

Tiny Url

Statement:-

A tiny url is a small url (handy url) for a large url. With this tiny url, share with other team members, typo error can be avoided.

Functional Requirement.

- 1) Create a short url for a given long url.
- 2) Re-direct/Return long url for the short url input.
- 3) If multiple user trying to ~~access~~ get the same long url, short url, for the given long url, system should give the same ~~long~~ short url.
- 4) Once short url is created for a long url, it is used for that same long url (forever) which means, that short url cannot be used again for different long url.

Non - Functional Requirement:-

- 1) System should be highly available.
- 2) Capture metrics & logs for analytics.
- 3) Auto scale.

Traffic & Capacity:

Let's assume 10M short url requests comes per month.

Also 1:200 is the write-read ratio, which means for every 1 writes there will be 200 reads. so it is a read heavy application.

Let's also assume size of the url will be 500 bytes.

$$\text{Write-requests per second} = \frac{10^M}{\cancel{30 \times 24 \times 60 \times 60}} \approx 4 \text{ URLs/sec}$$

$$\text{Read - request per second} = 4 \times 200 = 800 \text{ URLs/sec}$$

Size of the URL = 500 bytes.

Storage:-

For one month, number of URLs = 10^M ,

so total storage = $10^M \times 500 \text{ bytes} \approx 4.6 \text{ GB / per month.}$

Let's Assume system should be available for 100 years, then total storage capacity will be.

$$4.6 \times 100 \times 12 \approx 5.4 \text{ TB.}$$

Cache:-

Number of reads per second = 800 URLs

so for month = $800 \times 30 \times 24 \times 60 \times 60 = 2 \text{ billion requests}$

so cache limit will be to hold one month's URLs.

$$\text{Cache capacity} = 800 \times 30 \times 24 \times 60 \times 60 \times 500 \approx 1 \text{ GB.}$$

short-URL length:-

If it is decided to use Base-62 encoding,

Total number of URLs for 100 years = 120 Billion.

Total number of characters required in short URL = 7

$$62^5 \approx 1 \text{ billion}; 62^6 \approx \frac{56}{6} \text{ billion, but}$$

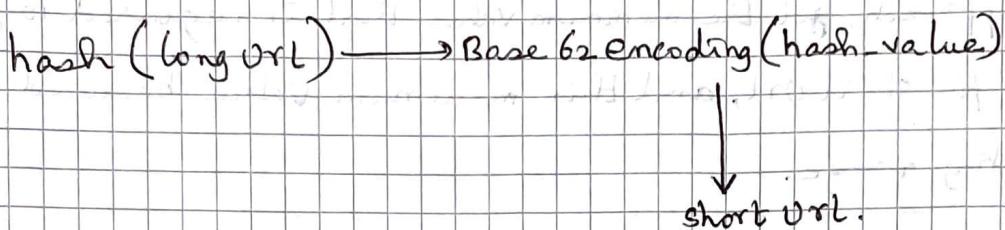
$$62^7 \approx 3.5 \text{ trillion.}$$

How base-62 works?

Base-62 encoding, performs hashing on ~~and~~ works on encoding the integers or alphanumerics.

~~for e.g.: 1 2 3~~

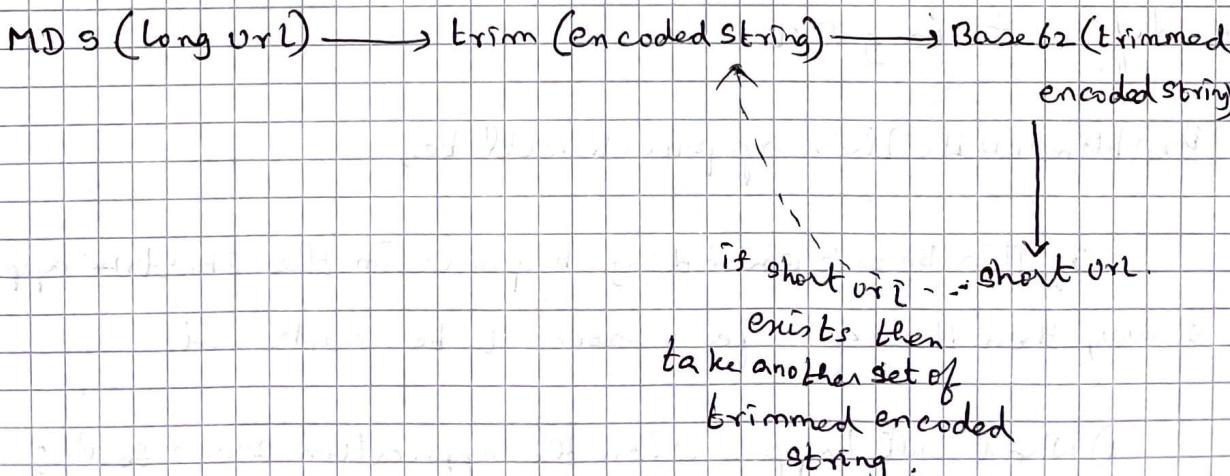
So convert the ~~long url~~ to decimal by hashing and then apply Base-62 encoding. So,



Drawback of this approach, with hashing chances of getting duplicate is more, and more i.e. 2 or more ~~hash~~ url's can return the same value, and so hashing with encoding is not possible.

