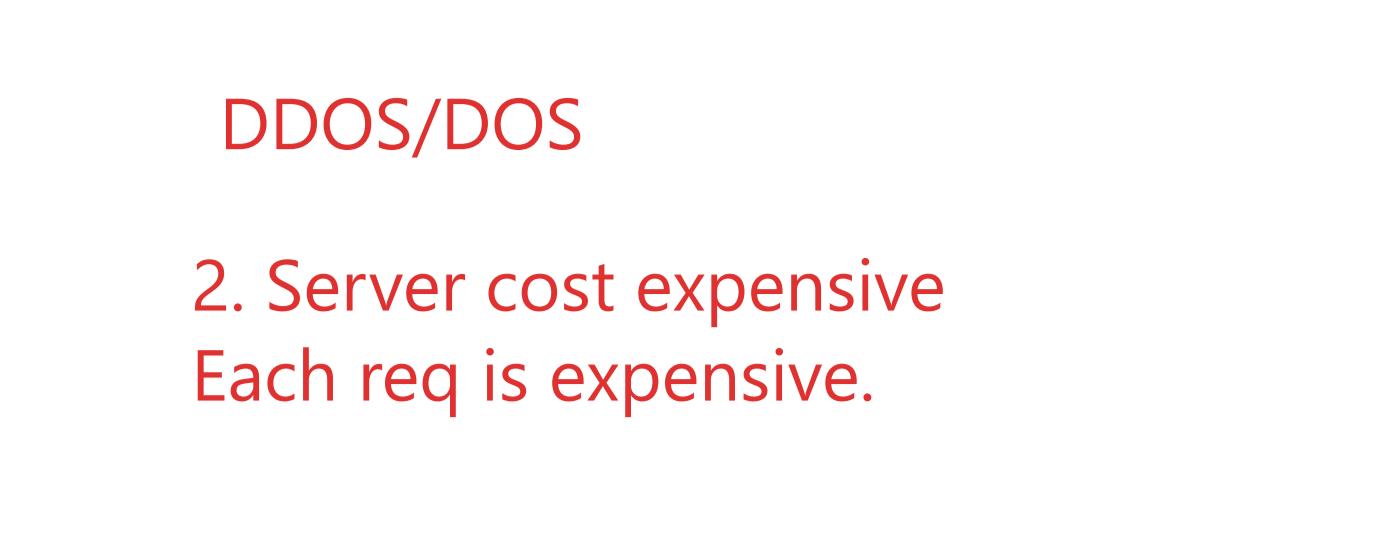


Rate limiting No.of times that a client can request a server in any timeframe
Caching Strategies

Rate limiting



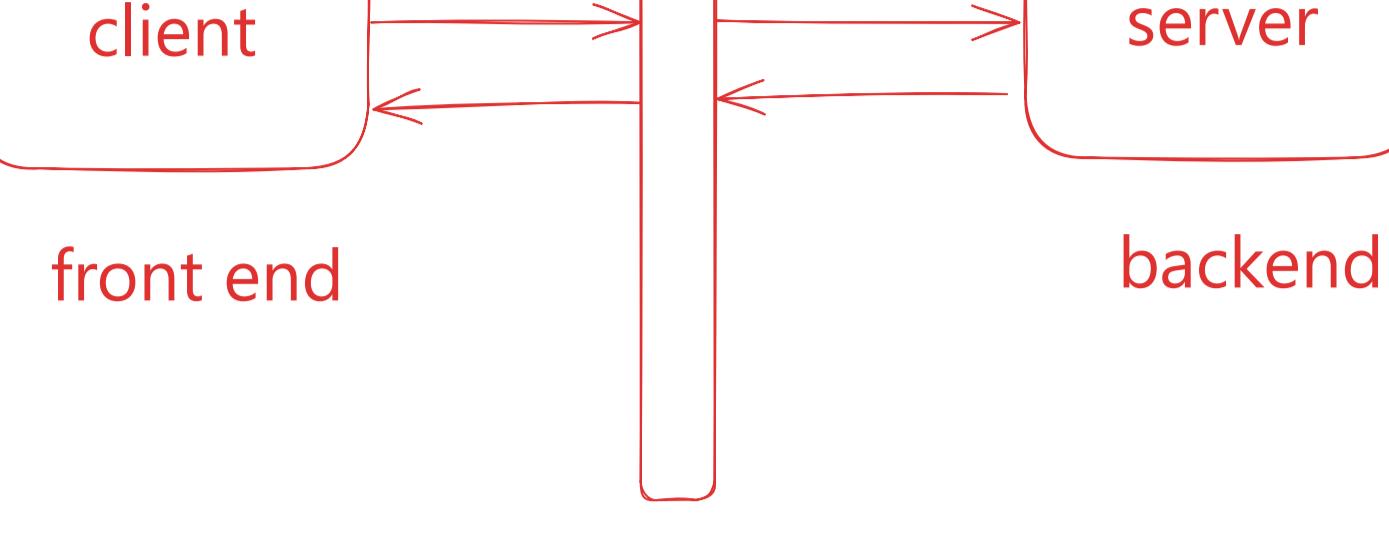
problems:

