# Agenda.

$\hookrightarrow$ Design **Rate Limiter**.

$\Rightarrow$ ChatGPT : 35 requests | 3hrs.

$\Rightarrow$ LinkedIn : Limit on no. of connections request we can send/accept per day

$\Rightarrow$ Twitter : View 10K tweets/day.

# Rate Limiter.

$\hookrightarrow$ Prevent (DDOS) attack

Distributed Denial of Service.

$\Rightarrow$ Our system can down.

$\Rightarrow$ fair usage of the services.

$\Rightarrow$ limit on the no. of API calls can be made from a particular user|IP.

M/c ≈≈≈≈ [≈≈≈] ——— 1M/sec ———→ Server

Ip.

API
Gw
+
LB

Rate
Limiter

Rate Limiting
Service

T/F

API
Gw
+
LB

⟹ Http Status Code.

└→ 429

└→ Too many request

# Functional Requirements (features)

→ Ability to limit the requests based on a particular criteria.

→ Configurable.

→ Different API's can have different limiting criterias.

API endpoint : 10 requests | user | min.

/creatPost

/ viewPost : 10K requests | user | Day.

# Non functional Requirements
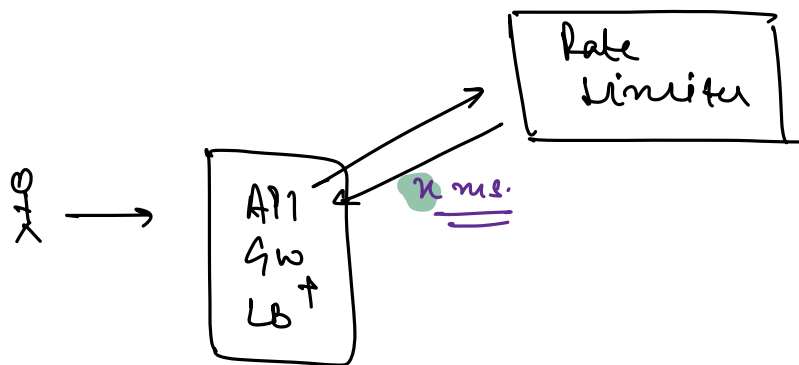(Trade Offs | Architecture)          CAP.

⇒ Consistency (vs) Availability ✓

High Availability.

⇒ If Rate limiting service is Unavailable then our entire service will be Unavailable. So, High availability is MUST of Rate limiter.

⇒ Consistency ⓥⓢ low latency.

⇒ Super low latency.



⇒ Rate limiter is going to be used for API calls, and if its taking more time then latency of all other API will increase.

# Back of the Envelope Calculations.

Twitter | facebook newsfeed.

→ CreatePost (≡)

→ ViewPost (≡)

$\Rightarrow$ 3B Total users.

DAU : 1 B

Avg No. of view post Api Call/user/Day = 50

No. of Api calls/Day = 50 B.
(view post)

No. of create post/day = (10% of 1 B) * 5
API

$= 100M \times 5$

$= 500M$

$= 0.5 B.$

Total # of API requests = 60B/Day.

$= \dfrac{60 \times 10^9}{10^5}$ q.ps.

86400 sec/Day $\longrightarrow 10^5$

$= 60 \times 10^4$ q.ps.

$= 600k$ qps

# Read Heavy (or) Write Heavy.

⇒ Both Read & Write Heavy.

# Storage.

api-endpoints.

| userId/ IP | API id | Count |
|---|---|---|
| 1042 | 1xyz | 100 |

8b

4B.

8B.

20B

| | id |
|---|---|

1 Billion users × 10 APIs × 20 Bytes.

$200 \times 10^9 \times$ Bytes.

200 GB

⇒ SHARDING ✗

⇒ No need to Persist Data.

⇒ Cache ✓

⇒ Distributed Rate Limiter

Rate Limiter

S1

S3

LB

S2

S4

U → →

API
GW
+
LB

# 1. Token Bucket Algorithm.

50 tokens/min ] refil rate.

1000

= 0

⇓

Decline

> 0, -1

⇓

Allow.

user-id = 1041

API ⇒ /accept Connection

⇓

50 req/min.

user-id + API ⇒ Count.

⇒ Bursty Traffic.

2. **Leaky Bucket Algo**

Refill rate
50 tokens/min.

