

Zuul - Netflix API Gateway that Handles 2M Requests Per Second

Susheel Aroskar

Cloud Gateway - Netflix

NETFLIX

Amazon AWS S3 outage is breaking things for a lot of websites and apps

Darrell Etherington @etherington / Feb 28, 2017



Comment

OPINION

Why Netflix didn't sink when Amazon S3 went down

When a typo caused an outage in part of Amazon Simple Storage Service's (S3) cloud infrastructure on February 28, many of their customers' websites went down but Netflix survived unblemished.



How??

How??



Zuul shifted all the traffic from the affected us-east-1 region to other healthy AWS regions



<https://github.com/Netflix/zuul>



Susheel Aroskar

Senior Software Engineer
Cloud Gateway

saroskar@netflix.com



github.com/raksoras

[@susheelaroskar](https://twitter.com/susheelaroskar)

Agenda

Agenda

- What is Zuul?

Agenda

- What is Zuul?
- What Zuul can do for you?

Agenda

- What is Zuul?
- What Zuul can do for you?
- How Zuul does it?

Agenda

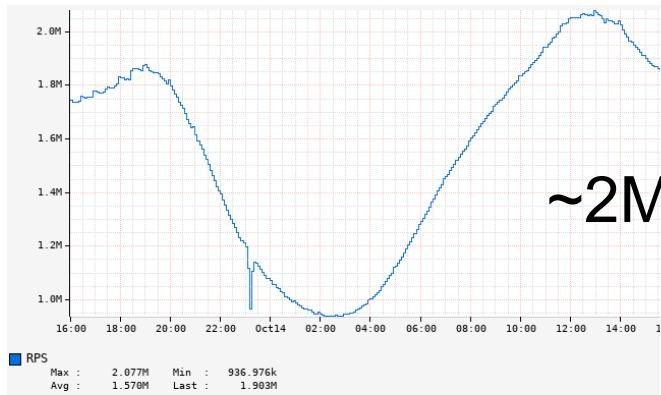
- What is Zuul?
- What Zuul can do for you?
- How Zuul does it?
- Future roadmap

What Is Zuul?

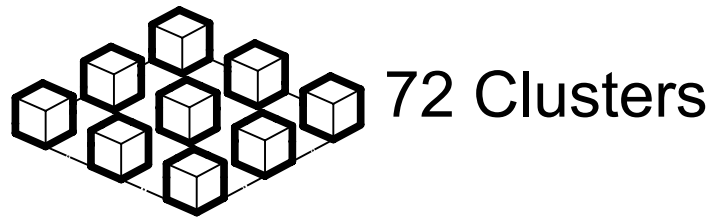
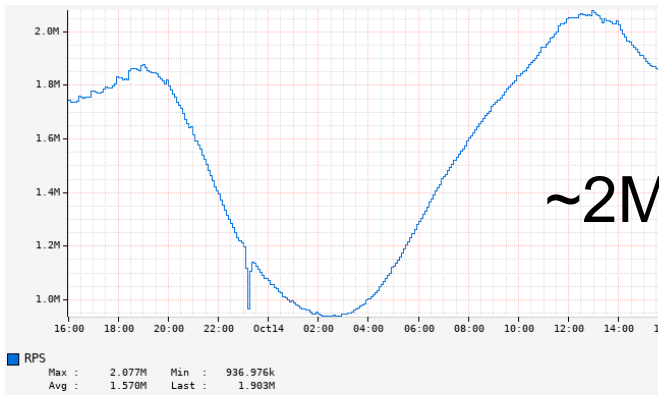
| Zuul Introduction



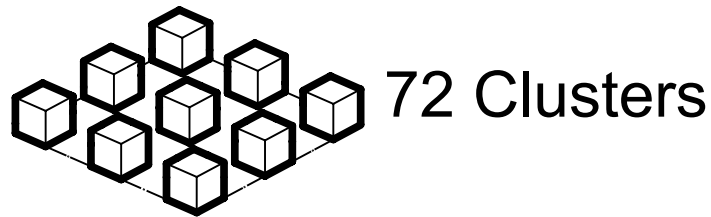
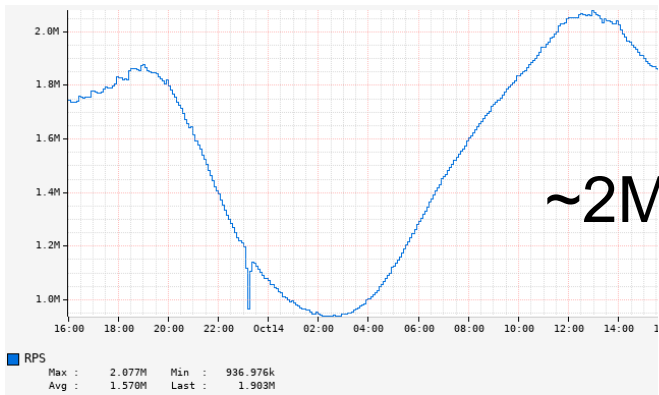
Zuul is Netflix's API Gateway



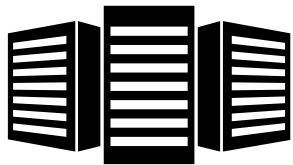
Zuul is Netflix's API Gateway



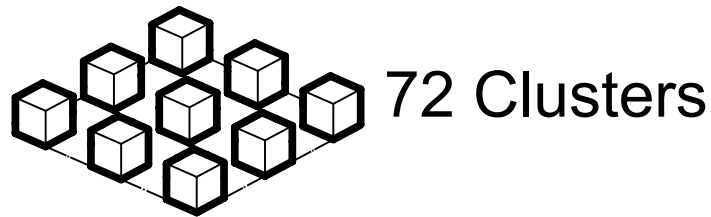
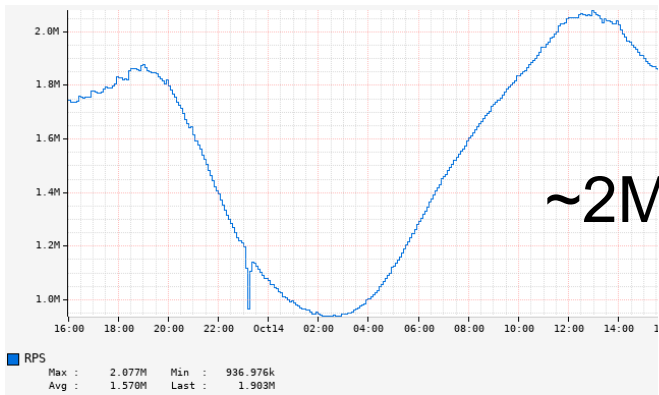
Zuul is Netflix's API Gateway



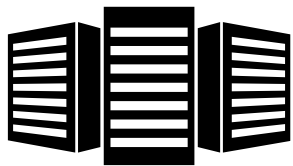
Zuul is Netflix's API Gateway



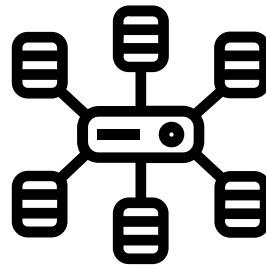
6000 servers



Zuul is Netflix's API Gateway



6000 servers



150 Backend
Services

What Can Zuul Do for you?

| Zuul Features



Core



Responsibilities



Route Traffic

Core



Responsibilities

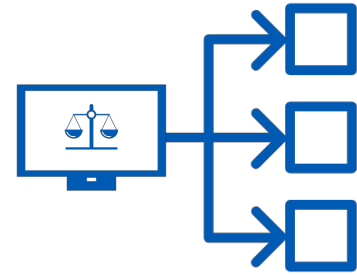


Route Traffic

Core



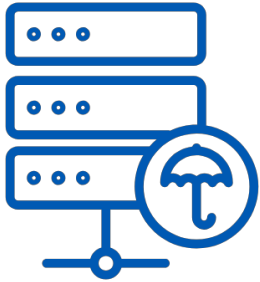
Responsibilities



Load Balance



Route Traffic

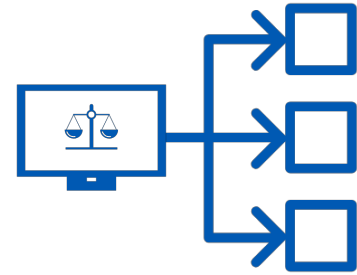


Protect Origins

Core



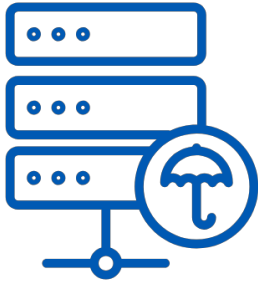
Responsibilities



Load Balance



Route Traffic

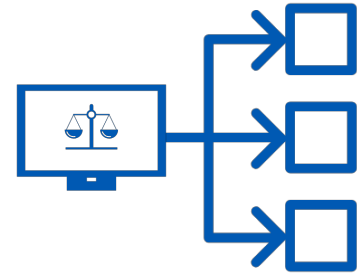


Protect Origins

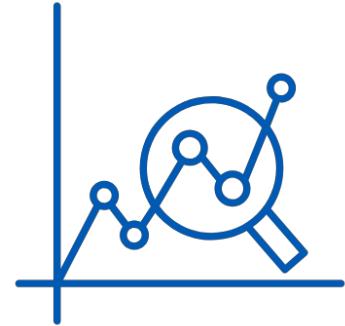
Core



Responsibilities



Load Balance



Metrics and Monitoring

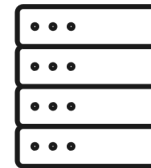
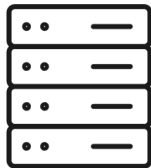
1. Route Traffic



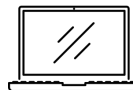
CORE
RESPONSIBILITIES

NETFLIX

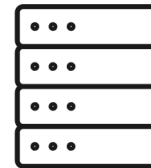
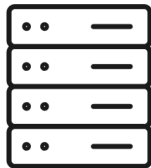
Login
Service



Playback
Service



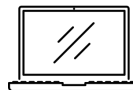
Login
Service



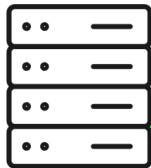
Playback
Service



Service
Registry
(Eureka)



Login
Service



Playback
Service



register

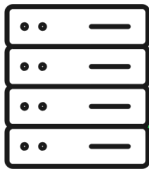
register



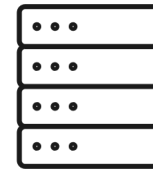
Service
Registry
(Eureka)



Login
Service



Playback
Service



register

register

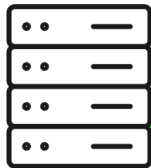
Service
Registry
(Eureka)



*Load
service
addresses*



Login
Service



Playback
Service



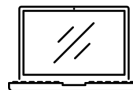
register

register

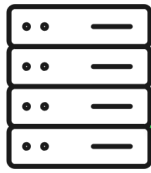
Service
Registry
(Eureka)



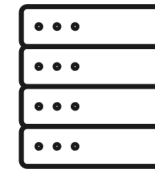
*Load
service
addresses*



Login
Service



Playback
Service



register

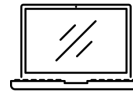
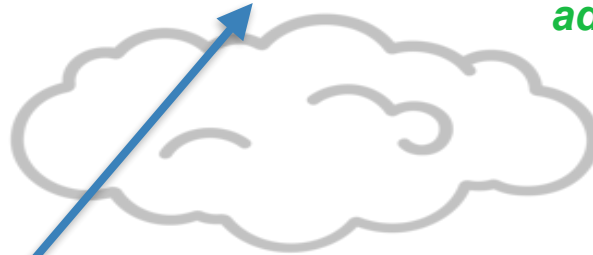
register

