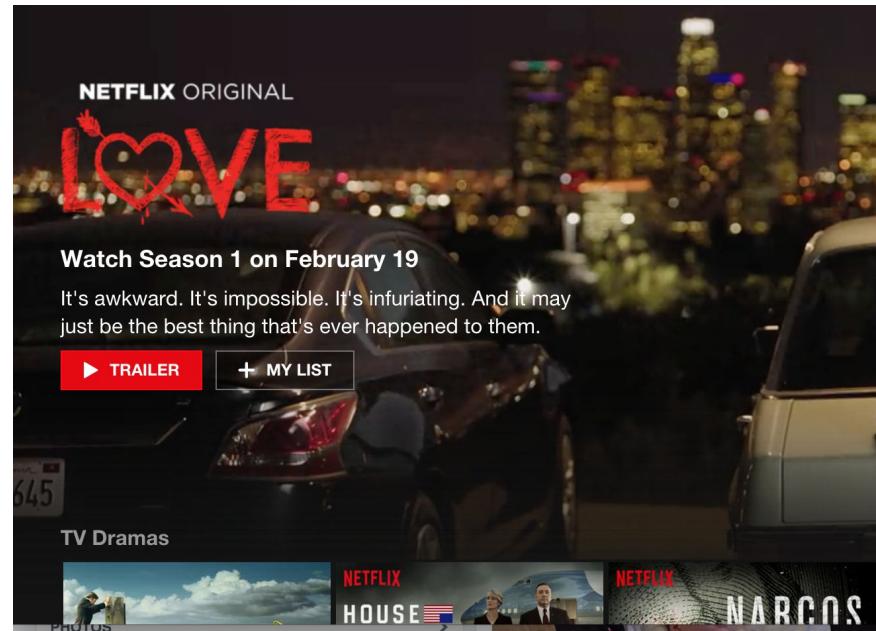


# The Netflix API for a global service

Katharina Probst  
Engineering Manager, API  
DevNexus, February 2016

# NETFLIX



# What is Netflix?

Stream TV shows and movies anywhere, any time.





# Global!

*(except China and where we can't operate for legal reasons)*

# Netflix Originals



A NETFLIX ORIGINAL FILM  
**BEASTS OF NO NATION**

★★★★★ 2015 NR 2h 17m

[Resume](#)

1 of 137m

A brutal war took a boy's family. A mercenary commander takes his youth. In this war, the demons come for everyone.

[MY LIST](#)



NETFLIX ORIGINAL  
**CHEF'S TABLE**

★★★★★ 2015 TV-14 1 Season



KEVIN SPACEY



**HOUSE of CARDS**  
BEGINNING FEBRUARY 1

AN ORIGINAL  
NETFLIX  
SERIES

NETFLIX ORIGINAL



Watch Season 1 on February 19

It's awkward. It's impossible. It's infuriating. And it may just be the best thing that's ever happened to them.

[► TRAILER](#)

[+ MY LIST](#)

TV Dramas



# Scale

- Peak downstream traffic in the US is 37%, upstream almost 7%.

Rank	Upstream		Downstream		Aggregate	
	Application	Share	Application	Share	Application	Share
1	BitTorrent	28.56%	Netflix	37.05%	Netflix	34.70%
2	Netflix	6.78%	YouTube	17.85%	YouTube	16.88%
3	HTTP	5.93%	HTTP	6.06%	HTTP	6.05%
4	Google Cloud	5.30%	Amazon Video	3.11%	BitTorrent	4.35%
5	YouTube	5.21%	iTunes	2.79%	Amazon Video	2.94%
6	SSL - OTHER	5.10%	BitTorrent	2.67%	iTunes	2.62%
7	iCloud	3.08%	Hulu	2.58%	Facebook	2.51%
8	FaceTime	2.55%	Facebook	2.53%	Hulu	2.48%
9	Facebook	2.25%	MPEG - OTHER	2.30%	MPEG	2.16%
10	Dropbox	1.18%	SSL - OTHER	1.73%	SSL - OTHER	1.99%
		65.95%		78.69%		76.68%

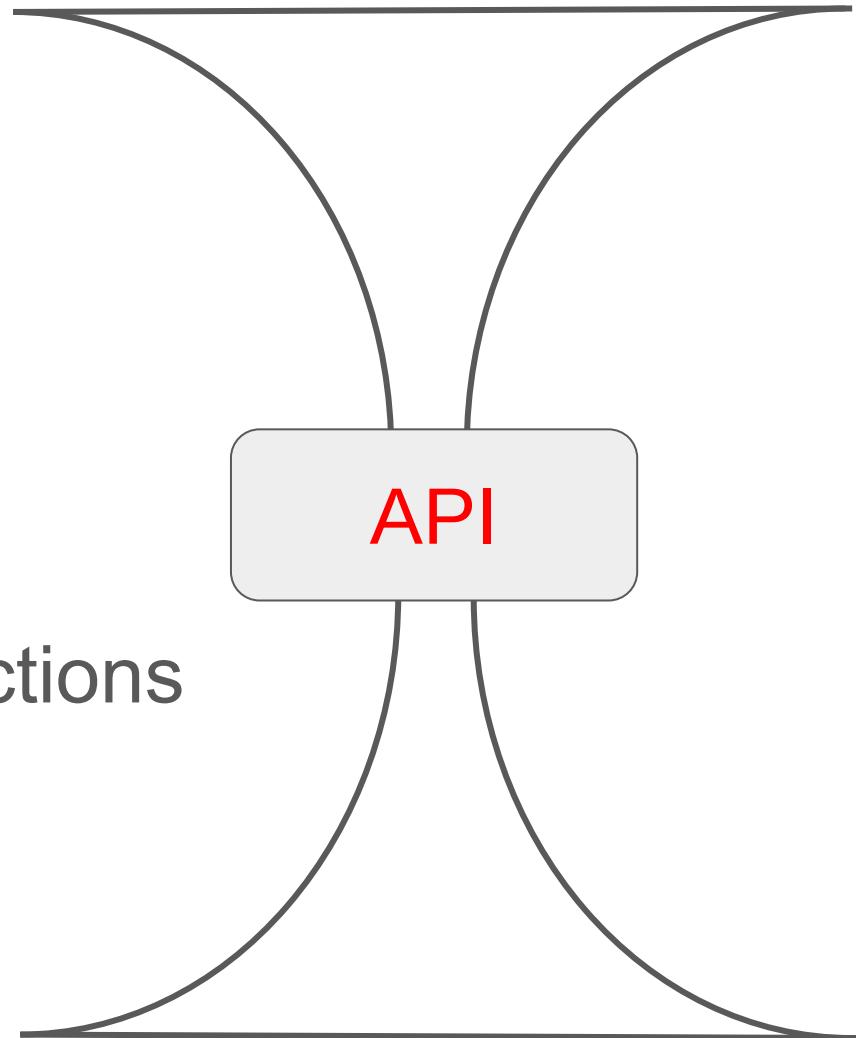


Source: [http://www.sandvine.com/news/global\\_broadband\\_trends.asp](http://www.sandvine.com/news/global_broadband_trends.asp)

- 75 Million subscribers worldwide and growing

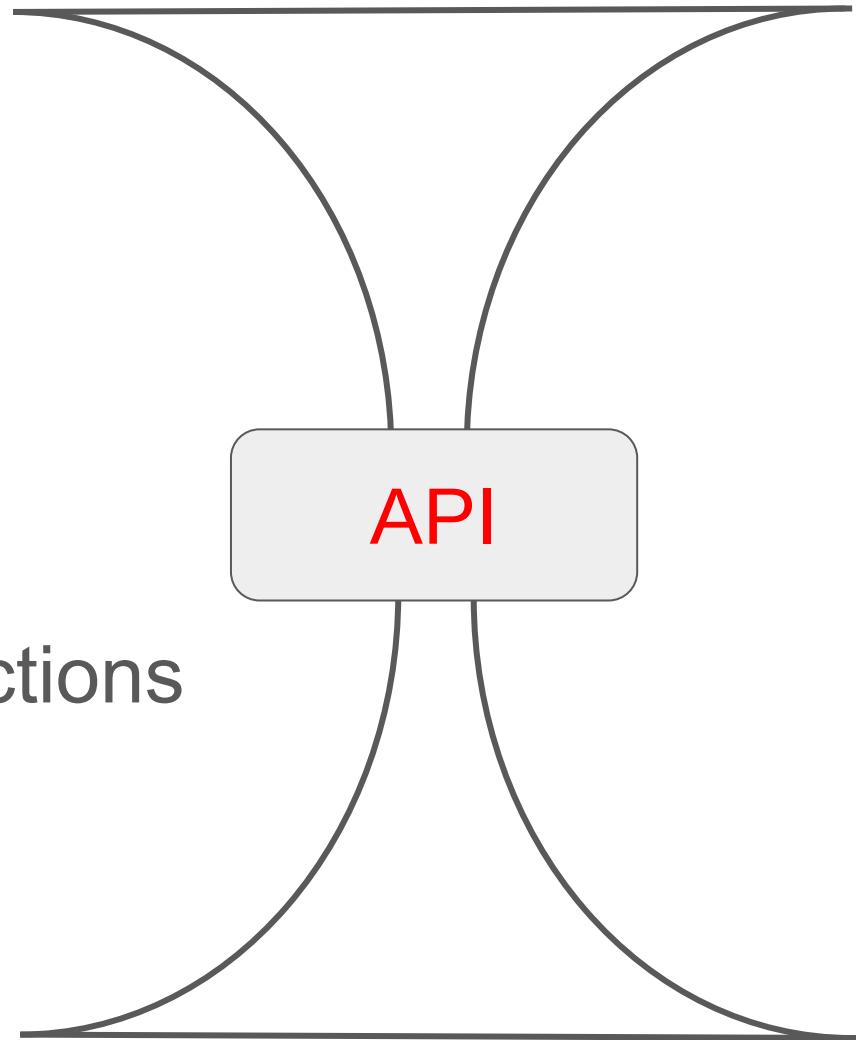
# Netflix API

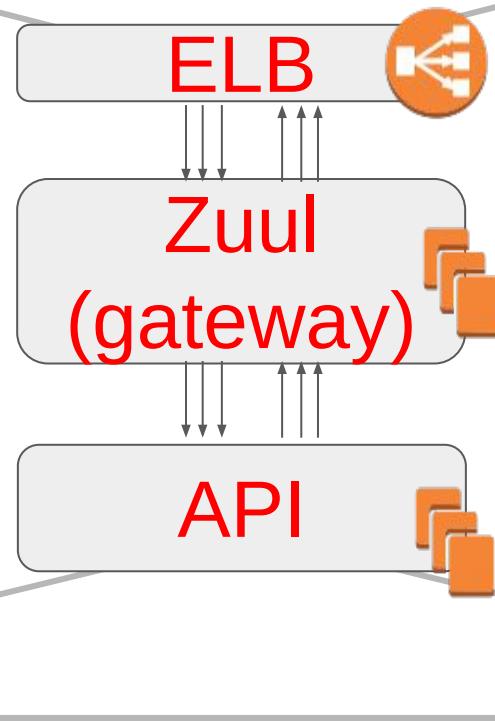
- ❑ Architecture
- ❑ Resiliency
- ❑ Developer velocity
- ❑ Tooling and DevOps
- ❑ Current and future directions



# Netflix API

- ❑ **Architecture**
- ❑ Resiliency
- ❑ Developer velocity
- ❑ Tooling and DevOps
- ❑ Current and future directions





Personalization  
Engine

User  
Info

....

