

Requirements

- 5 million new URL shortening requests / month
- 100 : 1 Read-write ratio
- Postgres DB.

Calculating storage :-

→ log URL → https://www.google.com
22 chars.

lets assume that max length of long URL be 2048 chars.
 $\Rightarrow 2 \times 2^{10} = \underline{\underline{2 \text{ KB}}}$

→ let's make short-url be of length 7 chars. $\Rightarrow \underline{\underline{7 \text{ B}}}$
https://www.chota-url.com/wX9yMpA
 7 chars.

mapping in database :-

short-url	long-url
wX9yMp	https://...google.com
⋮	⋮

- Created-at → 7 Bytes (7 chars)
- Expiry-at → 7 Bytes (7 chars)

Total storage = 2.024 KB / url shortening

so for 5 mn requests / month

$$\begin{aligned}\Rightarrow \text{storage} &= (5 \times 10^6) \times (2.021 \text{ KB}) / \text{month} \\ &= 10.105 \times 10^6 \text{ KB / month} \\ &= \underline{10.105 \text{ GB / month}} \\ &= 10.105 \times 12 \text{ GB / year} \\ &= 121.26 \text{ GB / year} \\ &\approx \underline{0.12 \text{ TB / year}}\end{aligned}$$

→ Generating 7 chars unique random id

Approach 1: Use B_{10} (base 10)
(0 to 9) = 10 unique values

$$\begin{aligned}\text{total combinations for 7 chars long string} &= \underbrace{\frac{10}{\text{char}} \frac{10}{\text{char}} \frac{10}{\text{char}} \dots \frac{10}{\text{char}}}_7 \\ &= 10^7 \\ &= 10 \text{ mn}\end{aligned}$$

∴ we have 5 mn req / month

we can only create unique short url for 2 months

⇒ NOT a good approach ⇒ Collisions after 2 months

Approach 2: Use B_{62} (base 62)

(0-9) = 10 values

(A-Z) = 26 values

(a-z) = 26 values

62 unique values

(3)

$$\begin{aligned} \text{total combinations} &= \underbrace{62 \ 62 \ \dots \ 62}_7 \\ &= 62^7 \\ &= 3.52 \text{ trillion} \end{aligned}$$

for 5 mn req / month
 \Rightarrow 60 mn req / year

$$\begin{aligned} \text{total unique combinations} &= \frac{3.52 \text{ trillion}}{60 \text{ mn}} \\ &= \frac{3.52 \times 10^{13}}{60 \times 10^6} = \frac{352 \times 10^{11}}{6 \times 10^4} \\ &= 58.6 \times 10^4 \text{ years} \\ &\approx 50,000 \text{ years} \end{aligned}$$

we can generate unique ids for 50,000 years using approach

Mapping short url & long url in DB :-

Generate 7 random numbers each can range from (0-61)

eg. 0 12 61 51 49 7 33

↓ ↓ ↓ ↓ ↓ ↓ ↓
 a l 9 z

⏟

"a l 9 z p q t" \rightarrow short url id

DB

short	long
a l 9 z p q t	\rightarrow google.com

for another turn, we can get same random nos for different long-url.
 if we simply insert in DB \Rightarrow COLLISION