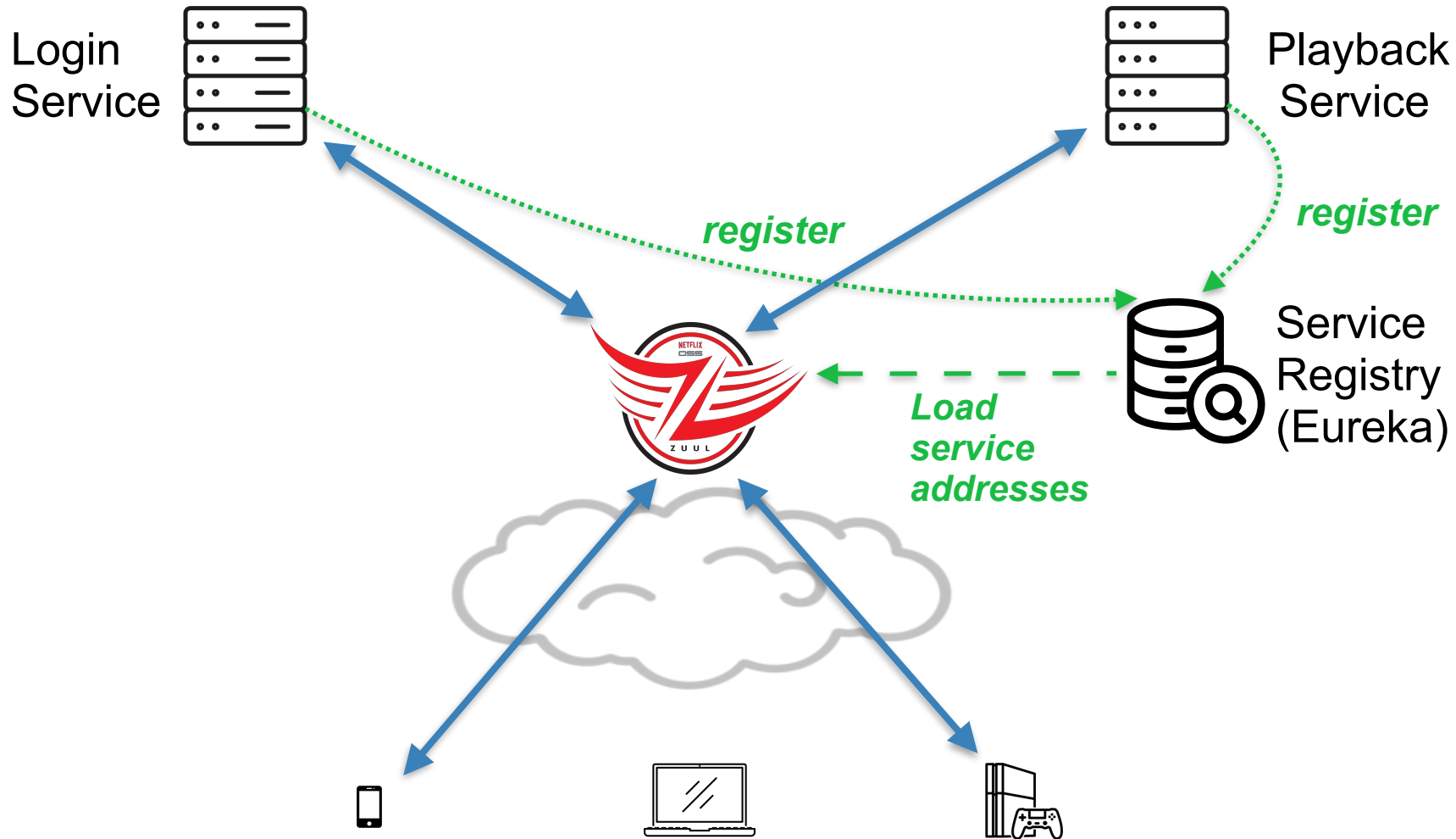
register

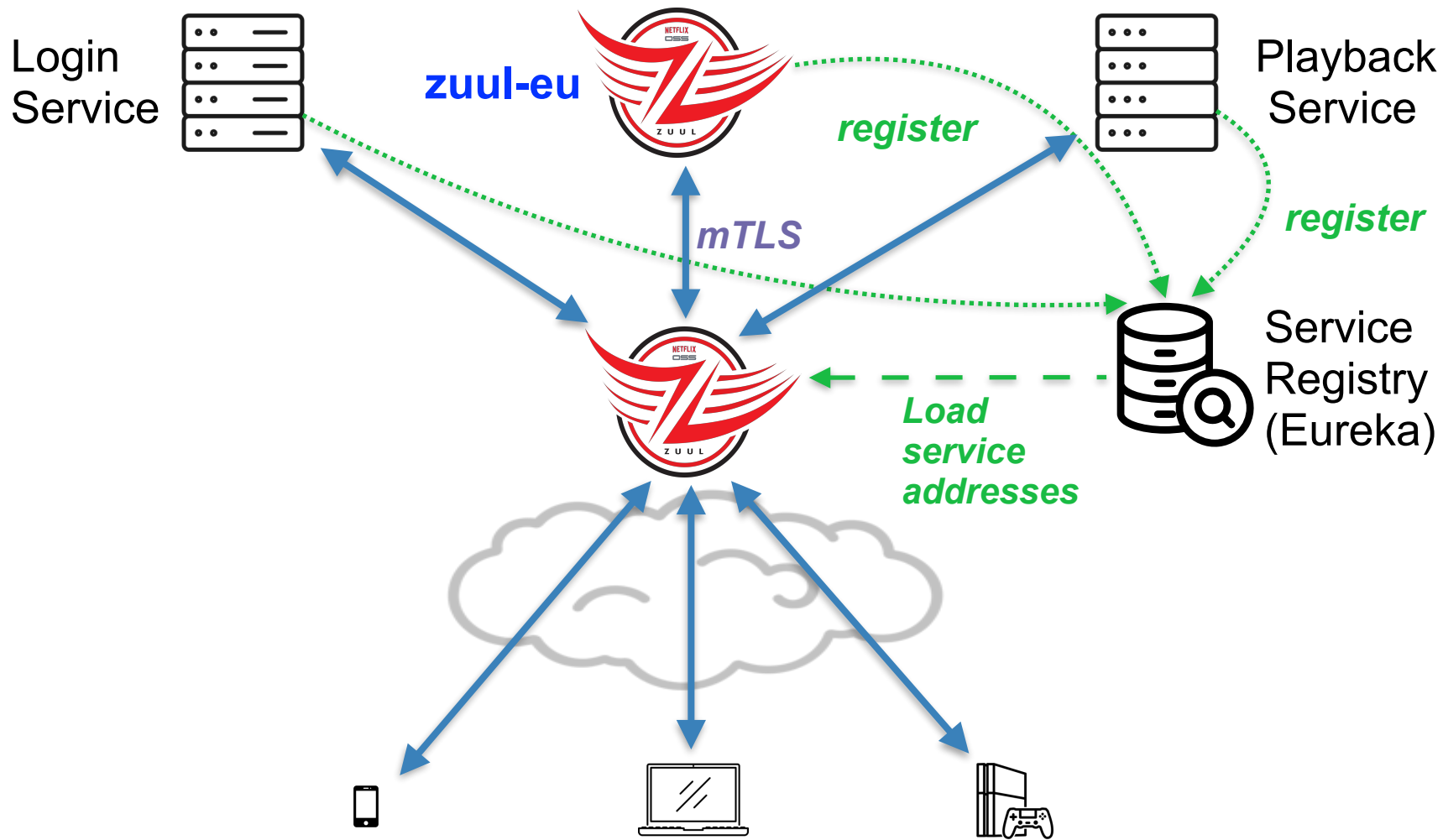
Service
Registry
(Eureka)



*Load
service
addresses*







You can change routes dynamically at runtime

You can change routes dynamically at runtime



Key	Operator	value
path	regexp	^(/api)?/(users logs)(/?.)*
region	==	all

Routing Action - Override VIP

vipOverride

Baseline App or VIP	Target
api-prod.netflix.net:7001	api-prod-rest.netflix.net:7001

You can change routes dynamically at runtime

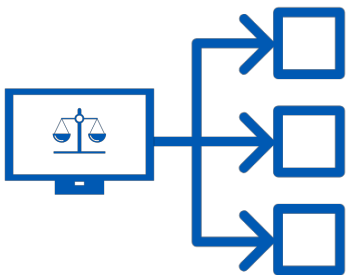
Key	Operator	value
path	regexp	^(/api)?/(users logs)(/?.)*
region	==	all

Routing Action - Override VIP

vipOverride

Baseline App or VIP	Target
api-prod.netflix.net:7001	api-prod-rest.netflix.net:7001

