Abhi tak hamare "Global Rate Limiting" dekh rahe
the matlab allover globe kahi se jaisne req di,
uspe rate limiting kar rahe the

lekin humlog Rate limiting lagate h user level pr
ya too, IP level pe.

① user ya IP level banane ki liye dikkat ye h ki,
mujhe haar ek ureid ya IP Address ke
corresponding mujhe bucket banana
padega.

Note:- Hum rate limiter API based bhi
banna sakte h,
matlab ek application mai let say
100 API, to mai chahoe to
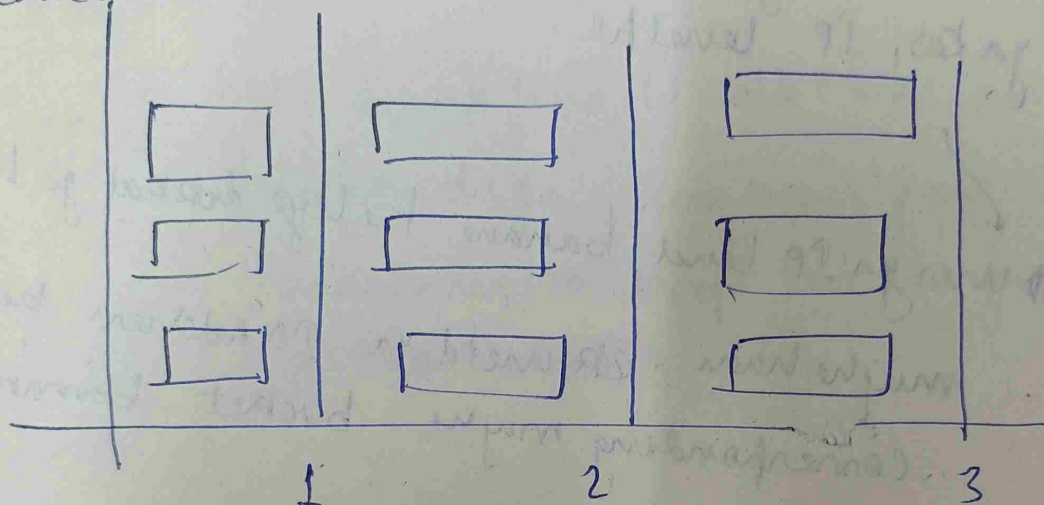haar ek ke liye alag Rate limiter
laga sakta hu
→ Reasou:- Saare API buch API
ka tank small ho sakta
jaldi process karle
aur mai uaka
rate big rakhna
type.

# Fixed Window Counter

① In a fixed interval (say a sec or a minute)
only a specific no. of req can come.

→ Isse bhi aap global ya user ya IP level
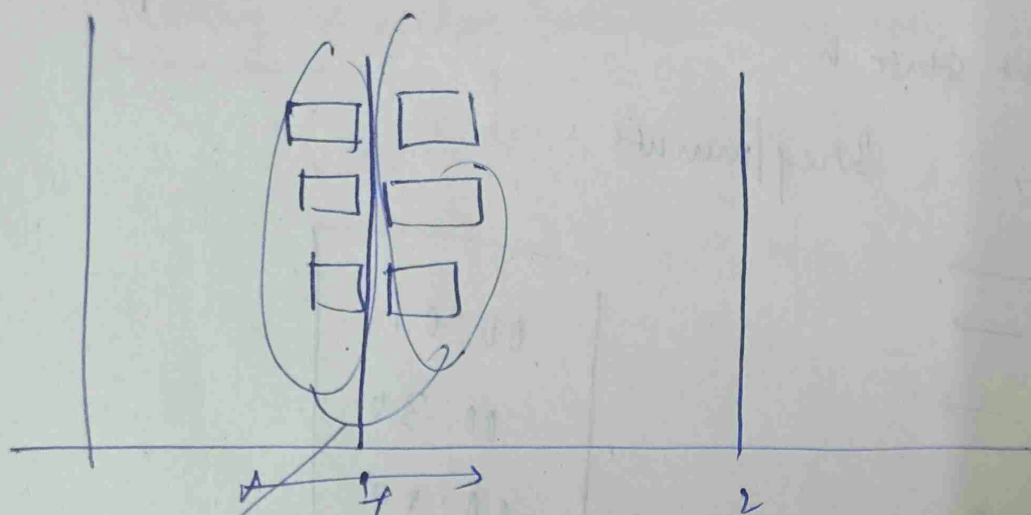pe implement kar sakte ho

I dead care m.