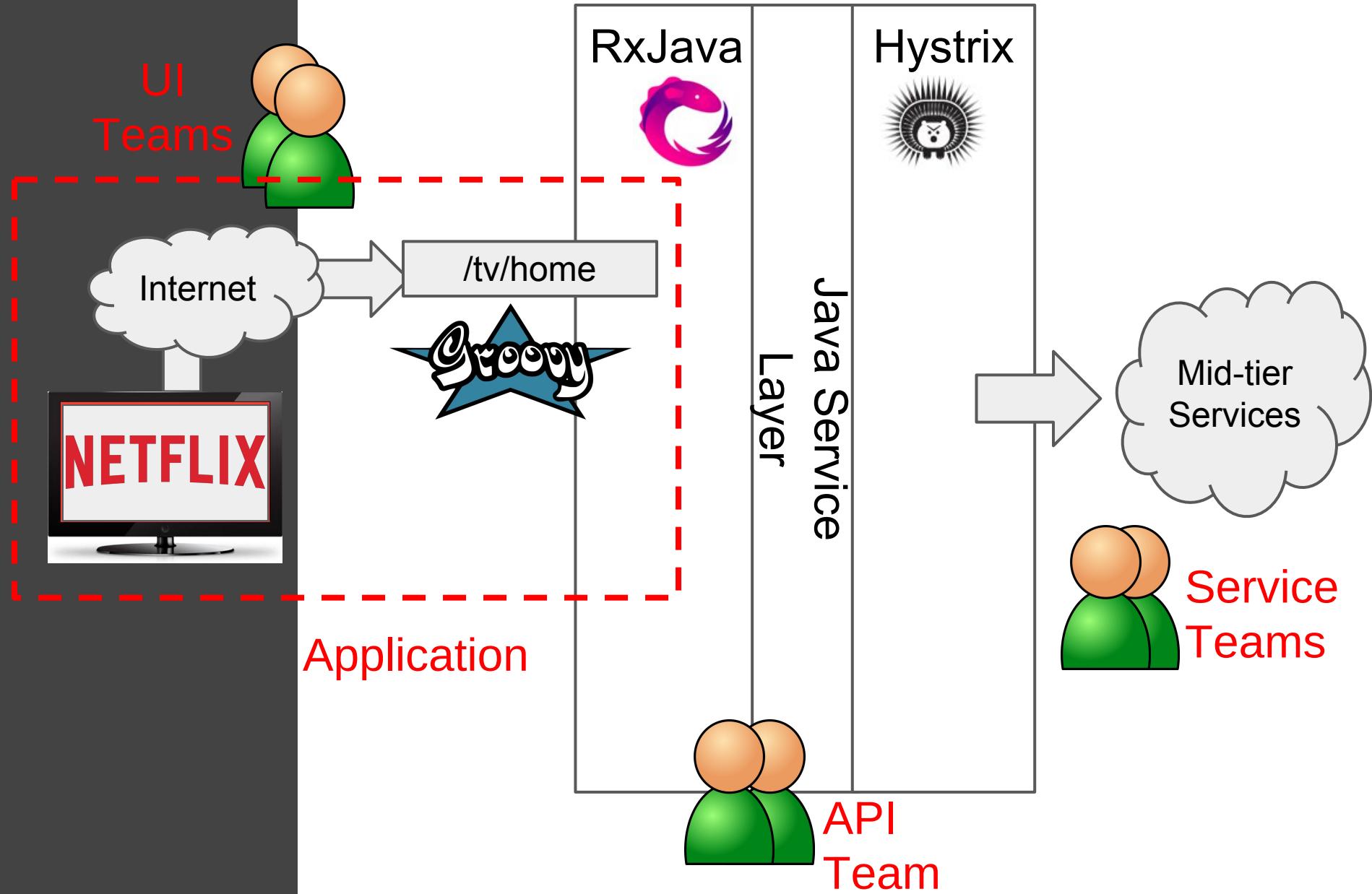
Ratings

Similar  
Movies

A/B Test  
Engine

# Client

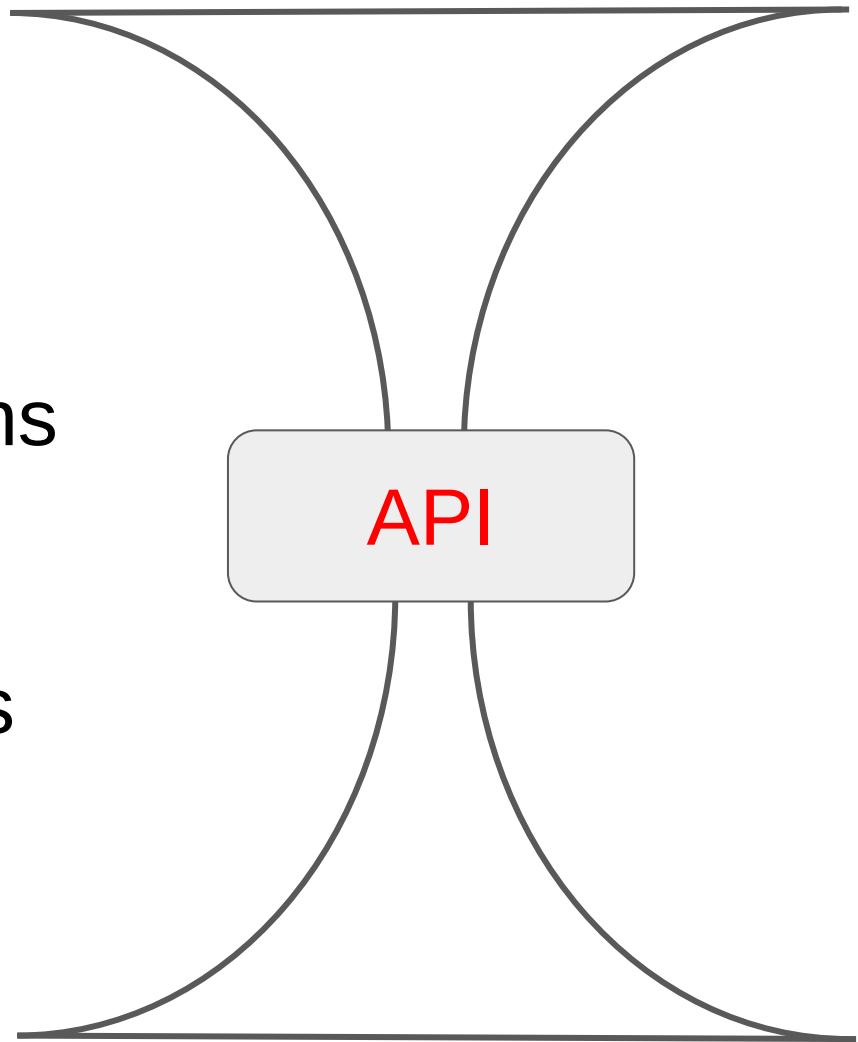
# Server



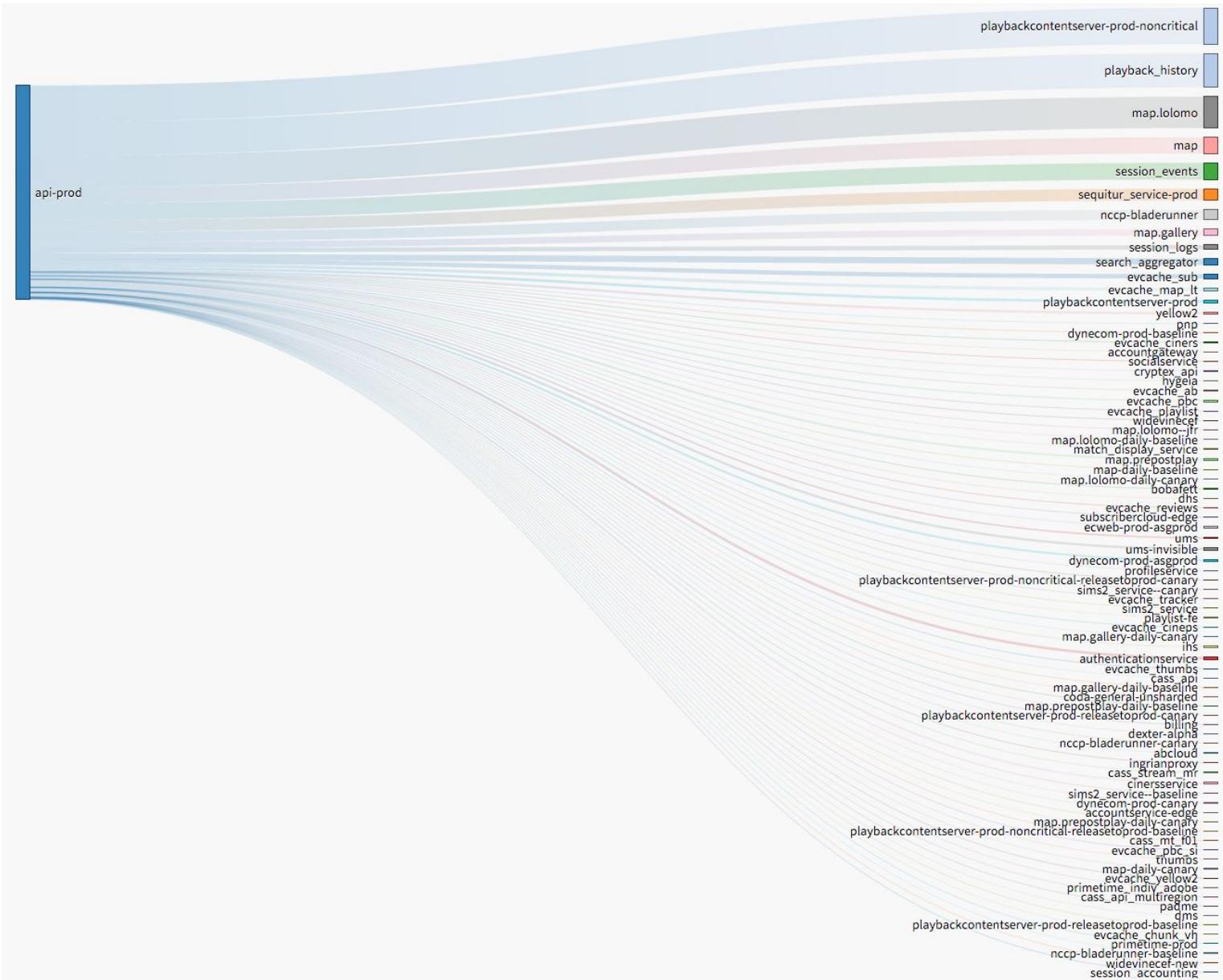
# What is the API used for?

Examples:

- Discovery
  - Recommendations
  - Move metadata
  - Ratings
- Sign-up and Profiles
- Playback
  - Bookmarks
  - DRM
- A/B testing

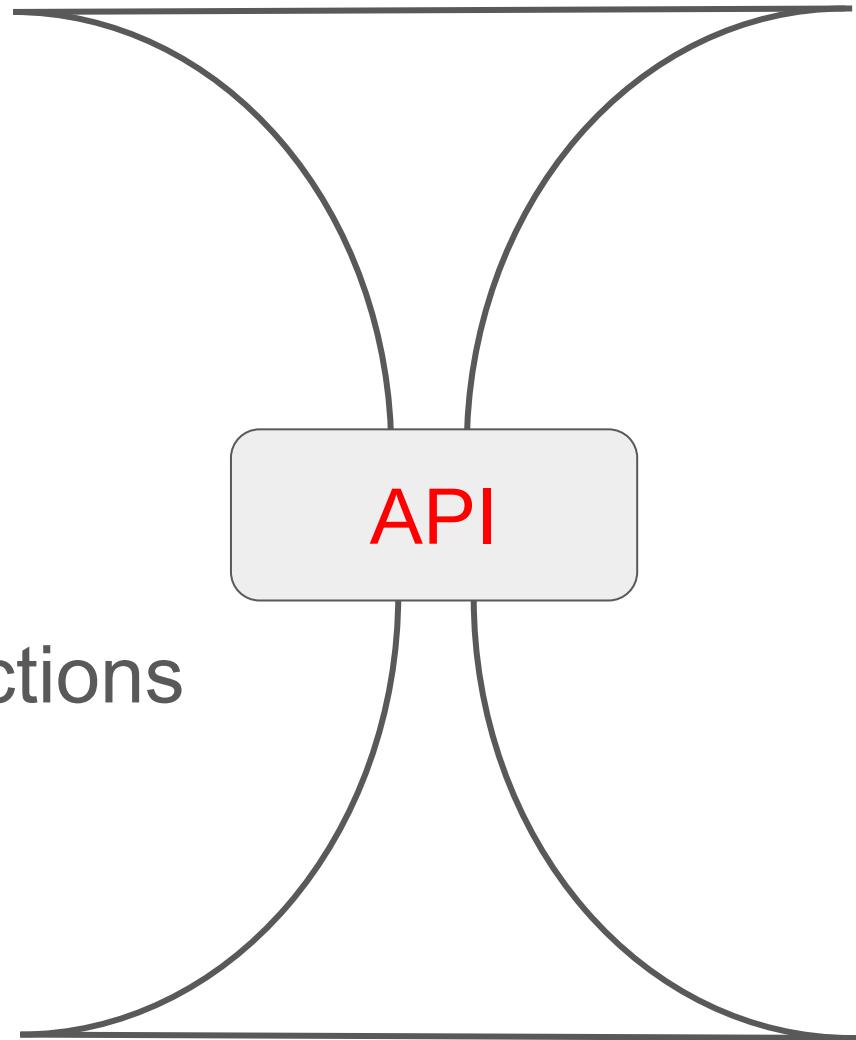


# Direct dependencies on other services

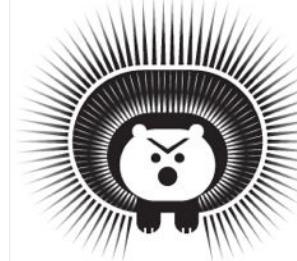


# Netflix API

- ❑ Architecture
- ❑ **Resiliency**
- ❑ Developer velocity
- ❑ Tooling and DevOps
- ❑ Current and future directions



