

Agenda

What & why of Rate Limiting

Proxy

Client side

Server side

Design a distributed Rate Limiter

Various Rate Limiting algorithms

Rate Limiter

limit the rate of requests / usage of a service / # of calls per hour
Safeguarding / DDoS attacks / QoS

We want to limit access to some resource

Why? API / Software service

① DDoS prevention (Distributed Denial of Service)

900 Million requests / sec

② fair usage (FUP limit)

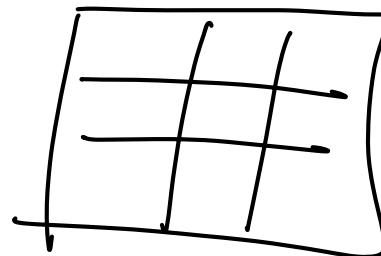
↳ fair usage policy.

We have some public resources - prevent some bad users from 'hogging' up all the resources.

③ Prevent Scraping. — Using a program to download & process data from a website.

Google

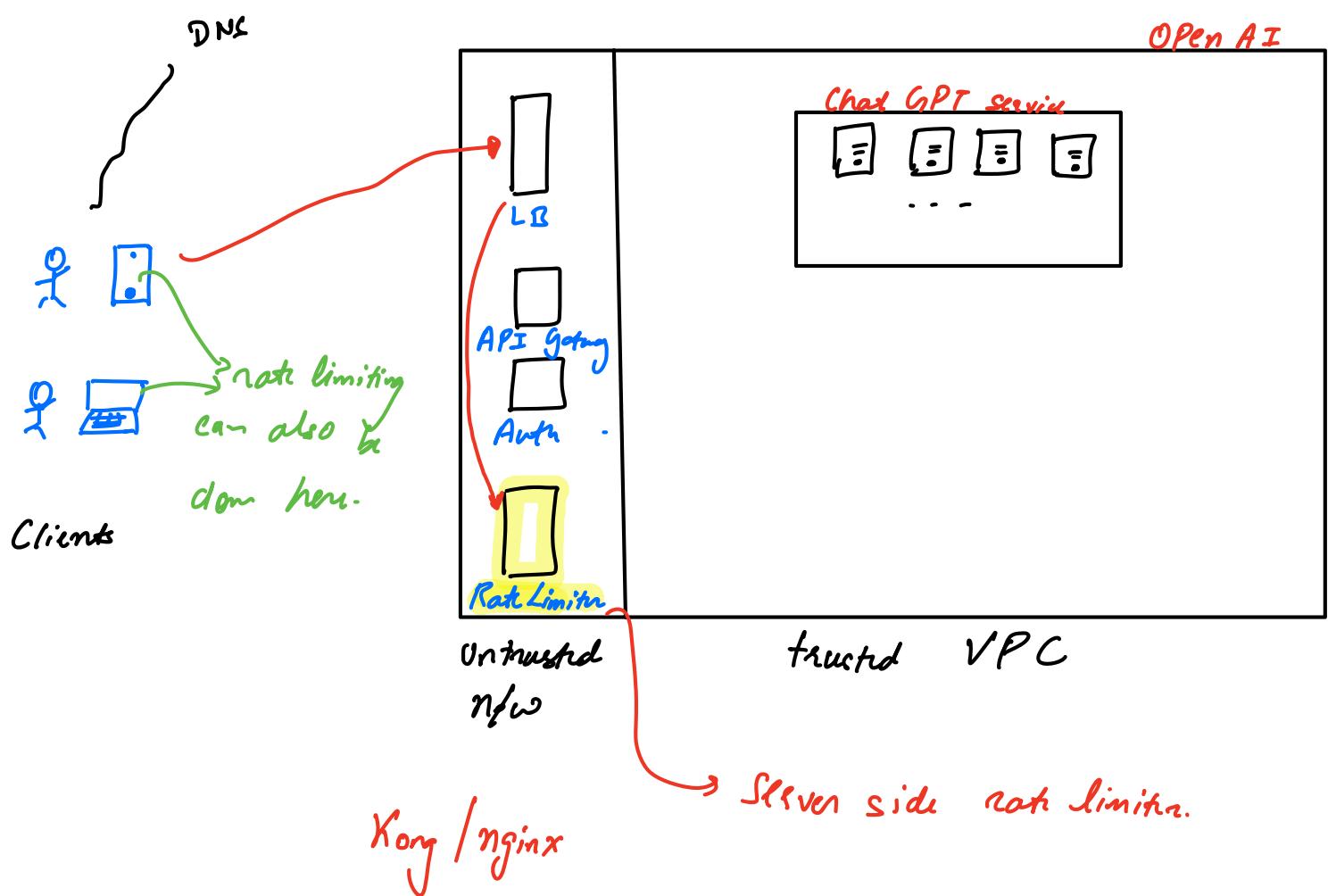
I'm not a robot



Bing

④ Regulatory Compliance — limit the no. of stock trades/second on some instruments.

