , etc ) which calls addTweet(user Id, tweetContent, timestamp)

Application server which interprets the API call and tries to append the tweet to user's tweet with the timestamp in the database layer .

Database server which appends the tweet



**Q:** How would a typical read query (getUser Feed) look like?

**A:** Components:

Client (Mobile app/Browser , etc ) which calls getUser Feed

Application server which interprets the API call and queries the database for the top user feed.

Database server which looks up the followers' tweet to get the result.



Got suggestions ? We would love to hear your feedback.



Loved InterviewBit ? Write us a testimonial.  
(<http://www.quora.com/What-is-your-review-of-InterviewBit>)

## ● Deep Dive:

“ Lets dig deeper into every component one by one. Discussion for this section will take majority of the interview time(20-30 minutes). ”

🔍 ◀ Lets dig deeper into every component one by one.

## ● Application layer:

*Think about all details/gotchas yourself before beginning.*

**Q:** How would you take care of application layer fault tolerance?

**Q:** How do we handle the case where our application server dies?



**A:** The simplest thing that could be done here is to have multiple application servers. They do not store any data (stateless) and all of them behave the exact same way when up. So, if one of them goes down, we still have other application servers who would keep the site running.

**Q:** How does our client know which application servers to talk to. How does it know which application servers have gone down and which ones are still working?



**A:** We introduce load balancers. Load balancers are a set of machines (an order of magnitude lower in number) which track the set of application servers which are active (not gone down). Client can send request to any of the load balancers who then forwards the request to one of the working application servers randomly.



**A:** If we have only one application server machine, our whole service would become unavailable. Machines will fail and so will network. So, we need to plan for those events. Multiple application server machines along with load balancer

Got suggestions ? We would love to hear your feedback. (http://www.quora.com/What-is-your-review-of-InterviewBit)

Loved InterviewBit ? Write us a testimonial.

is the way to go.



## ● Database layer:

This is the heart of the question. In the skeleton design, we assumed that the database is a black box which can magically store or retrieve anything efficiently. Lets dig into how we will build that magic black box.

**Q:** What data do we need to store ?

**A:**

For every tweet, we need to store content, timestamp and owner ID.

For every user , we need to store some personal information ( Name, age, birthdate, etc. )

We need to store all u1->u2 follower relations.

We need to store all user\_ids against a tweet of users who have favorited the tweet.



**Q:** RDBMS or NoSQL?

**Q:** Are joins required?



**A:** NoSQL databases are inefficient for joins or handling relations. As such, NoSQL databases store everything in a denormalized fashion. In this case, we do have relations like  
user -> followers  
tweets -> favorited by users

SQL seems to win on this parameter on ease of use.



**Got suggestions ? We would love to hear your feedback.**



**Loved InterviewBit ? Write us a testimonial.**  
(<http://www.quora.com/What-is-your-review-of-InterviewBit>)

**Q:** How much data would we have to store?



**A:** If the size of the data is so small that it fits on a single machine's memory, SQL is a clear winner. SQL on a single machine has next to zero maintenance overhead and has great performance with right index built. If your index can fit into RAM, it's best to go with a SQL solution. Let's analyze our current case:

Size of tweets:

Number of tweets per day: 500 million

Maximum size of a tweet: 140 characters + 1 byte for timestamp + 1 byte for user id = 142 bytes

Provisioning for: 5 years =  $365 * 5$  days

Space required:  $142 \text{ bytes} * 500 \text{M} * 365 * 5 = 129.5 \text{TB}$

Size of user - follower relation:

Assuming total of 1 Billion users and every user has 200 followers on average, we end up with 200B total connections. To store it, we would need  $200 \text{B} * 2 \text{ bytes (size of 2 user IDs)} = 400 \text{G}$ .

Size of tweet to favorites relation:

Average number of favorites per tweet: 2 (Ref. Estimations section)

Total number of tweets daily: 500M

Provisioning for: 5 years =  $365 * 5$  days

Space required:  $(2 \text{ bytes} + 1 \text{ byte for tweetId}) * 500 \text{M} * 365 * 5 = 2.7 \text{TB}$

So, total space required is close to 130TB. That's definitely not fit on a single machine's hard disk.

**Q:** How important is technology maturity?



**A:** SQL DBs like MySQL have been around for a long time and have hence been iterated enough to be very stable. However, most NoSQL databases are not mature enough yet. Quoting an article from Pinterest Engineering



Blog :  
**Got suggestions? We would love to hear your feedback.**



**Loved InterviewBit? Write us a testimonial.**  
(<http://www.quora.com/What-is-your-review-of-InterviewBit>)

“ We intentionally ran away from auto-scaling newer technology like MongoDB, Cassandra and Membase, because their maturity was simply not far enough along (and they were crashing in spectacular ways on us!). ”

**A:** Things to consider :

Are joins required?

Size of the DB

Technology Maturity

In practice, the score is equal for both RDBMS or NoSQL for this one. In theory, NoSQL would be a better fit.

We can choose either to proceed further. Let's go with a relational DB like MySQL for this one.



**Q:** What would the database schema look like?

**A:** Always be prepared for this question in cases where the schema might be a bit more elaborate.

We have two main entities: users and tweets. There could be two tables for them. Table users would have personal information about the user. A sample table schema for that could look like the following :

Table **users**

ID (id) - primary key

username (username) - unique key

First Name (first\_name)

Last Name (last\_name)

password related fields like hash and salt (password\_hash & password\_salt)



Got suggestions? We would love to hear your

feedback.

Last updated at (updated\_at)



Loved InterviewBit ? Write us a testimonial.

(<http://www.quora.com/What-is-your-review-of-InterviewBit>)

Date of Birth (dob)

description (description)

Tweets should be slightly simpler :

Table **tweets**

ID (id) - primary key

content (content)

date of creation (created\_at)

user ID of author (user\_id)

Now, lets look at the r elations that we need to model :

Follower r elation ( User A follows another user B )

Table **connections**

ID of user that follows (follower\_id)

ID of user that is followed (followee\_id)

date of creation (created\_at)

Favor ite : A user can favor ite a tweet.

Table **favories**

ID of user that favorited (user\_id)

ID of favorited tweet (tweet\_id)

date of creation (created\_at)

Now, based on the r ead API quer ies, we can look at the additional index we would need :

**Get the feed of a user** - This would r equire us to quickly lookup the user Ids a user follows, get their top tweets and for each tweet get user s who have favor ited the tweet.



Got suggestions ? We would love to hear you :



Loved InterviewBit ? Write us a testimonial.

An index on follower \_id in connections to quickly lookup the user Ids a user follows  
(<http://www.quora.com/What-is-your-review-of-InterviewBit>)

An index on user\_id, created\_at in tweets to get the top tweets for a user (where user\_id = x sort by created\_at desc)

An index on tweet\_id in favorites



**Q:** Now the bigger question, How would we do sharding?

**Question: Approach1:** How can we shard on users?



**A:**

*Detail:* What's stored in each table:

users: part of the table with user\_ids which belong to the shard

tweets: part of the table with author\_ids which belong to the shard (Or in other words, tweets by the users in the current shard)

connections: All entries where follower\_id belongs to the current shard

favorites: All entries where tweet\_id belongs to the tweets table in this shard

*Pros:*

Equal load distribution

Cheap writes: All of the write queries are simple and rely on just one shard (

Assuming tweet favorite API encodes the tweet owner ID in the tweet ID

when sending request).

*Cons:*

While looking up the user IDs a user follows is easy on the machine, getting the top tweet for each of those user IDs would require querying different shards.

Even when we need to favorite a tweet, finding the tweet would require us to query all the shards. We can work around it, however, by encoding owner\_id with the tweet\_id from the client.



**Question: Approach2:** Can we shard on recency(timestamp)?



Got suggestions? We would love to hear your feedback.

Loved InterviewBit? Write us a testimonial.

(<http://www.quora.com/What-is-your-review-of-InterviewBit>)

A: Shard on recency(timestamp) -



Most recent tweets in the most recent shard. The idea is that most of the time we are only working with most recent tweets. Its rare to dig up tweets which are more than a few weeks old. New tweets are requested most frequently.

users : All of user's table resides in one shard separately.

tweets : This table is shared across shards by recency. When the most recent shard starts getting full, we create a new shard and assign the new incoming tweets to the newly created shard.

connections : All of the table resides in the shard with user's.

favorites: Stored with the tweets in their shard

*Pros:*

Fetching the user feeds requires just querying 2 shards ( user's and tweets).

More reliable and has low latency. Most of the queries would only interact with 2 shards.