# Hystrix Primer

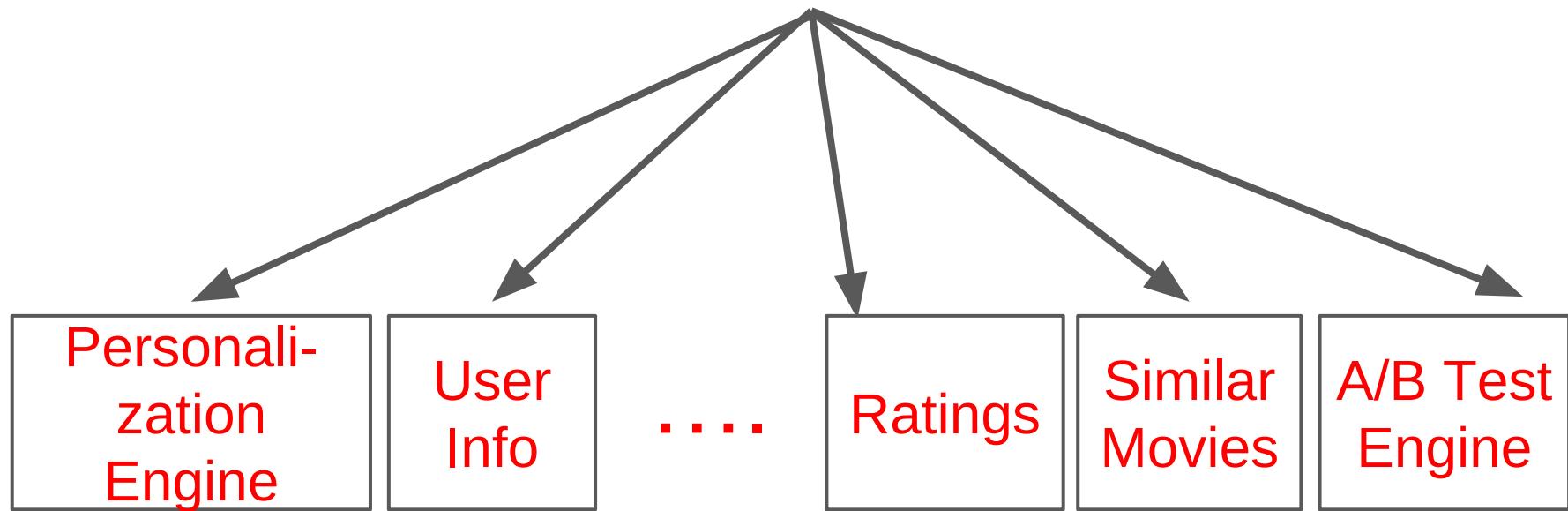


**HYSTRIX**  
DEFEND YOUR APP

- ❑ Protection from and control over latency and failure from dependencies
- ❑ Stop cascading failures in a complex distributed system
- ❑ Fall back and gracefully degrade
- ❑ Fail fast and rapidly recover

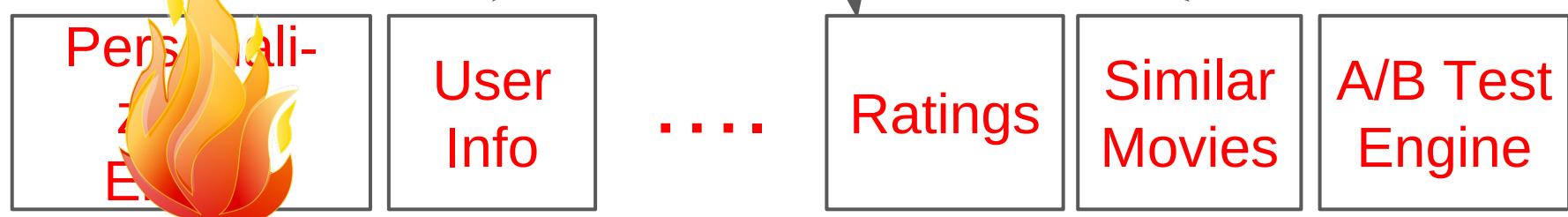
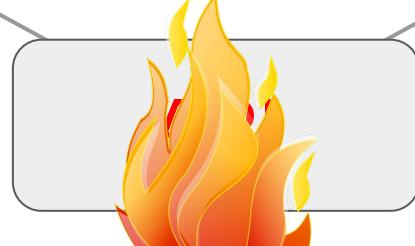


API



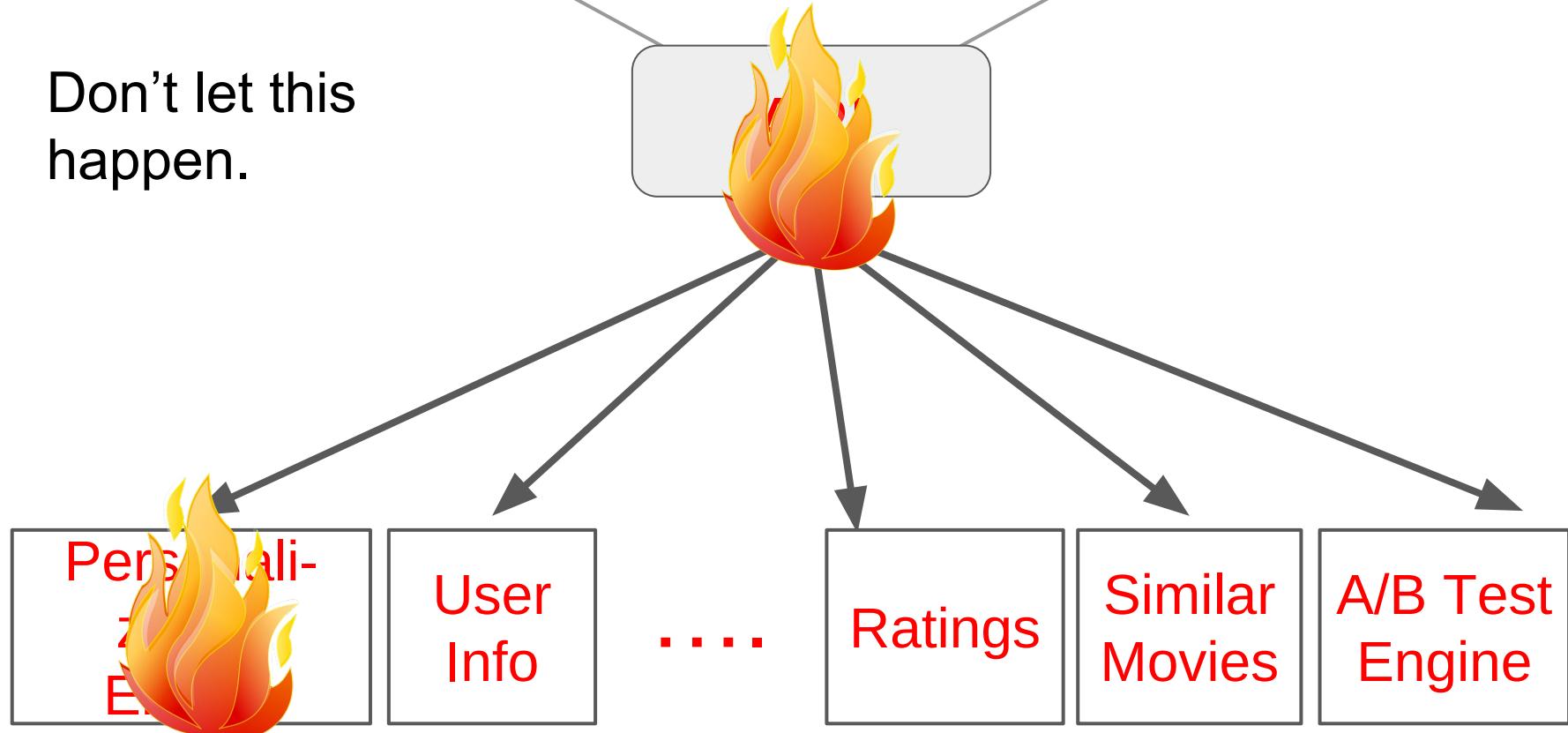


Don't let this  
happen.



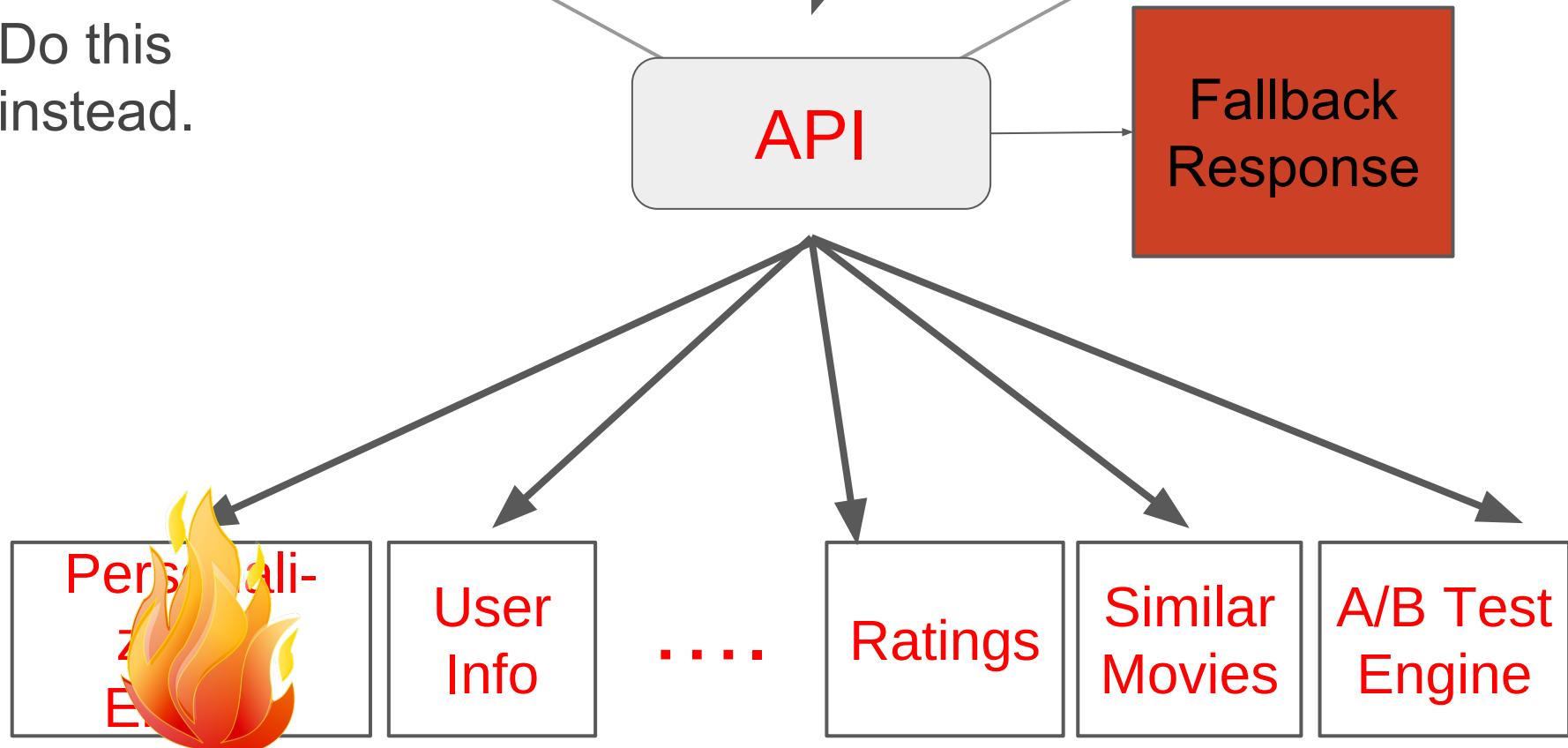


Don't let this  
happen.



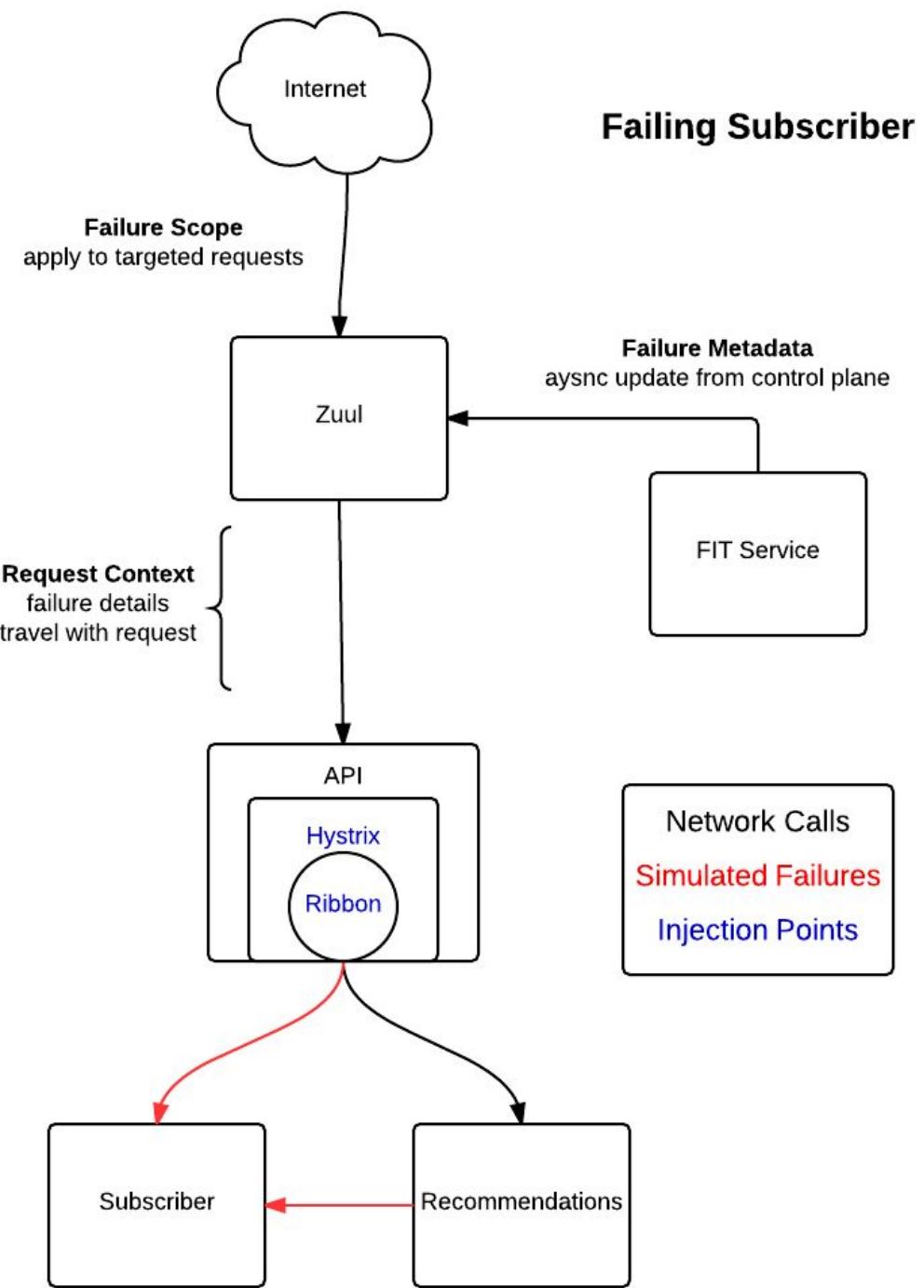


Do this  
instead.