So before inserting into DB,

a) check if short-url exists

b) If yes

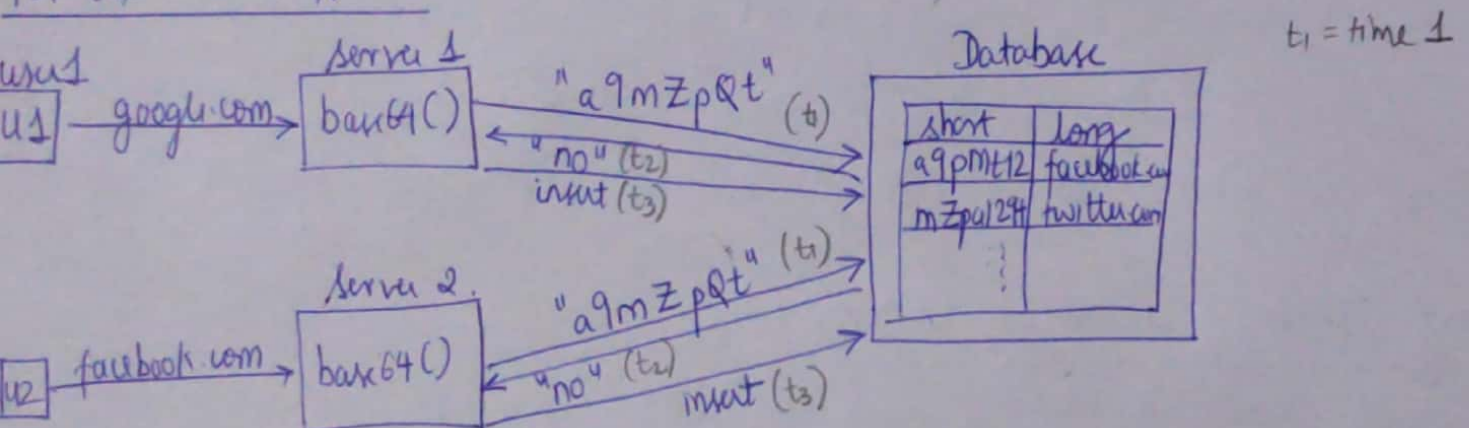
(i) generate another random nos.
do step (a)

c) else:

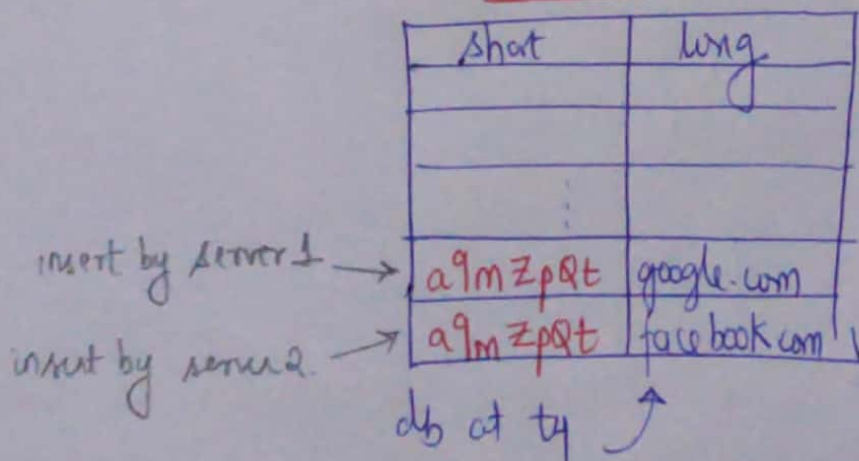
(i) insert into db.

→ This approach works well with centralized DB server, since it will execute write request serially one at a time.

For distributed server:-

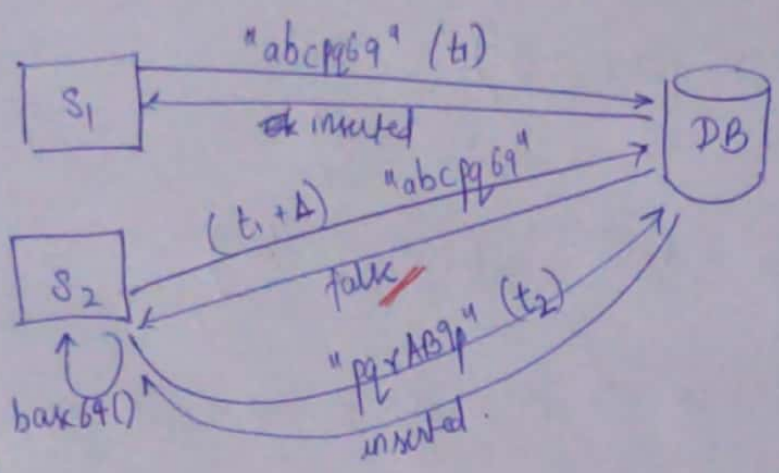


at time t₄ db will have two same short urls mapping different long urls
⇒ collision.



