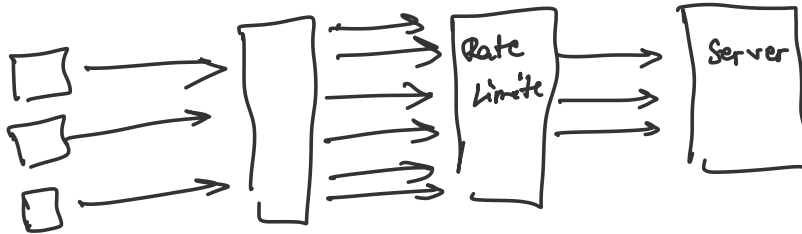# Rate limiter

Saturday, 15 November 2025    10:33 AM

Q. What is a Rate limiter?

→ A software component that can control the amount of requests coming to a system.



Q. why do we need a rate limiter?

→ To protect attack (DDOS)
   ↓        ↳ Denial of Service
   Distributed

→ Brute force login attempts

→ Reducing infra cost

→ Help making system available to users uniformly.

$256^8$

$(2^8)^8$

$\boxed{2^{64}}$

→ Reduce burst traffic spikes.

# Requirements

## Functional

→ Limit number of request to an API within a time limit.

{ → The rate limit should be considered across different servers.

## Non-functional

{ Component { separate component

→ Highly available

→ Low latency.

[ datastructure → map
  algo

( HTTP 429 → Too many requests )

# (Types)

→ Hard → requests cannot exceed the threshold.
  ↳ Banking.

→ Soft → requests can exceed a <u>certain percentage.</u>
  ↓ grace cushion

  (100) → 10% } → 110

→ Dynamic → • complex throttling
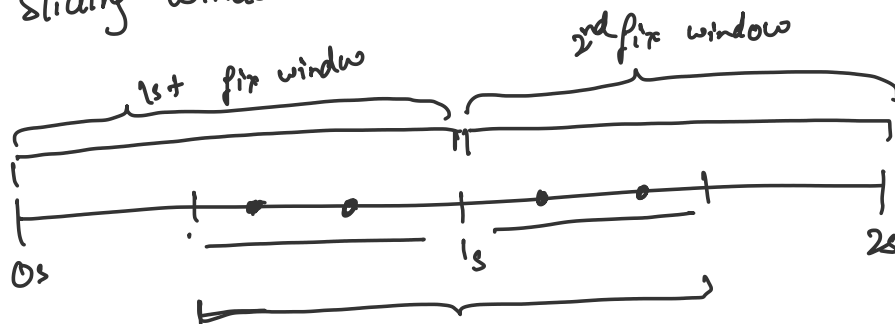  • when you want the system to adapt to user.

(1000)/m → (100/m)/customer
  (900)

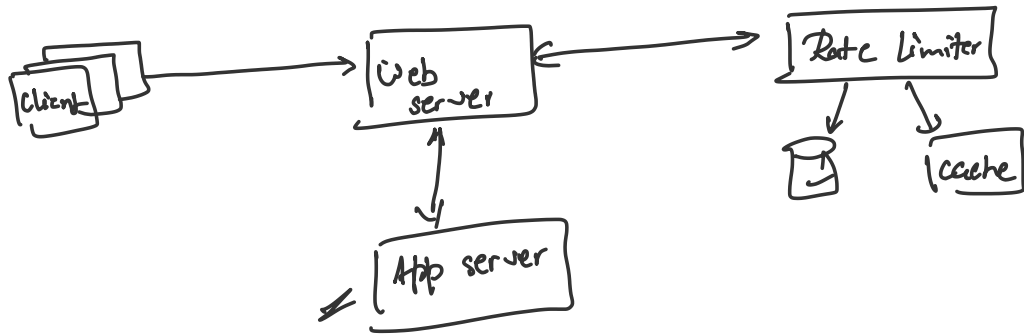(AWS S3) → 3500/s (write)
  5500/s (read)

# Algorithms

→ Fixed window ]→ limit after n-request in a fixed time window

→ Sliding window

# HLD



# System design

Fixed window

key → value.

user Id : { count, start Time }

Note: If
not p
address
can be

→ If user Id is not present in map, set count
(starttime) with user Id, allow the request
↳ normalize.

(los)　3:00:00
　　　　3:00:05

→ If user Id is present, & current Time - (start Ti
set current Time as start, count as 1
update the entry ) → (start a new t

→ If user Id is present & current - st

$\rightarrow$ If count < threshold, increment

$\rightarrow$ If count >= threshold, block.