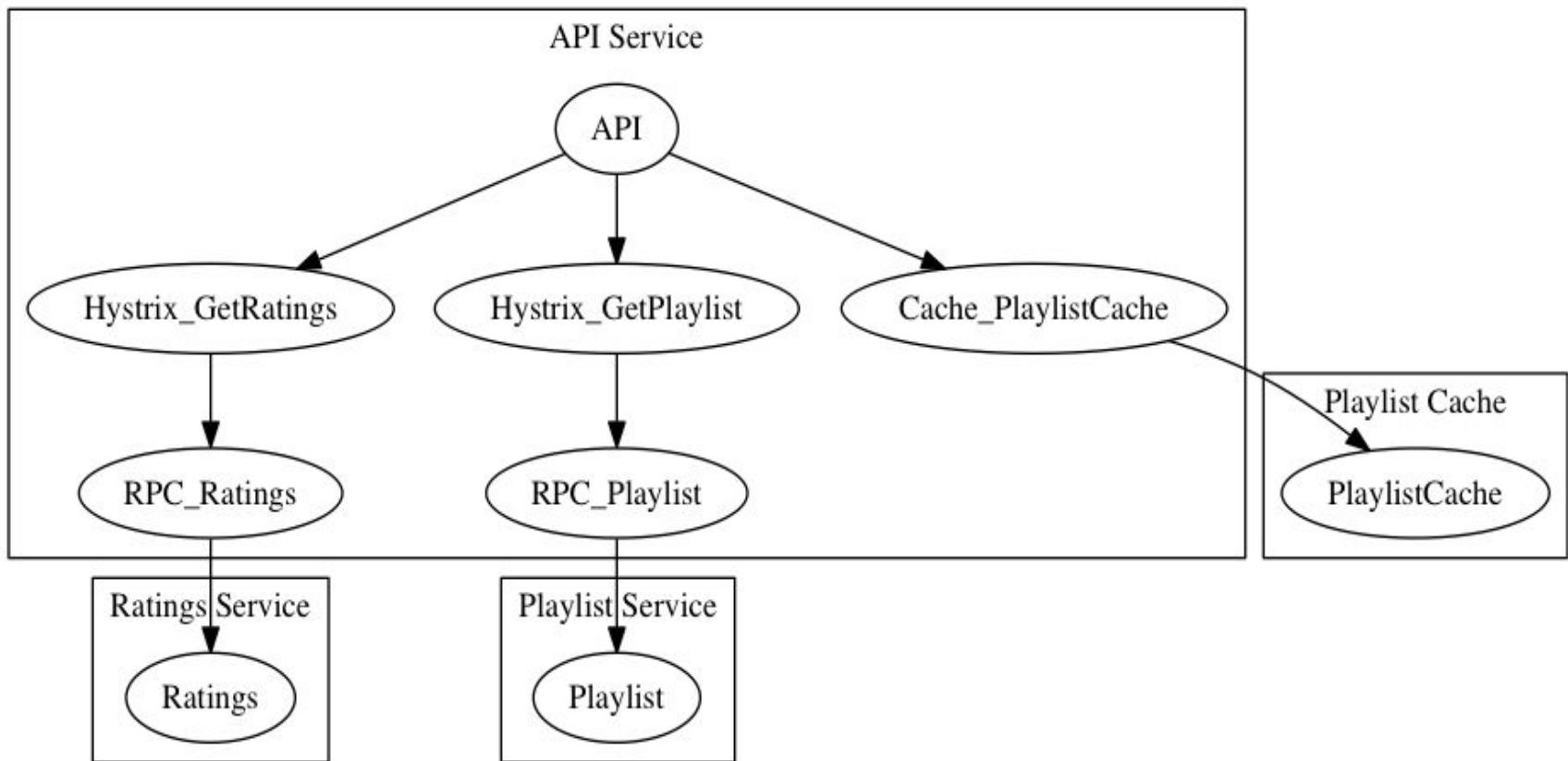
# Failure Injection Testing (FIT)

Goal: Study how the system behaves when a failures occur (e.g., backend service unreachable).



# More automated failure testing

Goal: Find groups of service calls that are needed for success.

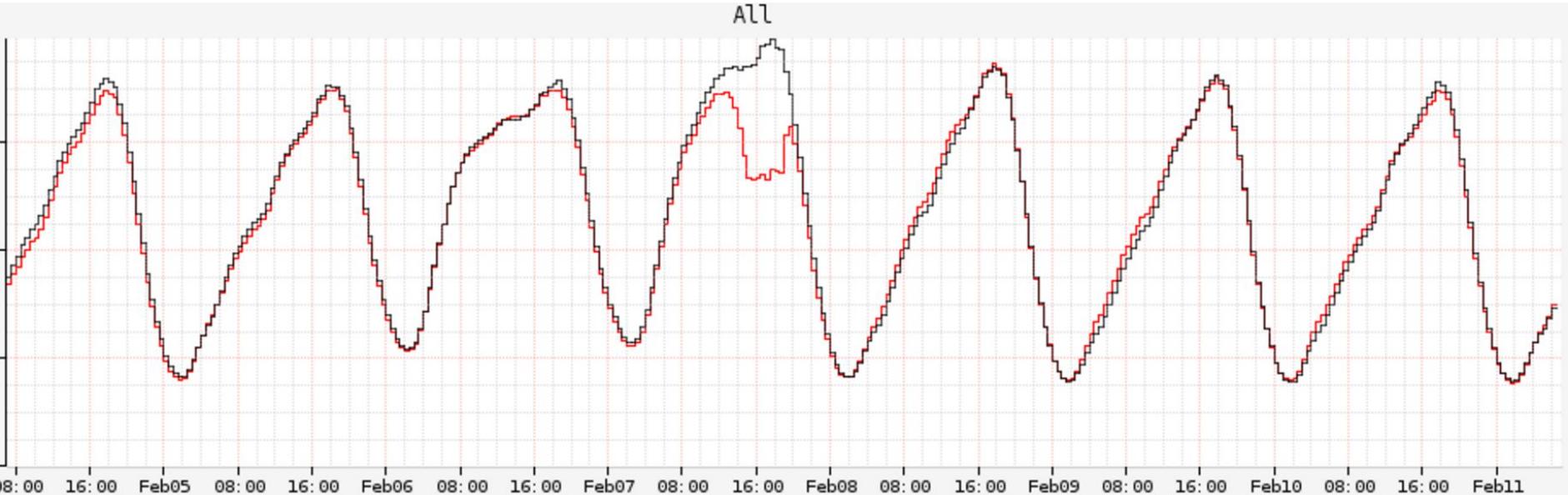


# Autoscaling & Capacity Management



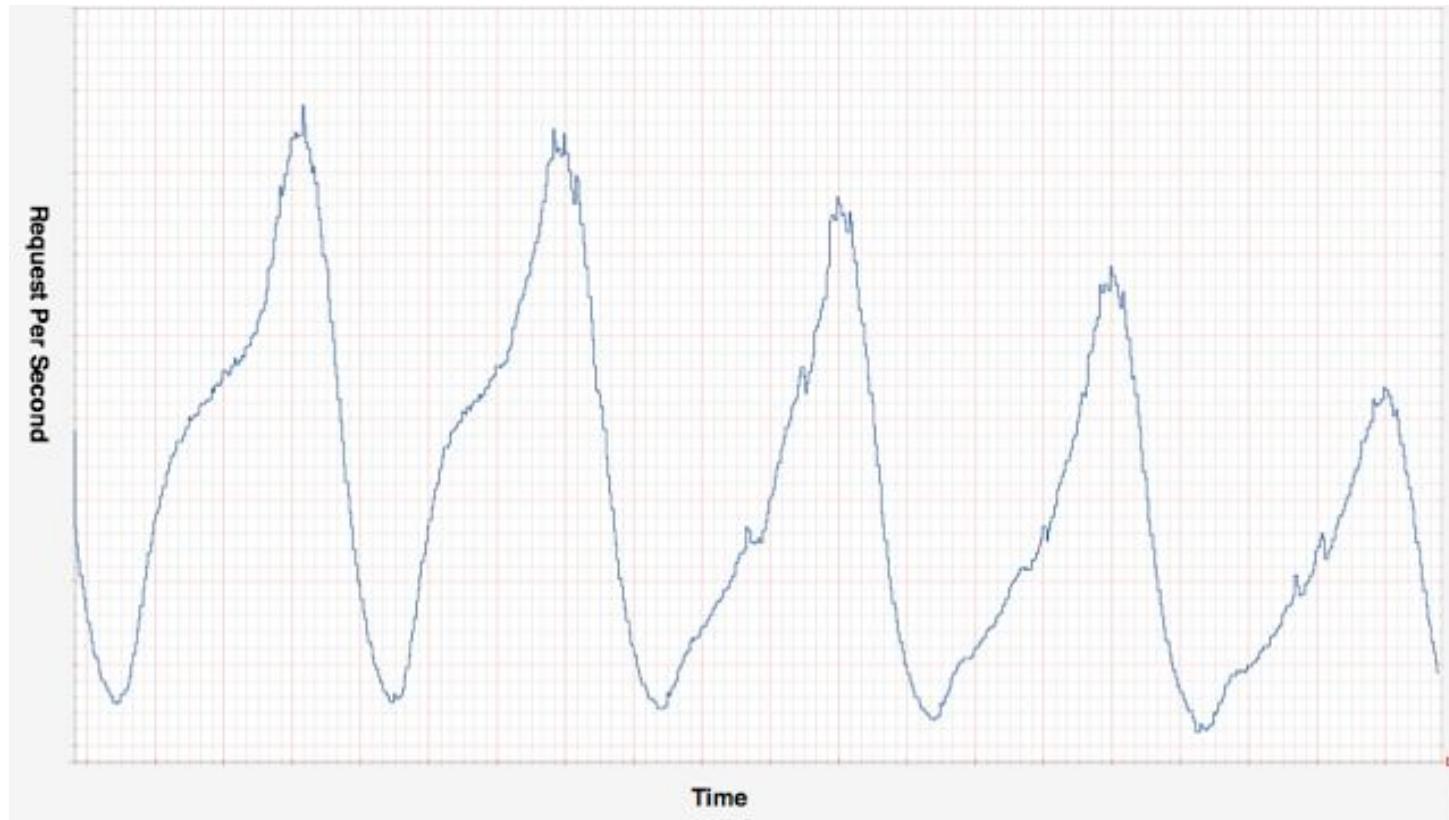
<http://nflx.it/1LvqLUi>

# Autoscaling & Capacity Management

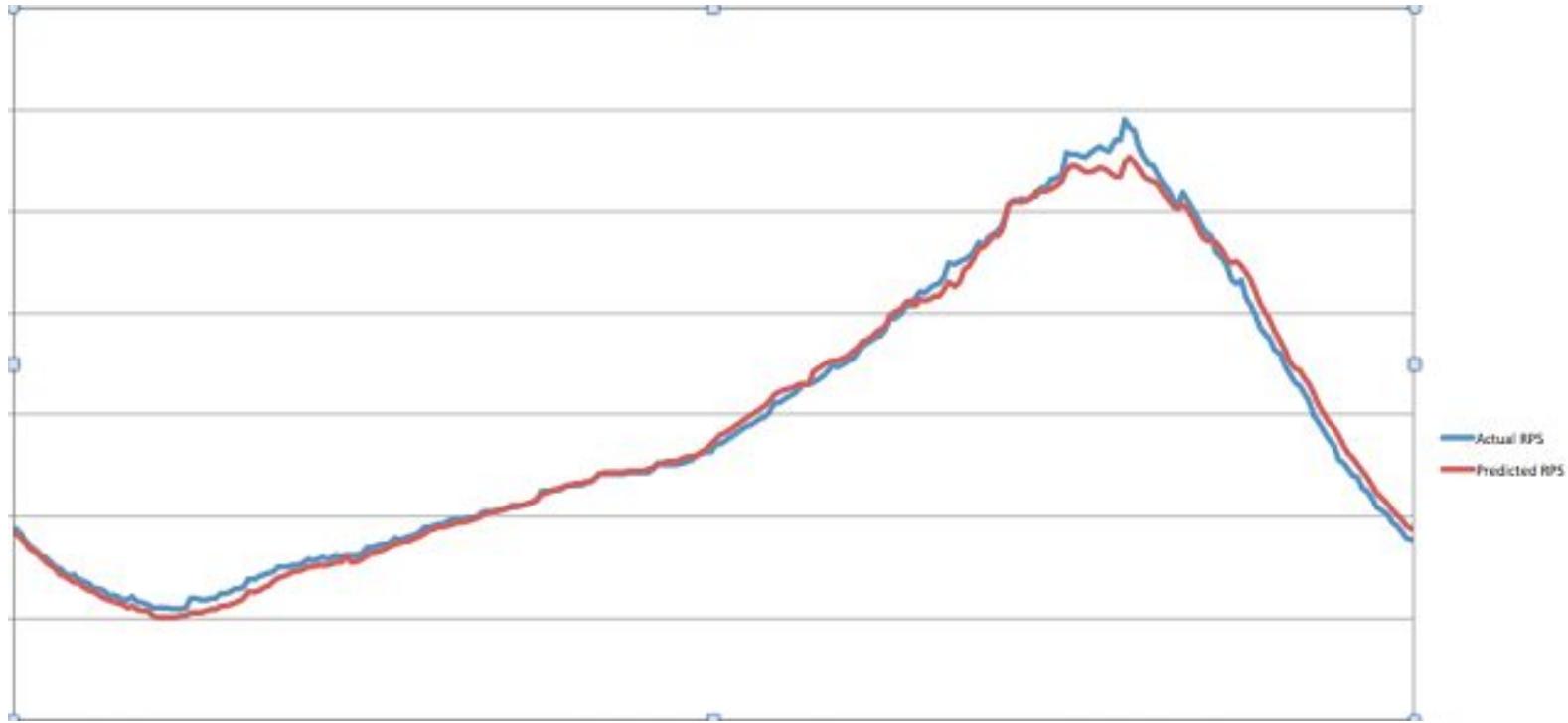


- ❑ Red: traffic for current week (x-axis)
- ❑ Black: traffic for previous week for comparison
- ❑ What happened on February 7? Superbowl!