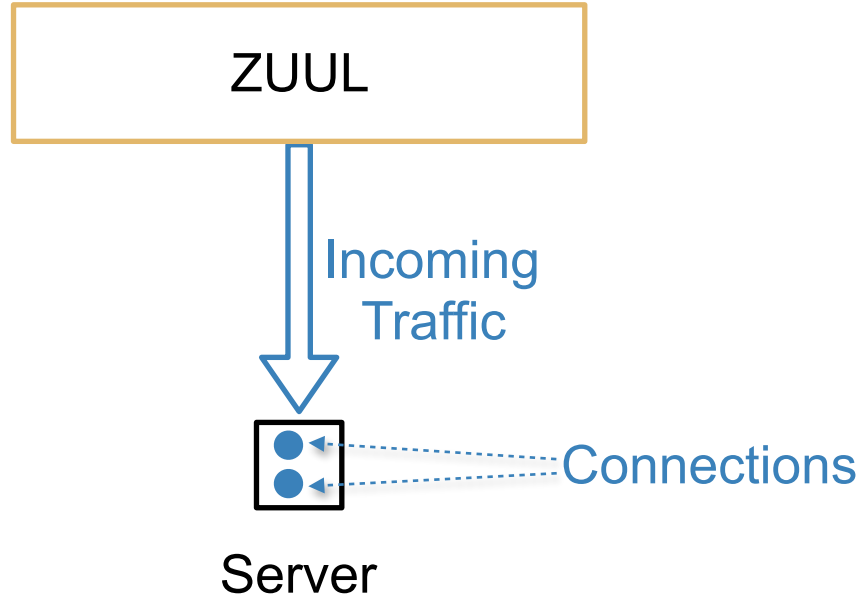
2. Load Balance



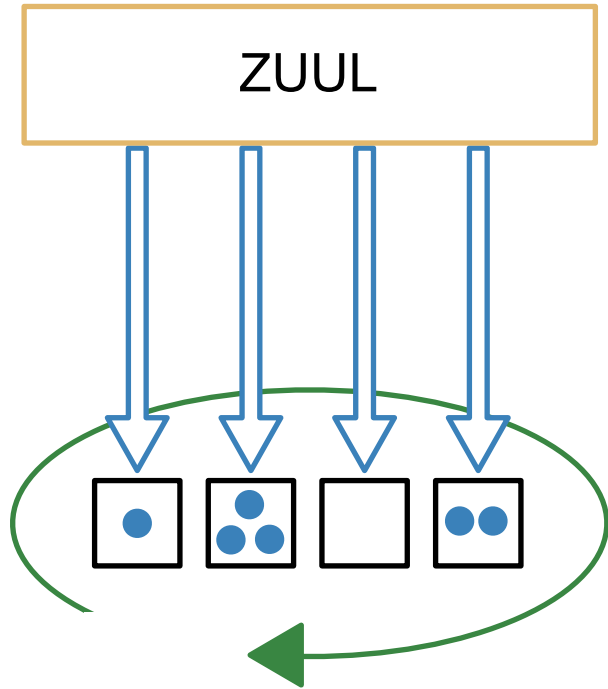
CORE
RESPONSIBILITIES

NETFLIX

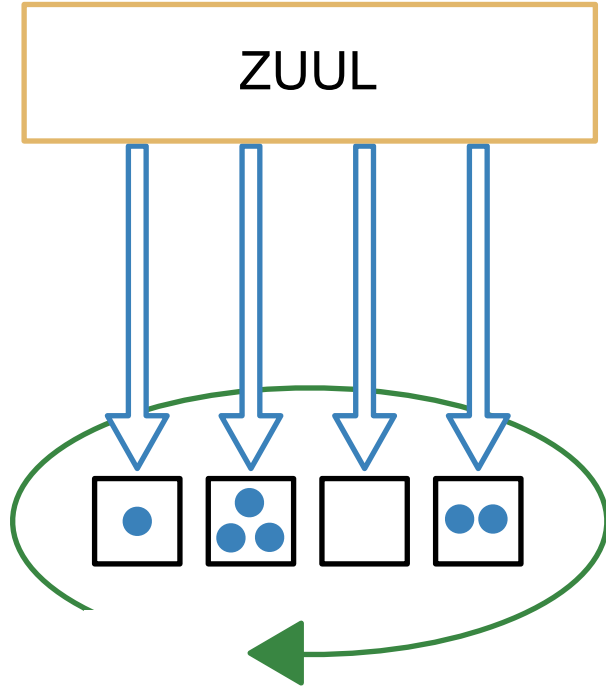
Load balancing notation



Round Robin

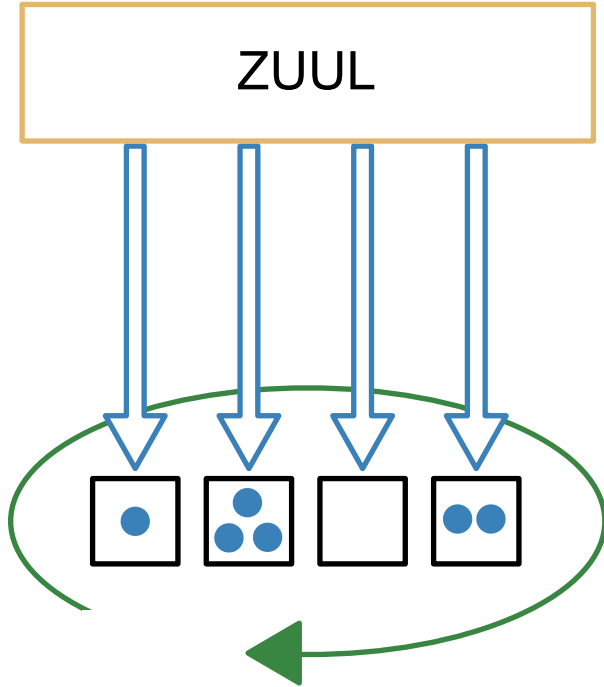


Round Robin



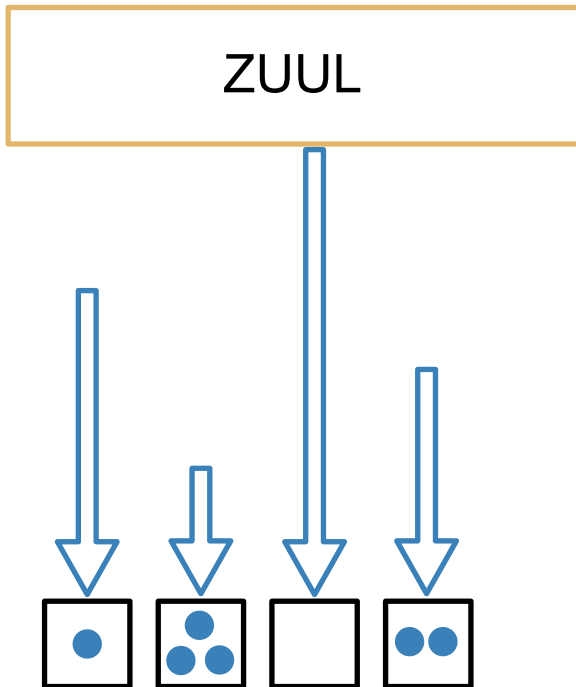
Easiest to implement,
No locking/ Synchronization

Round Robin

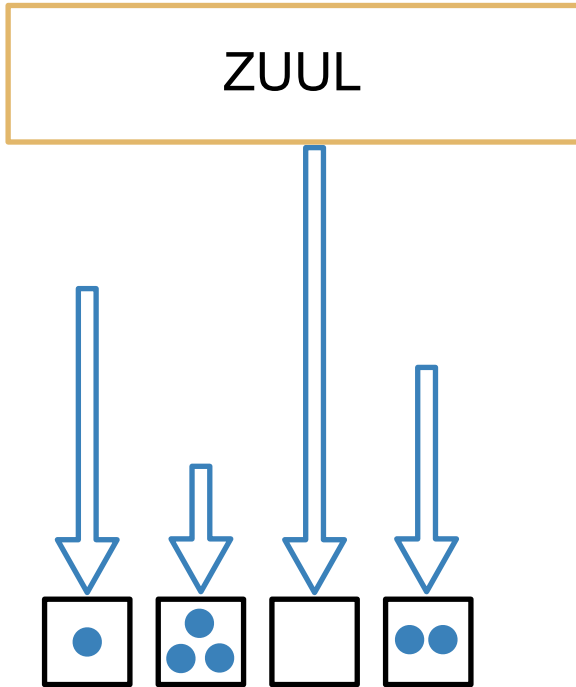


- 👍 Easiest to implement,
No locking/ Synchronization
- 👎 Unequal resource utilization

Least Connections

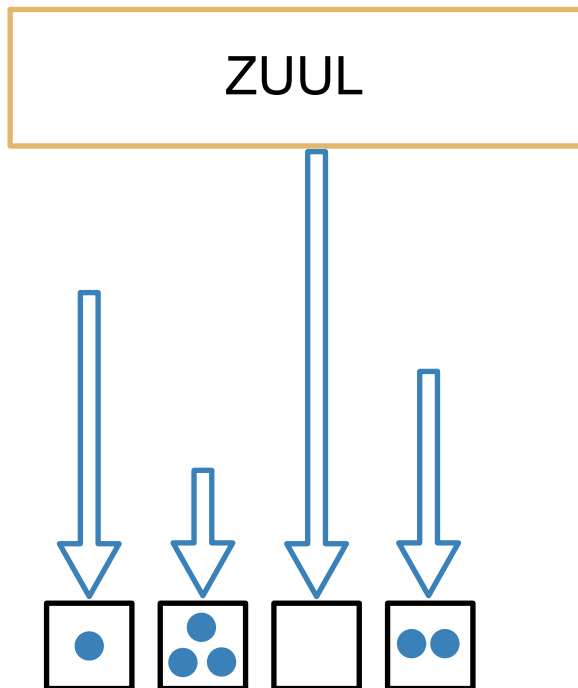


Least Connections



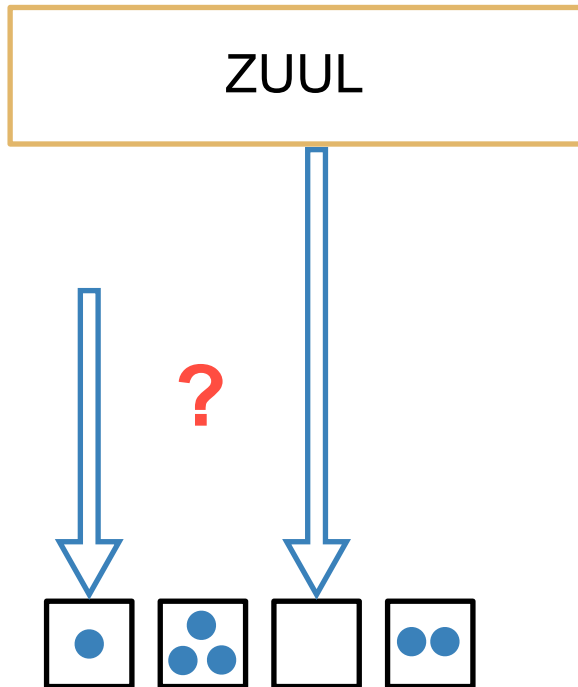
More balanced resource utilization

Least Connections

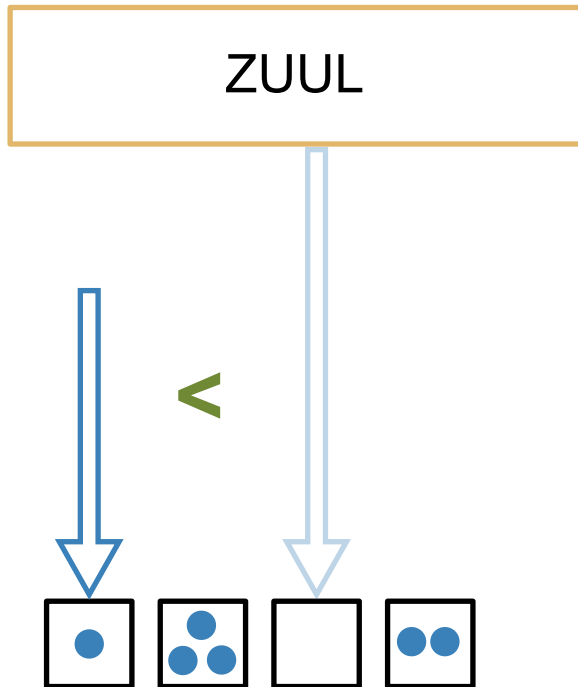


- 👍 More balanced resource utilization
- 👎 Selecting next server is expensive, $O(N)$ operation.