Memory

userId : { count, timestamp }

8　　　8　　　8　　=24 byte

(over head)
$\downarrow$ 20 bytes needed for overhead

Total size of 1 record = 2

for tracking 1 million users :

44 bytes × 1M = (44 MB) $\rightarrow$

Rate limit is of 10 req /s

Traffic on system $\Rightarrow$ 10 × 1m =

$\rightarrow$ we need to store in a distributed cache cluster

60 M queries

# Sliding window      (1min)

track requests in the last minute from request.

3: 0 1: 30
        ↳ see requests 'starting from

key → (value)
          ↳ sorted set of ti'

userId ⟹ $\{ t_1, t_2, t_3 \dots, t_n \}$

→ req →
   ↳ check the first item of set
   ↳ if it falls within time window
      & set size >= threshold
       (block the request)
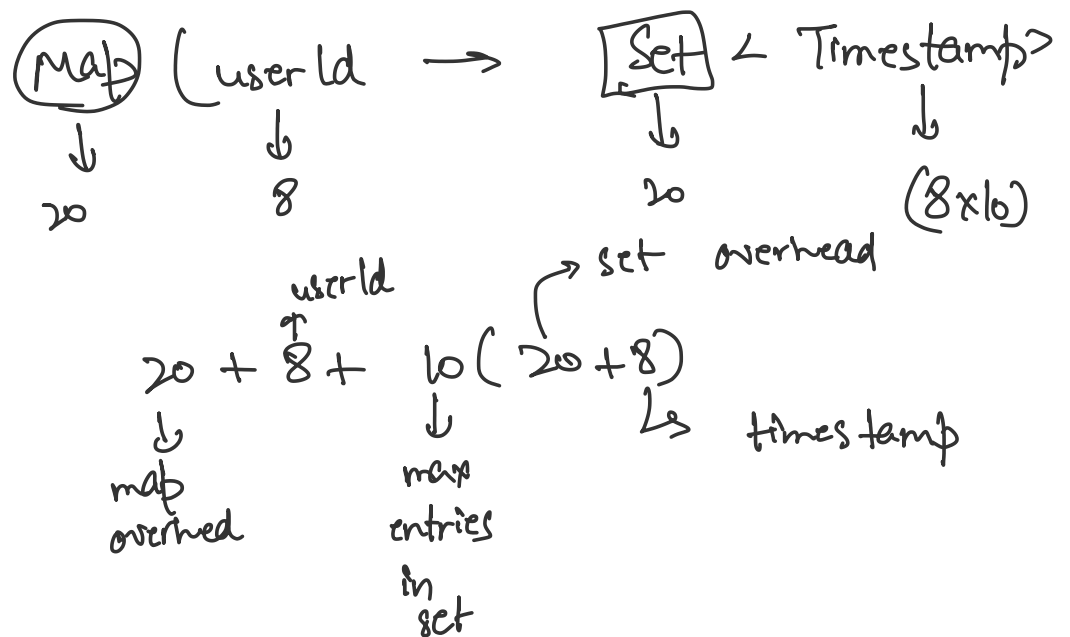

→ userId is not present.
       ↳ add userId with current time

→ userId is present

$\hookrightarrow$ remove from start till ...
than (current Time — time wi

$\hookrightarrow$ If size of set $\geq$ threshold,

$\hookrightarrow$ otherwise, add timestamp to

1 user $\Rightarrow$ <u>10 req /min</u>

(Map) (userId $\longrightarrow$ [Set] $<$ Timestamp$>$
$\downarrow$ ... $\downarrow$ ... $\downarrow$ ... $\downarrow$
20 ... 8 ... 20 ... $(8 \times 10)$

$$20 + \underset{\underset{userId}{\uparrow}}{8} + 10\overset{\rightarrow\ set\ overhead}{(20 + 8)}$$

map overhead

max entries in set

$\hookrightarrow$ timestamp

$$28 + 280 = 308 \text{ bytes.}$$

Tracking 1 million users would need.

$$1 M \times 308 = \boxed{308\ MB}$$

~7x increase.

11. Partitioning

# 

└→ key of the map

└→ range based.

# What to use as a key.

→ user → can limit requests per use
          multiple accounts can be

→ IP → can limit requests per de
       user can use multiple de

→ Hybrid → use combination of both.