Let's try with another approach, which is MD-5 algorithm.

In this approach, MD5 algorithm is applied directly on long url and the encoded string is long and so need to trim the encoded string ~~so~~ thereby to get 7-character short url. Due to this trimming duplicates are possible. ~~so there's a~~



~~KGS~~

The core problem of above approaches are getting the ~~same~~ duplicate value before encoding for the same URL. If that duplicate value for same URL is resolved, then Base62 encoding works flawlessly. To solve this duplicate issue can go with counter approach.

In this approach, get the minimum value that is required to generate 7-characters for short URL, and this minimum value be used as a starting range.

For example, let's assume 100000000 will be the starting value for generating 7-char short URL.

For the first time value will be 100000000 and increment by a fixed number, in this case let's increment by 1 & apply Base-62 encoding for the hash value.

Base62
Base (100000000) → P27AKLm (tiny URL).

Second time

Base62 (100000001) → AK79Q3t

so there won't be any duplicates i.e. no long short URL for 2 different long URL.

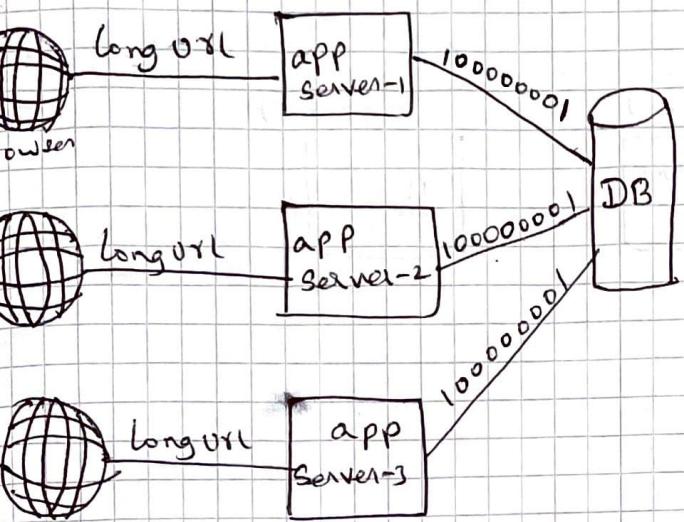
Problem with this approach will be,

If this Base62 encoding happens in the ~~shorter~~ application server, then there are few issues to be addressed.

1) What will happen, when the application server fails?

In this case, ~~range~~ counter value will be losted, and if there is only

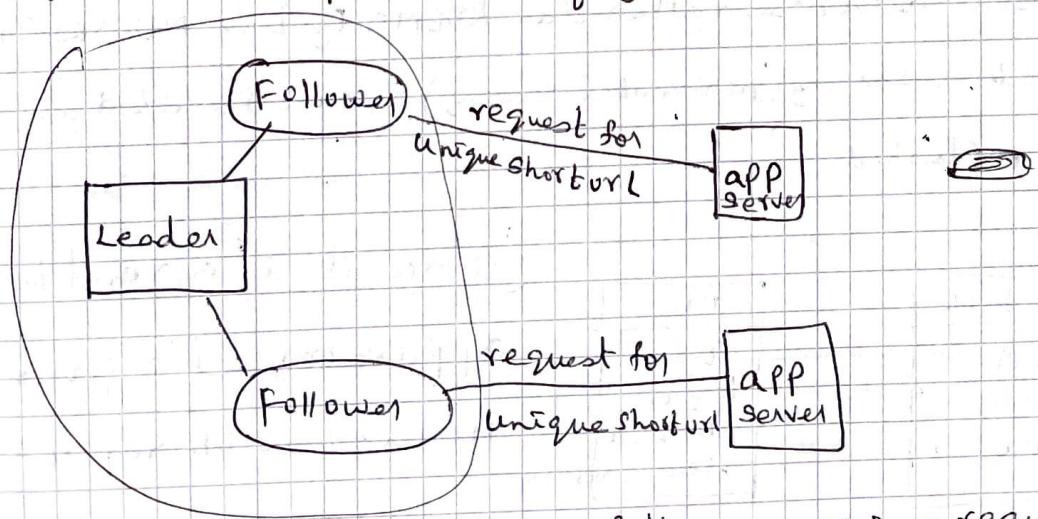
one application server, then it is a single point failure. If we have multiple application servers, and if its maintained the counter value is maintained in database, then parallel processing of request will be a problem, because of updating the same counter value from the app server to the database or different long URL, so duplicates can come up.