        1            2            3

Fixed Counter = ③ → ek particular fixed
                    time interval mai 3 se
                    jyada req ni aa sakte.

**Cons!** ~~Prost Server pe load n~~

Const-

If burst traffic comes at edges.
of window it may lead to
server crash, high latency.

( matlab PTO )

✓ ye teeno alag alag
interval mai h

→ lekin ins short interval
mai kaafi saare req aa gaye
lekin fir bhi accept karna
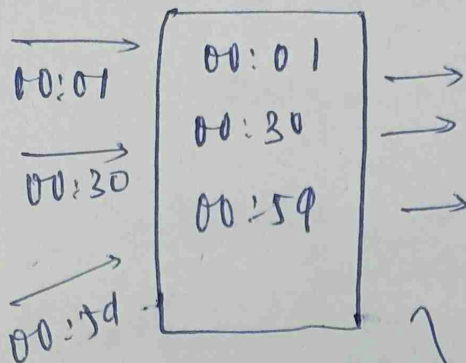padega.

## Sliding window log Algo

→ Best Algorithm for Rate limiter
→ Strict
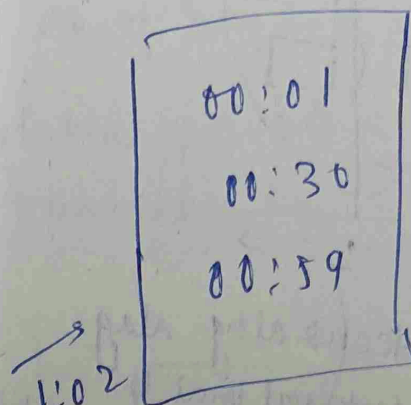→ Slow & memory consuming

### Algorithm:

(i) It will store each req in a log file.

(ii) Whenever a req comes, it will first remove all the outdated req from the file.

(iii) Then it will log the new req and check, if the counter limit reached, drop the req else send it to server.

lets say ye current state h.

3req/minute.



00:01 →
00:30 →
00:59 →
00:59 →

00:01
00:30
00:59

ablet
say ni req
aayi at 1:02

1:02

00:01
00:30
00:59

ye dekhege ki

1:02 - 00:01 > 1mn.

to 00:01 hatega.

an 1:02 aani

00:30
00:59
1:02

1:04

ab 1:04 ← 0:30 < 1 min to

ye req descard ho jayega.
ya menaging qnow man
daaldeg.

↳ Is procen p. mai memory jayda lagta h
an saan req ka log Store karke
rakhna parta h.
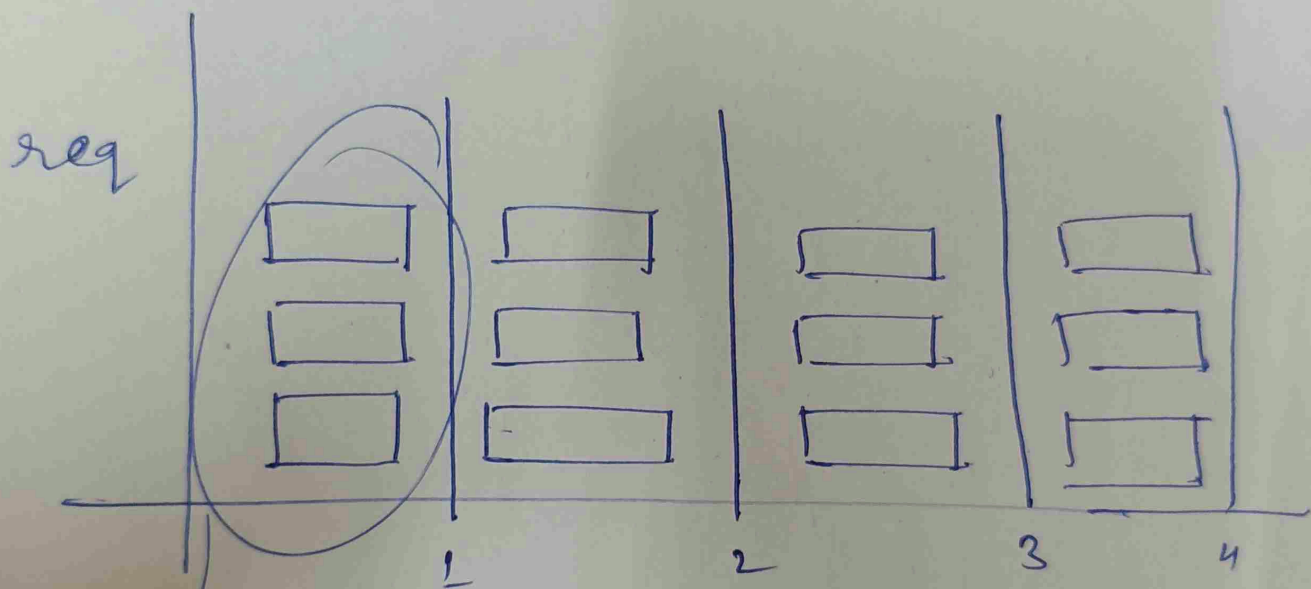
## Sliding Window Counter Algorithm