Power of 2

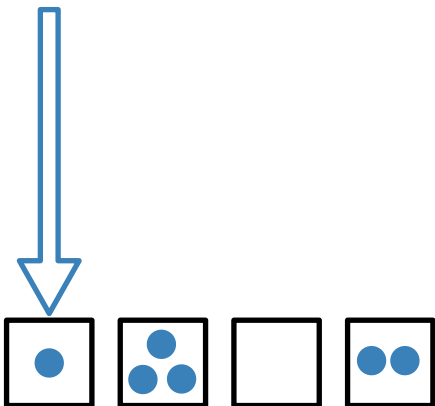


Power of 2



Power of 2

ZUUL



Power of 2

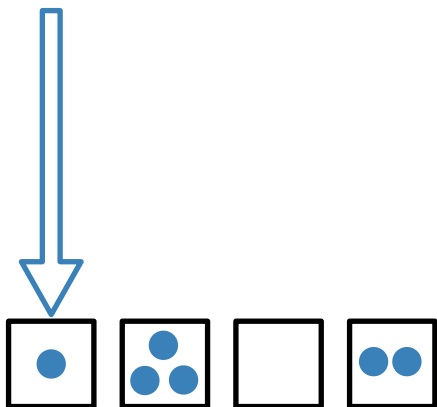
ZUUL



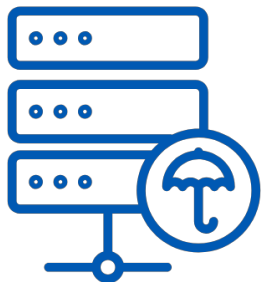
More balanced resource utilization



Inexpensive next server selection



3. Protect Origins



CORE
RESPONSIBILITIES

NETFLIX

Protection from slow clients - timeouts

Default = 100 ms



Protection from slow clients - timeouts



Default = 100 ms

└─> Origin *Billing* = 150 ms

Protection from slow clients - timeouts



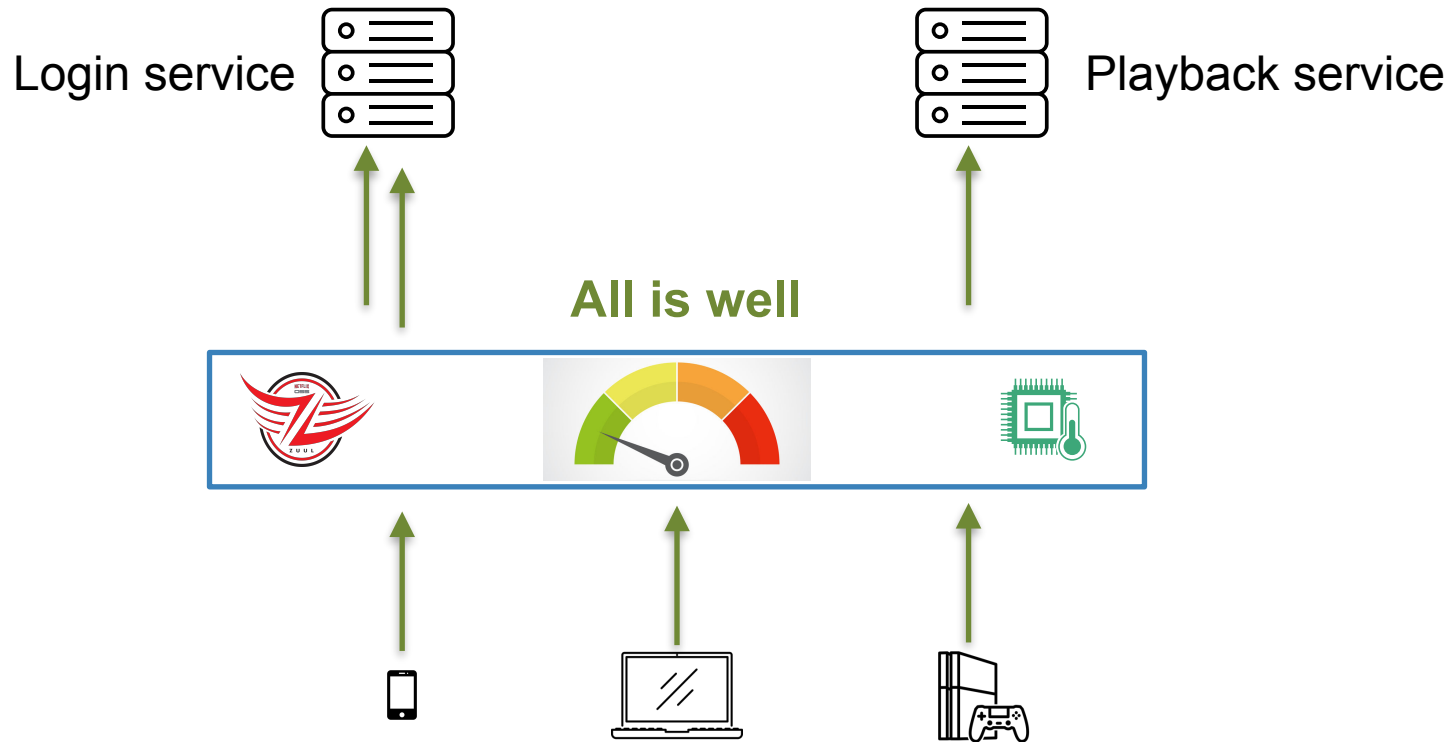
Default = 100 ms

└─> Origin *Billing* = 150 ms

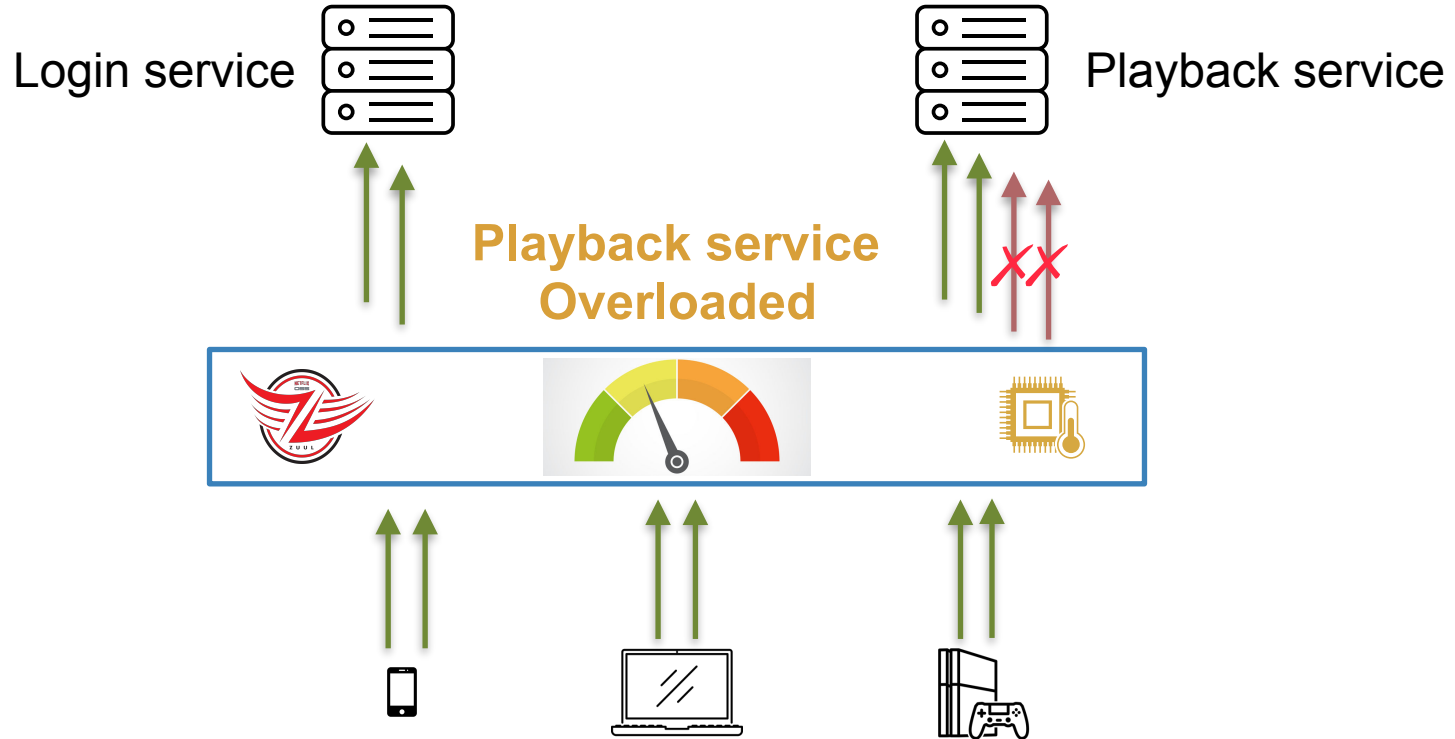
└─> Path */payment* = 200 ms

Surviving sudden surge in traffic - throttling and load shedding

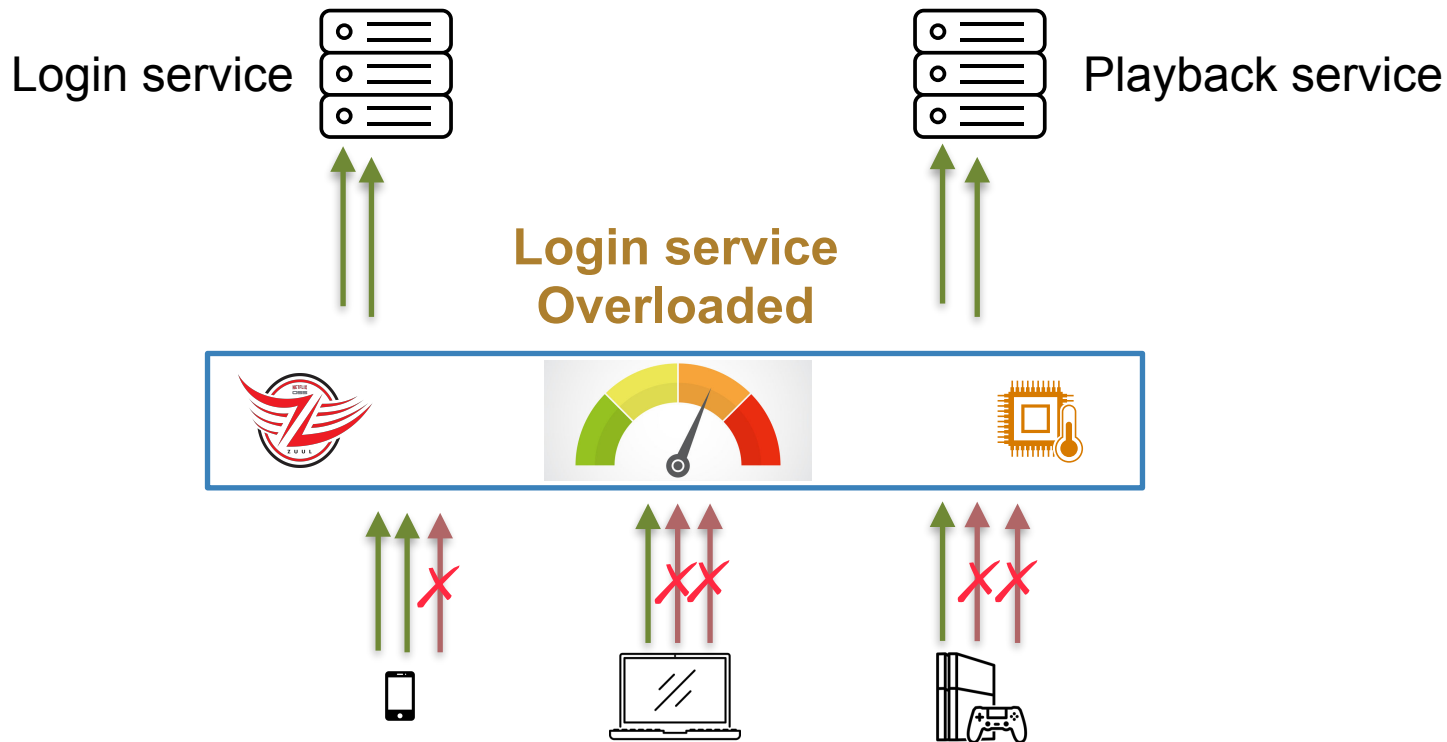
Surviving sudden surge in traffic - throttling and load shedding



Per origin concurrency throttling

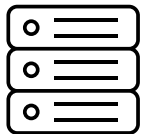


Global concurrency throttling at Zuul frontend

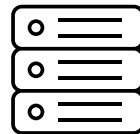


Zuul CPU load shedding

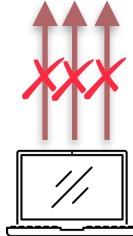
Login service



Playback service

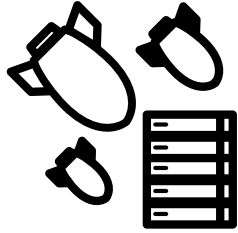


**Zuul CPU hits
the wall**



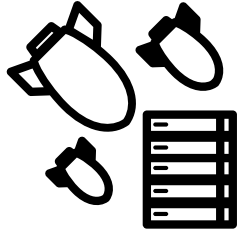
Other side of protection - security

Other side of protection - security



DDoS protection

Other side of protection - security

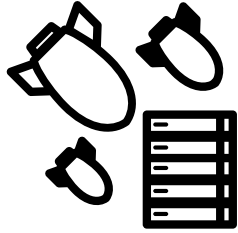


DDoS protection



Single sign-on

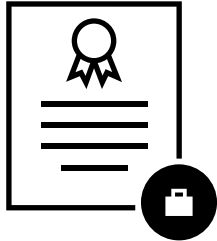
Other side of protection - security



DDoS protection

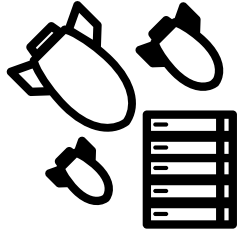


Single sign-on



Client auth certificates

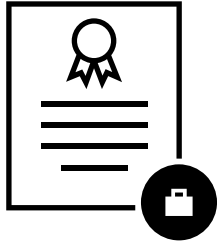
Other side of protection - security



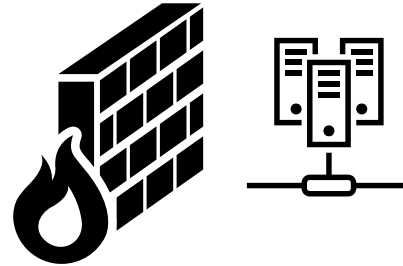
DDoS protection



Single sign-on

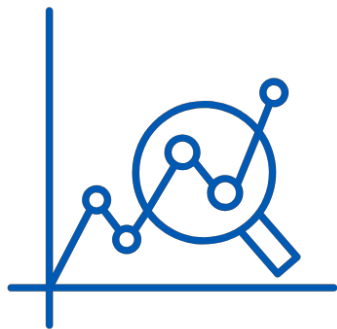


Client auth certificates



Whitelisted internal
CORP access

4. Metrics and Monitoring



CORE
RESPONSIBILITIES

NETFLIX

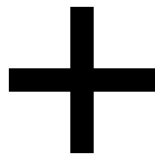
Historic and near real time aggregate metrics - Atlas