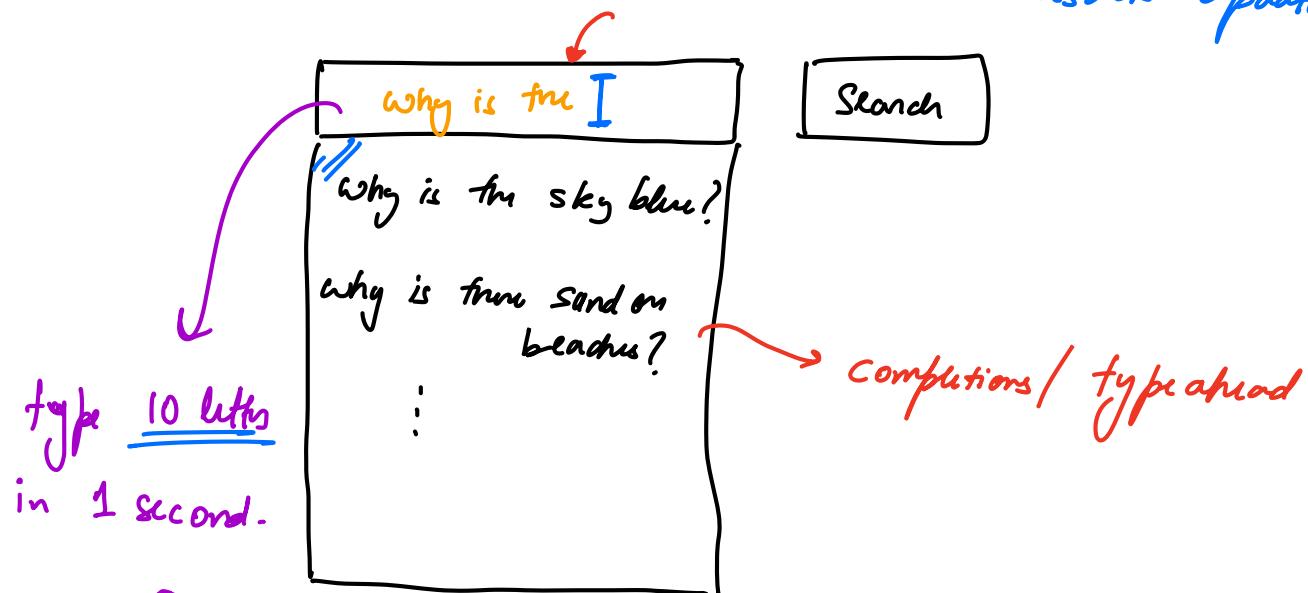
System Architecture



Client Side Rate Limiting

① Debouncing

→ every time you type a letter the auto complete results update.



```
#("input").keydown (e) => {  
    if (e.keyCode != ENTER) {  
        // make an API request  
        // Populate auto complete suggestions  
    }  
}
```

Only send the request once user stops typing!

① sending requests on user input

② auto saving

② Throttling → manually add a delay b/w any two requests



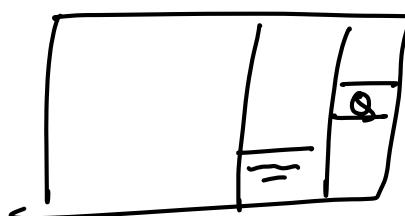
Exponential Backoff / Queuing / Caching / Batching

Purpose of all client side rate limiting → to reduce the no of API calls that are generated.

Issues with purely client Side rate limits

① It can always be bypassed.