# AWS Controls Reactive, does not scale up fast enough



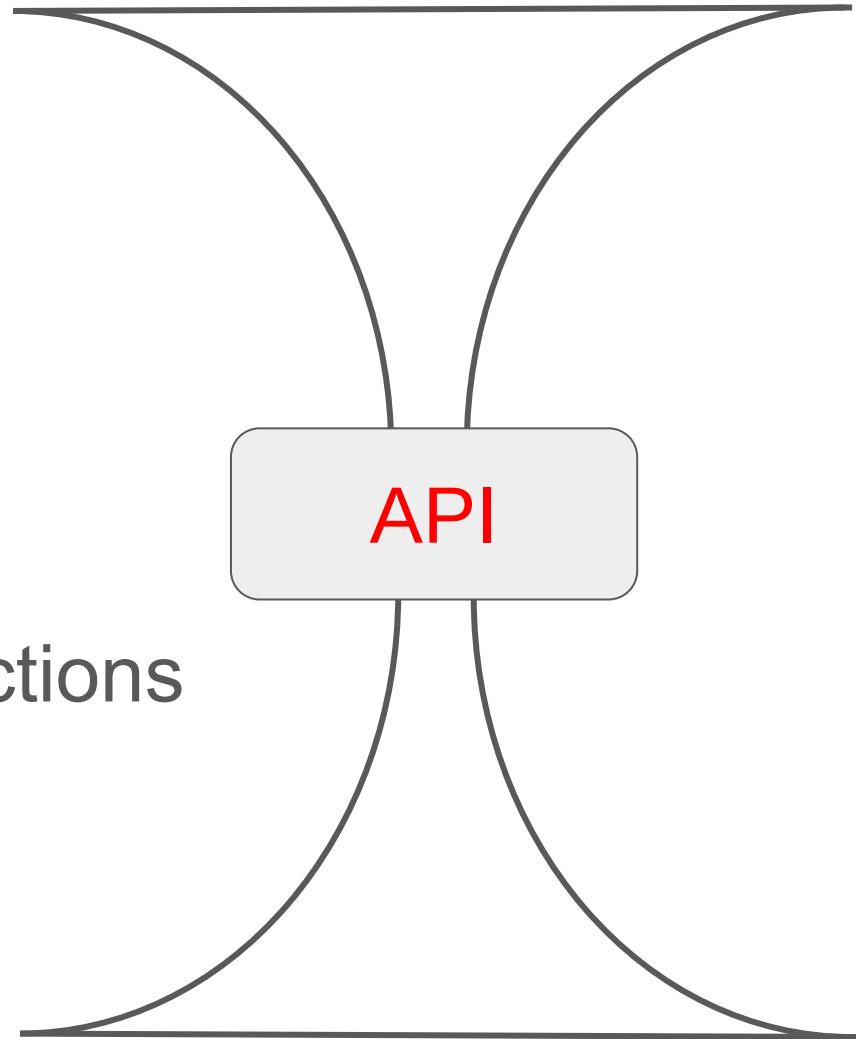
# Fine-grained Control with Scryer Complements AWS Controls



- ❑ Faster scale-up, improved cost
- ❑ Use reactive policy for organic scale down

# Netflix API

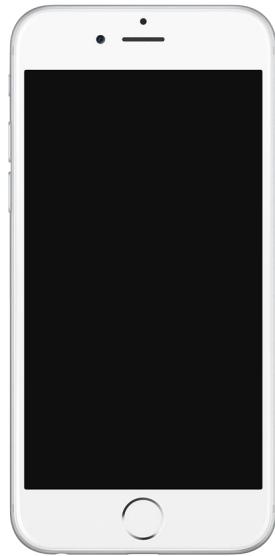
- ❑ Architecture
- ❑ Resiliency
- ❑ **Developer velocity**
- ❑ Tooling and DevOps
- ❑ Current and future directions



# Lots of devices, lots of variety

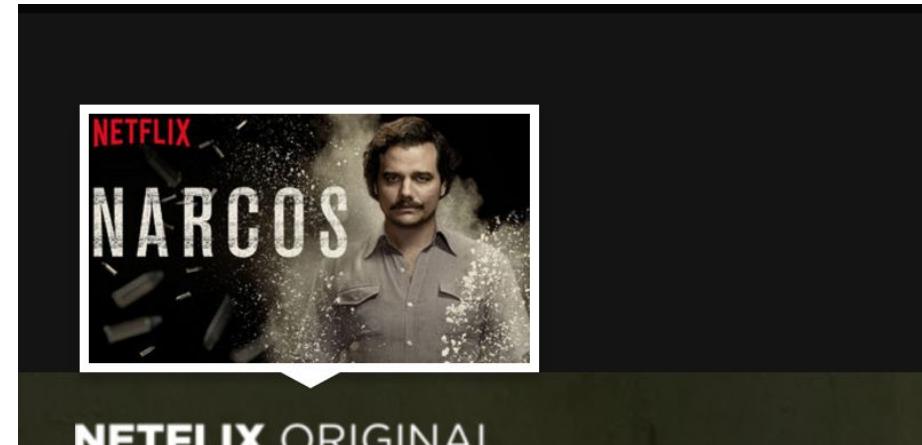


# Different interaction models



# And just to make things a little more interesting....

- A/B tests
- profiles
- localization



NETFLIX ORIGINAL  
**NARCOS**  
★★★★★ 2015 1 Staffel  
[Alle Folgen jetzt ansehen](#)  
Diese düstere neue Gangsterserie basiert auf der wahren Geschichte der berüchtigten, gewalttätigen und einflussreichen Drogenkartelle Kolumbiens.  
Mit: Wagner Moura, Boyd Holbrook, Pedro Pascal  
Genres: Serien, Action- und Abenteuerserien, Crime-/Krimiserien  
Diese Serie ist: Aufregend, Spannend, Rau, Düster

**MEINE LISTE**

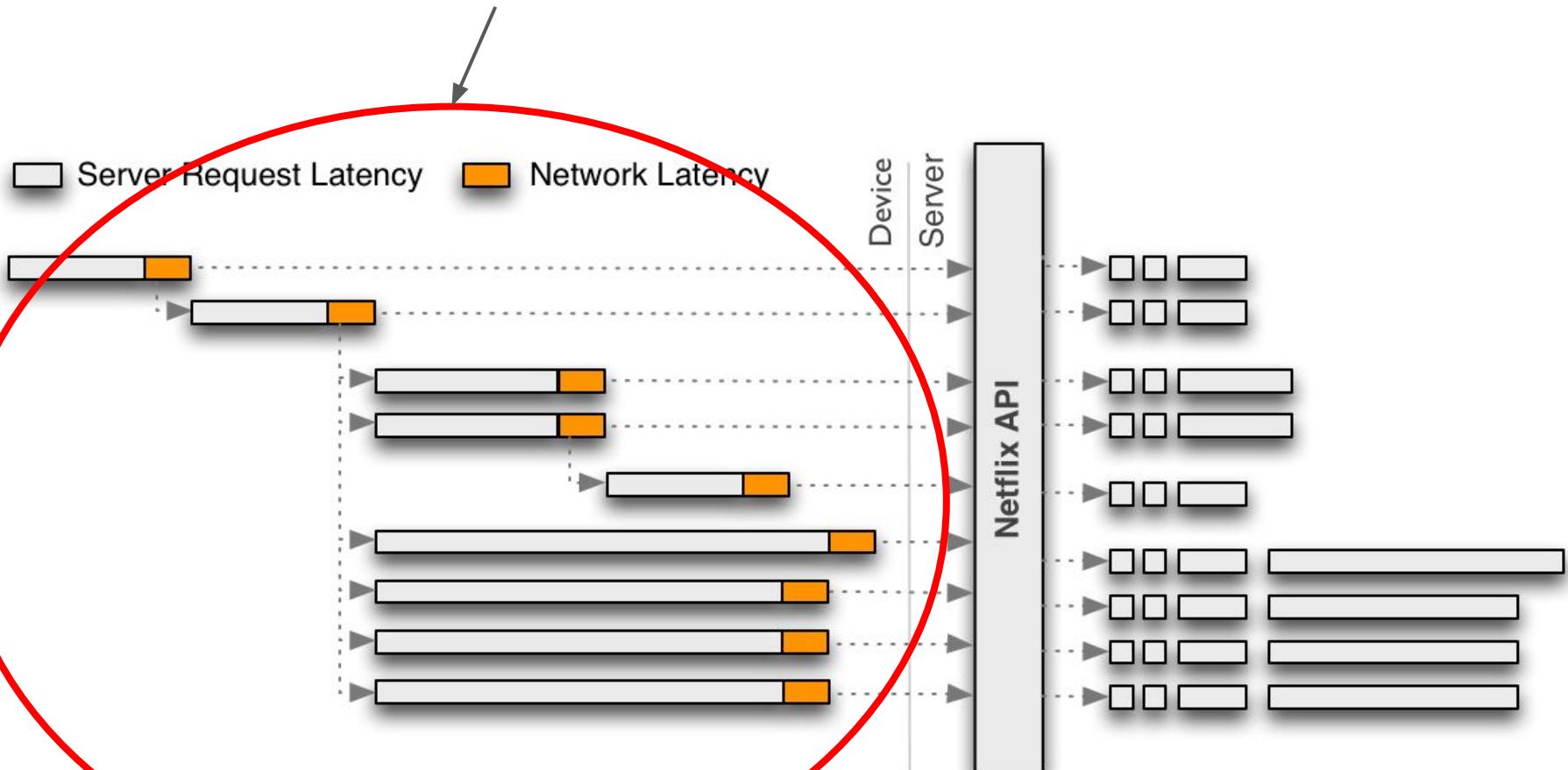
**ÜBERSICHT**

## Add server-side scripting capability

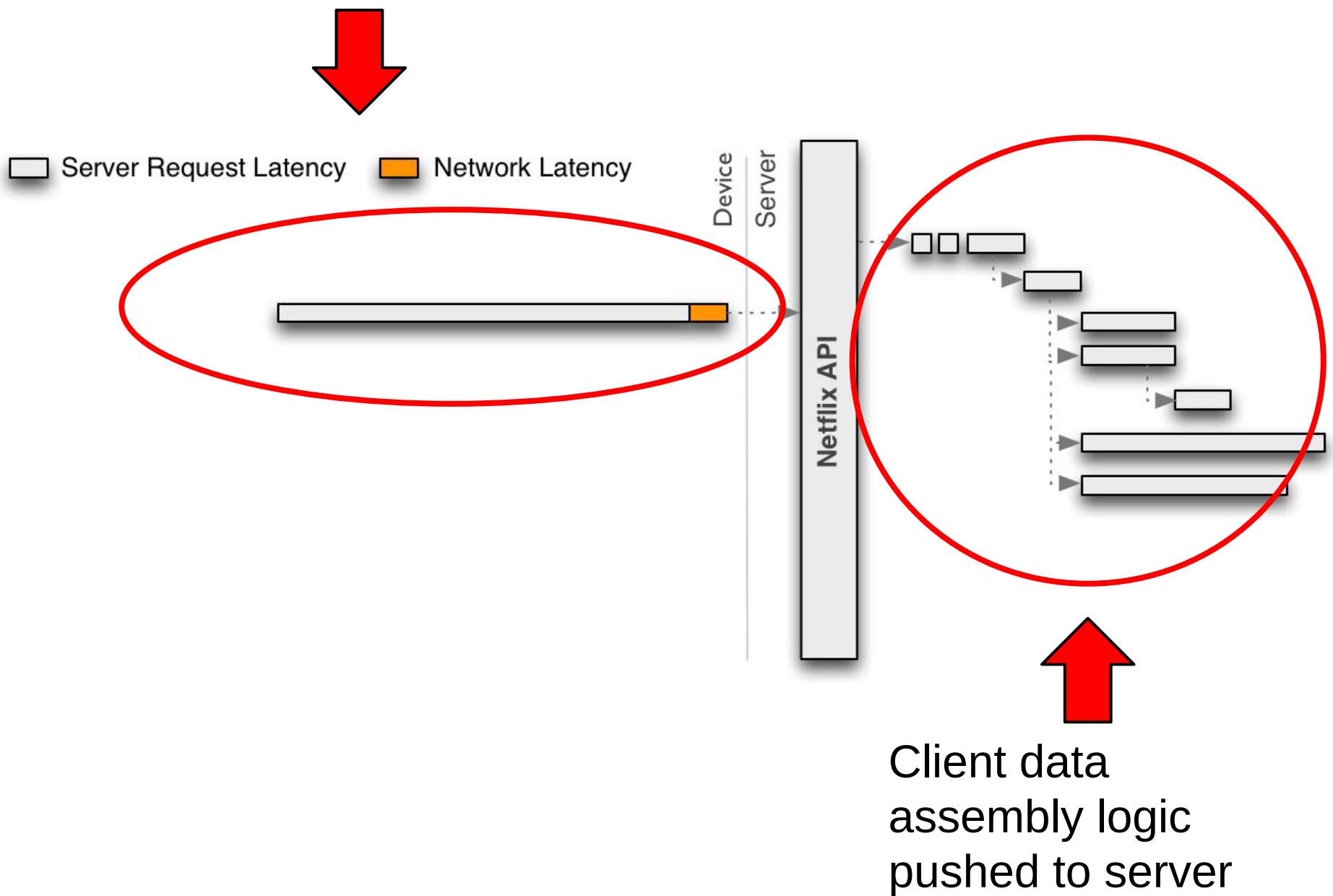


- ❑ Reduce network chattiness
- ❑ Support device optimizations
- ❑ Enable faster development for internal users