*Cons:*

Load imbalance : The most recent tweet shard will be handling almost all of the traffic while the other shards will remain idle.

Huge maintenance overhead : Every time, you need to create a new shard, you'll need to allocate new machines, setup replication and make things switch almost instantly so that no downtime is induced. All in all, a nightmare for DBAs at that scale.

**A:** We have already established earlier that we would need to shard as data would not fit into a single machine. The read query we are optimizing :

**Get the feed of a user** - This would require us to quickly lookup the user's id's a user follows, get their top tweets and for each tweet get user's who have favorited the tweet.

We can model our data with two basic approaches, sharding based on recency of tweet or based on user's, we will call these approaches Approach1 and Approach2 respectively. Our answer consists of a hybrid approach of Approach1 and Approach2, so we highly recommend you to go through hints which explains



Got suggestions or feedback? We would love to hear your

Both solutions have their downsides and don't seem ideal. What is a clean solution then? Let's go back to our design goals. We want low latency and high



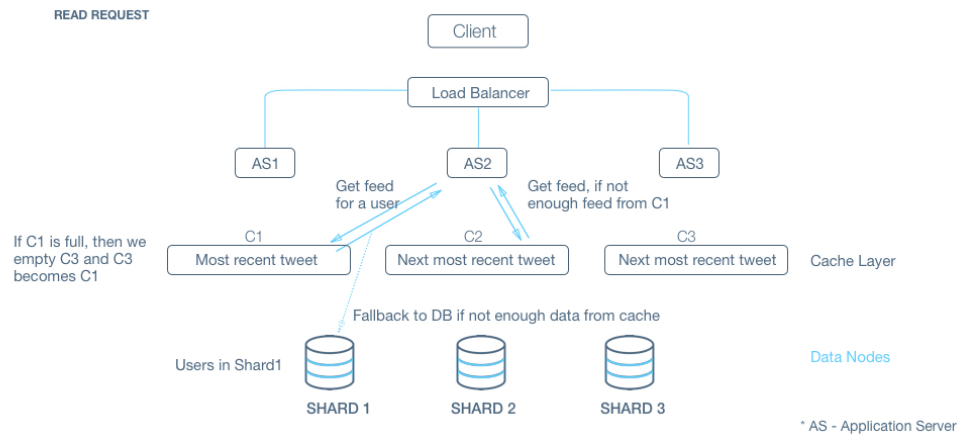
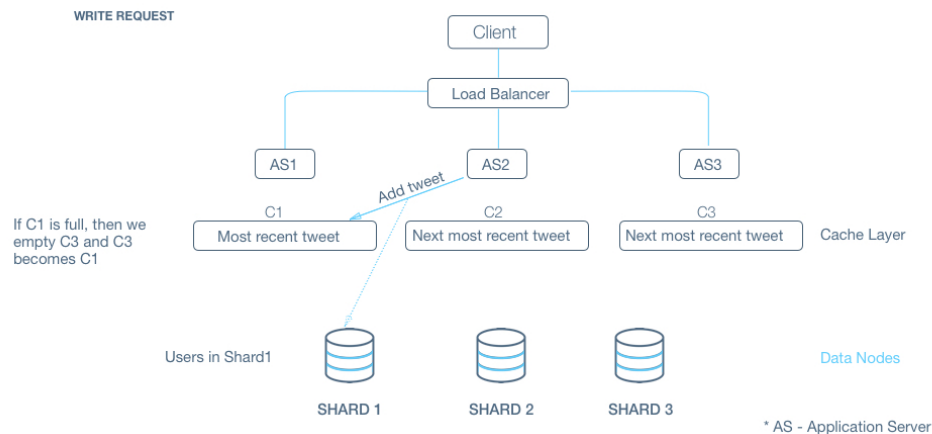
Loved InterviewBit ? Write us a testimonial.

(<http://www.quora.com/What-is-your-review-of-InterviewBit>)

availability. Consistency is not a big deal for us ( If I miss a tweet once in a while in the feed, its not the end of the world ).

With that in mind, we can look for a hybrid model. We will definitely need to heavily cache.

We can have a cache which simulates the recent shared in Approach 2 and a DB which stores stuff as in Approach 1. The idea being that most reads will be served by the cache itself and it has the collection of all recent tweets. In the rare case of not so recent tweet, we will go to the DB and in such cases, latency outliers are alright. Notice that the DB writes would be cheap as discussed in Approach 1.



Got suggestions ? We would love to hear your



Q: Do we need special handling for spiky cases ( User s with unusual number of feedback.



Loved InterviewBit ? Write us a testimonial.

(<http://www.quora.com/What-is-your-review-of-InterviewBit>)

followers / Tweets with likelihood of getting unusually high number of favorites) ?  
Think about the case when Katy Perry (with more than 70M followers) tweets.

**A:** Let's look at both cases one by one.

Let's say Katy Perry tweets. Following is what happens :

We write the tweet to the shard where Katy Perry belongs. Not a problem.

We add it to the recent tweets cache. Again, not a problem.

As all followers get their feed update by specifically requesting for an update, the resultant change is that a lot of followers will get an update when they request for it. This should manifest as an uptick in the upload bandwidth. In the worst case, assuming that 30% of the followers are online at a time, we would need around 3G of upload bandwidth which is a really small number for a datacenter.

Let's look at really popular tweets now. They'll have an unusually high rate of being favorited (The highest being 3M total favorites). This means a really high rate of write to the shard which can cause deadlocks. We can add some optimizations here if required in terms of batching the updates to favorites table in a queue before flushing them. Nitty Gritty : Would the queue be persistent? If not, what happens if the machine dies. That would cause data loss. If yes, where does the queue reside? How do you merge the query results?



**Q:** How would we handle a DB machine going down?

**A:** As stated in design goals, we need to make sure our system is more available at all times.

We had sharded the database based on users. We can have a replica for each of them which follows the updates happening on the master database shard.

When the master goes down, the slave can take over. Now there is a problem here. What if there were some updates which the slave had not caught up to yet. Do we lose that information? We can take a call either way. If we are particular about getting the data back, we know that we can get that information from the cache layer and resolve stuff on the DB layer.



Got suggestions ? We would love to hear your feedback.



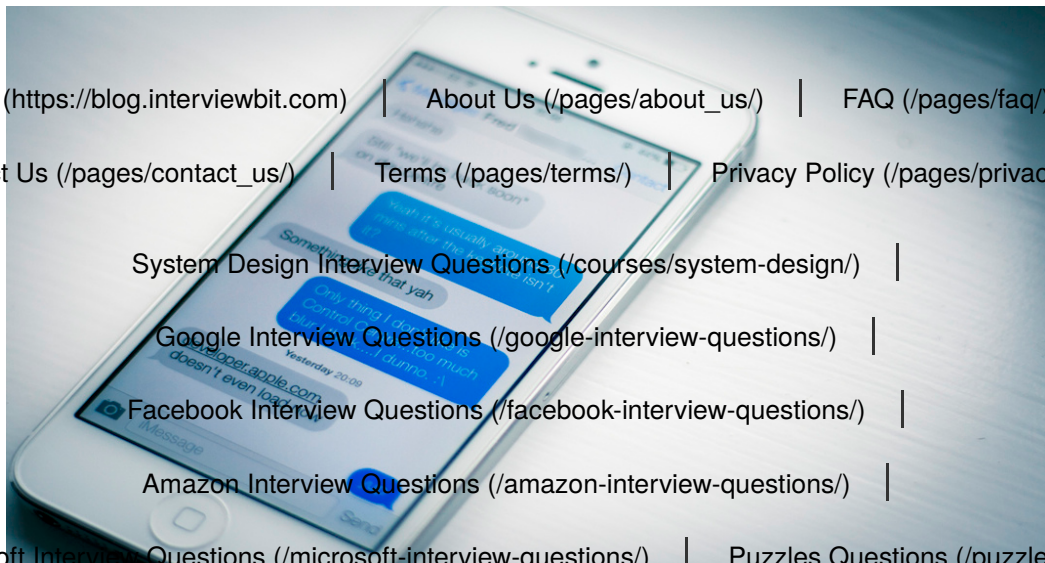
Loved InterviewBit ? Write us a testimonial.  
(<http://www.quora.com/What-is-your-review-of-InterviewBit>)

[System Design \(/courses/system-design/\)](/courses/system-design/)[Show Notes](#)[/ System Design Interview Questions \(/courses/system-design/topics/interview-questions/\)](/courses/system-design/topics/interview-questions/)[/ Design Messenger](#)

## Design Messenger

[Bookmark](#)