1. Client can send too many requests within a timeframe.

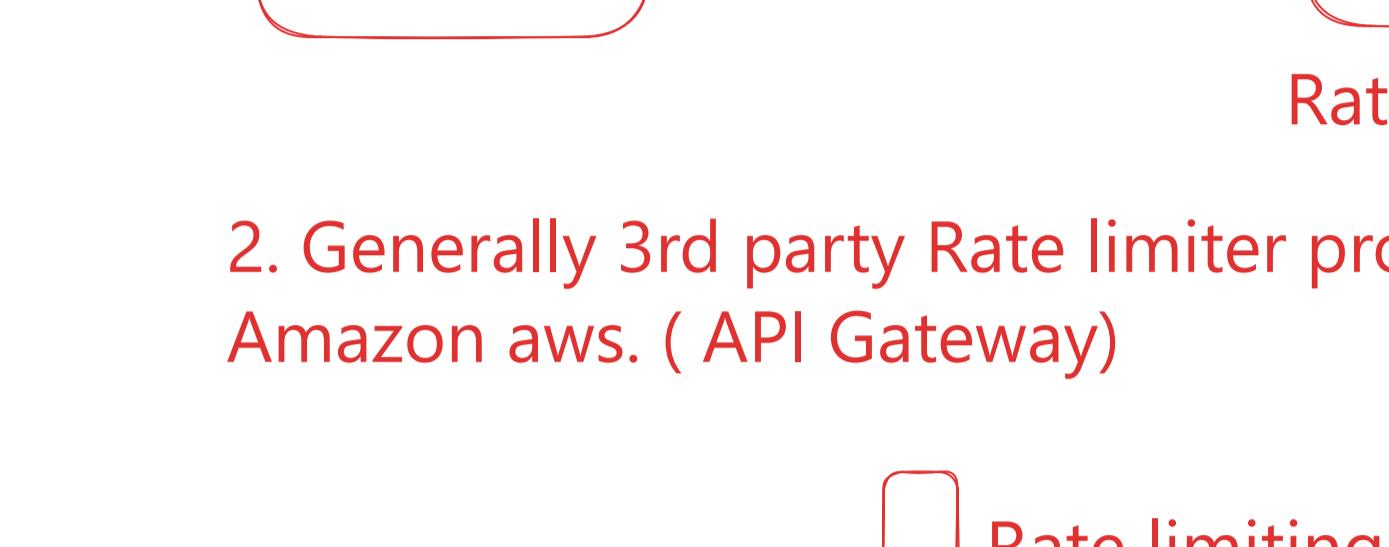
DDOS/DOS

2. Server cost expensive

Each req is expensive.



