# Data stores at Netflix

# Agenda

Introduction

High-level architecture

Workload characterization

Observability

A day in the life of EVCache

Incidents and takeaways

# EVCache

- Distributed, sharded, replicated key-value store
- Based on Memcached
- Tunable in-Region and global replication
- Resilient to failures
- Topology aware
- Linearly scalable
- Seamless deployments

# EVCache

Two variants of Memcached

- In-memory

- NVMe storage*

\* https://netflixtechblog.com/application-data-caching-using-ssds-5bf25df851ef

# EVCache footprint at Netflix

3 Regions

4 engineers

~160 clusters

~18,000 servers

~15,000,000 replications/sec

~350,000,000 ops/sec

~1,400,000,000,000 items

~14,000,000,000,000,000 bytes

# Agenda

Introduction
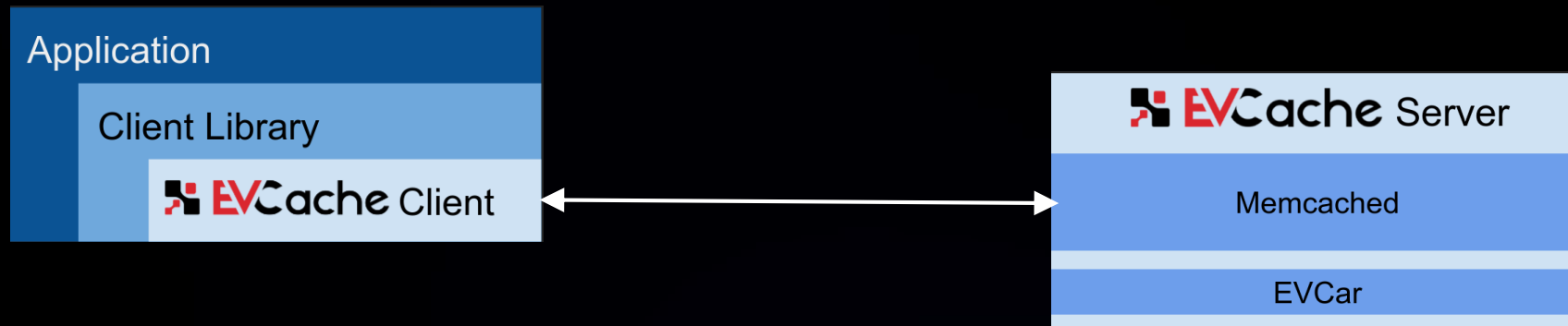
<span style="color:red">High-level architecture</span>