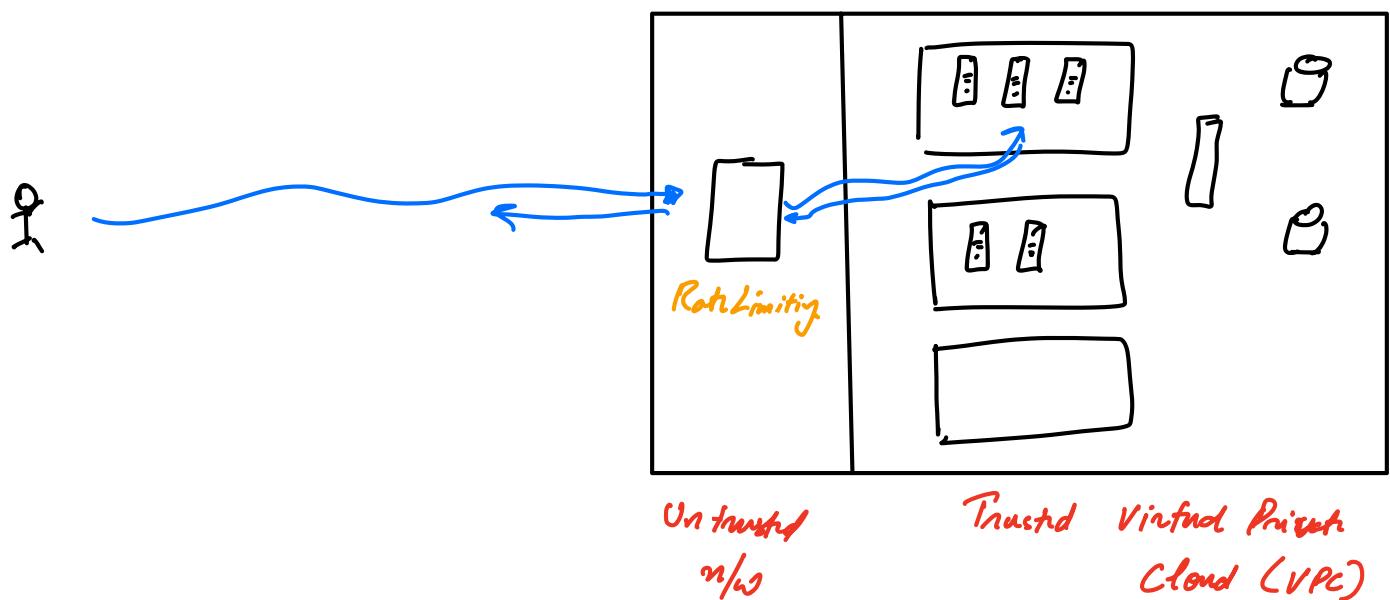
② If your clients have direct API access → no frontend.



Zerodha → stock broker

Rite API → to build your own stock trading platform.
→ to build your custom trading algorithm.

Always need Stream Side Rate Limiting as well!



If rate limits decides to not accept → Queue it for later drop/deny request.

Work flow

- ① Client request arrives at Load Balancer.
- ② Load Balancer checks with Rate Limiter about whether it should accept/reject it.
IP/van-id/geographic location/...
- ③ Rate Limiter will retrieve information from some 'data store'
to store ↗
how many req have we seen in past
- ④ Decision Making
 - Accept
 - Reject
 - Queued
 - Denied → Proper HTTP error code
↳ 429 : Too many Requests

⑤ Request Processing → let the request go through.

⑥ Logging & Monitoring.

↳ inspection of whether the request is accepted or denied

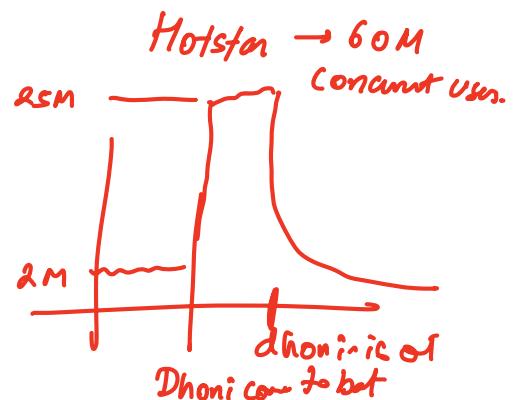
we should log it!!

↳ Metrics Collection : % of acc/rq / Top 10 users with most rq --

↳ Monitoring Tools

↳ Alerting

↳ Capacity Planning.