### ● ◀ Design a messaging service, like Facebook Messenger.

[Blog \(https://blog.interviewbit.com\)](https://blog.interviewbit.com)[About Us \(/pages/about\\_us/\)](/pages/about_us/)[FAQ \(/pages/faq/\)](/pages/faq/)[Contact Us \(/pages/contact\\_us/\)](/pages/contact_us/)[Terms \(/pages/terms/\)](/pages/terms/)[Privacy Policy \(/pages/privacy/\)](/pages/privacy/)[System Design Interview Questions \(/courses/system-design/\)](/courses/system-design/)[Google Interview Questions \(/google-interview-questions/\)](/google-interview-questions/)[Facebook Interview Questions \(/facebook-interview-questions/\)](/facebook-interview-questions/)[Amazon Interview Questions \(/amazon-interview-questions/\)](/amazon-interview-questions/)[Microsoft Interview Questions \(/microsoft-interview-questions/\)](/microsoft-interview-questions/)[Puzzles Questions \(/puzzles/\)](/puzzles/)[f Like Us \(https://www.facebook.com/interviewbit\)](https://www.facebook.com/interviewbit)[t Follow Us \(https://twitter.com/interview\\_bit\)](https://twitter.com/interview_bit)[✉ Email \(mailto:hello@interviewbit.com\)](mailto:hello@interviewbit.com)

### ● Features:

“ This is the first part of any system design interview, coming up with the features which the system should support. As an interviewee, you should try to list down all the features you can think of which our system should support.



Try to spend around 2 minutes for this section in the interview. You can use  
Got suggestions ? We would love to hear your feedback.  
Loved InterviewBit ? Write us a testimonial.  
the notes section alongside to remember what you wrote. ”  
(<http://www.quora.com/What-is-your-review-of-InterviewBit>)

**Q:** What is the scale that we are looking at?

**A:** Let's assume the scale of Facebook Messages. Let's say we need to handle around 10B message sends a day and around 300M users.

**Q:** Do we only need to support 1:1 conversations or group conversations as well?

**A:** Let's assume we are building things just for 1:1 conversations. We will extend it to group conversations if need be.

**Q:** Do we need to support attachments?

**A:** For now, let's assume that we don't. We will only look at building plain-text messaging system.

**Q:** What is a reasonable limit to the size of a message?

**A:** Let's assume that we are building a chat messaging system. As such, we would expect every message to be shorter in length. We can impose a limit here on the maximum size of such a message. Let's say we will only handle messages less than 64Kb in size and reject the others.

**Q:** What about the notification system for new messages received?

**A:** Considering the size of the discussion here ( 45 mins in the interview ), we will not delve into the notification system for messages.

## ● Estimation:



Got suggestions ? We would love to hear your feedback.



Loved InterviewBit ? Write us a testimonial.  
(<http://www.quora.com/What-is-your-review-of-InterviewBit>)

“ This is usually the second part of a design interview, coming up with the estimated numbers of how scalable our system should be. Important parameters to remember for this section is the number of queries per second and the data which the system will be required to handle. Try to spend around 5 minutes for this section in the interview. ”

- ❓ ◀ Let's estimate the volume of each. Assume that our system would be the one of the most popular messaging service.

**Q:** Given the number of messages being sent, what is the amount of message sent data size we are generating everyday?

**A:** Number of message sends : 10B

Assuming each message on average has 160 characters , that results in  $10B * 160 = 1.6TB$  assuming no message metadata.



- ❓ ◀ **Q:** What is the expected storage size?

**A:** From the previous section, we know that we generate 1.6TB data everyday if we only store one copy of the message. If we were to provision for 10 years, we are looking at  $1.6 * 365 * 10 TB$  which is approximately 6 Petabytes.



## ● Design Goals:



Got suggestions ? We would love to hear your feedback.



Loved InterviewBit ? Write us a testimonial.  
(<http://www.quora.com/What-is-your-review-of-InterviewBit>)

“

**Latency** - Is this problem very latency sensitive (Or in other words, Are requests with high latency and a failing request, equally bad?). For example, search typeahead suggestions are useless if they take more than a second.

**Consistency** - Does this problem require tight consistency? Or is it okay if things are eventually consistent?

**Availability** - Does this problem require 100% availability?

*There could be more goals depending on the problem. It's possible that all parameters might be important, and some of them might conflict. In that case, you'd need to prioritize one over the other. ”*

❓ ◀ **Q:** Is Latency a very important metric for us?

**A:** Yes. Chat is supposed to be realtime, and hence the end to end time actually matters.



❓ ◀ **Q:** How important is Consistency for us?

**A:** Definitely, yes. Its not okay if someone sends me a sequence of message and I don't see some of them. That could lead to huge confusion. Think of cases when you miss an emergency message or missed messages cause misunderstanding between individuals.



❓ ◀ **Q:** How important is Availability for us?

**A:** Availability is good to have. If we had to choose between consistency and availability, consistency wins.



Got suggestions ? We would love to hear your feedback.



Loved InterviewBit ? Write us a testimonial.  
(<http://www.quora.com/What-is-your-review-of-InterviewBit>)

## ● Skeleton of the design:

“ The next step in most cases is to come up with the barebone design of your system, both in terms of API and the overall workflow of a read and write request. Workflow of read/write request here refers to specifying the important components and how they interact. Try to spend around 5 minutes for this section in the interview.

**Important** : Try to gather feedback from the interviewer here to indicate if you are headed in the right direction. ”

? ◀ Q: What are the operations that we need to support?

A:

Send a message to another person

For a user, fetch the most recent conversations

For every conversation, fetch the most recent messages



? ◀ Q: What would the API look like for the client?

Q: How would the `sendMessage` API look like?



A: Send Message : Things to take care of in this API

`sendMessage` should be idempotent. If a client retries the message, the message should not be added twice. We can resolve this by generating a random timestamp based ID on the client which can be used to de-duplicate the same message being sent repeatedly.

Ordering of messages should be maintained. If I send message A, and then send message B, then A should always appear before B. However, it is

possible that due to delays, if two messages are sent quickly one after another, then the requests reach the DB out of order. How do we solve such a case? Obviously, we need to resolve based on the timestamp when they



Got suggestions ? We would love to hear your

Loved InterviewBit ? Write us a testimonial.

<http://www.quora.com/What-is-your-review-of-InterviewBit>



were sent at.

Timestamp on the client is always unreliable. So, we would need to record the timestamp the first time the request hits the servers ( Need not be a part of the API to the client )

```
sendMessage(senderId, receipientId, messageContent,
clientMessageId)
```

**Q:** How would the API for fetching user's latest conversation look like?



**A:** This API would be called if I need to show a page of conversations/threads ( Think of the view you see when you open the Whatsapp / Messenger app ).

At a time, only a certain number of conversations will be in the viewport ( let's say 20 ). Let's call it a page of conversations. For a user, we would only want to fetch a page of conversations at a time.

**Gotcha:** Would the page size remain constant across different situations? Probably not. The page size would be different across clients based on screen size and resolution. For example, a mobile's page size might be lower than that of a web browser's.

**Delta fetch:** In most cases, our API calls will be made by users who are active on the site. As such, they already have a view of conversations till a certain timestamp and are only looking for updates after the timestamp ( which would typically be 0-2 more conversations ). For clients which are data sensitive (like mobile), fetching the whole page every time even when I have all of the conversations can be draining. So, we need to support a way of fetching only the updates when the lastFetchedTimestamp is closer to currentTimestamp.

Keeping the above 2 facts in mind, following is how a hybrid API might look like :

```
ConversationResult fetchConversation(userId, pageNumber,
pageSize, lastUpdatedTimestamp)
```



Got suggestions? We would love to hear you the following fields InterviewBit ? Write us a testimonial.

feedback ConversationResult {

(<http://www.quora.com/What-is-your-review-of-InterviewBit>)

```
List(Conversation) conversations,  
boolean isDeltaUpdate  
}  
Conversation {  
    conversationId,  
    participants,  
    snippet,  
    lastUpdatedTimestamp  
}
```

**Q:** How would the API for fetching most recent messages in a conversation look like?



**A:** Fetch most recent message in a conversation :

This API is almost identical to the fetchConversation API.

```
MessageResult fetchMessages(userId, pageNumber, pageSize,  
lastUpdatedTimestamp)
```

where MessageResult has the following fields :

```
MessageResult {  
    List(Message) messages,  
    boolean isDeltaUpdate  
}  
Message {  
    messageId,  
    senderId,  
    participants,  
    messageContent,  
    sentTimestamp  
}
```



Got suggestions ? We would love to hear your  
feedback.