# Discrete HTTP requests pay network tax repeatedly



# Single, optimized request; pay network tax once



## Remote API

GET  
`/users/{user_id}/lists`

## Local Method

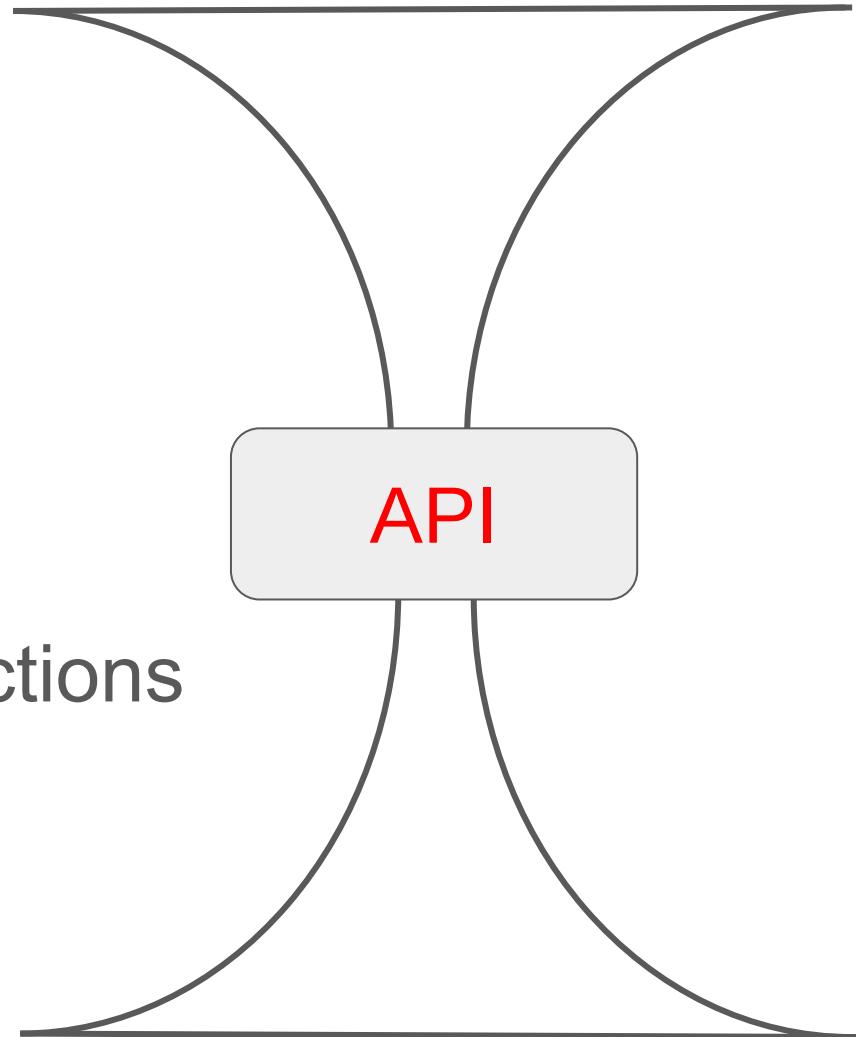
`getLists(userId)`

# Impact on velocity and collaboration

- ❑ UI (script) changes can happen independently
- ❑ Script changes can be pushed to running servers, so decoupled from API push schedule
- ❑ Decoupling leads to greater developer velocity

# Netflix API

- ❑ Architecture
- ❑ Resiliency
- ❑ Developer velocity
- ❑ **Tooling and DevOps**
- ❑ Current and future directions



Run 1% of your traffic on the new code and see how it does



# So you've run a canary. Now what?

Control

VS

Canary

- Errors: 2xx, 4xx, 5xx
- latency
- network
- busy threads
- load, memory consumption
- ...

# edge-server

Canary Score: 97.5% PASS

Report Date: [ ]  
Region: [ ]  
Canary: api-prod-autoprodpush-canary  
Type: cluster  
Duration: PT360M  
Version: [ ]

## Metric Analysis

Show  pass  low  high  nodata

Filter

Group name, metric

Group Name	metrics	Deviation	Score	
network		100%	<span style="background-color: #6aa84f; color: white; padding: 2px 10px;">PASS</span>	
logging		100%	<span style="background-color: #6aa84f; color: white; padding: 2px 10px;">PASS</span>	
system		100%	<span style="background-color: #6aa84f; color: white; padding: 2px 10px;">PASS</span>	
playback		100%	<span style="background-color: #6aa84f; color: white; padding: 2px 10px;">PASS</span>	
http-latency		100%	<span style="background-color: #6aa84f; color: white; padding: 2px 10px;">PASS</span>	
hystrix-latency		98%	<span style="background-color: #6aa84f; color: white; padding: 2px 10px;">PASS</span>	
gc		100%	<span style="background-color: #6aa84f; color: white; padding: 2px 10px;">PASS</span>	

Tag: system

Score: 100% PASS

## 10 Metrics

Group Name

processCpuLoad

freePhysicalMemorySize

Concurrent-Hystrix-Threads

tomcat-busy-threads

openFileDescriptorCount

atlas-metrics

systemCpuLoad

Concurrent-Http-Requests

systemLoadAverage

totalCompilationTime

# Successful canary



# red/black push

# Continuous Delivery with Spinnaker

api 

PIPELINES

CLUSTERS

LOAD BALANCERS

SECURITY GROUPS

PROPERTIES

TASKS

CONFIG

Filters  Clear All

SEARCH 

api-staging

ACCOUNT

- prod
- test

REGION



STACK

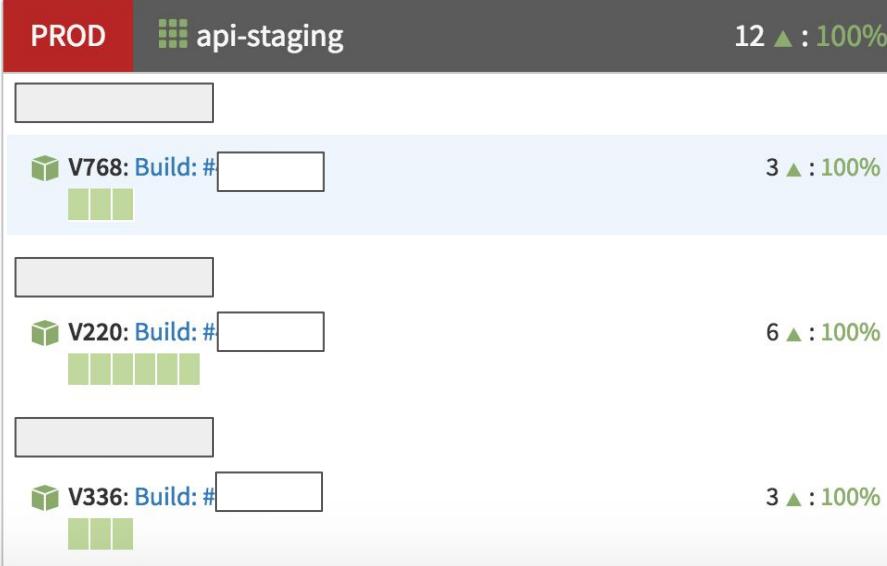
- ci
- content
- debug

## Clusters

Show  Instances  with details

 Create Server Group

Filtered by: SEARCH: api-staging  ACCOUNT: prod  ACCOUNT: test  Clear All



 api-staging-v

Server Group Actions 

Insight 

## SERVER GROUP INFORMATION

Created:

In:

VPC:

Availability ...

# Quickly see status of all clusters

**Filters** [Clear All](#) [Unpin](#) [+](#)

**SEARCH** [api-test](#)

**ACCOUNT**  prod  test

**REGION**

**STACK**  ci  content

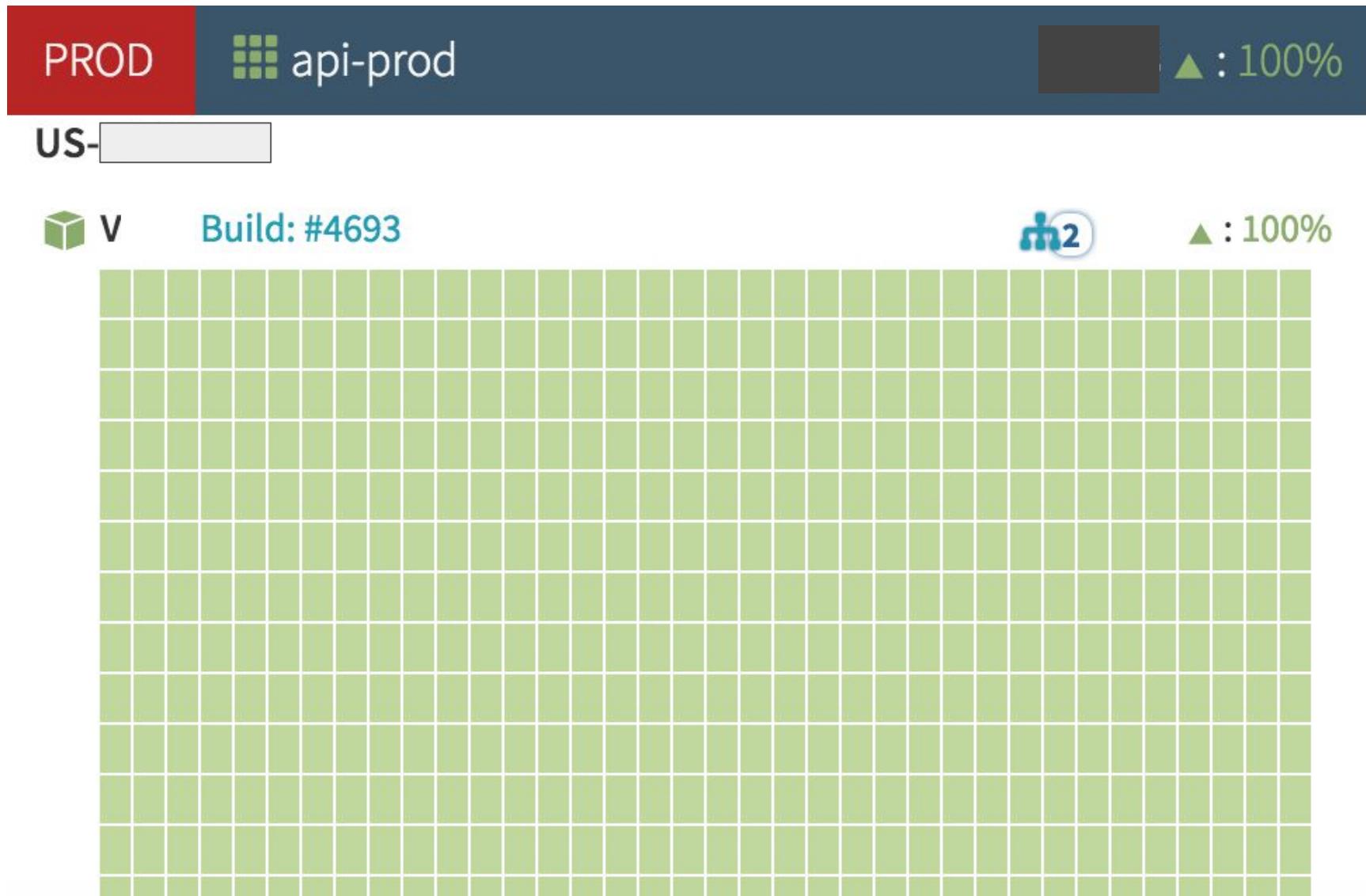
**Clusters** [Show Instances](#) [with details](#)

Filtered by: SEARCH: api-test ACCOUNT: test [Clear All](#)

TEST	api-test	31 ▲ : 100%
V495: Build: #	3 ▲ : 100%	
V039: Build: #	25 ▲ : 100%	
V867: Build: #		3 ▲ : 100%