To avoid this issue, we can either go with

- 1) Server synchronization (Apache zookeeper)
- 2) Range counter approach.

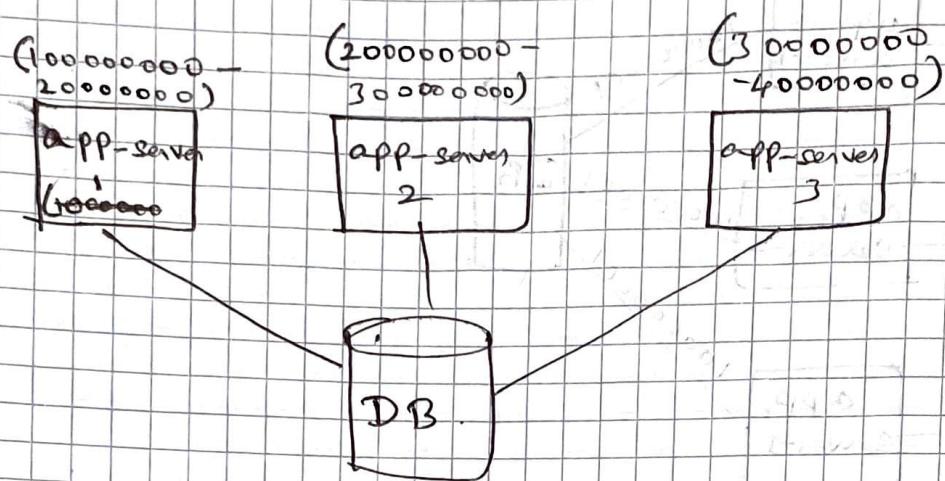
With server synchronization, we can use zookeepers, which will take the responsibility of generating the keys using counters.



When follower receives request, leader it will be forwarded to leader & leader will generate counter incrementing URL.

A Second approach is range based, where a set of ranges of counter value is assigned allocated to each application server, so that uniqueness will be achieved.

problem with this approach is, if one application server dies, then range allocated to that server will be freed, in order to avoid this, use a db to keep track of used range values & un-used ones.



With these information, in mind let's discuss the design,

Let's start with DB, ~~first~~

Since it is write-read heavy application, and database record entities have no complex relationships, so it is better to use no-sql database rather than RDBMS. By this NO-SQL DB, even ~~auto~~-scaling (horizontal scaling) is more feasible & is easily accomplished.