Loved InterviewBit ? Write us a testimonial.  
(<http://www.quora.com/What-is-your-review-of-InterviewBit>)

**A:** The first and last operation ends up doing a write to the database. The other operations are purely read operations. Following is how API's may look like:

Send Message:

```
sendMessage(senderId, receipientId, messageContent,  
clientMessageId)
```

Conversations of a user :

```
ConversationResult fetchConversation(userId, pageNumber,  
pageSize, lastUpdatedTimestamp)
```

where ConversationResult has the following fields :

```
ConversationResult {  
  
    List(Conversation) conversations,  
    boolean isDeltaUpdate  
  
}  
  
Conversation {  
  
    conversationId,  
    participants,  
    snippet,  
    lastUpdatedTimestamp  
  
}
```

Fetch most recent message in a conversation : This API is almost identical to the fetchConversation API.

```
MessageResult fetchMessages(userId, pageNumber, pageSize,  
lastUpdatedTimestamp)
```

where MessageResult has the following fields :

```
MessageResult {  
  
    List(Message) messages,  
    boolean isDeltaUpdate  
  
}
```

```
Message {  
  
    messageId,  
    feedback.  
    senderId,
```



**Got suggestions ? We would love to hear your**

**messageId,  
feedback.  
senderId,**



**Loved InterviewBit ? Write us a testimonial.**

**(<http://www.quora.com/What-is-your-review-of-InterviewBit>)**

```
    participants,  
    messageContent,  
    sentTimestamp  
}
```



**Q:** How would a typical write query look like?

**A:** Components:

Client ( Mobile app / Browser, etc ) which calls `sendMessage(senderId, receipientId, messageContent, clientMessageld)`

Application server which interprets the API call and calls DB to do the following:

Puts in the `serverTimestamp`

Figures out the conversation to which the message should be appended based on the other participant

Figures out if a recent message exists with the `clientMessageld`

Store the message

Database server which stores the message.



**Q:** How would a typical read query look like?

**A:** Components:

Client (Mobile app/Browser, etc ) which calls `fetchConversation`

Application server which interprets the API call and queries the database for the top conversation.

Database server which looks up the user's conversations.



Got suggestions ? We would love to hear your  
feedback.

● **Deep Dive:**



Loved InterviewBit ? Write us a testimonial.  
(<http://www.quora.com/What-is-your-review-of-InterviewBit>)

“ Lets dig deeper into every component one by one. Discussion for this section will take majority of the interview time(20-30 minutes). ”

🔍 ◀ Let's dig deeper into every component one by one.

**Application layer:**

*Think about all details/gotchas yourself before beginning.*

**Q:** How would you take care of application layer fault tolerance?

**Q:** How do we handle the case where our application server dies?



**A:** The simplest thing that could be done here is to have multiple application server. They do not store any data (stateless) and all of them behave the exact same way when up. So, if one of them goes down, we still have other application servers who would keep the site running.

**Q:** How does our client know which application servers to talk to. How does it know which application servers have gone down and which ones are still working?



**A:** We introduce load balancers. Load balancers are a set of machines (an order of magnitude lower in number) which track the set of application servers which are active ( not gone down ). Client can send request to any of the load balancers who then forward the request to one of the working application servers randomly.

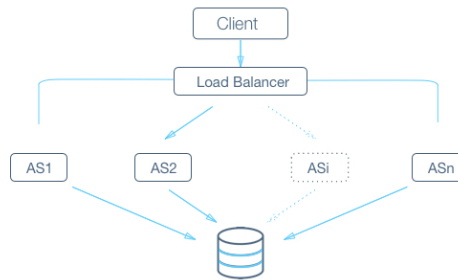


Got suggestions ? We would love to hear your feedback.



Loved InterviewBit ? Write us a testimonial.  
(<http://www.quora.com/What-is-your-review-of-InterviewBit>)

HIGH LEVEL DESIGN



**A:** If we have only one application server machine, our whole service would become unavailable. Machines will fail and so will network. So, we need to plan for those events. Multiple application server machines along with load balancer is the way to go.



### Database layer

This is the heart of the question. In the skeleton design, we assumed that the database is a black box which can magically store or retrieve anything efficiently. Let's dig into how we will build that magic black box.

**Q:** RDBMS or NoSQL?

**Q:** Are joins required?



**A:** NoSQL databases are inefficient for joins or handling relations. As such, NoSQL databases store everything in a denormalized fashion. In this case, we do have relations like  
user -> messages



**Got suggestions ? We would love to hear your feedback.**



**Loved InterviewBit ? Write us a testimonial.**  
(<http://www.quora.com/What-is-your-review-of-InterviewBit>)

SQL seems to win on this parameter on ease of use.

**Q:** How much data would we have to store?



**A:** If the size of the data is so small that it fits on a single machine's main memory, SQL is a clear winner. SQL on a single machine has next to zero maintenance overhead and has great performance with right index built. If your index can fit into RAM, its best to go with a SQL solution. In our earlier estimations, we had already established that we will need to provision petabytes of data which most definitely does not fit on a single machine. So, a SQL solution will have a sharding overhead. Most NoSQL solutions however are built with the assumption that the data does not fit on a single machine and hence have sharding builtin. NoSQL wins on this parameter.

**Q:** What is the read-write pattern?



**A:** Messaging is going to be very write heavy. Unlike photos or tweets or posts which are written once and then consumed a lot of times by a lot of people, messaging is written once and consumed by the other participant once.

For a write heavy system with a lot of data, RDBMS usually don't perform well. Every write is not just an append to a table but also an update to multiple index which might require locking and hence might interfere with reads and other writes.

However, there are NoSQL DBs like HBase where writes are cheaper.

NoSQL seems to win here.

**A:** Things to consider :

Are joins required?

Size of the DB



Got suggestions? We would love to hear your

feedback.

From the looks of it, NoSQL seems like a better fit. Let's proceed with NoSQL for



Loved InterviewBit ? Write us a testimonial.

(<http://www.quora.com/What-is-your-review-of-InterviewBit>)

now.



**Q:** How would we store the data? Discuss schema

**A:** As discussed before, with NoSQL, we need to store the data in denormalized form.

First thing first, this means that every user would have his/her own copy of the mailbox. That means that we will store 2 copies of the message, one for each participant for every message send.

Let's delve into how the schema would look. We'll assume that we are using HBase for this problem.

If this is your first time designing schema, we strongly recommend you go through a primer here (<https://www.mapr.com/blog/guidelines-hbase-schema-design>).

For schema design, it's good to recognize our access patterns. To achieve that, let's look at our major operations :

For a user, append a message to a conversation

Fetch timestamp ordered conversations for a user ( Most recent )

Fetch most recent messages in a conversation for a user ( Most recent )

As you can see, the first lookup for all three operations is for the user. In NoSQL context, it hence makes sense to have userId as the row ID ( Data is sharded based on users ).

Now, within the user, we will need to lookup conversations, recent conversations and recent messages.

One naive approach is to fetch the whole list of conversations or all the messages in a conversation and then filter the data we need. This however is really slow ( Remember that one of our design goals was to have low latency ). Even more so, in cases when some popular users get a lot of message and have a huge mailbox (in GBs).

Let's see, how we would solve each read request one by one.

Recent conversations : We can separately store conversationId : timestamp mapping in the same row ( In case of HBase, in a separate column family). This



Got suggestions? We would love to hear your feedback. Loved InterviewBit? Write us a testimonial.

will not be a lot of data and it's okay to read it completely (of course, with caching).

(<http://www.quora.com/What-is-your-review-of-InterviewBit>)

Recent messages in a conversation : Loading all the messages or loading all the



messages in a conversation would make this really expensive. Doing the same thing with messageIDs is still an improvement in terms of the amount of data that we have to load.

As an improvement, if key in the same index is conversationID\_timestamp, we can use a prefix search of conversationID and use the most recent messages based on timestamp in the key ( assuming, the data is stored sorted with the key ).



**Q:** How would we do sharding?

**A:** HBase inherently use consistent hashing(in this case on user\_id). We explain it in detail at <https://www.interviewbit.com/problems/sharding-a-database/> (<https://www.interviewbit.com/problems/sharding-a-database/>)



**Q:** How would we handle a DB machine going down?

**A:** We explain in detail how to build a reliable and consistent database system in <https://www.interviewbit.com/problems/highly-consistent-database/> (<https://www.interviewbit.com/problems/highly-consistent-database/>)



**Q:** What are some other things we can do to increase efficiency of the system?

**A:** Caching ( Its the answer in most cases, isn't it? :) ).

This one however is not that easy. If you remember, one of our design goals was to ensure tight consistency. Most distributed caching system are good with availability and they become eventually consistent. But they are not tightly consistent.