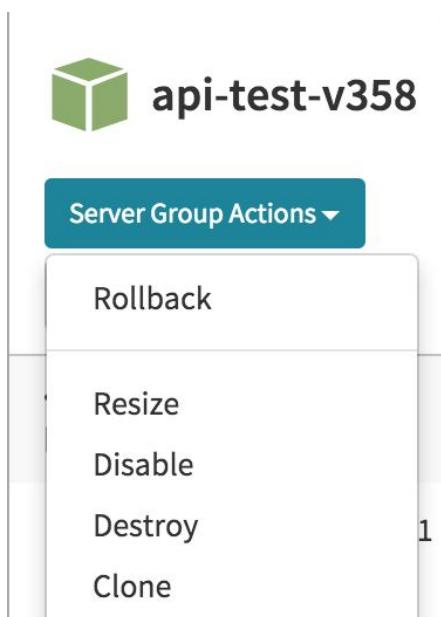
[Permalink](#)

# Prod is a little different....

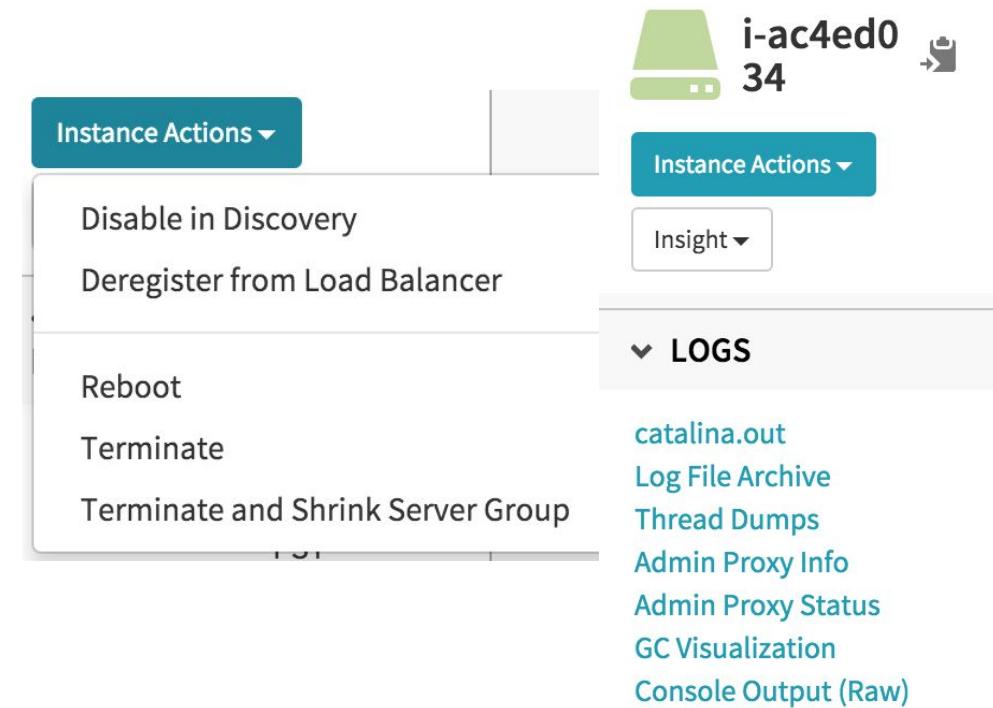


# The things you can do

... with server groups



... with instances



## ENDPOINTS

## Overview

## Manage

## Health

## Activity

## RESOURCES

## .Next SDK

## .Next Javadoc

## .Next User Guide

## .Next REPL (PROD)

## Deployments

## Environments

Env

PROD

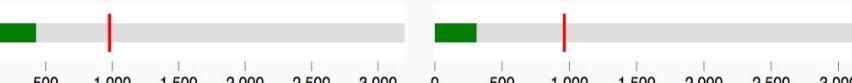
STAGING

INT

## Endpoints

## Usage

Active



Endpoints,

Total

Endpoints,

Active Limit

## Recent

## Activity

Show/Hide

2014-06-20 21:42:13 PDT

/mobile/s

Revision 9 activated.

2014-06-20 17:32:26 PDT

/apps/

Revision 22 activated.

2014-06-20 21:42:13 PDT

/mobile/

Revision 9 activated.

2014-06-20 21:42:03 PDT

/mobile/staging

Revision 12 activated.

2014-06-20 17:04:19 PDT

/tvui/lgm/-

Was deactivated.

2014-06-20 21:42:03 PDT

/mobile

Revision 12 activated.

/android/samurai/config

Revision 55 activated.

/ecapi/acco

Revision 4 activated.

/android/

Revision 55 activated.

2014-06-20 15:16:56 PDT

/mobile/

Revision 24 activated.

2014-06-20 15:55:03 PDT

/ecapi/acc

Revision 10 activated.

2014-06-20 15:16:56 PDT

/mobile

Revision 24 activated.

/mobile/ti

Revision 26 activated.

/ecapi/acc

Revision 11 activated.

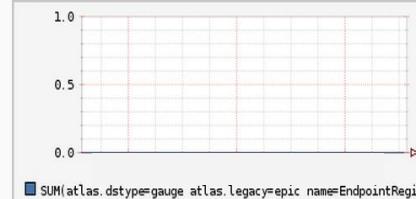
/mobile

Revision 26 activated.

## Compilation

## Failures

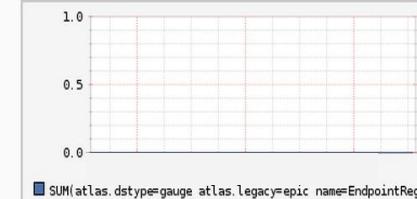
## Atlas Link



## Atlas Link



## Atlas Link



## ENDPOINTS

Overview

Manage

Health

Activity

RESOURCES

.Next SDK

.Next Javadoc

.Next User Guide

.Next REPL (PROD)

Deployments

Environments

## Endpoint Groups



Group	# Endpoints	# Active	Usage	Health	# Requests (e-2w)	Bytecode Size	Owner
tvui	459	69	76%			184.4 MB	tflix.com
mobile	53	26	52%			20.9 MB	e@netflix.com
android	25	16	80%			27.8 MB	flix.com
desktop	73	63	70%			6.0 MB	s@netflix.com
win	35	34	68%			52.1 MB	i@netflix.com
users	5	5	50%			8.9 MB	s@netflix.com
atv	31	6	15%			9.8 MB	tflix.com
halo	1	1	20%			2.2 MB	Dnetflix.com
streaming	2	2	40%			58.9 kB	Dnetflix.com
shakti	154	72	72%			2.9 MB	x.com
cbp	10	10	83%			3.3 MB	s@netflix.com
apps	4	3	60%			405.1 kB	Dnetflix.com
ecapi	46	40	66%			1.1 MB	s@netflix.com
ios	2	2	10%			1.3 MB	Dnetflix.com
account	5	5	50%			69.3 kB	s@netflix.com
internal	23	20	66%			1.4 MB	s@netflix.com
signup	4	4	80%			204.5 kB	Dnetflix.com
jellyvision	13	12	60%			463.4 kB	tflix.com
utility	9	9	60%			747.8 kB	Dnetflix.com
dse	1	1	20%			74.3 kB	Dnetflix.com
social	6	6	60%			265.1 kB	Dnetflix.com

# Operations

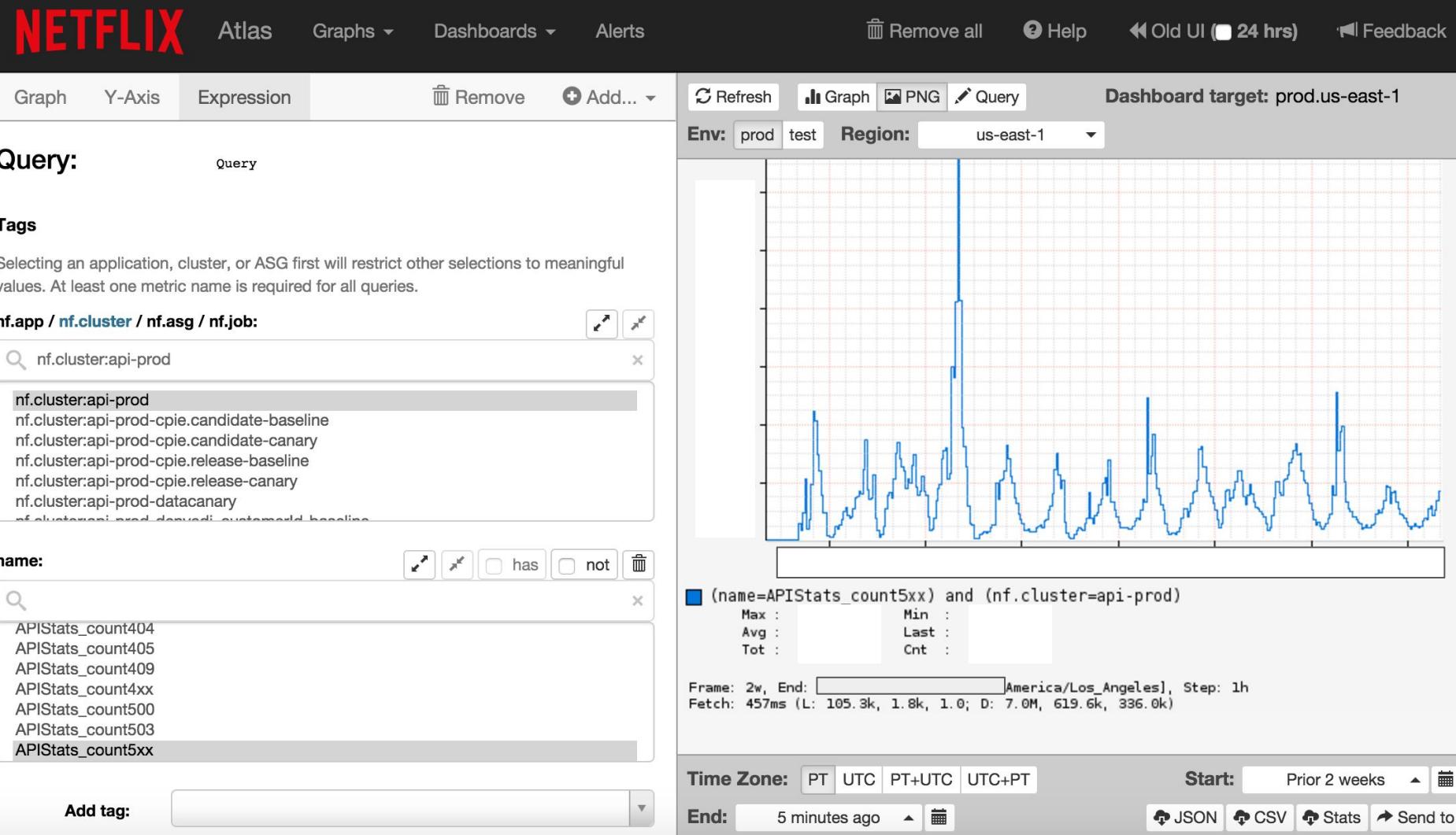
EDGE CR DASH INFRASTRUCTURE SERVER CALLS SERVICES THREAD POOLS WORLD MAP ZUUL

No Facets Selected +

CLUSTER STATUS (33) SORT: RPS Filter by...

Cluster\ASG	RPS Trend (5 min)	RPS	Error %	Per Instance RPS	Load Average	Requests Throttled	Latency (ms)	Discovery Instances	AWS Status
api-prod api - us-east-1					5.1	0.06			
api-prod-v EDGE-Master-Family-Build #					5.1	0.06			

# Operations



# Operations

DASHBOARDS    New    List    About    Help

## Edge Services Dashboard

VIEW

EDIT

Target

Search



API Overview

Refresh All

Auto Refresh

Show Legend

Logarithmic

Start

Last 3 hours



End

Minus 1 min



Shift

None

Step

Auto

Time Zone

US/Pacific

Zuul Overview

API Overview

NCCP Overview

Zuul Overviews per Cluster

Zuul Load per Cluster

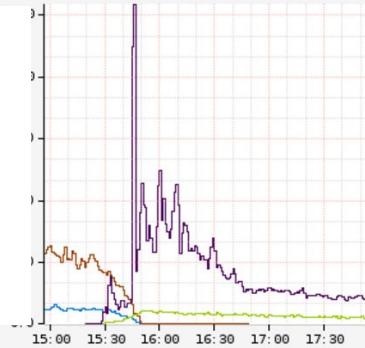
Zuul Latencies per Cluster

Zuul Requests by Origin

API Throttling

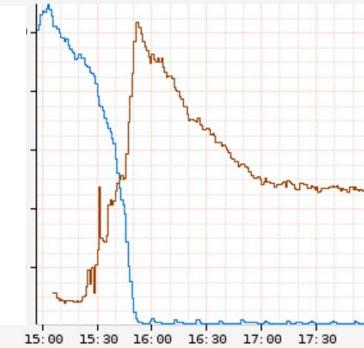
PROD:GLOBAL

4XX, 5XX Status by ASG: eu...