(will do at end)

# Create our own Rate Limiter

→ Aapko khud ka Rate limiter banane ke liye kin cheero ke jarurat padne wale h?

→ Which algo to implement?

→ let's say for impliary → Fixed window counter.

req



|   |   |   |   |
|---|---|---|---|
| 1 | 2 | 3 | 4 |

Assume kar rahe critically ki uniformly aa raha req aur server pe load ni pad raha.

Excepted Counter = 3
Current Count = [0 to 3]

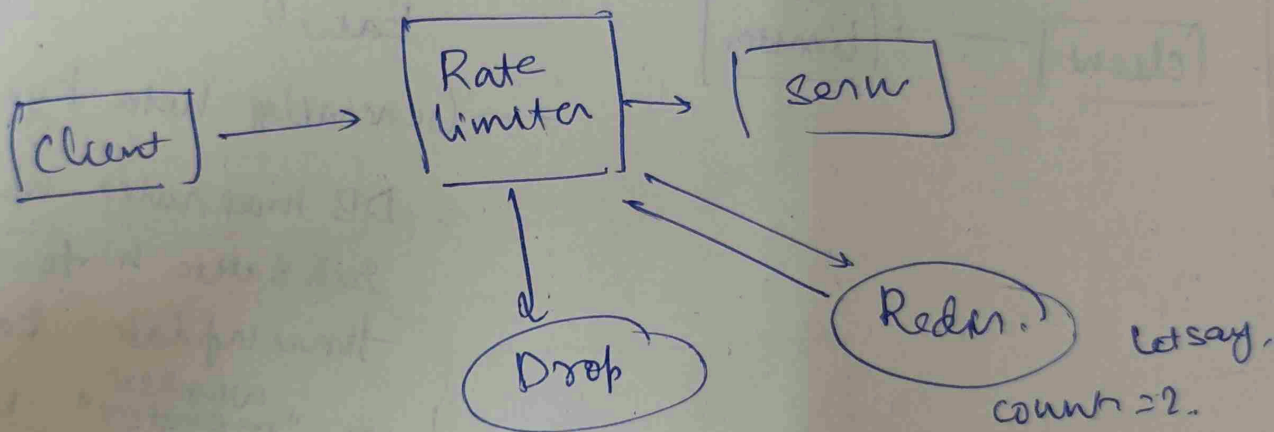} to in values ko kaha store karenge

to in values ko kaha store karenge

→ cache or Redis kyuki fast chahiye

```
[client]                    [Server]
```