< 0
Decline

> 0, -1
Accept.
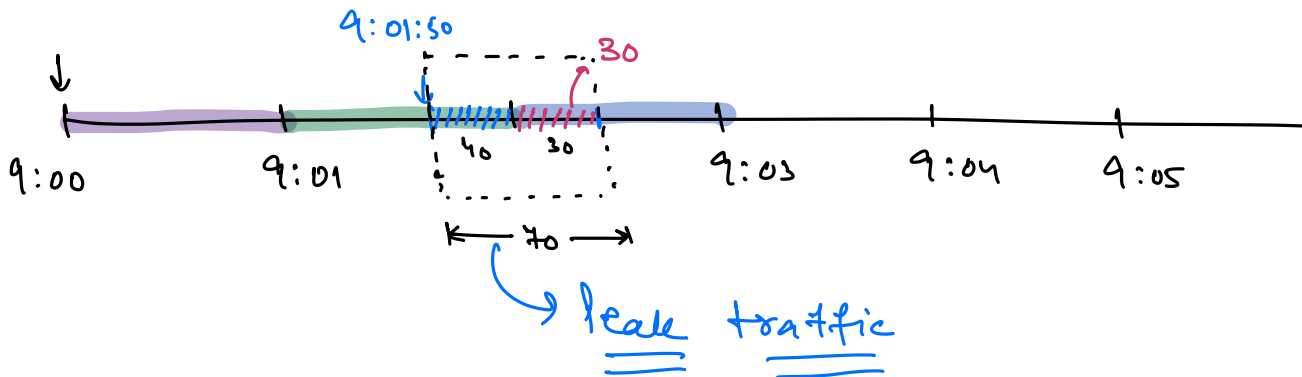
leak rate

3.

50 tokens/min

Queue ⇒

0

50

⇒ Resource extensive algo.

⇒ Unfair to the new requests as the older requests will get executed first.
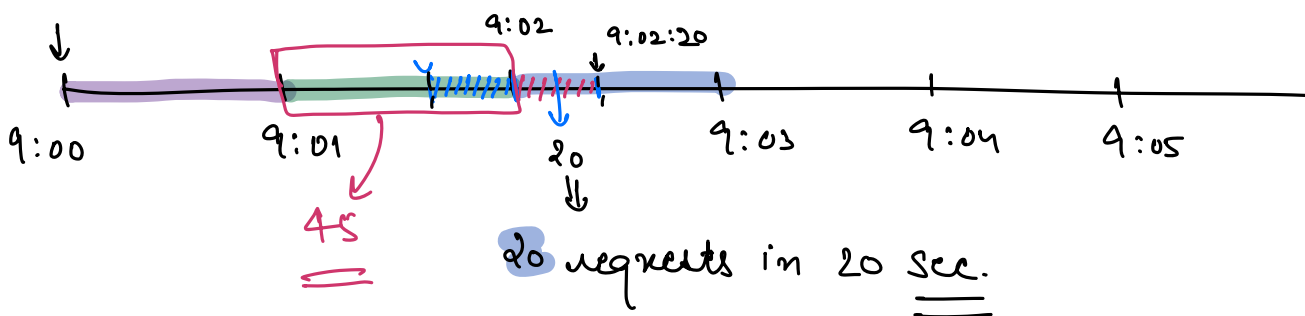
# 4. Fixed Window Counter

50 requests | min.

At the start of every window, reset the counter.



9:01:50    30
9:00   9:01   40   30   9:03   9:04   9:05

← 70 →

→ Peak traffic

# 5. Sliding Window Counter

Fixed Window Counter + Solve peak traffic



9:02   9:02:20
9:00   9:01   9:03   9:04   9:05

45

20
⇓
20 requests in 20 sec.

60 sec → 45 req
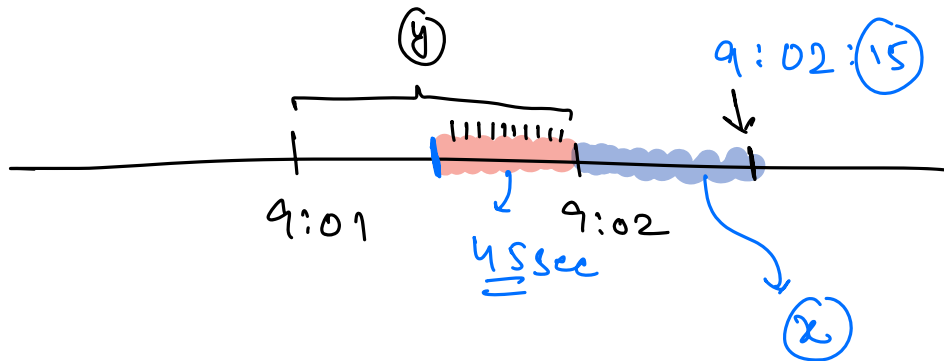
40 sec → $\dfrac{40}{60} \times 45$   3

(30)

\# of requests in the current min

$$= 20 + \frac{40}{60} \times 45$$

$$= \underline{50}$$



\# of req $= x + \frac{\overset{3}{48}}{\underset{4}{60}} \times y$

$\Rightarrow$ Approximation.