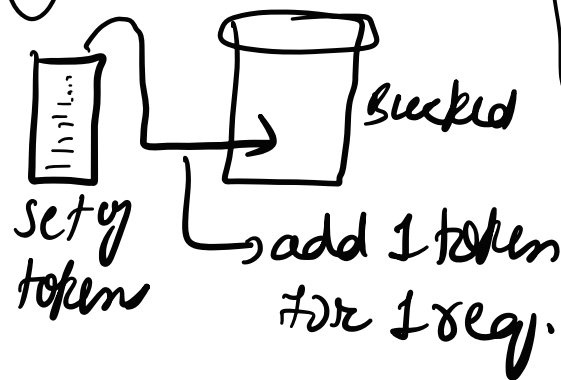# Rate limitter (RL)

- locks users request
  or
- allowes certain no. of requests in a certain time

→ prevent DDOS
— prevent cost
— Prevent/manage load of system

— ~~lesser~~ Serverside RL ⌐

## Technique to RL

1. Token bucket



set of token

add 1 token for 1 req.

Spikes in traffic can let ~~good user~~

→ IP adders (most common)
  └ — always unique
  └ — narrow down set of IPs which are bad one
  └ Send 429 status code to inform user that they are throttled
  └ Logging mechanism on our end to analyze the traffic

② <u>Fixed window system</u>

gives a time window you have fix
no. of ~~responses~~ req.
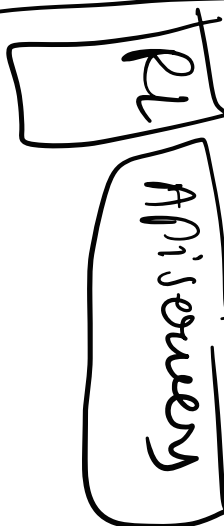no other req is entertained until that window

↳ downside → may req come towards end of window & start of window

↳ variation → Sliding window
→ tweak window size as needed

---

Server side why
Safe as DB hour control own it

<u>Components</u>
server side

RL

RL API server

Rule Engine
- algo
- rules for RL

Cache high throughput
→ writes in cache
as req comes if
not compute if
it's already
there

↳ token

# I system

Clients

API RL

success

failure 429 1771 call

Rules cache ( faster access )

Rule Engine

success

success

API web server

Cache

Logging mechanism

Long term storage of data

( Cache to store the flagged req for future )

{ long term storage to analyze

Server side blackboard on 429 1771 on blocking log mechanism

Distributed env

Client

Load Balancer

Success

API KV

429 failure Status code

Rule Engine

Cache Rules

Success

API web server

Cache

Logging mechanism

Long term storage

Scale horizon for cache to make it constant avg across SaaS? prevent any DDoS

constant to keep rules same as per region or country &

If you have different caches, then you will rate limit per datacenter

Here VS Rate limiting across globe

Left config for cache can be used
Like 1 Read & 1W cache
or 1 RW cache
across all datacenter

→ Rule cache can have multiple entries as Rules won't change over a period of time

→ IT is notgood to keep different RL based on geography as Apl can respons RL based on N