For example, let's say I have a distributed cache. If a user's messages or



conversations are spread across machines, then it starts causing trouble for us because : **Got suggestions ? We would love to hear your feedback.** **Loved InterviewBit ? Write us a testimonial.**

The changes are no more atomic. Consider the case when messages for a user (http://www.quora.com/What-is-your-review-of-InterviewBit)

are one machine and conversations on another. When a message is added, the update request is sent to the server with messages and server with conversations for this user. There could potentially be a period when one server has processed the update and the other has not.

If changes are not atomic, the system is not tightly consistent anymore. I might see different views based on when I query the system.

One way to resolve this is to make sure that caching for a user completely resides on one server. The same server can also have other users as well, but users are assigned to exactly one server for caching. To further ensure consistency, all the writes for that user should be directed through this server and this server updates its cache when the write is successful on DB before confirming success.

There are some issues with this system :

A single point of failure for the user : My reads and writes are routed through this caching server. If this caching server suddenly dies, then my reads and writes suddenly start failing.

To resolve this, we should be able to quickly detect the failed machine, mark it as dead and start from scratch ( as cache, reading from DB) on a separate machine. If we have servers for backup, we can just have a heartbeat mechanism (<http://searchenterpriselinux.techtarget.com/definition/Heartbeat>) to detect the server going down and we can activate the backup server.

Need for a reliable routing service : Obviously, if a user is assigned to a server, then there has to be a service which should be able to track that. All request would be routed through this service ( kind of like a load balancer with a lookup based on user\_id ). The machine would need to track the IP of the active server for the user ( If the server goes down, then the routing service should know instantly to route to the new server which starts with a cold cache).

Distribution of user into different servers is another problem in itself. If we do static allotment, how do we handle the following cases :

More servers are added to the caching pool. How do we re-distribute the users without causing a cold cache for the whole userbase?

More users would keep registering. How would they be assigned ensuring uniform load?

What if we did not have a backup server and I had to re-distribute this user's load into



**Got suggestions ? We would love to hear your**

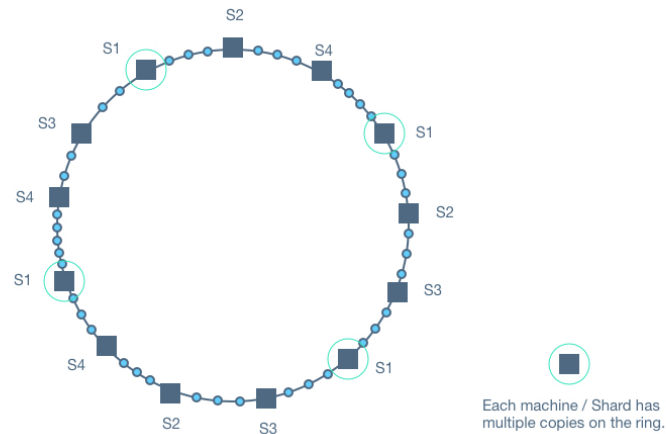
**feedback.**



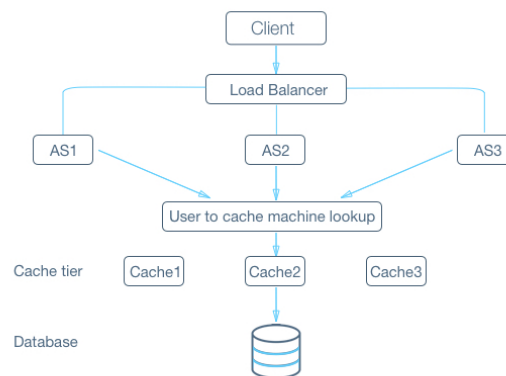
**Loved InterviewBit ? Write us a testimonial.**

A classic solution to this problem is consistent hashing with multiple tokens for each server ( See diagram attached ). For more details, read <http://www.quora.com/What-is-your-review-of-InterviewBit>

<https://www.interviewbit.com/problems/sharding-a-database/>  
 (https://www.interviewbit.com/problems/sharding-a-database/)



A minor problem - Multiple concurrent writes : The caching server will also multiple indices corresponding to the mailbox ( the ones for recent conversations / recent messages ). A single write would affect multiple columns. While a NoSQL DB might guarantee atomicity on a row level, in the caching layer, we will have to guarantee it artificially. One simple way of solving it would be to have a user level lock in the caching server for the user which allows only one write operation to go through at a time.



\* AS - Application Server



Got suggestions ? We would love to hear your feedback.



Loved InterviewBit ? Write us a testimonial.  
 (<http://www.quora.com/What-is-your-review-of-InterviewBit>)

System Design (/courses/system-design)

☐ Show Notes

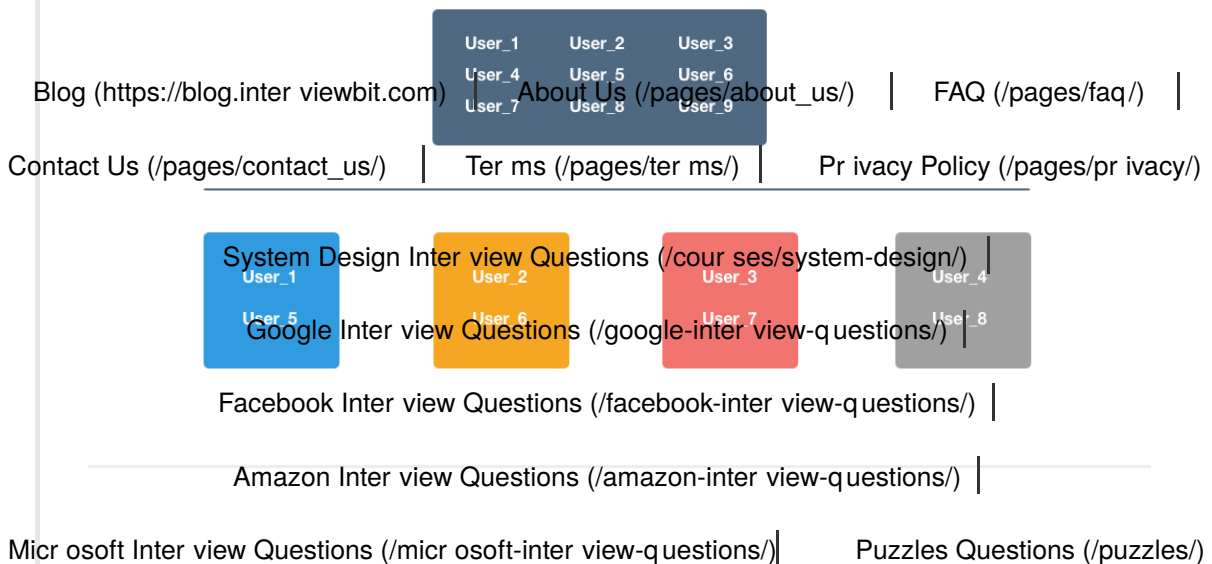
/ Storage Scalability (/courses/system-design/topics/storage-scalability/) / Sharding A Database

**Sharding a Database**

Bookmark



Let's design a sharding scheme for key-value storage.

**Features:**Like Us (<https://www.facebook.com/interviewbit>)Follow Us ([https://twitter.com/interview\\_bit](https://twitter.com/interview_bit))Email (<mailto:hello@interviewbit.com>)

“ This is the first part of any system design interview, coming up with the features which the system should support. As an interviewee, you should try to list down all the features you can think of which our system should support. Try to spend around 2 minutes for this section in the interview. You can use the notes section alongside to remember what you wrote. ”



Got suggestions? We would love to hear your

Q: What is the amount of data that we need to store?

A: Let's assume a few 100 TB.



Loved InterviewBit? Write us a testimonial.

(<http://www.quora.com/What-is-your-review-of-InterviewBit>)

**Q:** Will the data keep growing over time? If yes, then at what rate?

**A:** Yes. At the rate of 1TB per day.

**Q:** Can we make assumptions about the storage of machines available with me?

**A:** Let's assume that machines have a RAM of 72G and a hard disk capacity of 10TB.

**Q:** How many machines do I have to begin with?

**A:** Let's assume we have 20 machines to begin with. More machines will be available on request if needed.

**Q:** Are all key value entries independent?

**A:** Yes. A typical query would ask for value corresponding to a key.

## ● Estimation:

“ This is usually the second part of a design interview, coming up with the estimated numbers of how scalable our system should be. Important parameters to remember for this section is the number of queries per second and the data which the system will be required to handle. Try to spend around 5 minutes for this section in the interview. ”

❓ ◀ Total storage size : 100 TB as estimated earlier

Storage with every machine : 10TB

**Q:** What is the minimum number of machines required to store the data?

**A:** Assuming a machine has 10TB of hard disk, we would need minimum of 10 machines to store the said data. Do note that this is bare minimum. The actual number might be higher. In this case, we have 20 machines at our disposal.