Prometheus / Grafana / Load Balancing / API Gateway / DNS. --

High Level Design / System Design → how to scale a system.

Google → Staff Engineer (Bengaluru 1.5 - 2 Cr / annum)

• simple given a list of strings → sort them alphabetically.

• difficult twist: 50 Petabytes of data

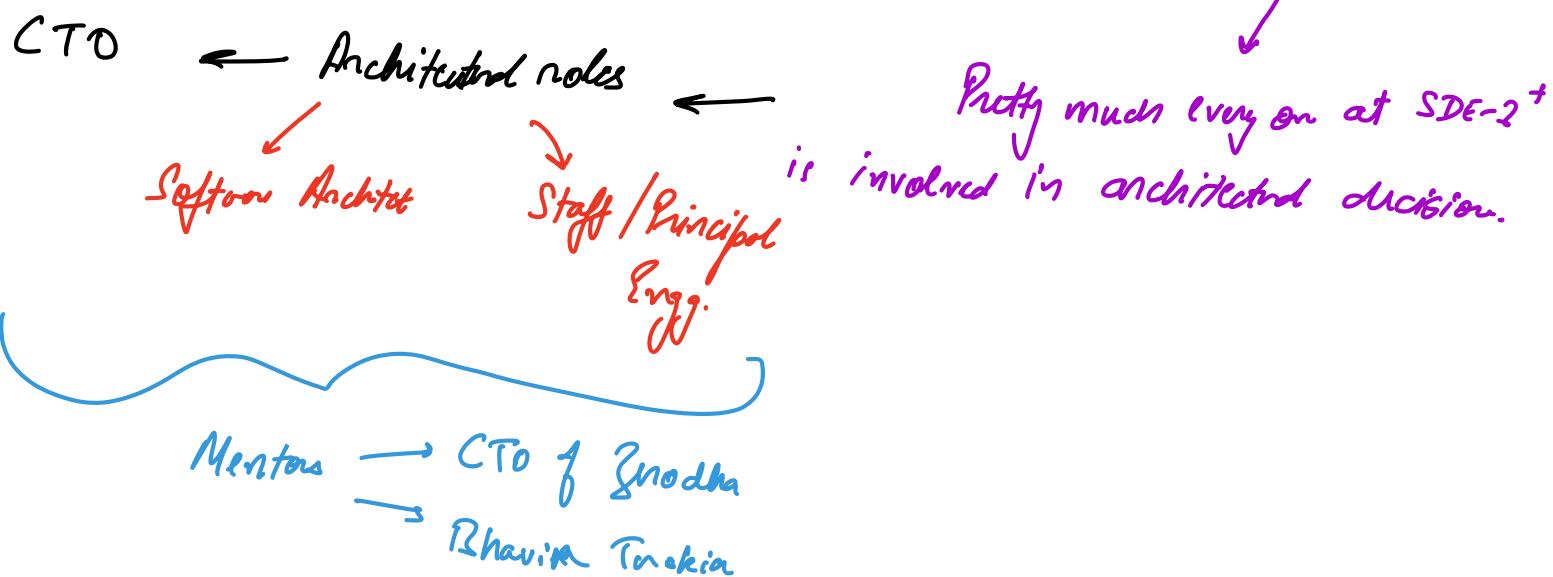
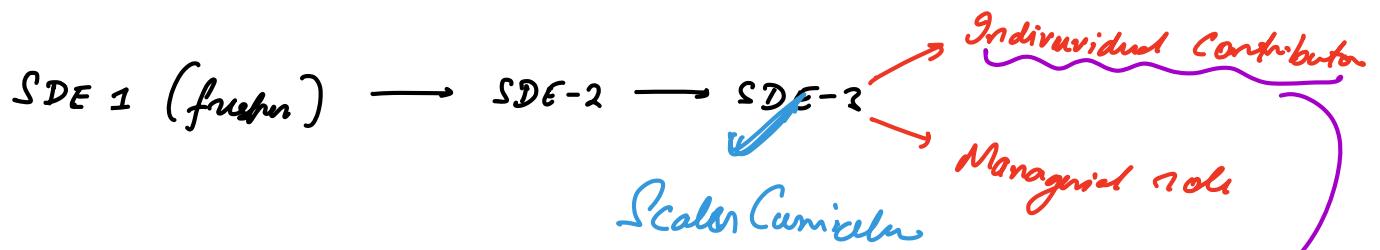
KB MB GB
TB PB