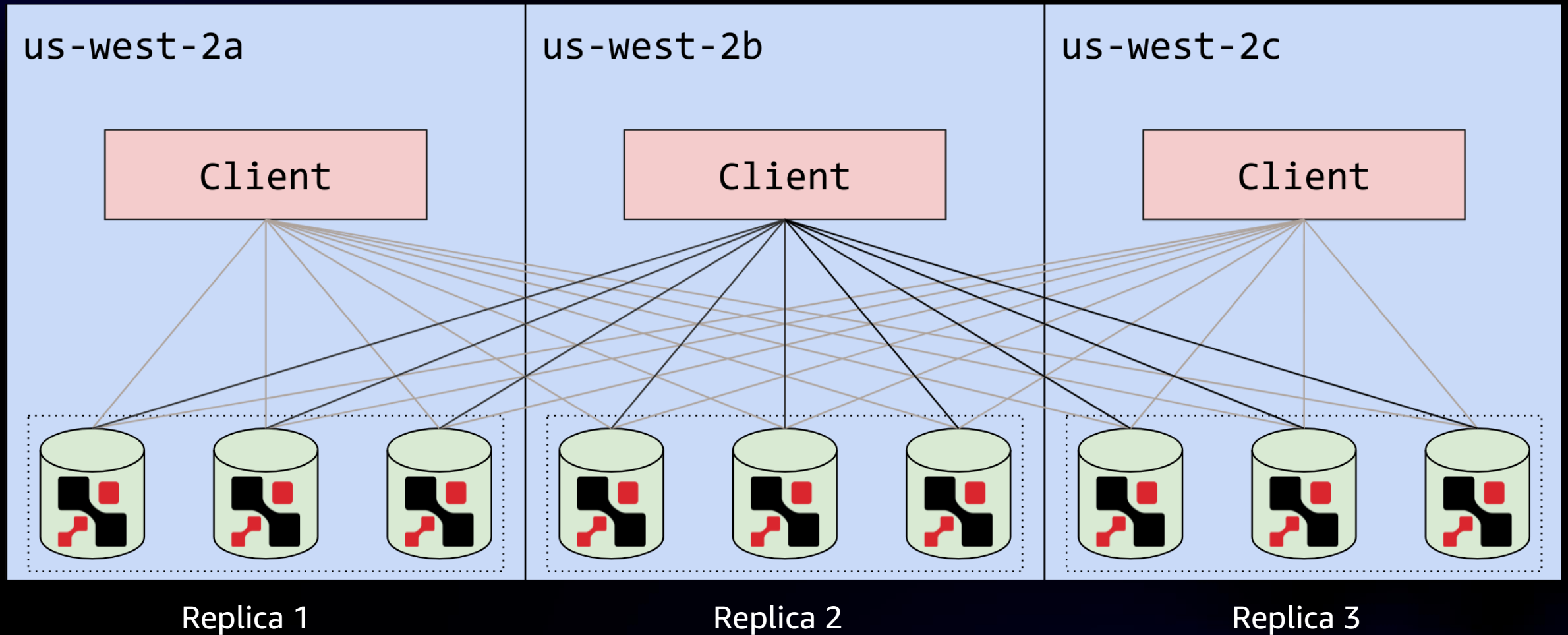Workload characterization

Observability

A day in the life of EVCache

Incidents and takeaways

# High-level architecture

# High-level architecture

# High-level architecture