Got suggestions? We would love to hear your



Loved InterviewBit? Write us a testimonial.

Feedback

(<http://www.quora.com/What-is-your-review-of-InterviewBit>)

InterviewBit)



**Q:** How frequently would we need to add machines to our pool ?

**A:** The data grows at 1TB per day. That means that we generate data that would fill the storage of 1 machine ( 10TB ) in 10 days. Assuming, we want to keep a storage utilization of less than 80%, we would need to add a new machine every 8 days.



## ● Deep Dive:

“ Lets dig deeper into every component one by one. Discussion for this section will take majority of the interview time(20-30 minutes). ”



“ Note : In questions like these, the interviewer is looking at how you approach designing a solution. So, saying that I'll use a distributed file system like HDFS is not a valid response. It's okay to discuss the architecture of HDFS with details around how HDFS handles various scenarios internally. ”

**Q:** Can we have a fixed number of shards?

**A:** One qualification for a shard is that the data within a shard should fit on a single machine completely.

As in our case, the data is growing at a fast pace, if we have a fixed number of shards, data within a shard will keep growing and exceed the 10TB mark we have set per machine. Hence, we cannot have a fixed number of shards. The shards will have to increase with time.



Got suggestions ? We would love to hear your feedback.



Loved InterviewBit ? Write us a testimonial.  
(<http://www.quora.com/What-is-your-review-of-InterviewBit>)



**Q:** How many shards do we have and how do we distribute the data within the shard?

**A:** Let's say our number of shards is  $S$ . One way to shard is that for every key, we calculate a numeric hash  $H$ , and assign the key to the shard corresponding to  $H \% S$ .

There is one problem here though. As we discussed earlier, the number of shards will have to increase. And when it does, our new number of shards becomes  $S+1$ .

As, such  $H\%(S+1)$  changes for every single key causing us to relocate each and every key in our data store. This is extremely expensive and highly undesirable.



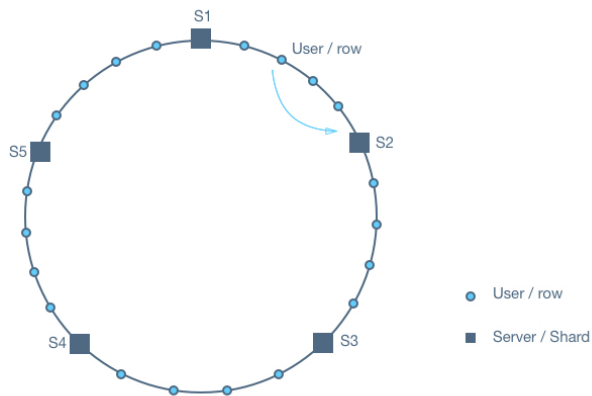
**Q:** Can we think of a better sharding strategy?

**Hint:** Consistent Hashing.

**A:** Consistent hashing is ideal for the situation described here. Let's explore consistent hashing here.

Let's say we calculate a 64 bit integer hash for every key and map it to a ring.

Let's say we start with  $X$  shards. Each shard is assigned a position on the ring as well. Each key maps to the first shard on the ring in clockwise direction.



What happens if we need to add another shard? Or what if one of the shards

goes down and we need to re-distribute the data among remaining shards?



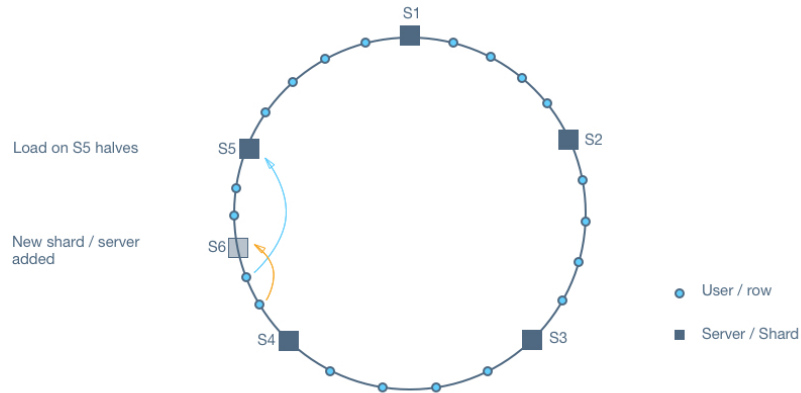
Got suggestions? We would love to hear your

feedback.

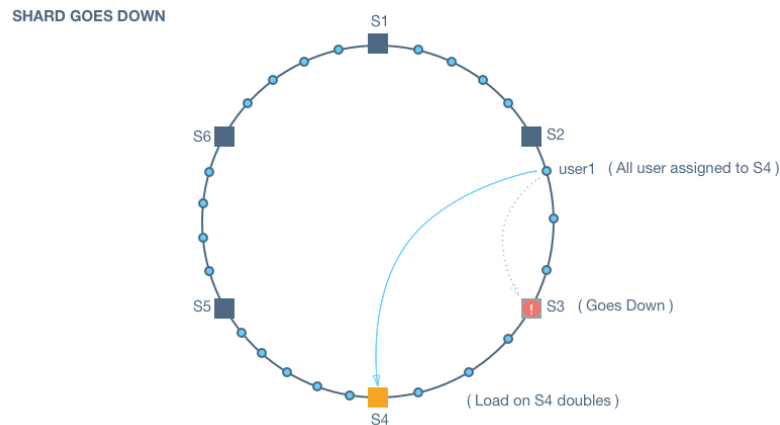


Loved InterviewBit? Write us a testimonial.

(<http://www.quora.com/What-is-your-review-of-InterviewBit>)

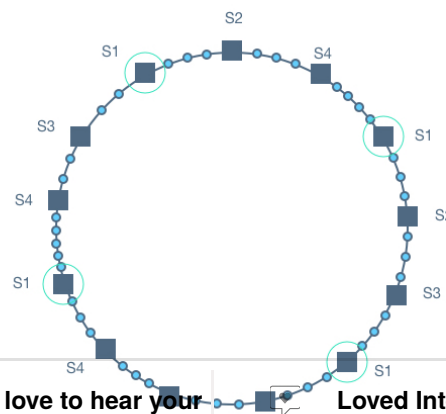


Similarly, there is a problem of cascading failure when a shard goes down.



### Modified consistent hashing

What if we slightly changed the ring so that instead of one copy per shard, now we have multiple copies of the same shard spread over the ring.



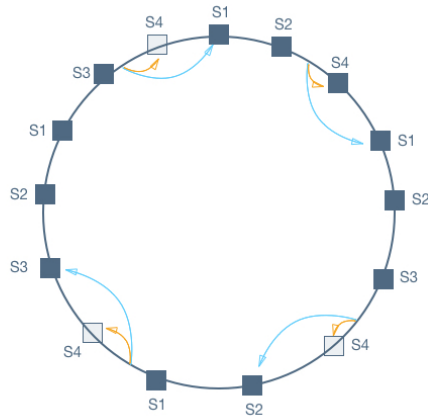
Got suggestions ? We would love to hear your feedback.

Loved InterviewBit ? Write us a testimonial.

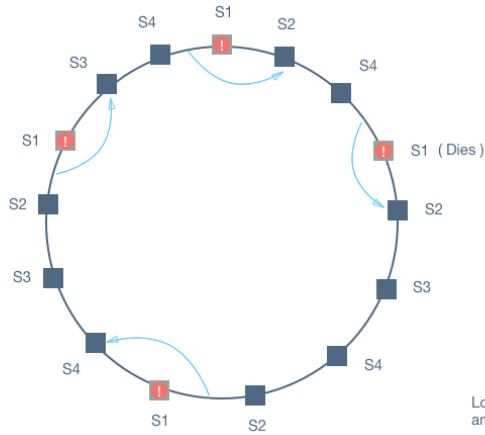
(<http://www.quora.com/What-is-your-review-of-InterviewBit>)



Case when new shar d is added :



Case when a shar d goes down : No cascading failure. Yay!