50,000,000 GB of data

① Scale → 1.5 months on High Level Design

→ 2-5 months on Low Level Design

9:02 → 9:15



Various Rate Limiting Algorithms

① Sliding Window Counter

limit to 1000 requests / hour
Max count

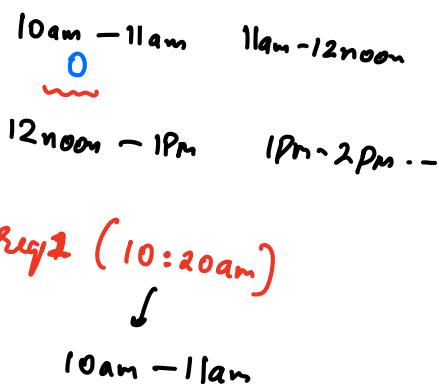
Windows

Product team
Infra team
Stakeholders

- for every window of T time we will store an intg count

- when a new req comes

↳ find the window that it falls under.
↳ retrieve fm count for window
↳ if count is within the limit (1000)
 ↳ accept req
 ↳ Increment count
↳ else
 ↳ reject request.



Data Structure

- On what 'entity' am I putting the rate limit on?

↳ user-id → simple limiting
↳ access-token/API Key → exposed APIs
↳ IP address → anonymous users
↳ geographic → DDoS prevention.

Key

Attempt 1 : Start a count for each key
 ↳ # of key for this key for this window.

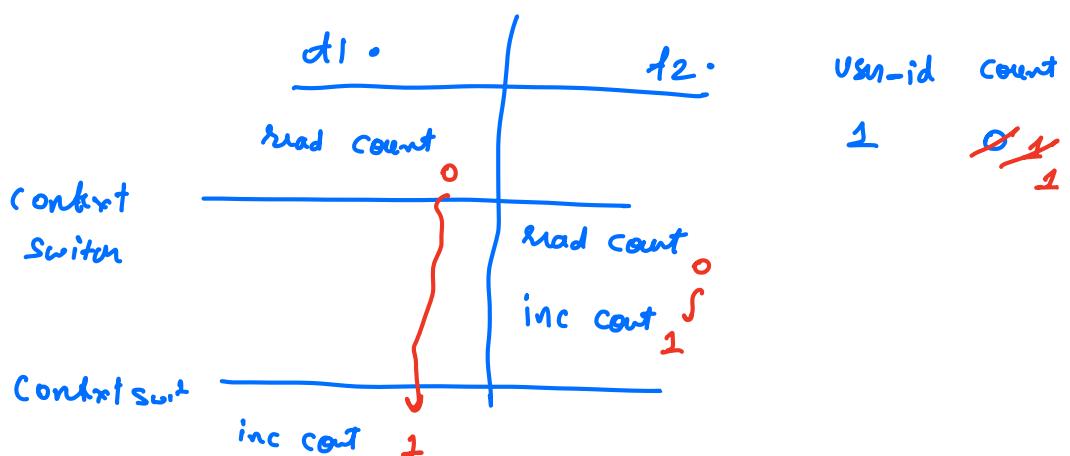
```

Map < Key, Integer > counts;
Map < Key, Integer > limits;
Void process (Request r) {
    Key k = r.key;
    int count = counts.get(k);
    if (count < limit) {
        // accept
        counts.set(k, count + 1);
    } else {
        // reject
    }
}
    
```

when will our count true?
 (in RAM of Rate Limiter system)
 (what if we have a dist system?)

global variable (same limit for everyone)
 loaded from some config (custom limits for each key)

'fetch then set' pattern



Topics for later

- ① Data Store
- ② Config Store
- ③ Mutual Exclusion

How to reset the count at start of next window?

- ① every 1 hour, go through all counts & set them to 0.

Eager

= `for (Key k : counts)`
= `counts.set(k, 0);`



2 Billion users.

0.1 ms / Key

$$2 \times 10^9 \times 0.1 \times 10^{-3} \text{ seconds}$$

$$= 2 \times 10^5 \text{ seconds}$$

≈ 60 hours

- ② Map < Key, Pair < Timestamp, Count >>

{
 userid/
 IP address/
 access token -
 start of current window
 count of current window.