Metrics

Total request counts

Latency

Geo location



Atlas

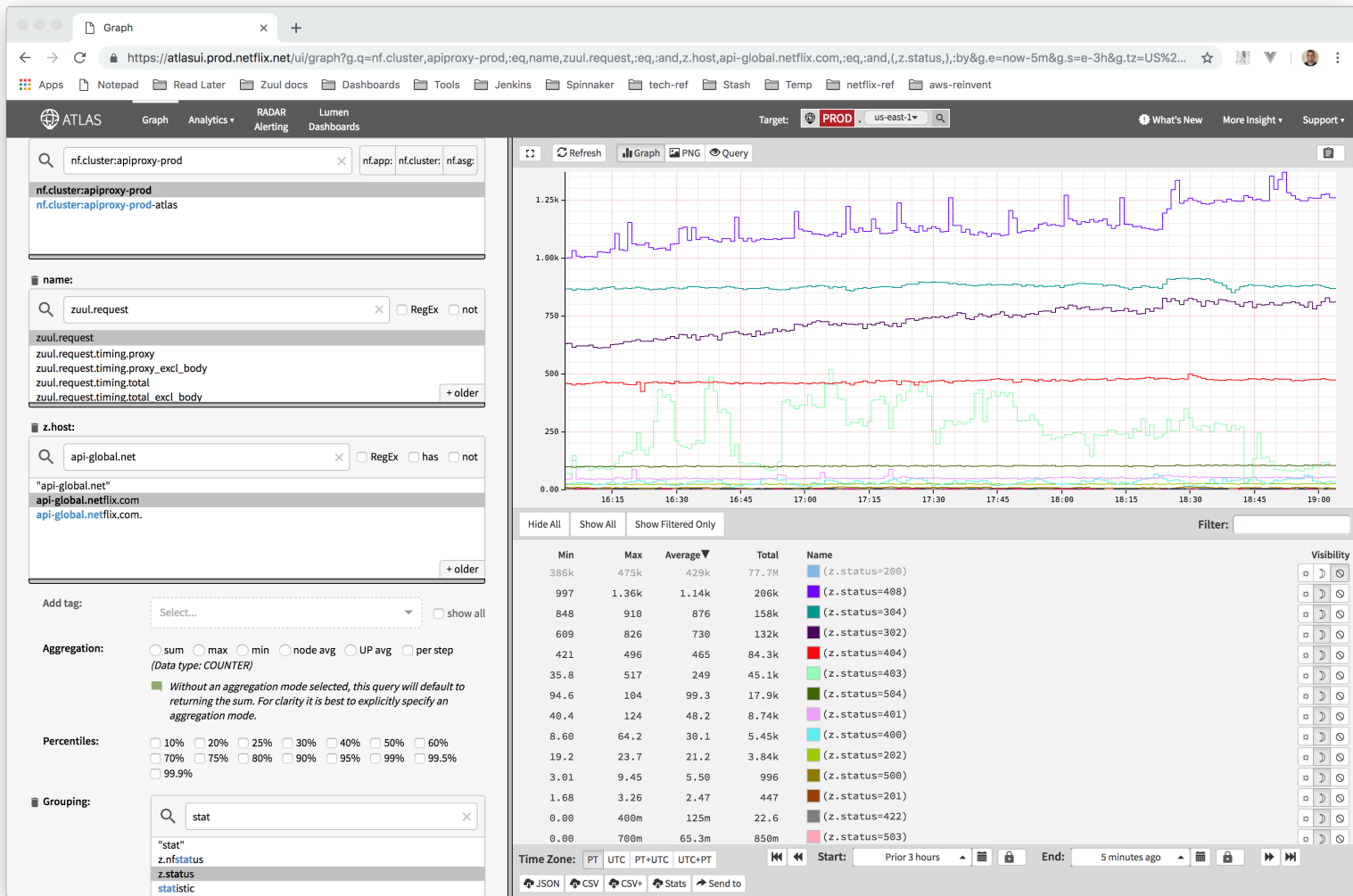
Tags

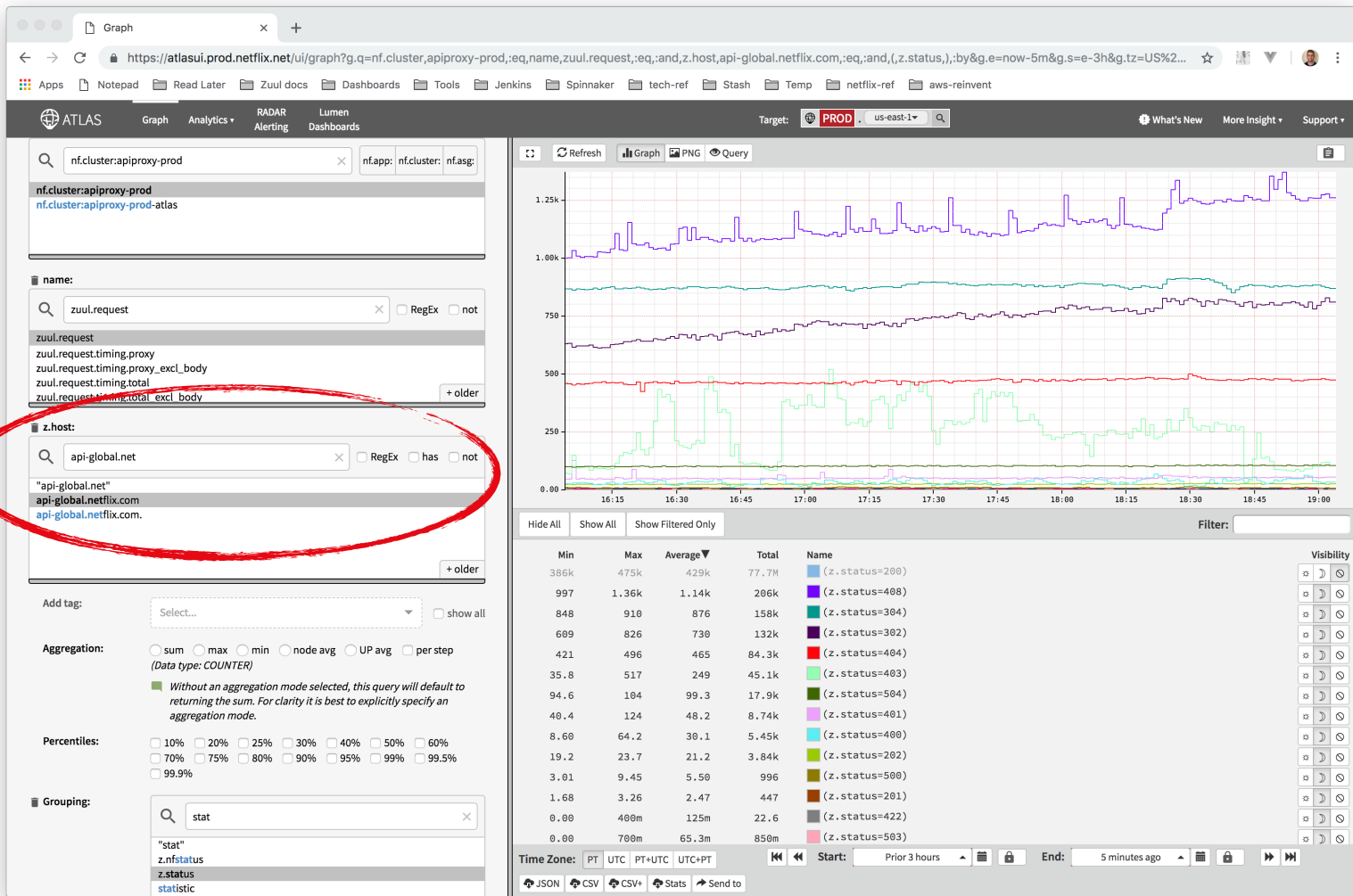
Status

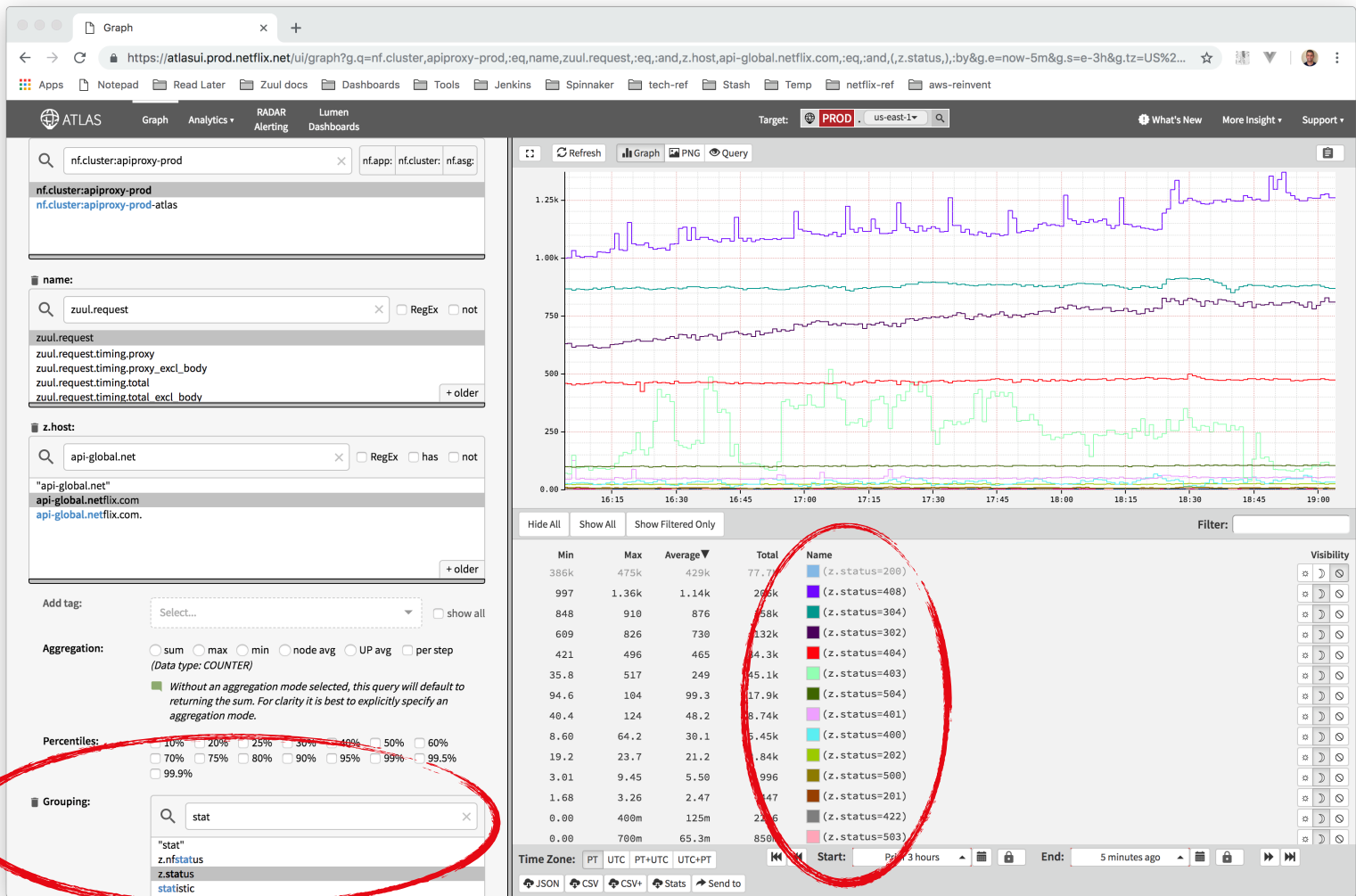
Host name

Device type





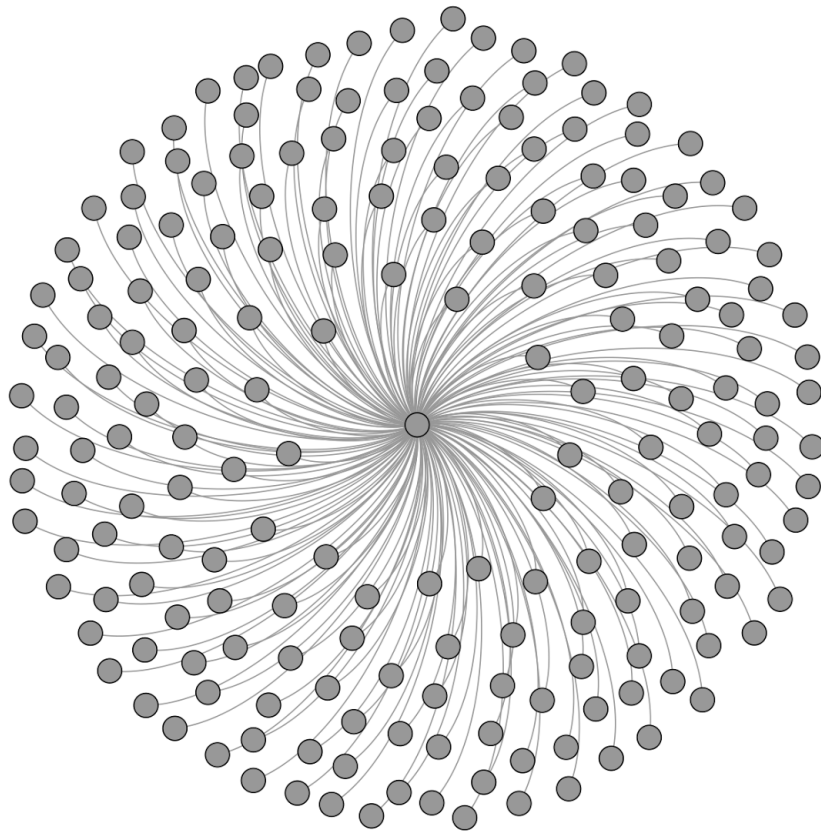




How Zuul does it?