Rate limiting ko kaha implement karo?
→ Mai ban raha middle wou.

```
[Client] ——→  ┌─────────┐  ——→  [Seru]
              │ Rate    │
              │ limiter │  ————→  (Reder)   Let say.
              └─────────┘ ←————             count = 2.
                  │
                  ↓
               (Drop)
```
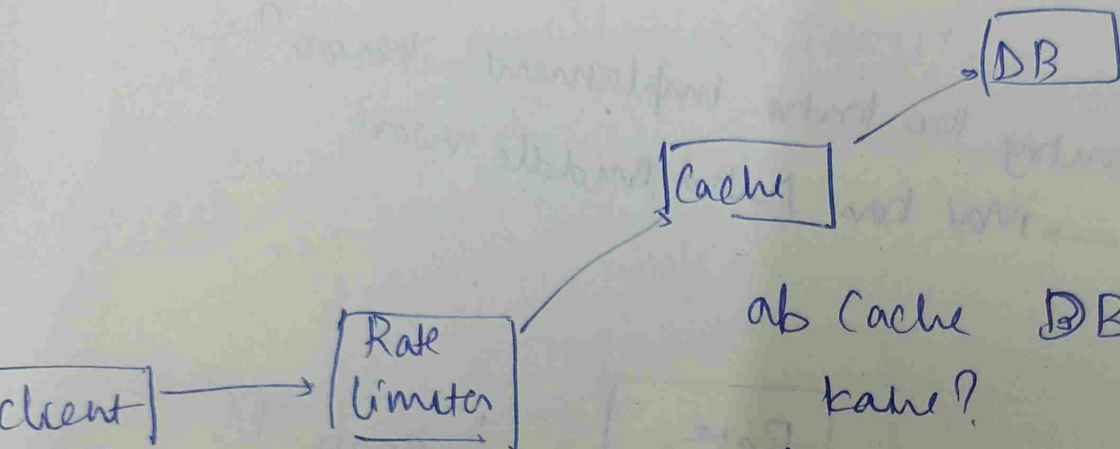
Jabhi req aayega RateLimiter to
mo Reder se puchega ki me allow
karu ya mi,
aur kuch time baad sab reset

Is Rate limiter ke kuch rules honge →Jo mai DB mai
store karunga

lekin DB se access kaafi slow ho jayega,
isiliye cache bhi use kareng.

$$\rightarrow \boxed{DB}$$

$$\boxed{Cache}$$

client $\longrightarrow$ $\boxed{\begin{array}{l}Rate \\ limiter\end{array}}$

ab Cache DB se data lega
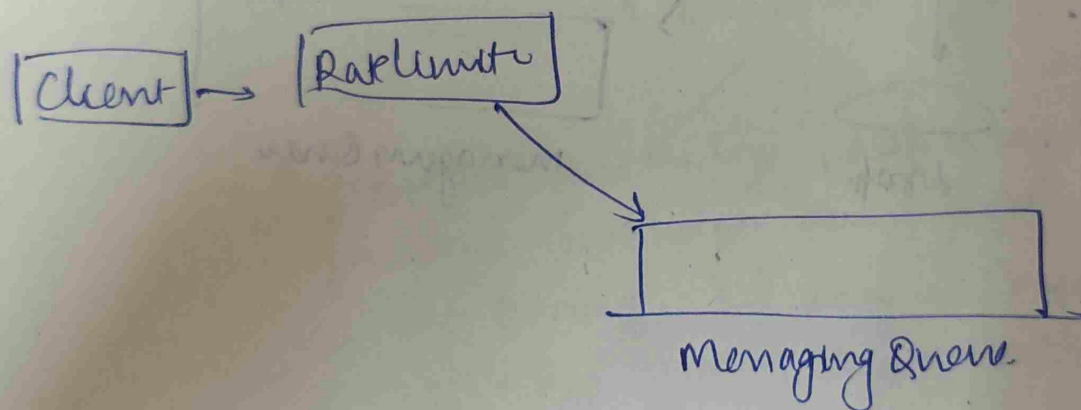kahe?

$\longrightarrow$ Generally hota kya hki

DB mai rules update hote
reh sakte h' to usse
~~time~~ update karne ke liye
hum "~~monitor~~ worker" install
karenge

ye bas ek code type h jo
haar ek particular
interval pe DB se
latest data fetch
karke Cache mai
daal dega

```
                                              [DB]
                                               ↑
                          [Cache] ← [worker]

[Client] ⟶ [Rate
           limiter] ⟶ Cache
```

ab mai saare req loojo limit exceed kare unhe.
Drop ni karunga, kuch imp req ko.
menaging Queue mai daaldiya karung

```
[Client] ⟶ [RateLimiter]
                    ↘
                     ⟶ ┌──────────────┐
                        │              │
                        └──────────────┘
                         Menaging Queue.
```

Ab mensaging Queue se req server pe chale jayeg
            by "workers"
      "workers" dekhenge ki serv free h.
          aur menagingQueue mai req
            pending to uno unre utha
             ke server ko de dega.

# whole architectur

[DB] rule.

↑

[Cache] ← [workers]

[client] → [Rate limiter] → [Server]

429
too
many
request

drop.

redis

worker

messaging Quer