Lazy

when req comes & check cur time.
if cur time is in same window
then just add
else & update the window & rest count.

Properties / Behavior of Fixed Window Counter.

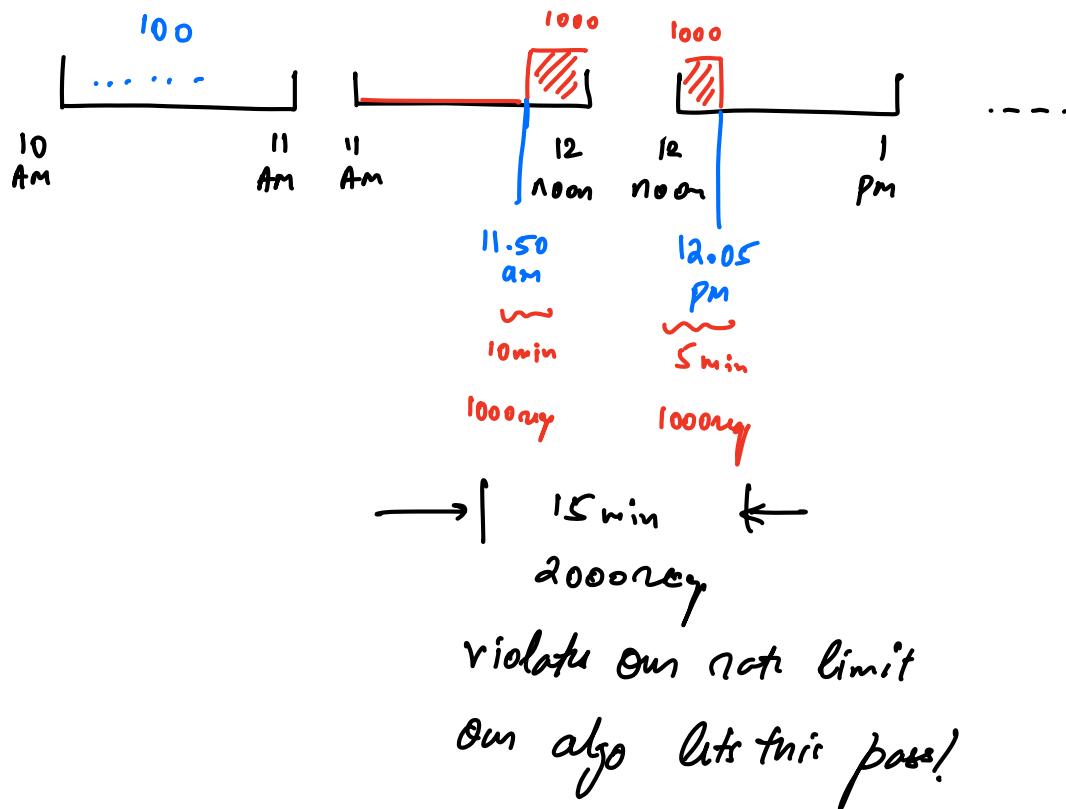
- ① v. simple to implement
- ② allows short bursts
- ③ successfully prevents long bursts.
- ④ allows high throughput across windows

Burstiness
if suddenly a lot of
req come in a short
period (within rate
limit)

long burst } continued
 | high req
 | above the rate limit.

short burst }
 | 1000 req / hour.
 | in a 5 sec window user
 | makes 999 req.

Boundary Case $\rightarrow 1000 \text{ req/hour}$



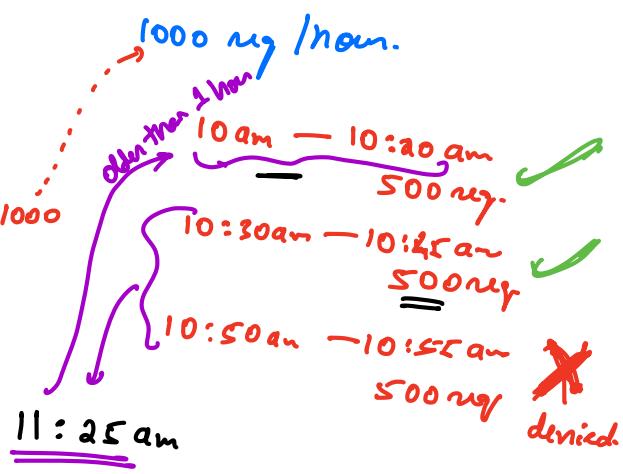
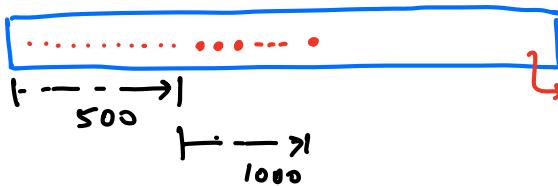
Sliding Window Log