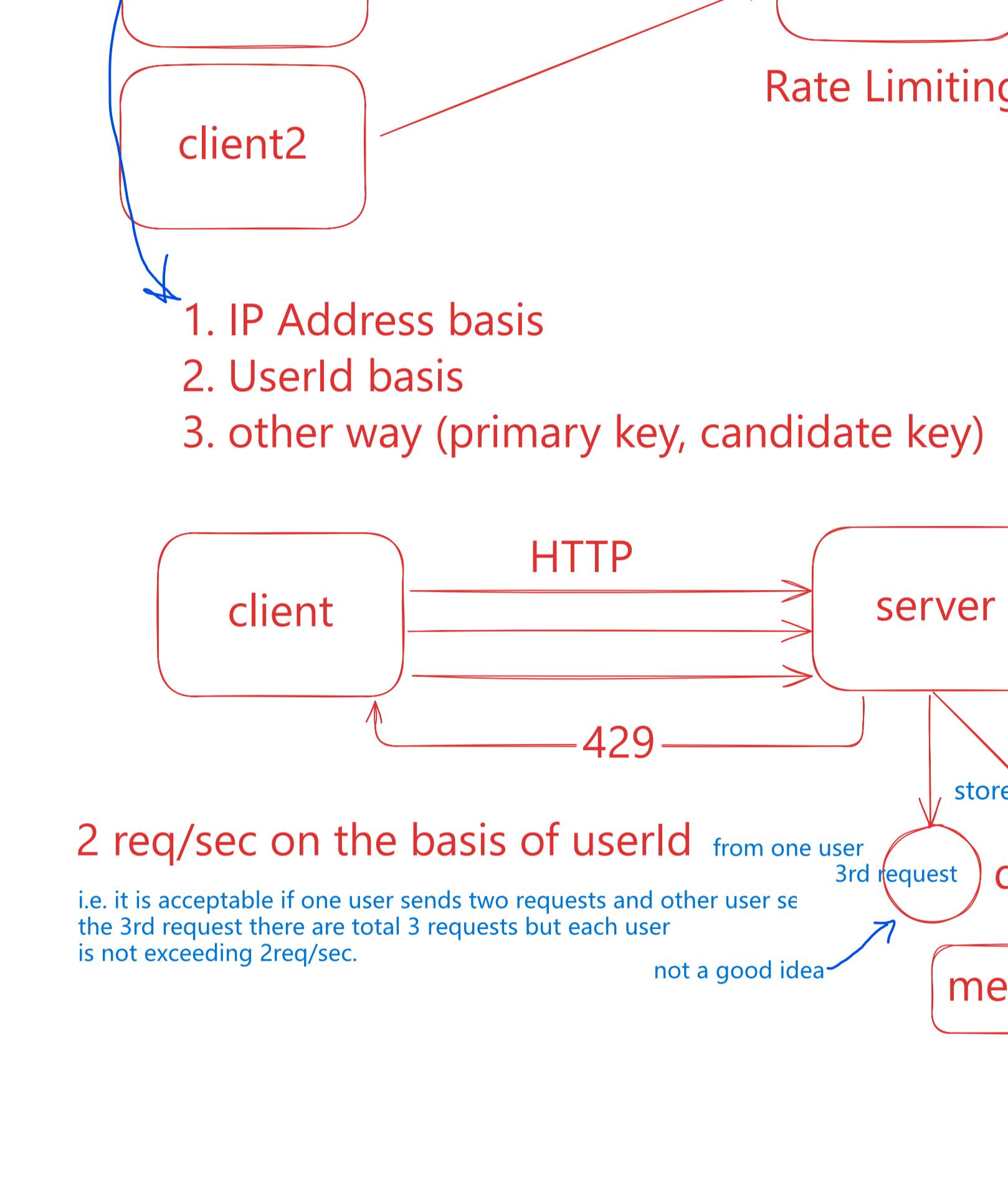
How many req client can send in any timeframe



