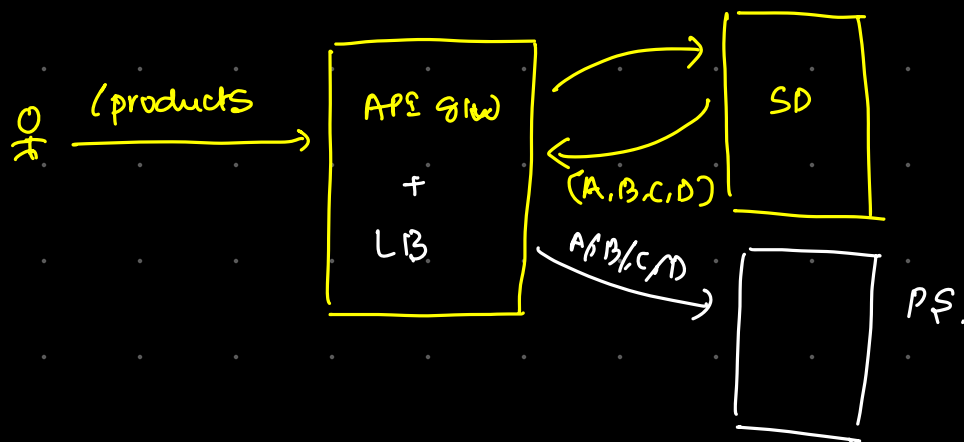


Two options.

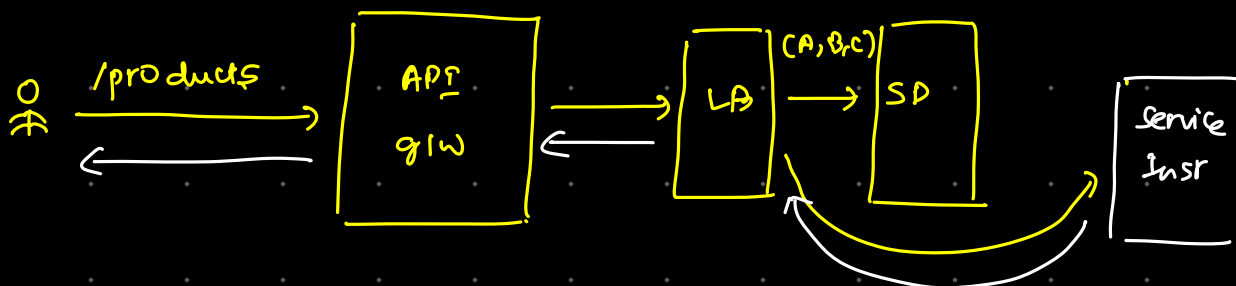
(Client side LB)

- API gw. → service discovery → LB @ API gw → service Inst
[A, B, C, D, E] [C]



Server side LB.

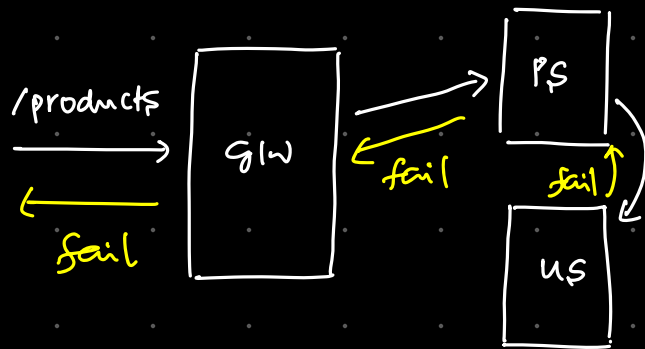
- API gw → LB of service → service discovery → service Inst
[A, B, C, D, E]



Qm.

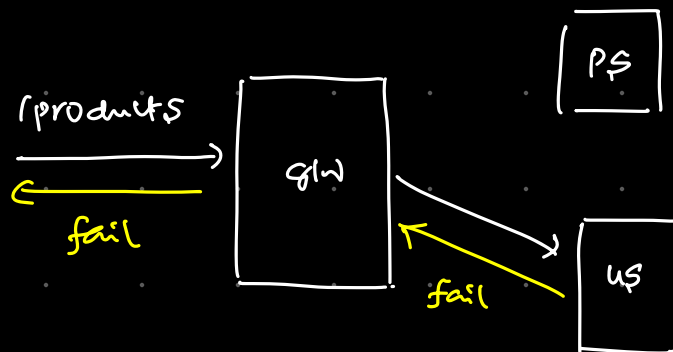
1. Can API gw take care of authentication.

①. API gw not doing auth.



auth after API gw.

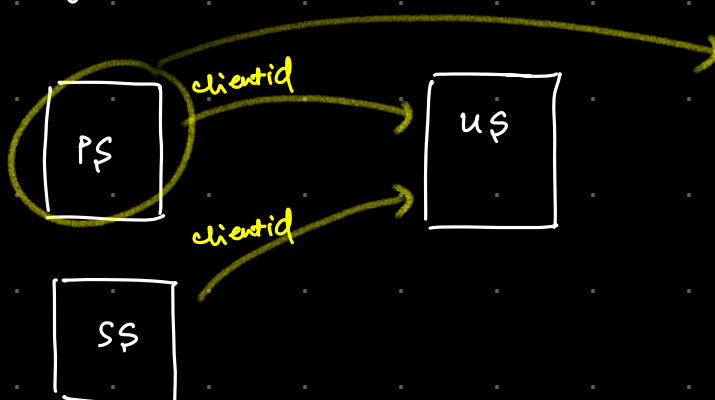
②. API gw not doing auth



We're filtering the requests
@ gw itself.

Throttling

limiting the no. of requests.



> 100 req/second.
Throttling.

API gw can also be used for throttling.

APE gw in action

Monitoring

100 instances in each service

around 20-30 services.

You would want to monitor:-

1. CPU utilization
2. Disk utilization
3. RAM utilization
4. latency of API'S [How much time API takes to respond]
5. Requests per second / RPS / QPS.

Actuator

It exposes some endpoints where you see multiple details about resources and other things.

localhost:8080 / actuator

- ① / metrics : memory usage, CPU usage etc
- ② / health : health info about the appⁿ.

Prometheus

This will hit the metric api's periodically and capture the data in a time series format.

Grafana

Queries prometheus to receive metrics and visualizes it through customisable dashboards.