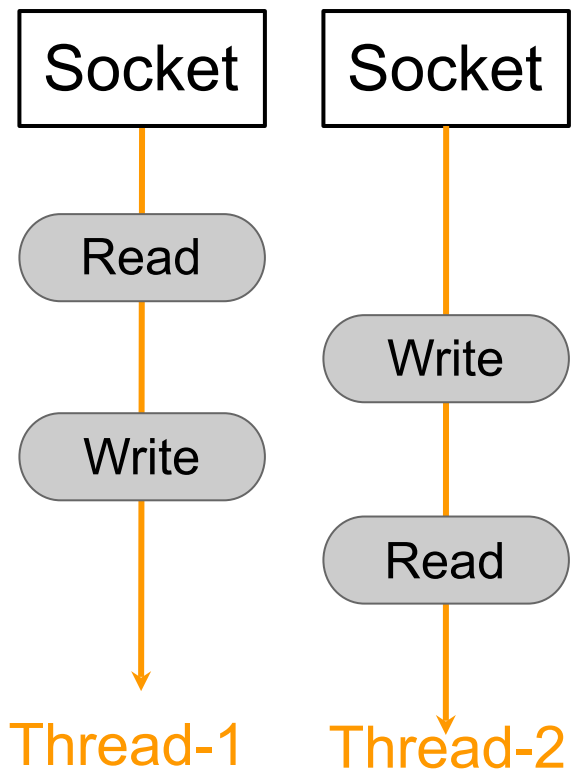
| Zuul Architecture



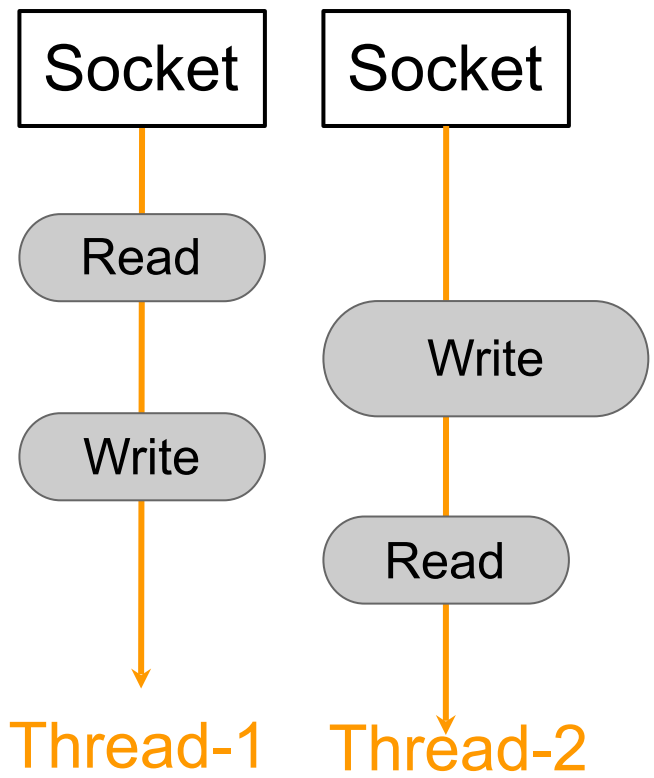


C10K challenge

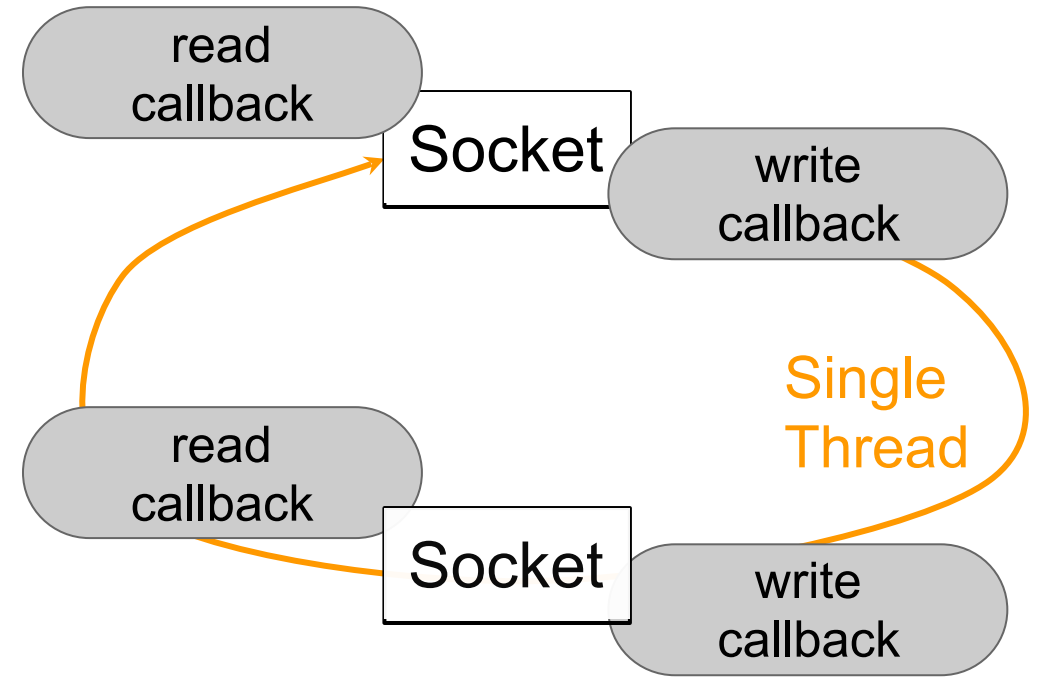
Thread per Connection



Thread per Connection



Async I/O



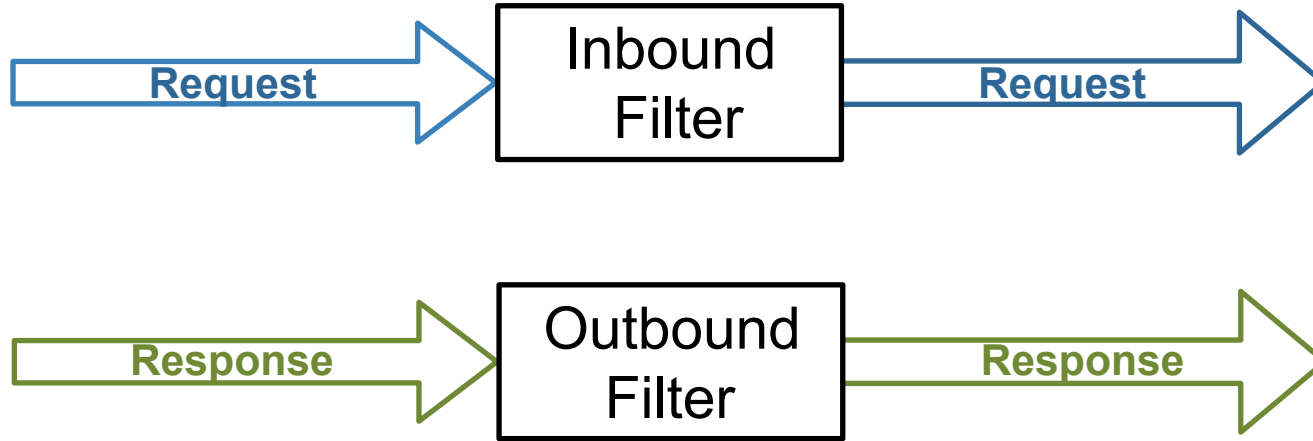
Core abstraction at the heart of the Zuul - Filter



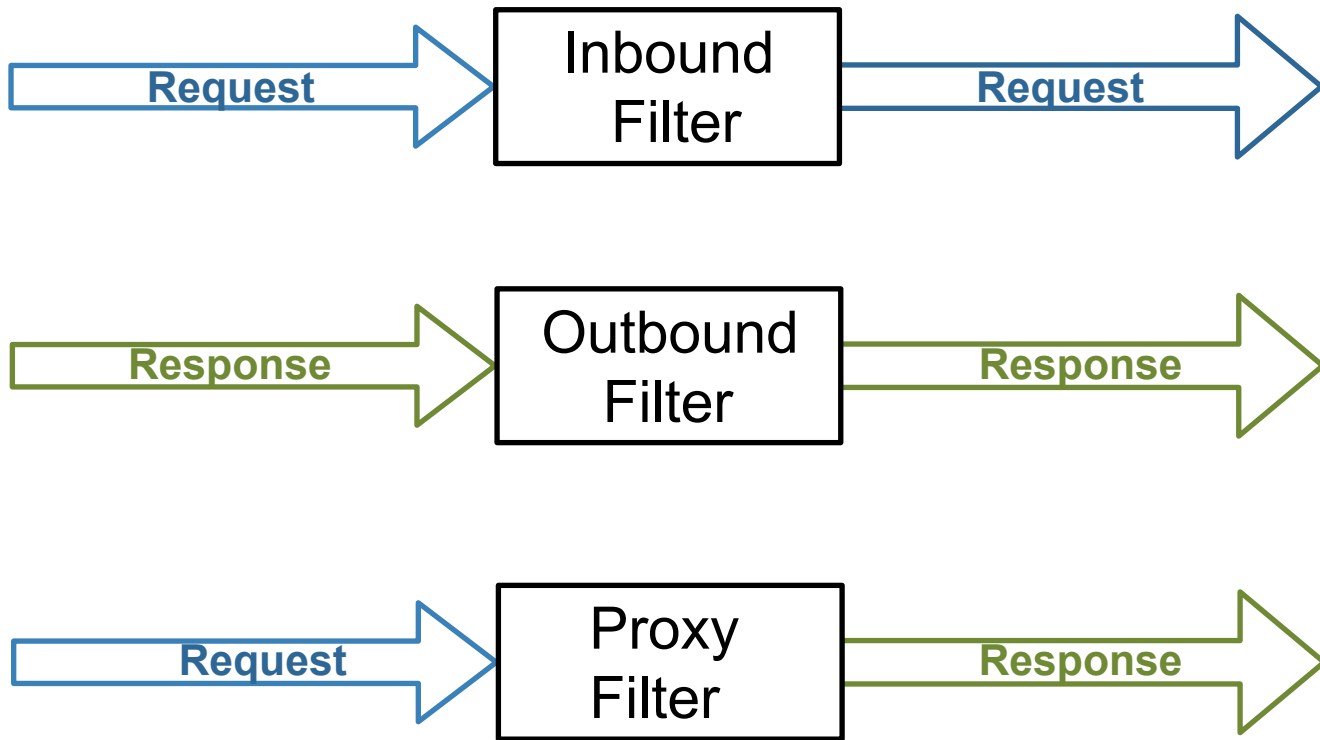
Zuul is made up of different types of filters



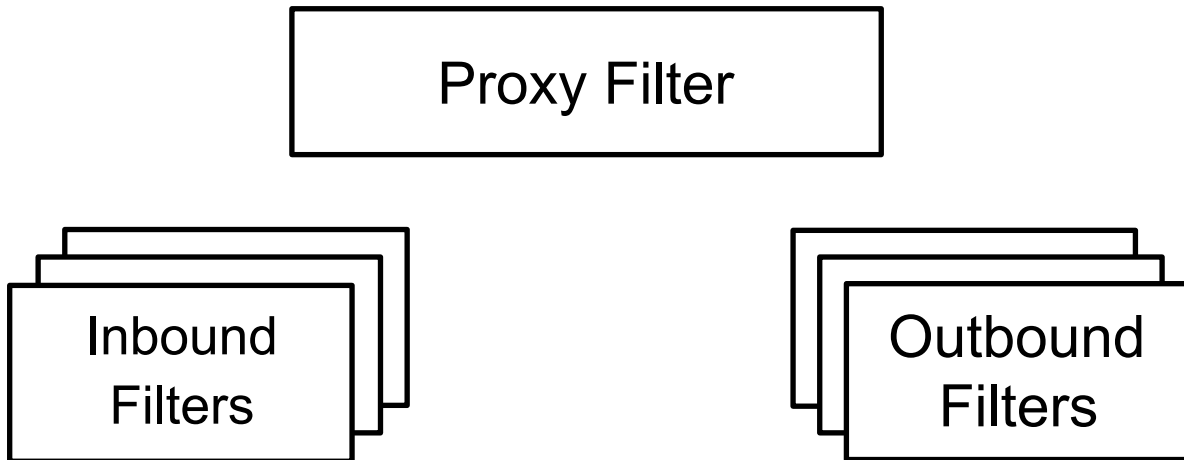
Zuul is made up of different types of filters