Rate Limiting

Rate limiting can be implemented

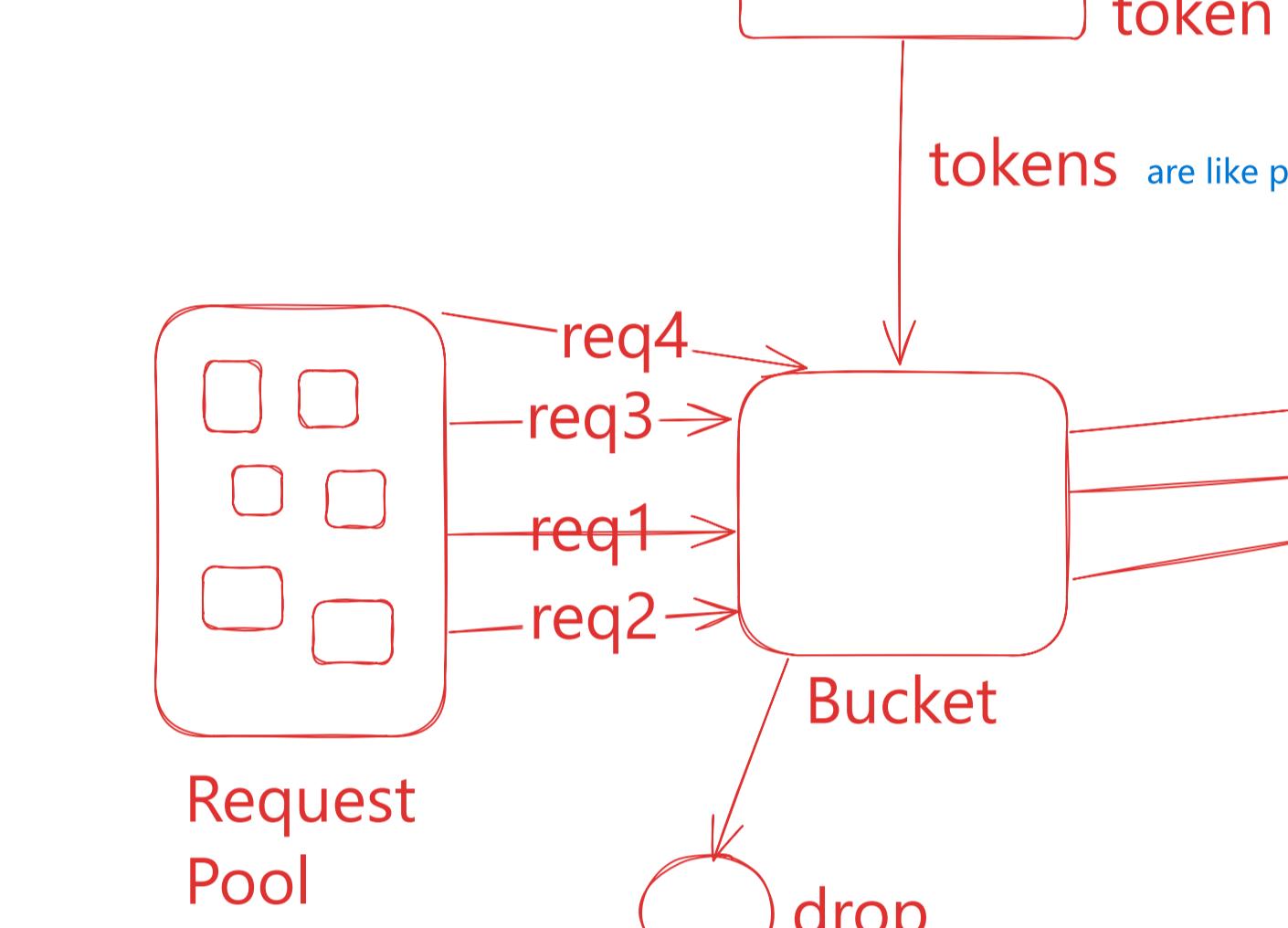
1. Client side : worst client is publicly available
2. Server side
3. Middleware



Application use case:

1. Server side language(Java, C++, javascript, Python) should be fast

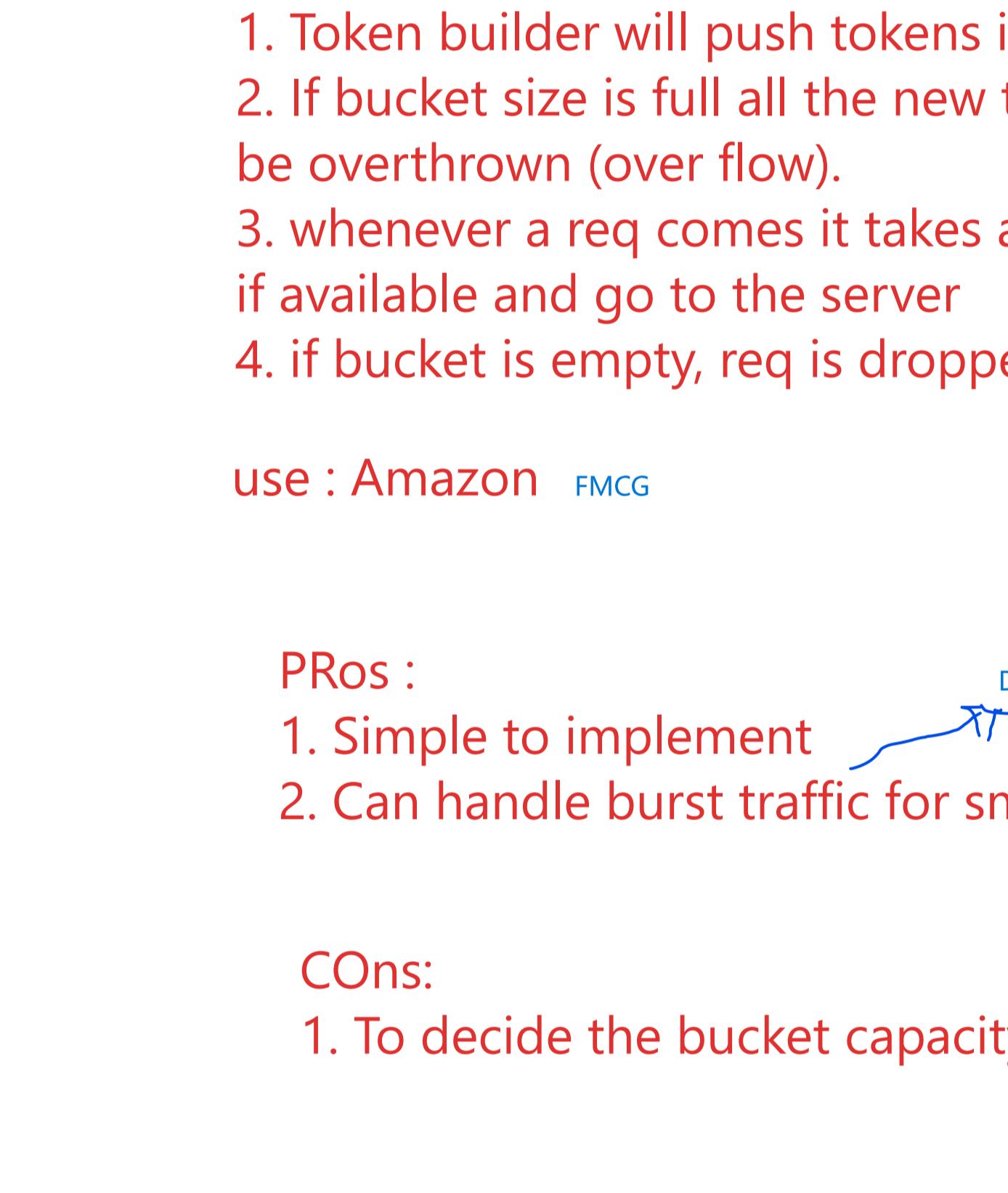
Server side



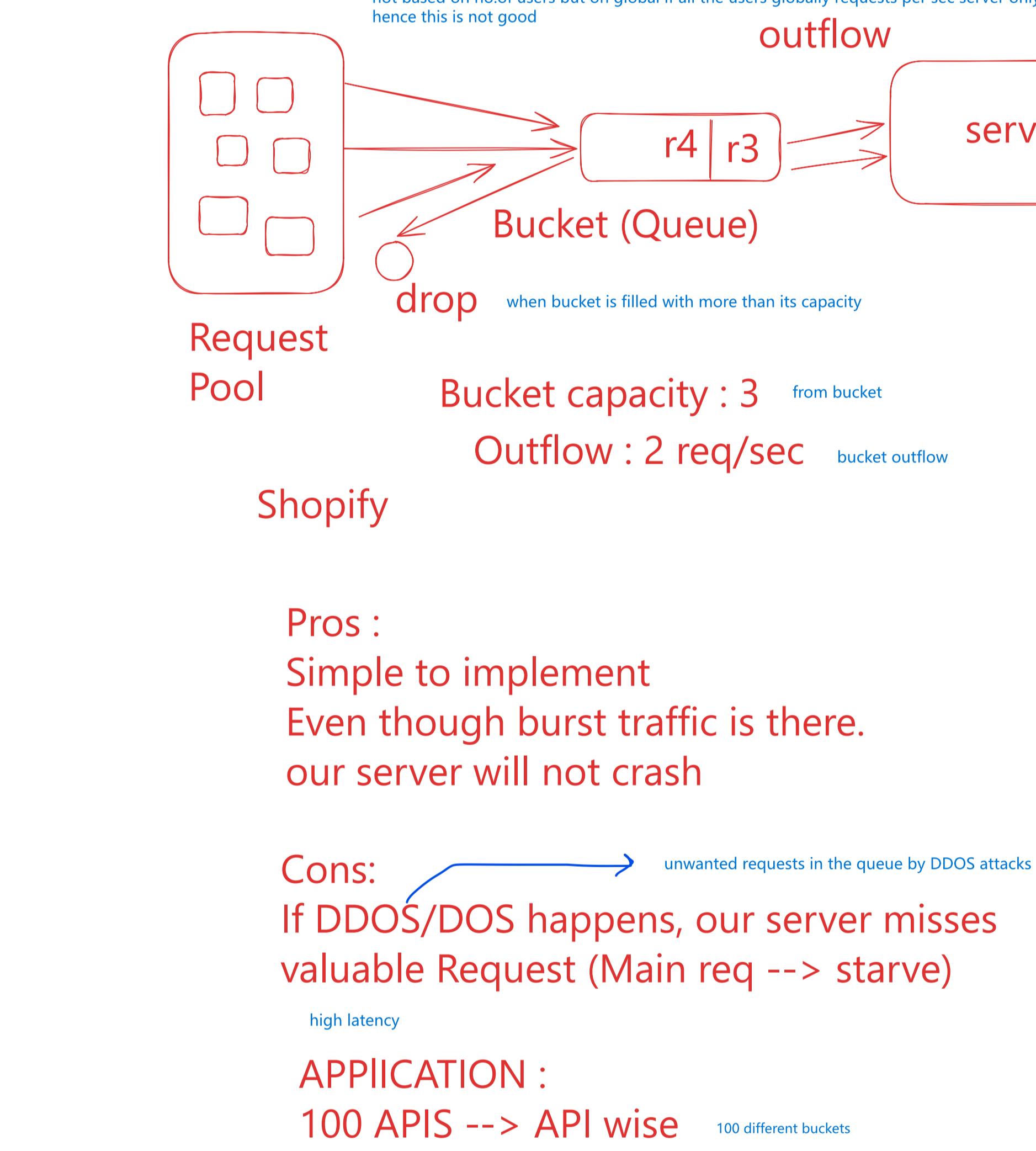
Rate Limiting

2. Generally 3rd party Rate limiter providers.

Amazon aws. (API Gateway)



Rate limiting on What ??



1. IP Address basis

2. Userid basis

3. other way (primary key, candidate key)

Response codes:

200 : Ok

429 : Too many Requests

Headers:

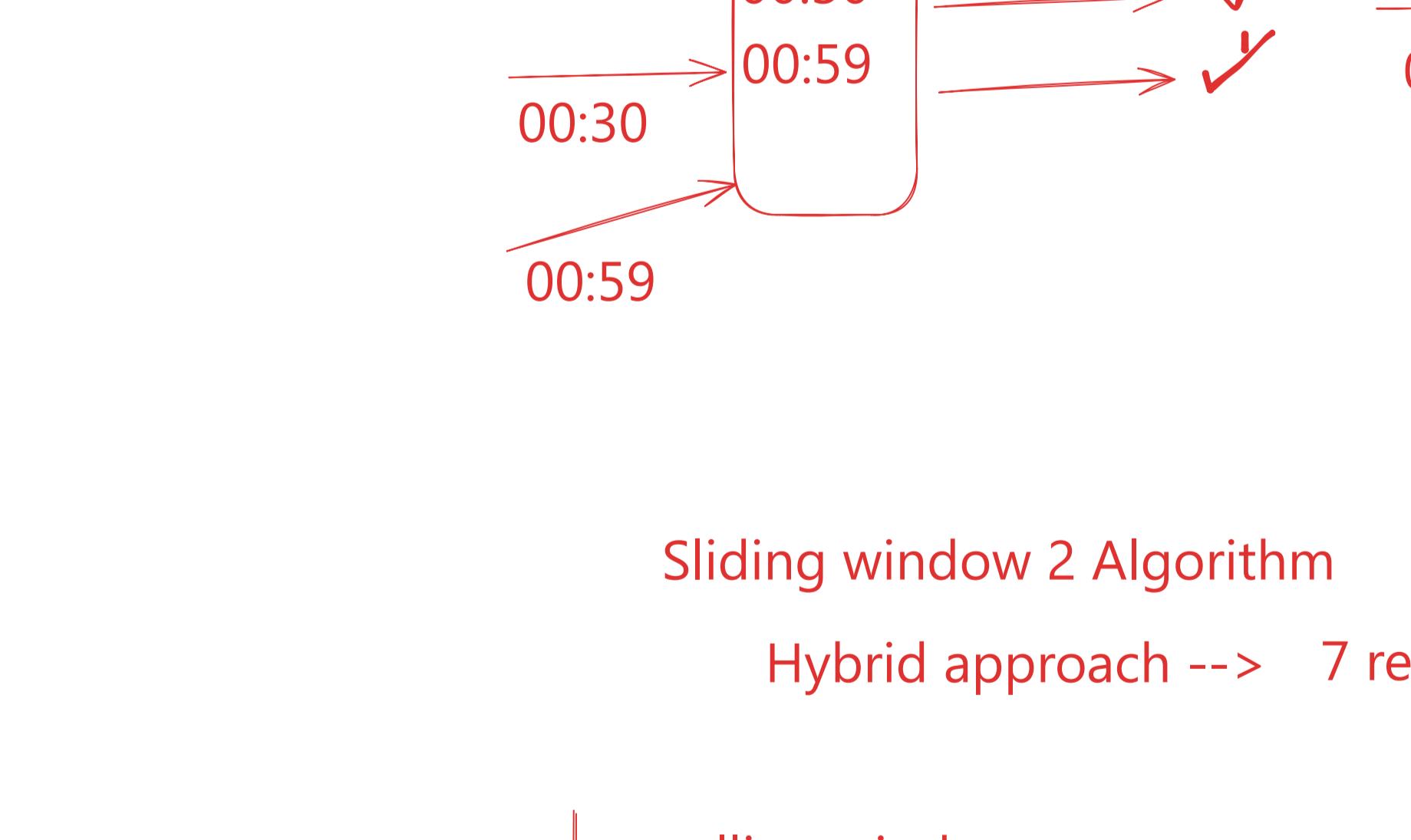
X-api-limit : 2

x-duration : seconds

Algorithms of Rate Limiting

1. Token Bucket
2. Leaky Bucket
3. Fixed window counter
4. Sliding window log
5. sliding window counter

Token Bucket Algorithm



1. Bucket size : 20
2. Inflow : 10 tokens/sec

1. Token builder will push tokens in fixed intervals
2. If bucket size is full all the new tokens will be overwritten (over flow).

3. whenever a req comes it takes a token from bucket if available and go to the server

4. if bucket is empty, req is dropped.

use : Amazon FMCG

PRos :

1. Simple to implement

2. Can handle burst traffic for small duration

COns:

1. To decide the bucket capacity and inflow

Leaky Bucket (Global Rate limiting) same for token based bucket its token capacity is high hence it can handle burst

when you want your application to not handle burst traffic then you can use leaky not based on no.of users but on global if all the users globally request per sec server only handles 3 buckets per sec hence this is not good

outflow

Bucket capacity : 3 from bucket

Outflow : 2 req/sec bucket outflow

App : 1000 users i.e. you want server to serve the requests at its own server outflow cap

IP level : every unique limit if based on userids then you have to have bucket for each user which can be a lot of buckets (same for token based buckets)

Shopify

Pros :

Simple to implement

Even though burst traffic is there, our server will not crash

Cons: unwanted requests in the queue by DDOS attacks

If DDOS/DOS happens, our server misses valuable Request (Main req --> starve)

high latency

APPLICATION :

100 APIs --> API wise 100 different buckets

FIXED WINDOW COUNTER

1. In a fixed interval (Say a sec, or a minute) Only a specific no. of req can come.

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

Sliding window Log Algorithm

Best Algorithm for Rate Limiting

strict

slow & memory consuming

Algorithm:

1. It will store each req in a log file.
2. Whenever a req comes, it will first remove all the outdated req from the file.

3. Then it will log the new req and check, if the counter limit reached, drop the req else

send it to server.

3 req/sec

Sliding window 2 Algorithm

Hybrid approach --> 7 req/ minute

window size and window capacity depends on the server

Formula : No. of req in curr interval

+ No. of req in prev interval

* 70% assumption as per the example

$$= 4 + 5 * 0.7 = 7.5 > 7.5 \text{ (hybrid approach) then drops the req}$$

Create our own Rate Limiter

Which algo to implement ?

--> Fixed window counter

Excepted Counter = 3

current counter = 3

Db ? SQL / NoSQL

Db : It is slow.

Cache : In-Memory Cache : Redis

DB rules for rate limit

fetch rules from DB from time to time and send to cache i.e. workers refreshes the cache

workers

Redis helps with the count value

high priority queues

Messaging Queue --> Kafka

429 Too many Req