DB Design:-

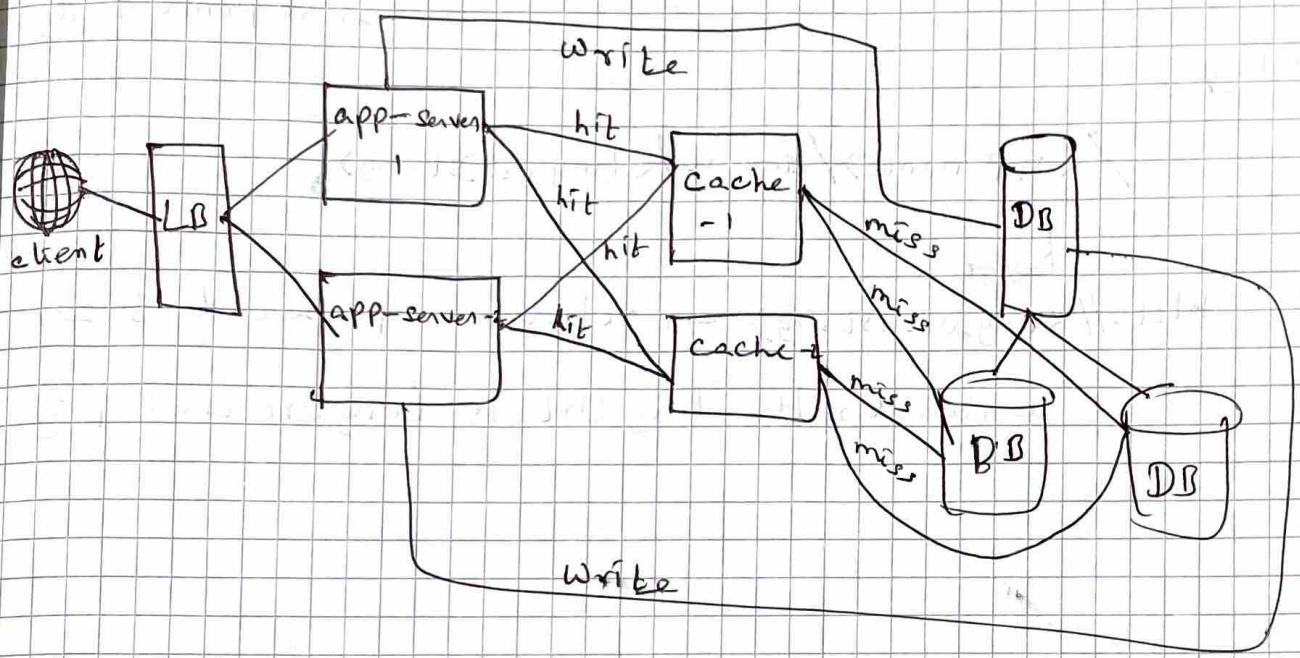
User-table
uid
name
age

URL
vid
long-vid
short-vid
create-timestamp

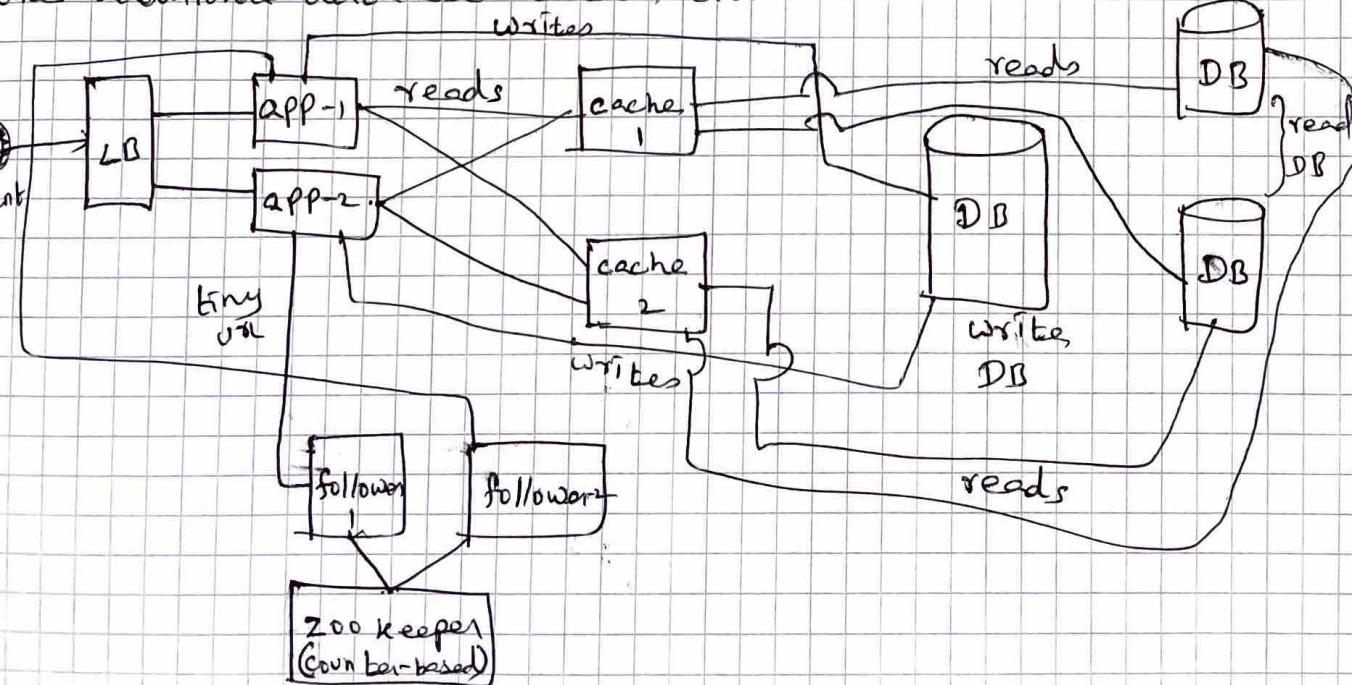
(Since it is a read-heavy application, let's go for Master-slave DB replica).

Now let's discuss about cache,

As discussed earlier, per month there will be 2 billion requests & cache storage is decided to have 1 GB. It is required to store the short URL for a long URL first in database, then in cache, because, in-distributed system, especially for this concurrent request of tiny URL, duplicate short-URL may be possible, if eventual consistency is followed over consistency. So consistency is on priority & so let's go for write-around cache.



Let's consider the tiny URL service, can go with Zookeeper Service for tiny URL generation, as range based approach needs more resource like additional database servers, etc.



Rest API calls:- (HTTP request)

Post Request:-

/<app-name>/~~tinyUrl~~ → URL
body → {
 User-id : <user-name>
 Long-URL : <URL>
}

If successful returns response code 200 with tiny url
in response body

Get Request:-

/<app-name>/tinyurl/<tinyurl-string>
http://<tinyurl-string> → returns response code 302
redirects the tiny url to long url web page