

URL Shortener

Saturday, 26 September 2020

4:48 PM

APPROACH

- ① Gather Requirements.
- ② Estimate Scale
- ③ Design the schema
- ④ Write high level APIs
- ⑤ Scale the system.

Design URL shortener

- Step 1: Gather Requirement
- X - Analytics required? how many times opened?
 - X - support customised URL?
 - ↳ /scaler
 - ↳ persist for how long? Expiry?
 - shortening capability.

Step 2: Traffic & data

- ① How much Data (single or multiple) sharding?

- ② Read heavy or write heavy
 - [QPS] ↗ Read
 - ↘ Write

- ③ No. of machines

- ④ Data access pattern.

→ LinkedIn will have much higher scale.

- 10M URL shortened every day.
- Every URL is clicked 10-50 times.
- HOT URLs → caching (Modi ji)

$$10M \times 365 \times 10 = 36.5B \text{ URL in year}$$

- ⑤ Length of shortened URL = $\geq 2^6 > 36.5B$

⑥ Size of 1 entry: $(7, 200) \sim 200B$ $\sim 7-8 \text{ chars.}$

$$\text{size in 1 year} = 200 \times 36.5 \times 10^9$$

not in a single machine $= 7.3 \text{ TB}$

⑦ Max data in a single m/c = 1TB

- ⑧ Every URL is read 10-50 times

⑨ Estimate QPS = write = 10 M / day

$$= \frac{10}{24 \times 3600} \text{ sec} \approx 1000 \text{ QPS}$$

$$\text{read} = \text{write} * 50$$

$$= 50000$$

⑩ Read Heavy system

Step 3: Design Goals

- Consistency X
- Availability ✓
- Latency very low
- Resiliency Multi-AZ

Step 4: Single system

API : a) create(url, expiry)

↳ return shortened-url.
cons.
- expiry-time
- analytics
- save-space

option 1 : Store next 10M strings in an array. Generate a random number between 0 to N. Replace element at random pos to N. and then generate from 0 to N-1.

owner-id | orig-url | short-url | expiry-url.

↳ Delete(short-url): Auth-header

c) // Analytics APIs

d) Read(short-url)

→ read from cache

- if not in cache then fetch from DB.



→ multiple machines as reg. by scale determined by step 2.

Step 1: Search, Typeahead, Autocomplete

- alphanumeric characters
- suggest term/phrases
- personalisation X
- 5 suggestion
- suggestions based on frequency.
- spell checking X
- suggest only after 3 chars.

- step 2 - 1 Billion search queries

$$- \text{autocomplete queries} : 5 \times 1B = 5B \text{ reads}$$

- writes? : 1 B writes

$$- QPS = \frac{6B}{86400} \approx 70K \text{ (100-200 machines)}$$

- Data size: ~ 120TB

10% new query Avg. 750 chars 90% old query count

(create) (update) $150 \times 100PA = 15B = 15GB/day$

$$= (5 \times 365 \times 25) GB$$

$$= 135TB$$

- Read > write

Step 3: available

- very low latency

- geo specific.

2²

2² × 2⁶

2⁶ × 2⁶ × 2⁶

2⁶⁴