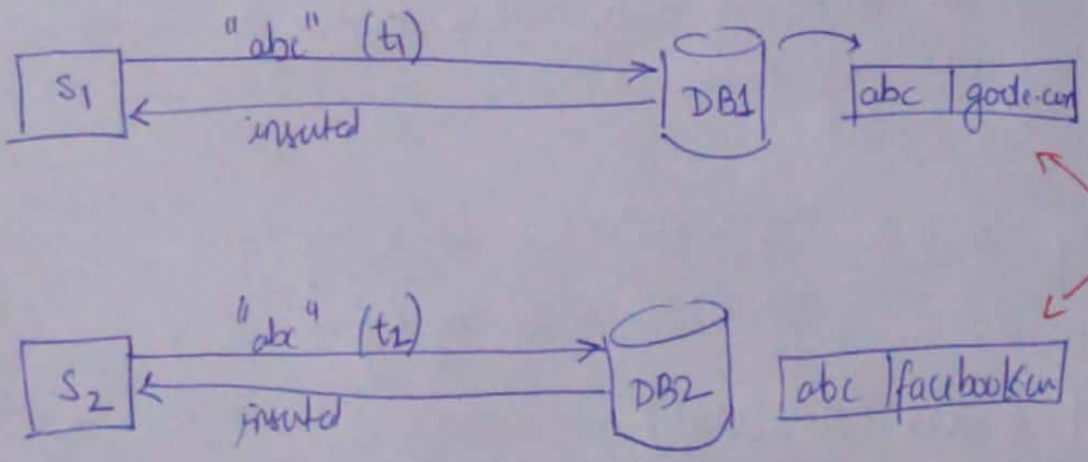
To avoid above collision problem, we can use in distributed servers

RDBMS function : "Insert if doesn't exist"



Problem 1

The above solution works well when we have CENTRALIZED DB
if we implement distributed DB, chances of data inconsistency can arise



collisions due to distributed dbs
↓
data inconsistency

Problem 2

The above solution works well when no. of ~~request~~ entries in db is less.
As the entries increase, the shortcut might already present in db & server will again generate base64.random id.
This can be time taking if no of db request increase (ie to many collisions)

To solve above 2 problems:-

- We can assign a counter to each of the server
- Each server will have a value which be reflected in counter & each server will have counter range of 1mn.

e.g. server2: value = 1 (server number - 1)
counter range = $[1,000,000 - 2,000,000)$
 $= [value * 1mn \quad to (1+value) * 1m)$

Now this counter will increase its value by 1 on each WRITE request

Calculating time required for each server to end its counter range

assumption: - We have 10+ servers in distb server.

- write req are uniformly distb across servers by load balancer
- 100:1 read-write ratio.
- 5 mn req / month.

$$\Rightarrow \frac{5 \times 10^6}{100} \text{ write req / month} \Rightarrow 5 \times 10^4 \times 12 \text{ write req / year}$$
$$\Rightarrow 6 \times 10^5 \text{ write req / year.}$$

if we have M servers (distributed)

then each server will have $\frac{6 \times 10^5}{M}$ write req / year.

$$\text{Time to complete counter range by each server} = \frac{1 \times 10^6}{\frac{6 \times 10^5}{M}} = 1.66 \times M \text{ years}$$

if no. of servers = 10 $\Rightarrow M=10$

\Rightarrow each server ~~is~~ will take 16 year to end its counter

This solution suits best for our application. \therefore No collision

⑦