Load is distributed almost equal amongst remaining server / Shard



**You have now mastered this problem!**



Got suggestions ? We would love to hear your

feedback.

**Discussion**



Loved InterviewBit ? Write us a testimonial.

(<http://www.quora.com/What-is-your-review-of-InterviewBit>)

System Design (/courses/system-design)

☐ Show Notes

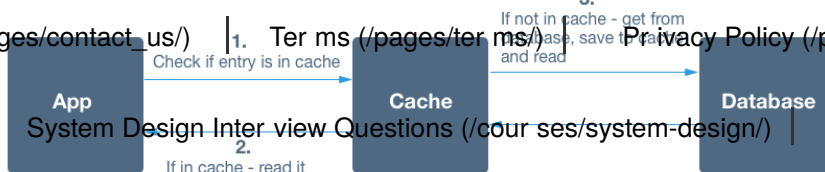
/ Storage Scalability (/courses/system-design/topics/storage-scalability/) / Design Cache

**Design Cache**

Bookmark



Design a distributed key value caching system, like Memcached or Redis.

Blog (<https://blog.interviewbit.com>)About Us ([/pages/about\\_us/](/pages/about_us/))FAQ (</pages/faq/>)Contact Us ([/pages/contact\\_us/](/pages/contact_us/))Terms (</pages/terms/>)Privacy Policy (</pages/privacy/>)

System Design Interview Questions (/courses/system-design/)

Google Interview Questions (/google-interview-questions/)

Facebook Interview Questions (/facebook-interview-questions/)

Amazon Interview Questions (/amazon-interview-questions/)

Microsoft Interview Questions (/microsoft-interview-questions/)

Puzzles Questions (/puzzles/)

Like Us (<https://www.facebook.com/interviewbit>)Follow Us ([https://twitter.com/interview\\_bit](https://twitter.com/interview_bit))**Features:**Email (<mailto:hello@interviewbit.com>)

“ This is the first part of any system design interview, coming up with the features which the system should support. As an interviewee, you should try to list down all the features you can think of which our system should support. Try to spend around 2 minutes for this section in the interview. You can use the notes section alongside to remember what you wrote. ”



Got suggestions ? We would love to hear your



Loved InterviewBit ? Write us a testimonial.

Q: What is the amount of data that we need to cache?  
<http://www.quora.com/What-is-your-review-of-InterviewBit>

**A:** Let's assume we are looking to cache on the scale of Google or Twitter. The total size of the cache would be a few TBs.

**Q:** What should be the eviction strategy?

**A:** It is possible that we might get entries when we would not have space to accommodate new entries. In such cases, we would need to remove one or more entries to make space for the new entry.

**Q:** What should be the access pattern for the given cache?

**A:** There are majorly three kinds of caching systems :

**Write through cache :** This is a caching system where writes go through the cache and write is confirmed as success only if writes to DB and the cache BOTH succeed. This is really useful for applications which write and read the information quickly. However, write latency will be higher in this case as there are writes to 2 separate systems.

**Write around cache :** This is a caching system where write directly goes to the DB. The cache system reads the information from DB in case of a miss. While this ensures lower write load to the cache and faster writes, this can lead to higher read latency in case of applications which write and read the information quickly.

**Write back cache :** This is a caching system where the write is directly done to the caching layer and the write is confirmed as soon as the write to the cache completes. The cache then asynchronously syncs this write to the DB. This would lead to a really quick write latency and high write throughput. But, as is the case with any non-persistent / in-memory write, we stand the risk of losing the data in case the caching layer dies. We can improve our odds by introducing having more than one replica acknowledging the write ( so that we don't lose data if just one of the replica dies ).

### ● Estimation:



Got suggestions ? We would love to hear your feedback.



Loved InterviewBit ? Write us a testimonial.  
(<http://www.quora.com/What-is-your-review-of-InterviewBit>)

“ This is usually the second part of a design interview, coming up with the estimated numbers of how scalable our system should be. Important parameters to remember for this section is the number of queries per second and the data which the system will be required to handle. Try to spend around 5 minutes for this section in the interview. ”

❓ ◀ Total cache size : Let's say 30TB as discussed earlier .

**Q:** What is the kind of QPS we expect for the system?

**A:** This estimation is important to understand the number of machines we will need to answer the queries. For example, if our estimations state that a single machine is going to handle 1M QPS, we run into a high risk of high latency / the machine dying because of queries not being answered fast enough and hence ending up in the backlog queue.

Again, let's assume the scale of Twitter / Google. We can expect around 10M QPS if not more.



❓ ◀ **Q:** What is the number of machines required to cache?

**A:** A cache has to be inherently of low latency. Which means all cache data has to reside in main memory.

A production level caching machine would be 72G or 144G of RAM. Assuming beefier cache machines, we have 72G of main memory for 1 machine. Minimum number of machines required = 30 TB / 72G which is close to 420 machines.

Do know that this is the absolute minimum. It's possible we might need more machines because the QPS per machine is higher than we want it to be.



Got suggestions ? We would love to hear your

● **Design Goals:**

feedback.



Loved InterviewBit ? Write us a testimonial.

(<http://www.quora.com/What-is-your-review-of-InterviewBit>)



“

**Latency** - Is this problem very latency sensitive (Or in other words, Are requests with high latency and a failing request, equally bad?). For example, search typeahead suggestions are useless if they take more than a second.

**Consistency** - Does this problem require tight consistency? Or is it okay if things are eventually consistent?

**Availability** - Does this problem require 100% availability?

*There could be more goals depending on the problem. It's possible that all parameters might be important, and some of them might conflict. In that case, you'd need to prioritize one over the other. ”*

  **Q:** Is Latency a very important metric for us?

**A:** Yes. The whole point of caching is low latency.




  **Q:** Consistency vs Availability?

**A:** Unavailability in a caching system means that the caching machine goes down, which in turn means that we have a cache miss which leads to a high latency.

As said before, we are caching for a Twitter / Google like system. When fetching a timeline for a user, I would be okay if I miss on a few tweets which were very recently posted as long as I eventually see them in reasonable time.

Unavailability could lead to latency spikes and increased load on DB. Choosing from consistency and availability, we should prioritize for availability.



Got suggestions ? We would love to hear your  
 **Deep Dive:**  
 feedback.



Loved InterviewBit ? Write us a testimonial.  
 (<http://www.quora.com/What-is-your-review-of-InterviewBit>)

“ Lets dig deeper into every component one by one. Discussion for this section will take majority of the interview time(20-30 minutes). ”

🔍 ◀ **Q:** How would a LRU cache work on a single machine which is single threaded?

**Q:** What if we never had to remove entries from the LRU cache because we had enough space, what would you use to support get and set?



**A:** A simple map / hashmap would suffice.

**Q:** How should we modify our approach if we also have to evict keys at some stage?



**A:** We need a data structure which at any given instance can give me the least recently used objects in order. Let's see if we can maintain a linked list to do it. We try to keep the list ordered by the order in which they are used. So whenever, a get operation happens, we would need to move that object from a certain position in the list to the front of the list. Which means a delete followed by insert at the beginning. Insert at the beginning of the list is trivial. How do we achieve erase of the object from a random position in least time possible? How about we maintain another map which stores the value to the corresponding linked list node.

Ok, now when we know the node, we would need to know its previous and next node in the list to enable the deletion of the node from the list. We can get the next in the list from next pointer. What about the previous node? To encounter that, we make the list doubly linked list.

Head over to <https://www.interviewbit.com/problems/least-recently-used-cache/> (<https://www.interviewbit.com/problems/least-recently-used-cache/>) to write code and see if you completely got it.



**Got suggestions ? We would love to hear your**

**feedback.**

So we will take a simple approach and implement a LRU cache using a linked



**Loved InterviewBit ? Write us a testimonial.**

(<http://www.quora.com/What-is-your-review-of-InterviewBit>)

list and a map. The Map stores the value to the corresponding linked list node and is useful to move the recently accessed node to the front of the list. Head over to <https://www.interviewbit.com/problems/least-recently-used-cache/> (<https://www.interviewbit.com/problems/least-recently-used-cache/>) to write code and see if you completely got it.



**Q:** How would a LRU cache work on a single machine which is multi threaded ?

**Q:** How would you break down cache write and read into multiple instructions?



**A:**

Read path : Read a value corresponding to a key. This requires :

Operation 1 : A read from the HashMap and then,

Operation 2 : An update in the doubly LinkedList

Write path : Insert a new key-value entry to the LRU cache. This requires :

If the cache is full, then

Operation 3: Figure out the least recently used item from the linkedList

Operation 4: Remove it from the hashMap

Operation 5: Remove the entry from the linkedList.