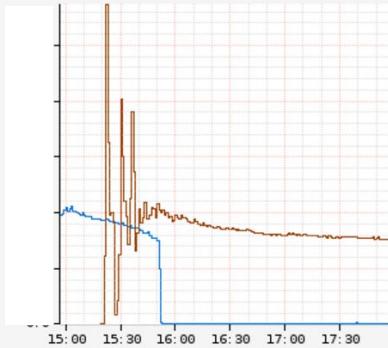
PROD:GLOBAL

Load Avg by ASG: eu-west-1



PROD:GLOBAL

Latency 90pct by ASG: eu-west-1



# Real-time analysis

EDX<sup>beta</sup>

Insights



kathrin

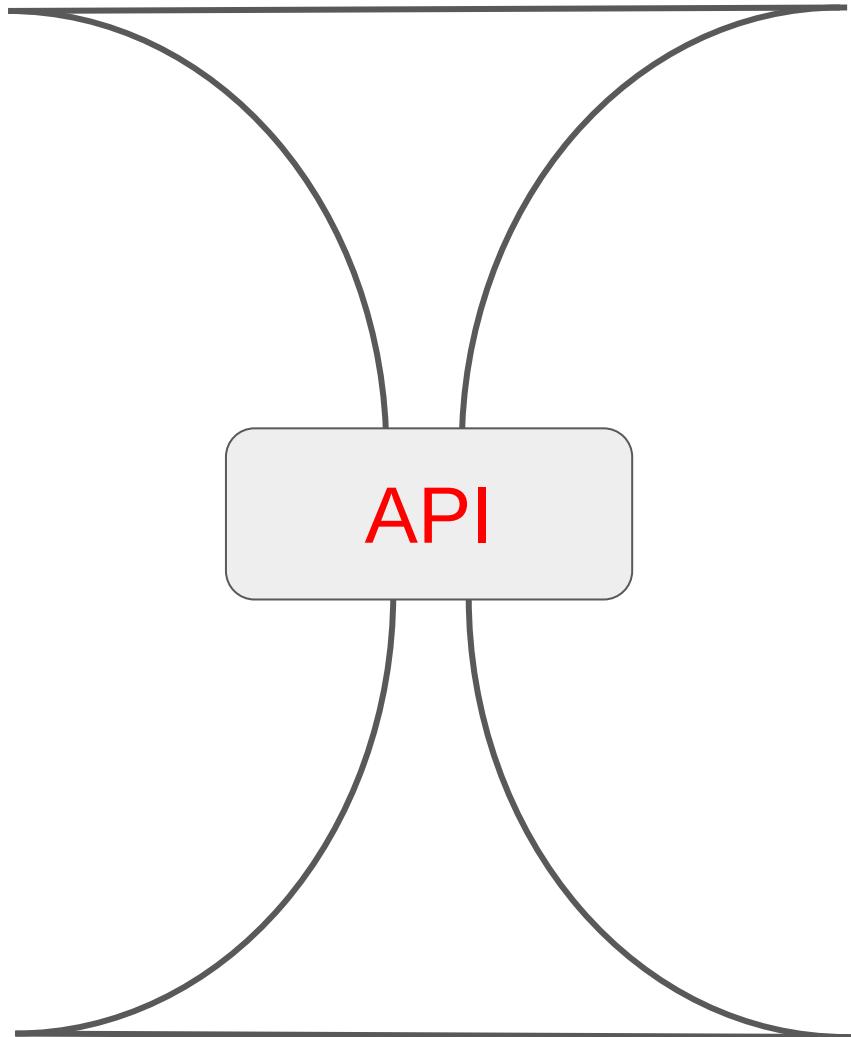
```
{"currentTime":1446582396522,"path":"/shakti/738b97b7/pathEvalua
{"currentTime":1446582395943,"path":"/shakti/738b97b7/dynamicMes
 {"currentTime":1446582395307,"path":"/shakti/738b97b7/billboardi
 {"currentTime":1446582393388,"path":"/shakti/738b97b7/pathEvalua
 {"currentTime":1446582392453,"path":"/shakti/738b97b7/pathEvalua
 {"currentTime":1446582390916,"path":"/shakti/738b97b7/pathEvalua
 {"currentTime":1446582384252,"path":"/shakti/738b97b7/dynamicMes
 {"currentTime":1446582383364,"path":"/shakti/738b97b7/dynamicMes
 {"currentTime":1446582382383,"path":"/shakti/738b97b7/dynamicMes
 {"currentTime":1446582381639,"path":"/shakti/738b97b7/dynamicMes
 {"currentTime":1446582380789,"path":"/shakti/738b97b7/billboardi
 {"currentTime":1446582379352,"path":"/shakti/738b97b7/dynamicMes
 {"currentTime":1446582378354,"path":"/shakti/738b97b7/dynamicMes
 {"currentTime":1446582377811,"path":"/shakti/738b97b7/signupCont
 {"currentTime":1446582375772,"path":"/shakti/738b97b7/fakiraHome
 {"currentTime":1446582374406,"path":"/shakti/738b97b7/dynamicMes
 {"currentTime":1446582361156,"path":"/shakti/738b97b7/pathEvalua
 {"currentTime":1446582360906,"path":"/shakti/738b97b7/pathEvalua

{
  "currentTime": 1446582396522,
  "path": "/shakti/738b97b7/pathEvaluator",
  "esn": "<redacted>",
  "method": "POST",
  "duration": 55,
  "status": 200,
  "region": "<redacted>",
  "index": 16225,
  "fullPayload": "<debug info redacted>"
}
```

Submit a query, see requests in real time.

# Netflix API

- ❑ Architecture
- ❑ Resiliency
- ❑ Developer velocity
- ❑ Tooling and DevOps
- ❑ **Current and future directions**



# What we've grown to

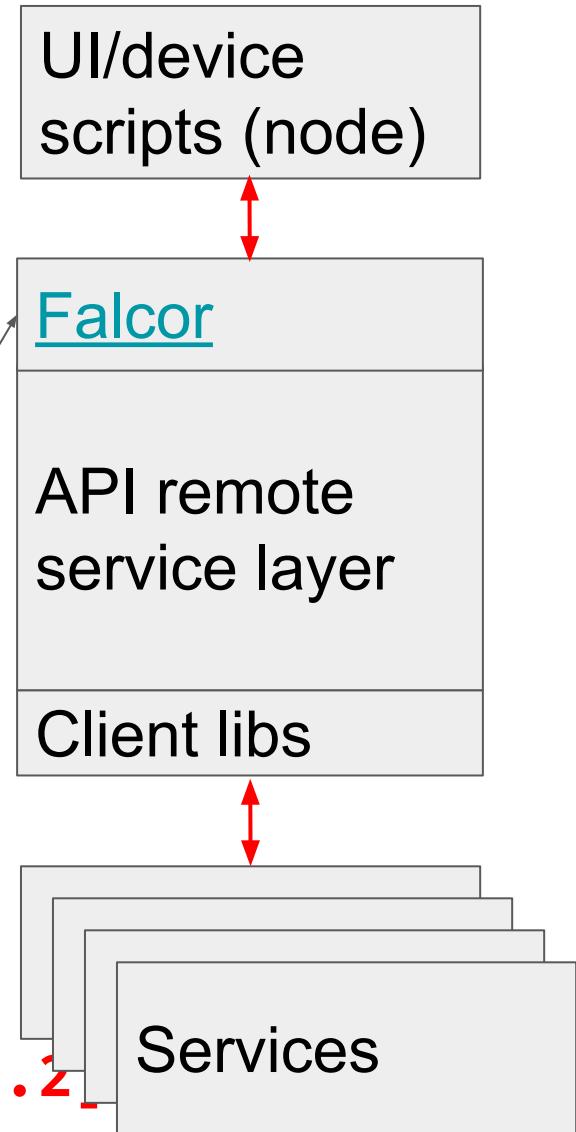
- > 900 active endpoints
- ~60 direct dependencies
- 78 thread pools
- 1000+ threads
- high memory usage



# Script isolation & node

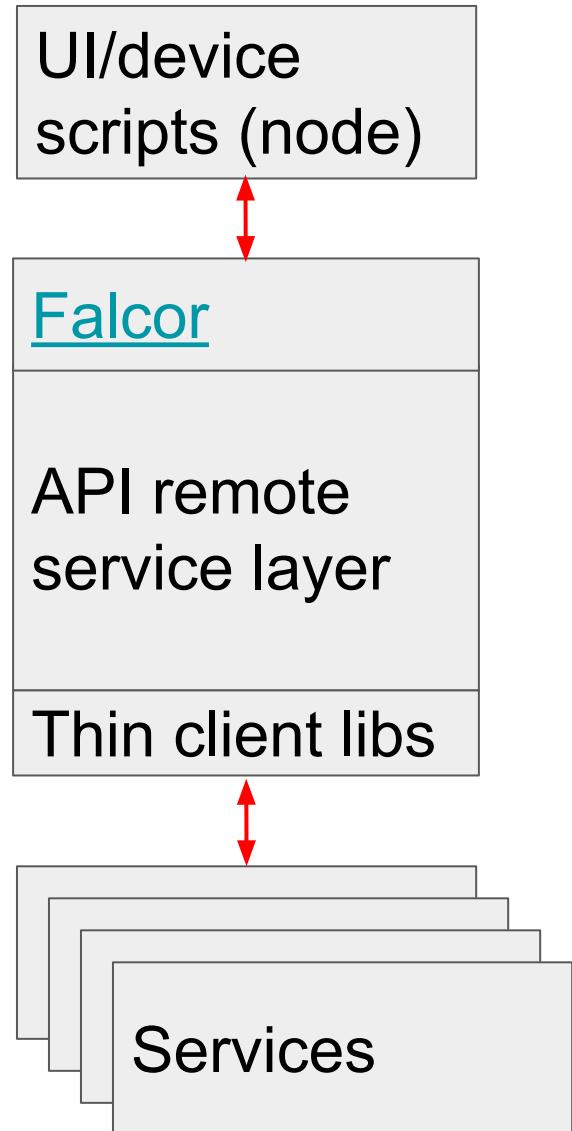
- ❑ Groovy scripts run as part of API process
- ❑ UI teams would like to use other languages (in particular node.js)

```
var response = model.get("todos[0..2]  
['name', 'done']");
```



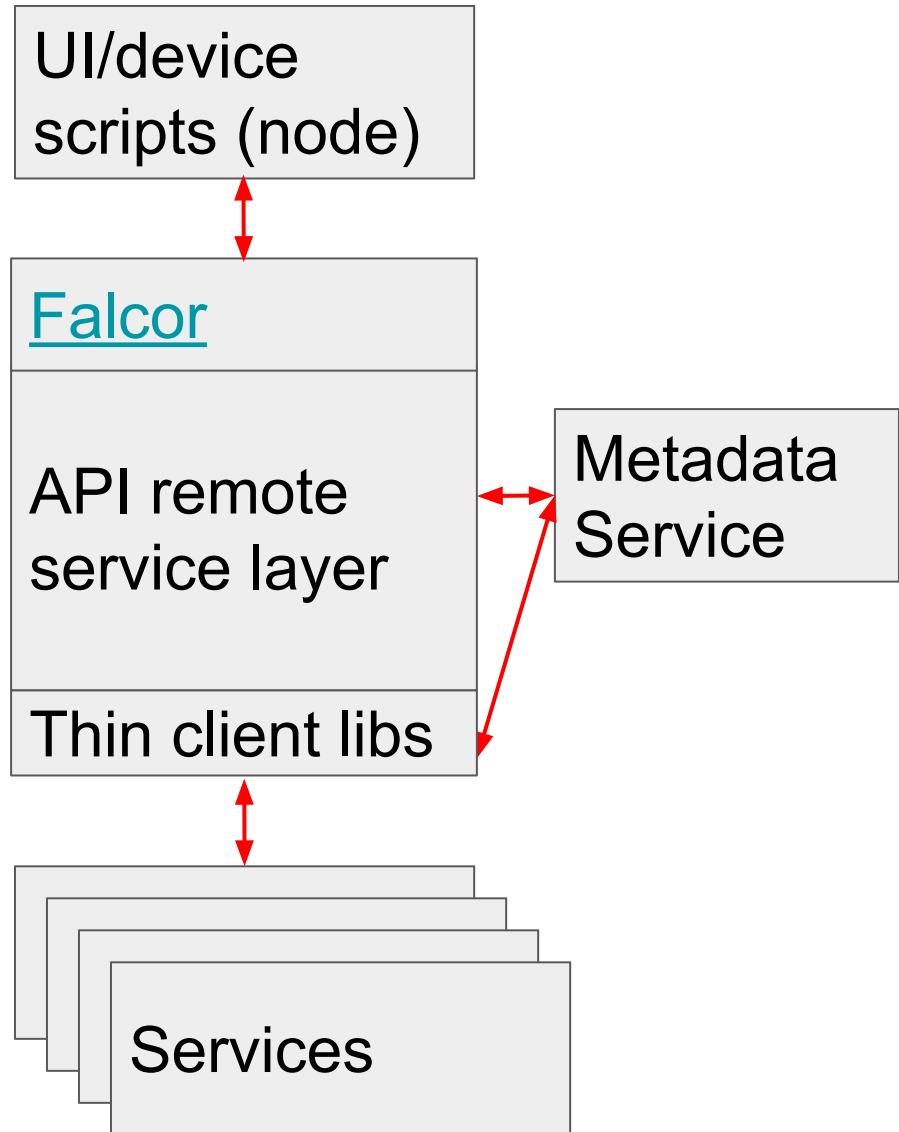
# Thin client libraries

- ❑ Fat client libraries
  - ❑ business logic and have
    - ❑ multiple dependencies
  - ❑ Move business logic and dependencies to services



# Remove metadata from API servers

- ❑ Metadata takes up significant memory in API servers
- ❑ Challenge: reduce chattiness to metadata



Thank you