

Api Rate Limiter

Monday, March 25, 2024 8:53 PM

FUNCTIONAL REQUIREMENT

1. limits no. of requests for an entity and blocks requests after threshold
2. provide API rate limiting across diff servers

NON FUNCTIONAL REQUIREMENTS

1. Availability over consistency (one or more requests going forward is okay)
2. Reliability
3. Speed (low latency)

TYPES OF THROTTLING

Hard throttling: strictly limits

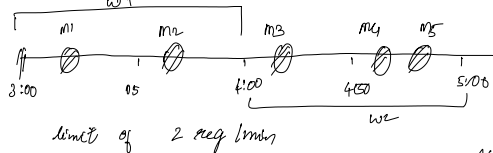
(3) min

Soft throttling: 500 reqs/hr + 10% exceed limit
550

Elastic throttling: allow additional requests if system has resources

ALGORITHMS:

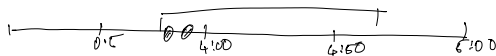
1) FIXED WINDOW ALGORITHM



W1 - m1, m2 are served.
W2 - m3, m4 are served
m5 is rejected.

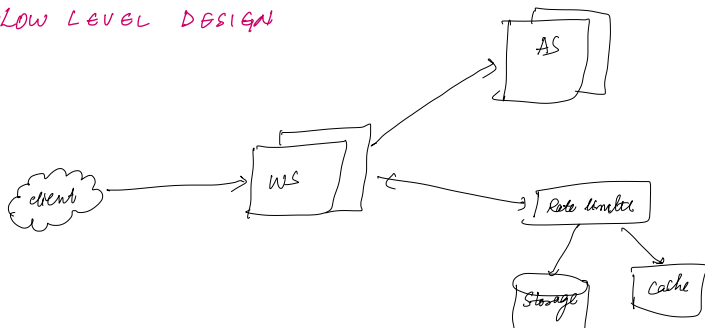
If 2 reqs are from second half of window and the other from 1st half -

2) SLIDING WINDOW ALGORITHM

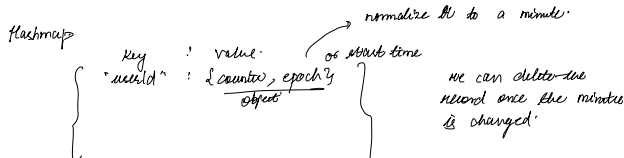


window slides every time a request comes in

LOW LEVEL DESIGN



1) FIXED WINDOW



500 reqs/hr
floor: sliding window

Ex: 10 req/hr
3 req/min

Window: 8 bytes
Epoch: 4 bytes
Counter: 2 bytes
500 reqs/hr
20 bytes (redis cache)

8 + (4+2+20) x 60 minutes

8 + 1560 + 20 = 1588

for middleware used 15

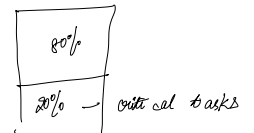
SHARDING AND PARTITIONING

19 x nodes
company / safe

MAC address - unique

Cache Eviction: LRU

Types of rate limiters:
1. Request rate limit
2. Concurrent rate limit
3. Token bucket



4. Work utilization

- 1) Critical
- 2) POST
- 3) GET
- 4) Testing

We'd love your feedback!

We have just two questions for you.

example: $\{ \text{user } 14308 : \{ 1, 14274008 \} \}$

assume: 3 req/min.

memory requirements:

userid: 8 bytes
count: 2 bytes
epoch: 4 bytes (32-bit number)

Total: $(8+2+4) + 18 \text{ (bytes overhead of hashmap)} \Rightarrow 32 \text{ bytes}$

1 million users active.

$32 \text{ bytes} \times 1 \times 10^6 \Rightarrow 32 \text{ MB}$

SLIDING WINDOW

HashMap
key: or OR
value: sorted set < >

3 req/min are allowed.

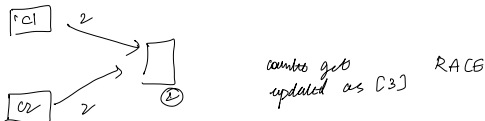
user 14308: 14308742, 14308744, 14308748, 14308750
14308774
user 14308: [14308742, 14308744, 14308748, 14308750]
14308774
if > above 1 min.

userid: 8 bytes
epoch: 4 bytes
500 req/min
20 bytes overhead sorted set

$8 + (4+20) \times 500 + 20 = 8 + 12000 + 20$
 $= 12028 \text{ Bytes}$
↑
overhead of hashmap

$= 12 \text{ kb}$

1 million users $\rightarrow 12 \text{ kb} \times 1 \times 10^6 \Rightarrow 12 \text{ GB}$



LOCKING MECHANISM

IN memory, cache - redis

mutex lock

Trade off for both is hybrid

Sliding window with the help of counters