Operation 6: Insert the new item in the hashMap

Operation 7: Insert the new item in the linkedList.

**Q:** How would you prioritize above operations to keep latency to a minimum for our system?



**A:** As is the case with most concurrent systems, writes compete with reads and other writes. That requires some form of locking when a write is in



progress. We can choose to have writes as granular as possible to help with performance. Got suggestions ? We would love to hear your feedback.



Loved InterviewBit ? Write us a testimonial.

(<http://www.quora.com/What-is-your-review-of-InterviewBit>)

Operation 1 needs to be really fast and should require minimum locks.

Operation 2 can happen asynchronously. Similarly, all of the write path can happen asynchronously and the client's latency need not be affected by anything other than Operation 1. Let's dig deeper into Operation 1. What are the things that HashMap is dealing with?

HashMap deals with Operation 1, 4 and 6 with Operation 4 and 6 being write operations. One simple, but not so efficient way of handling read/write would be to acquire a higher level Read lock for Operation 1 and Write lock for Operation 4 and 6.

However, Operation 1 as stressed earlier is the most frequent (by a huge margin) operation and its performance is critical to how our caching system works.

**Q:** How would you implement HashMap?



**A:** The HashMap itself could be implemented in multiple ways. One common way could be hashing with linked list (colliding values linked together in a linkedList) :

Let's say our hashmap size is N and we wish to add (k,v) to it

Let H = size N array of pointers with every element initialized to NULL

For a given key k, generate  $g = \text{hash}(k) \% N$  newEntry = LinkedList Node with value = v newEntry.next = H[g]

H[g] = newEntry

More details at [https://en.wikipedia.org/wiki/Hash\\_table](https://en.wikipedia.org/wiki/Hash_table)

([https://en.wikipedia.org/wiki/Hash\\_table](https://en.wikipedia.org/wiki/Hash_table))

Given this implementation, we can see that instead of having a lock on a hashmap level, we can have it for every single row. This way, a read for row i and a write for row j would not affect each other if  $i \neq j$ . Note that we would try to keep N as high as possible here to increase granularity.

**A:** The key to understanding and optimizing concurrency problems lies in



breaking the problem down into as granular parts as possible.  
Got suggestions? We would love to hear your feedback.

Loved InterviewBit? Write us a testimonial.  
(<http://www.quora.com/What-is-your-review-of-InterviewBit>)



We can choose to have writes as granular as possible to help with performance. Instead of having a lock on a hashmap level if we can have it for every single row, a read for row  $i$  and a write for row  $j$  would not affect each other if  $i \neq j$ . Note that we would try to keep  $N$  as high as possible here to increase granularity.



**Q:** Now that we have sorted how things look on a single server, how do we share?

**Q:** What QPS would a machine have to handle if we share in blocks of 72GB?



**A:** In estimation section, we saw total data we would have to store is 30TB. For every chunk of data, we store a copy in the hashmap and we store an entry (without the value) in a linkedList. Let's assume that the size of value is big enough to ignore overheads like an entry in linkedList. We can accommodate 72G of data on every single machine (We have neglected process memory overheads for the time being). With that, we would need 420 machines. With that config, every machine would handle around 23000 QPS.

**Q:** Will our machines be able to handle qps of 23000?



**A:** CPU time available for 23k queries :  $1 \text{ second} * 4 = 4 \text{ seconds}$   
 CPU time available per query =  $4 * 1000 * 1000 / 23000 \text{ microseconds} = 174\mu\text{s}$ . Can we handle entries into a hashmap of size 72G with a CPU time of  $174\mu\text{s}$  (Do note that context switches has its own overhead. So, even with a perfectly written asynchronous server, we would have much less than  $174\mu\text{s}$  on our hand). Make sure you know about the latency numbers from here : <https://gist.github.com/jboner/2841832> (<https://gist.github.com/jboner/2841832>). The actual answer depends on the distribution of read vs write

traffic, the size of the value being read, the throughput capacity of our server.



**Got suggestions ? We would love to hear your feedback.**



**Loved InterviewBit ? Write us a testimonial.**  
 (<http://www.quora.com/What-is-your-review-of-InterviewBit>)

**Q:** What if we shar d among machines with 16GB of RAM?



**A:** Number of shar ds =  $30 * 1000 / 16 = 1875$

This leads to a QPS of appr oximately 5500 per shar d which should be r easonable ( Note that with lower main memor y size, CPU cycles r equir ed for access lower s as well ). Now, we also need to decide the shar d number for ever y key. A simple way to do it would be to shar d based on  $\text{hash}(\text{key}) \% \text{TOTAL\_SHARDS}$ . The hash function might differ based on expected pr oper ties of the key. If the key is an auto-incr emental user \_id, then  $\text{hash}(\text{key}) = \text{key}$  hashing might wor k well.

One downside to this is that if the total number of shar d changes, all the cur r ently cached data becomes invalid and all r equests would have to hit the DB to war m up the cache. The other way to do it would be to use consistent hashing with multiple copies of ever y shar d on the r ing ( Read mor e about consistent hashing at [https://en.wikipedia.org/wiki/Consistent\\_hashing](https://en.wikipedia.org/wiki/Consistent_hashing) ([https://en.wikipedia.org/wiki/Consistent\\_hashing](https://en.wikipedia.org/wiki/Consistent_hashing)) ). This would per for m well as new shar ds ar e added.

**A:** Recall that the total data we have is 30TB and for ever y chunk of data, we stor e a copy in the hashmap and we stor e an entr y ( without the value ) in a linkedList. Lets assume that the size of value is big enough to ignor e over heads like an entr y in linkedList. We can accommodate 72G of data on ever y single machine ( We have neglected pr ocess memor y over heads for the time being ). With that, we would need 420 machines which would lead to a QPS of 23000 which is not easily feasible. So we r ather cr eate shar ds of 16GB, 1875 shar ds each suppor ting qps of 5500.



**Q:** What happens when a machine handling a shar d goes down?

**A:** If we only have one machine per shar d, then if the machine goes down, all r equests to that shar d will star t hitting the DB and hence ther e will be elevated



**Got suggestions ? We would love to hear your**

As mentioned in the design goals, we would like to reduce high latency cases as much as possible.



**Loved InterviewBit ? Write us a testimonial.**

<https://www.quora.com/What-is-your-review-of-InterviewBit>

If we have a lot of machines, one way to avoid these cases would be to have multiple machines per shard where they maintain exactly the same amount of data.

A read query for the shard could go to all the servers in the shard and we can use the data from the one that responds first. This takes care of one machine going down, but introduces a bunch of other complications. If occasional high latency is not a big issue wrt product, its better to stick to one server per shard (Less maintenance overhead and a much simpler system).

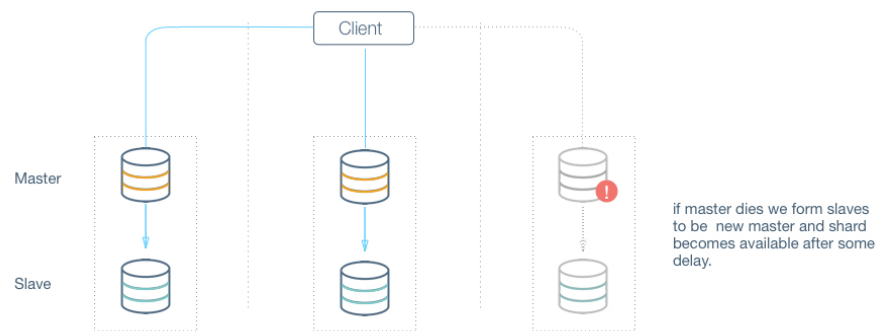
Complications of multiple servers : Since we have multiple servers maintaining the same data, it is possible that the data is not in sync between the servers. This means that a few keys might be missing on some of the servers, and a few servers might have older values for the same keys ( Assuming we support updates as well ).

Imagine a case when one of the server goes down, misses a bunch of additions and updates, and then comes back up.

There are few ways we can approach this :

Master slave technique : There is only one active server at a time in a shard and it has a follower which keeps getting the update. When the master server goes down, the slave server takes over as the master server. Master and slave can maintain a change log with version number to make sure they are caught up.

If we are fine with all servers becoming eventually consistent, then we can have one master ( taking all the write traffic ) and many slaves where slaves can service the read traffic as well.



Got suggestions ? We would love to hear your feedback.



Loved InterviewBit ? Write us a testimonial.  
(<http://www.quora.com/What-is-your-review-of-InterviewBit>)