Zuul is made up of different types of filters

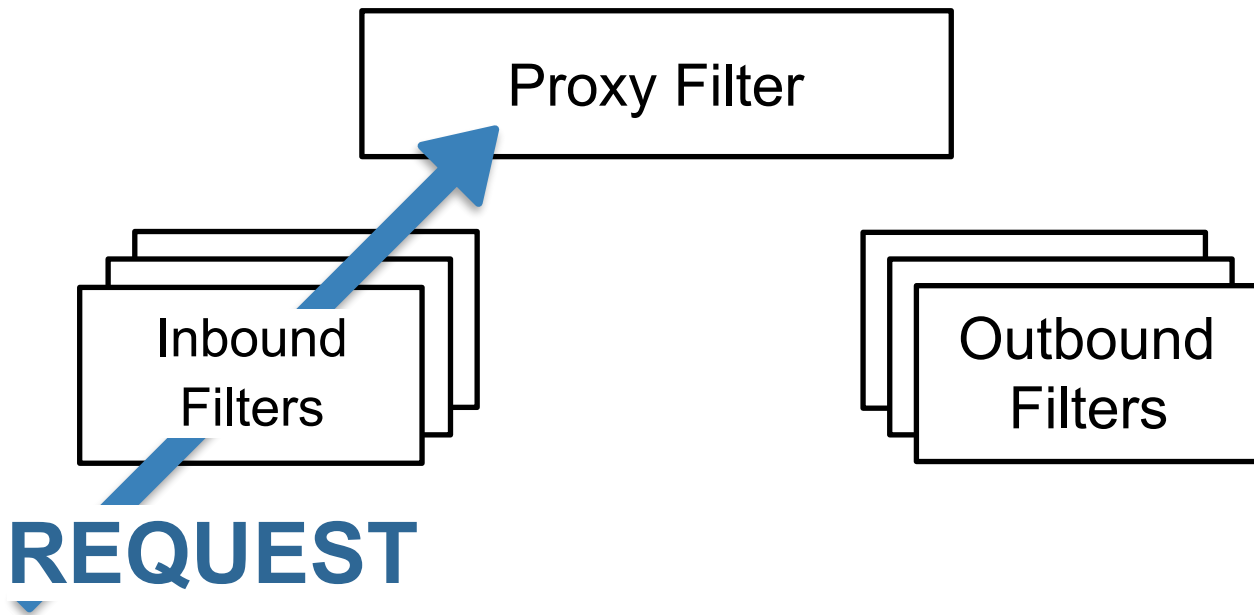


Simple, extensible Zuul core architecture

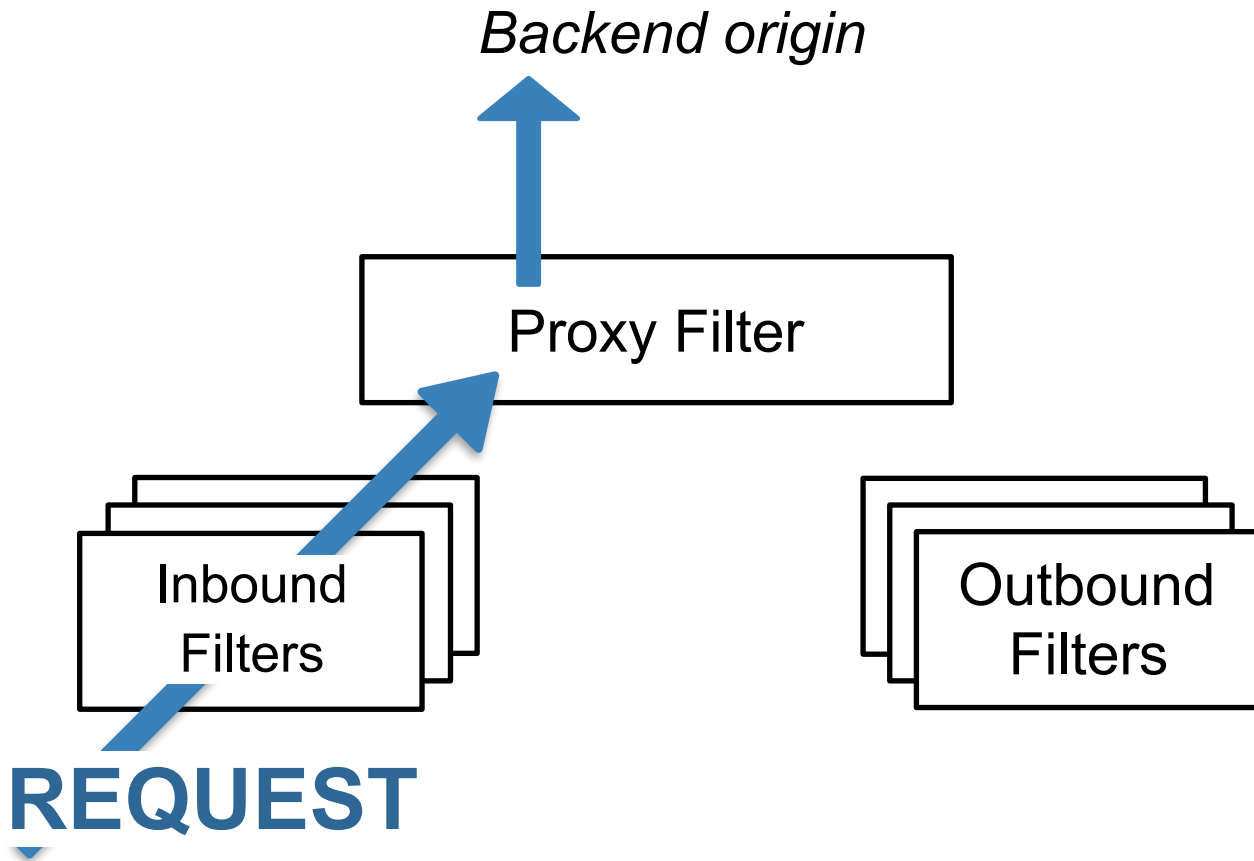


Simple, extensible Zuul core architecture

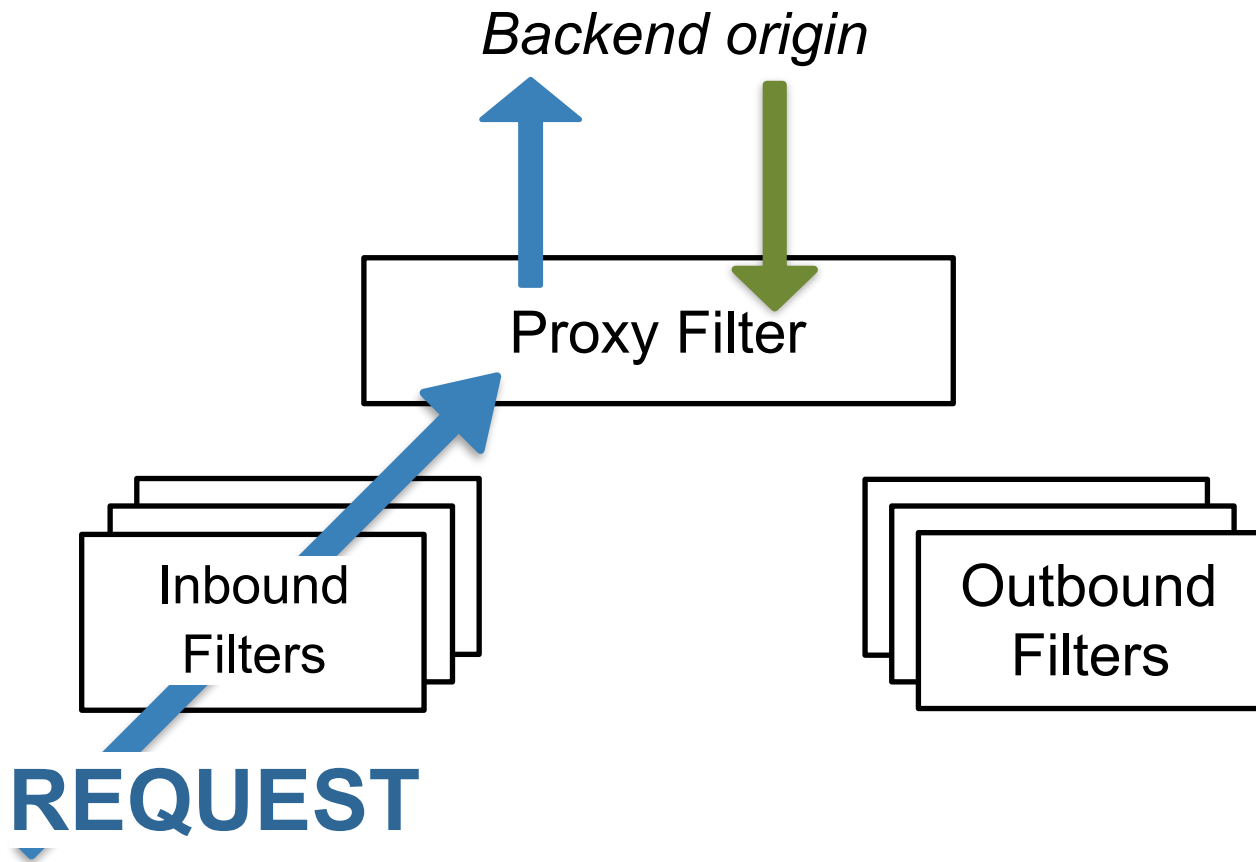
Backend origin



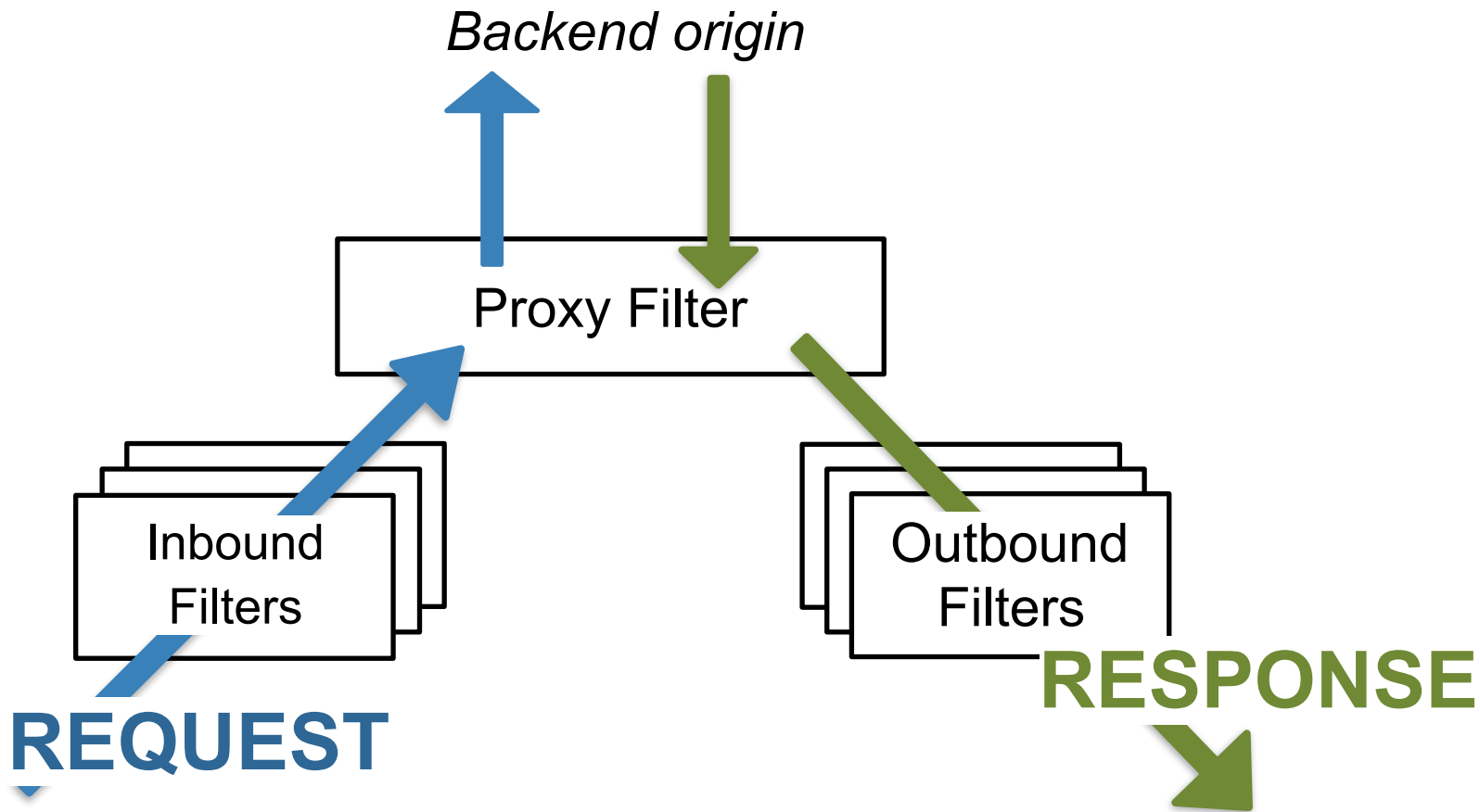
Simple, extensible Zuul core architecture