- Store / Log every req that comes
- Aim : "Strictly" enforce the rate limit

not a msg Queue \rightarrow it is a log queue

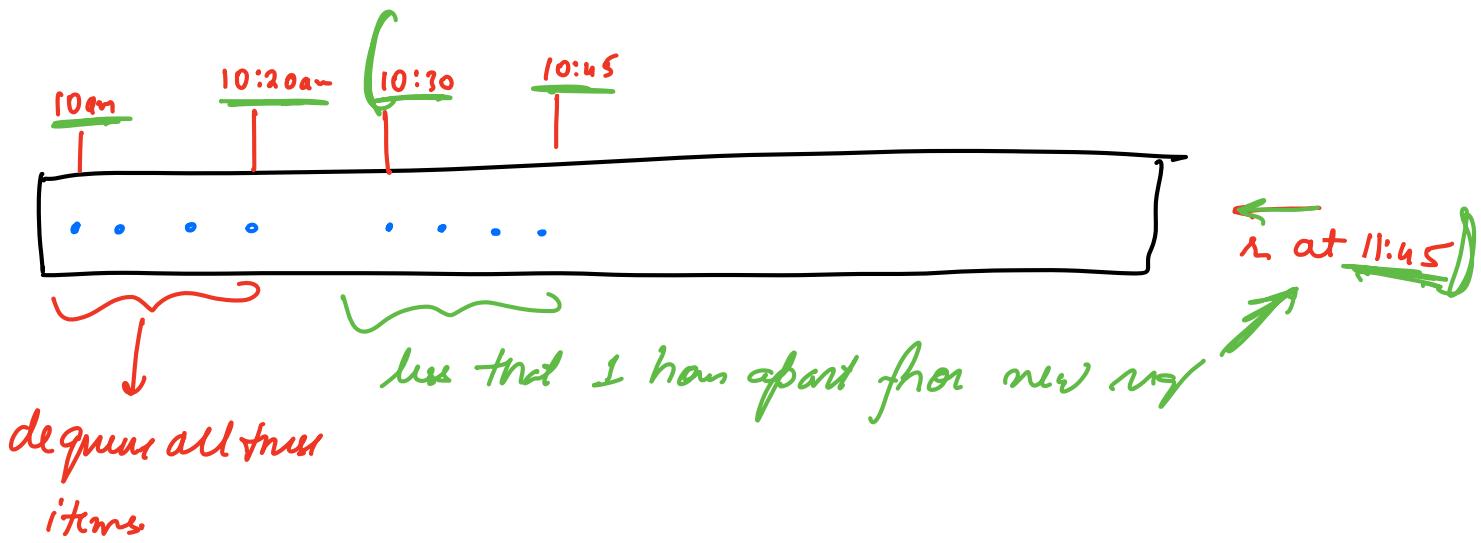
Map < Key, Queue < Timestamp > >

Q1's
Queue



① when a new avg comes.

- remove all items from queue that are outside curr window.
- check how many items are in log Queue
 \because all these are in the same hour.
- if # items \leq rate limit
 - ↳ accept & add to Queue
 - ↳ use
 - ↳ reject