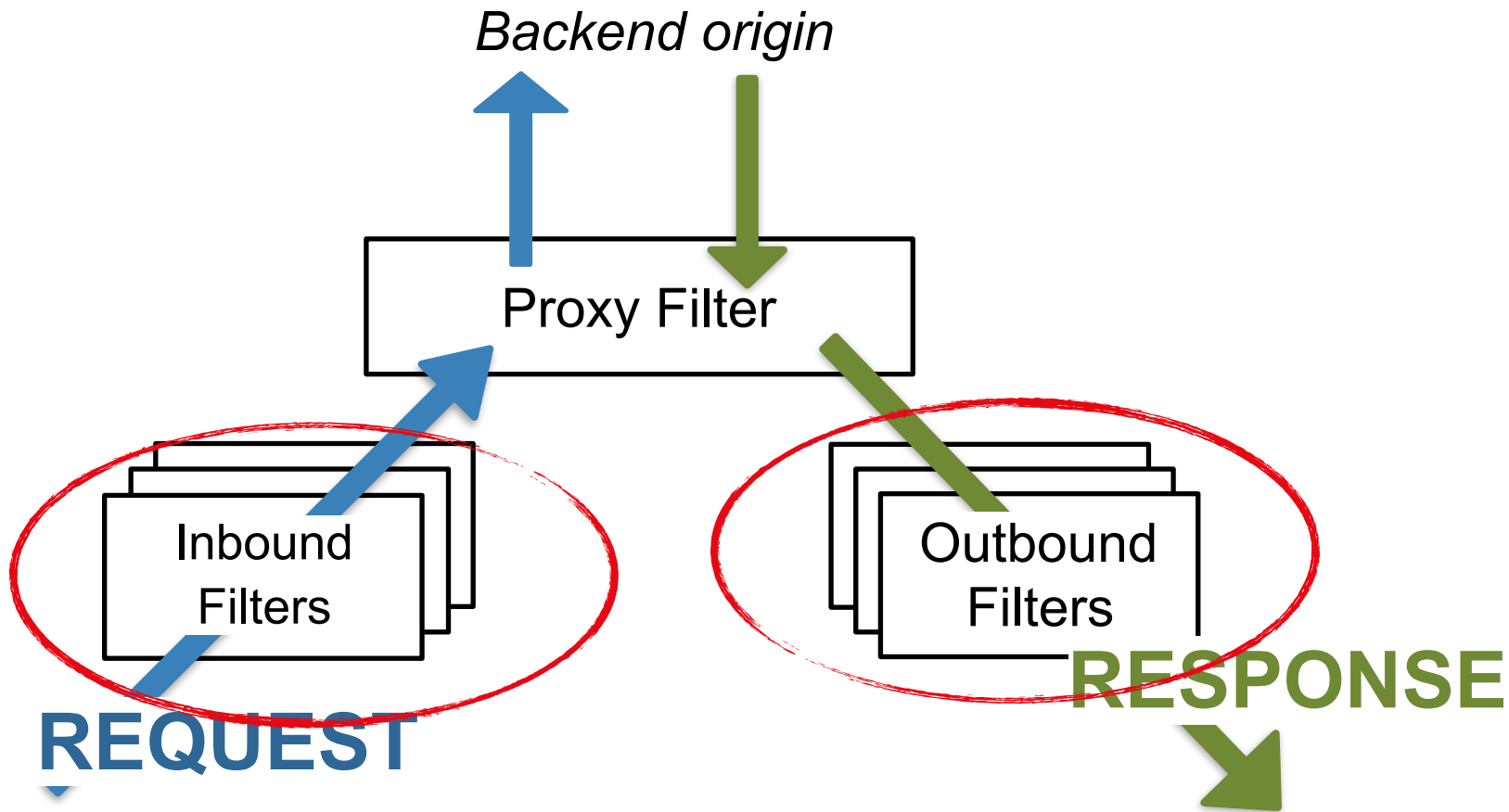
Simple, extensible Zuul core architecture



Simple, extensible Zuul core architecture



Simple, extensible Zuul core architecture



Future roadmap

| Plans for future



Planned features

- End to end HTTP/2 proxying

Planned features

- End to end HTTP/2 proxying
- WebSocket proxying

Planned features

- End to end HTTP/2 proxying
- WebSocket proxying
- gRPC proxying

Planned features

- End to end HTTP/2 proxying
- WebSocket proxying
- gRPC proxying
- Edge caching

In conclusion



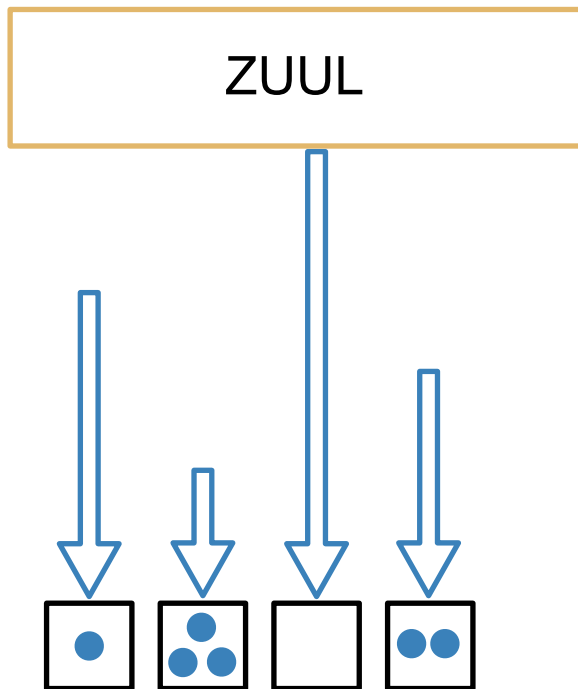
You can deploy Zuul at the edge of your network to gain higher **availability** and better **visibility** with unparalleled **flexibility** in handling your traffic.

Thank you.

NETFLIX



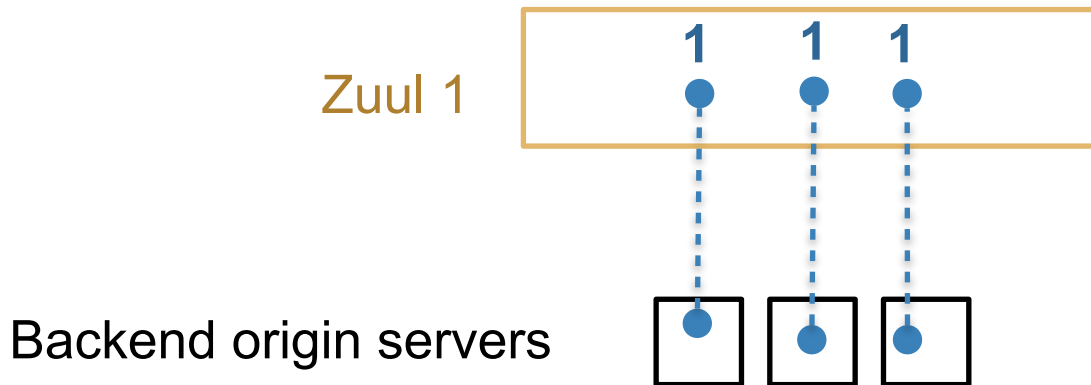
Least Connections



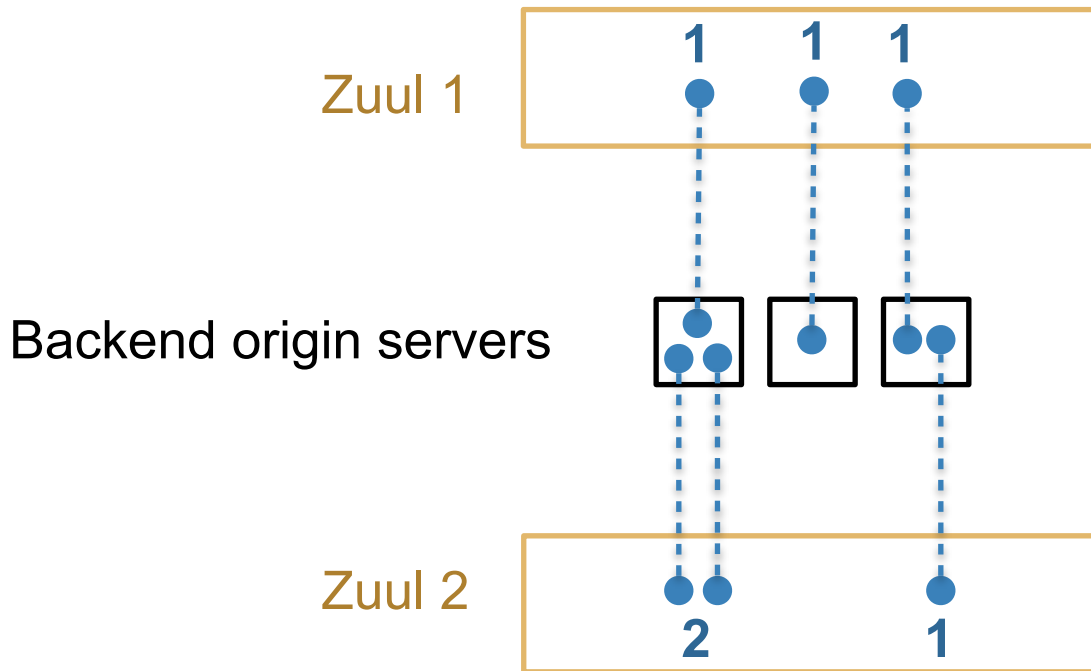
👍 More balanced resource utilization

👎 Tricky to track number of connections when multiple load balancing nodes are involved

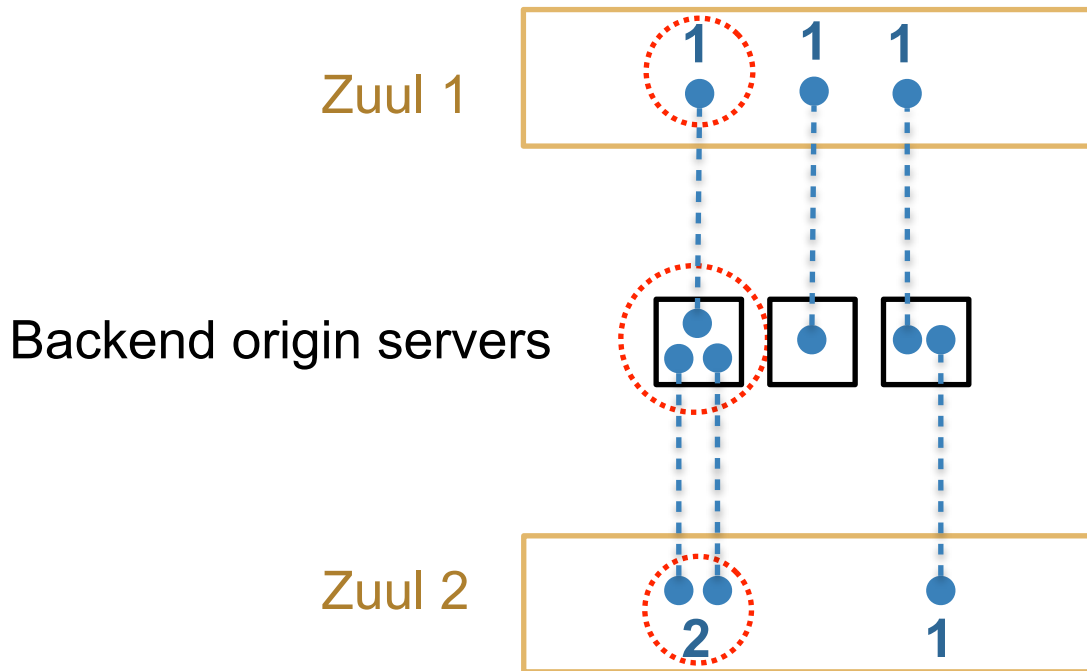
Determining number of connections with more than one LB



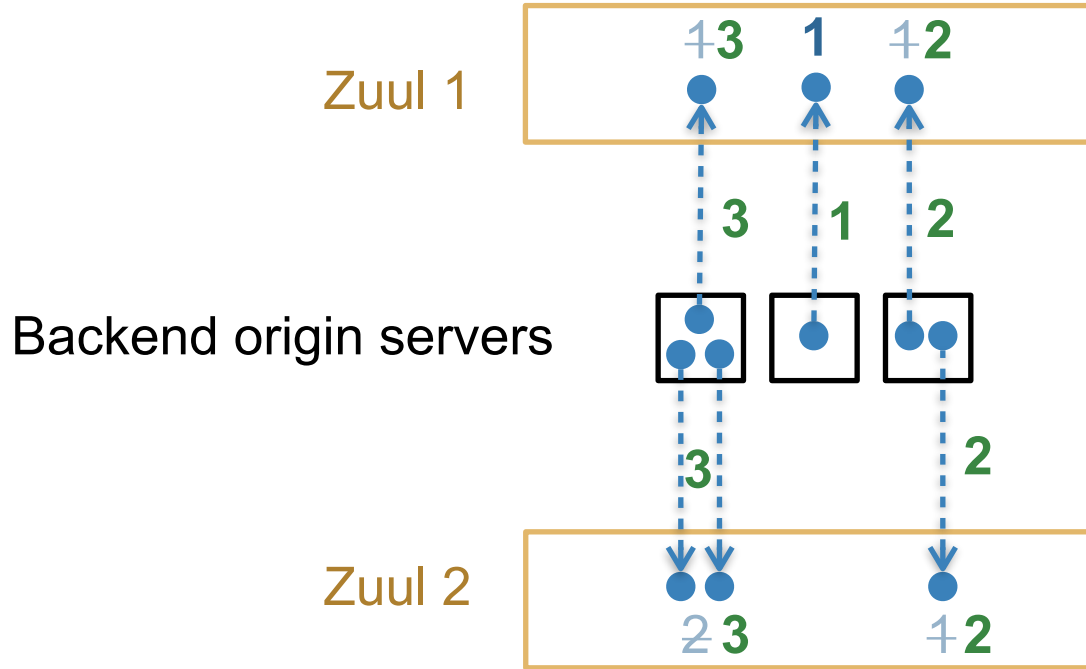
Determining number of connections with more than one LB



Determining number of connections with more than one LB



Zuul's improvement to least connections algorithm



Zuul's improvement