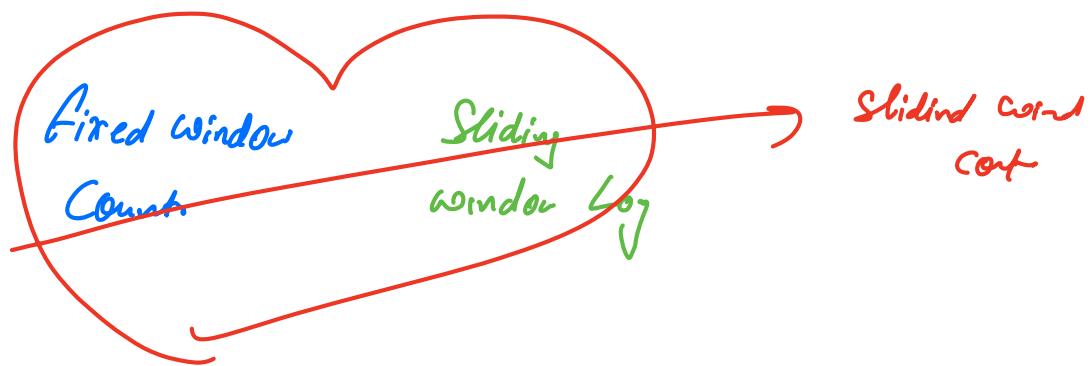


Pros ① absolute correctness

Cons ① takes more memory \because Queue

 ② for some avg the latency is v. high

Sliding Window Counter → Approximate algo



Token Bucket

Leaky Bucket

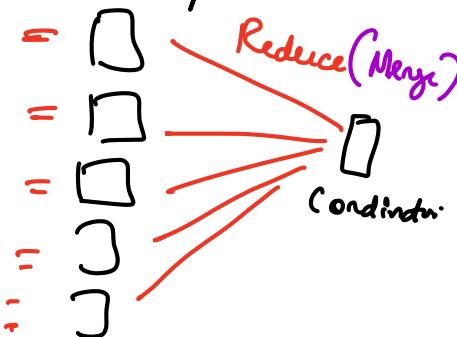
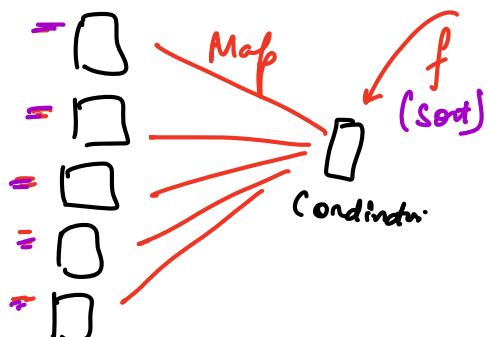
-
- ① How is data stored → Replication
Resiliency
 - ② How is data being read → CAP theorem
Read Replicas
Caching
 - ③ How do we sort it?

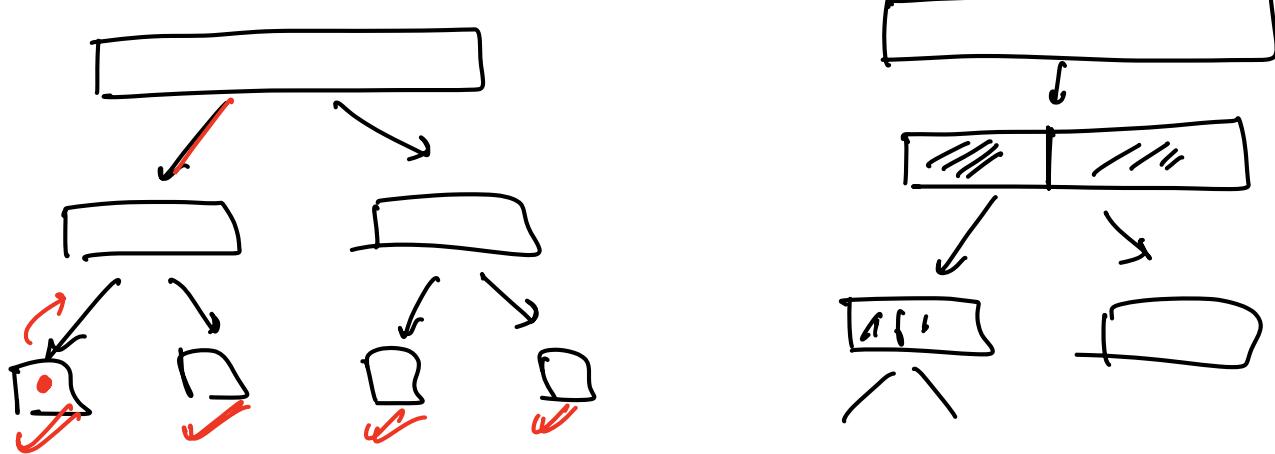
↳ usecase: archiving → ---

↳ usecase: analytics → Map - Reduce

Apache Hadoop

Apache Spark





AnyCast