# Netflix home page

# Request breakdown

# Agenda

Introduction

High-level architecture

Workload characterization

Observability
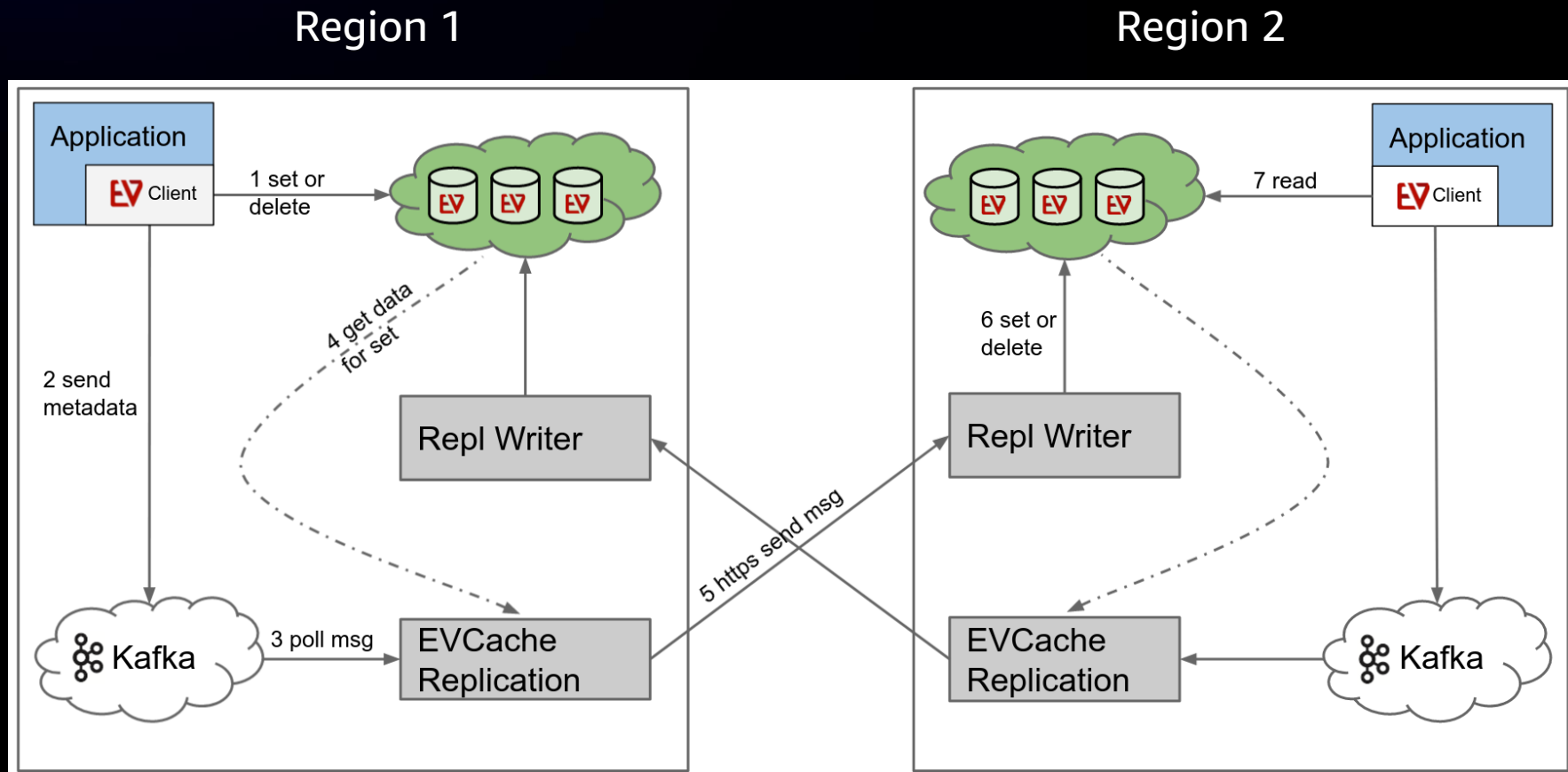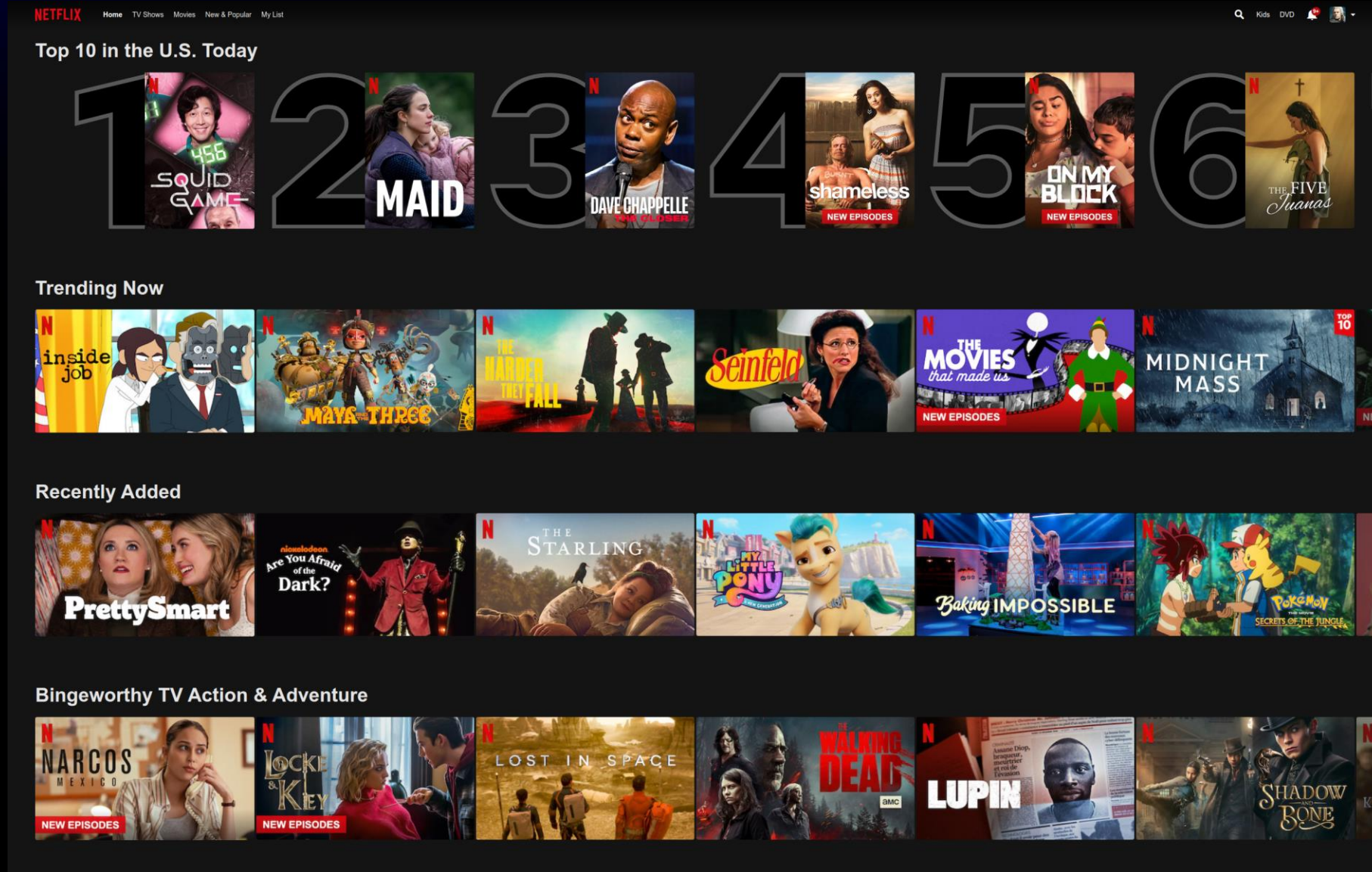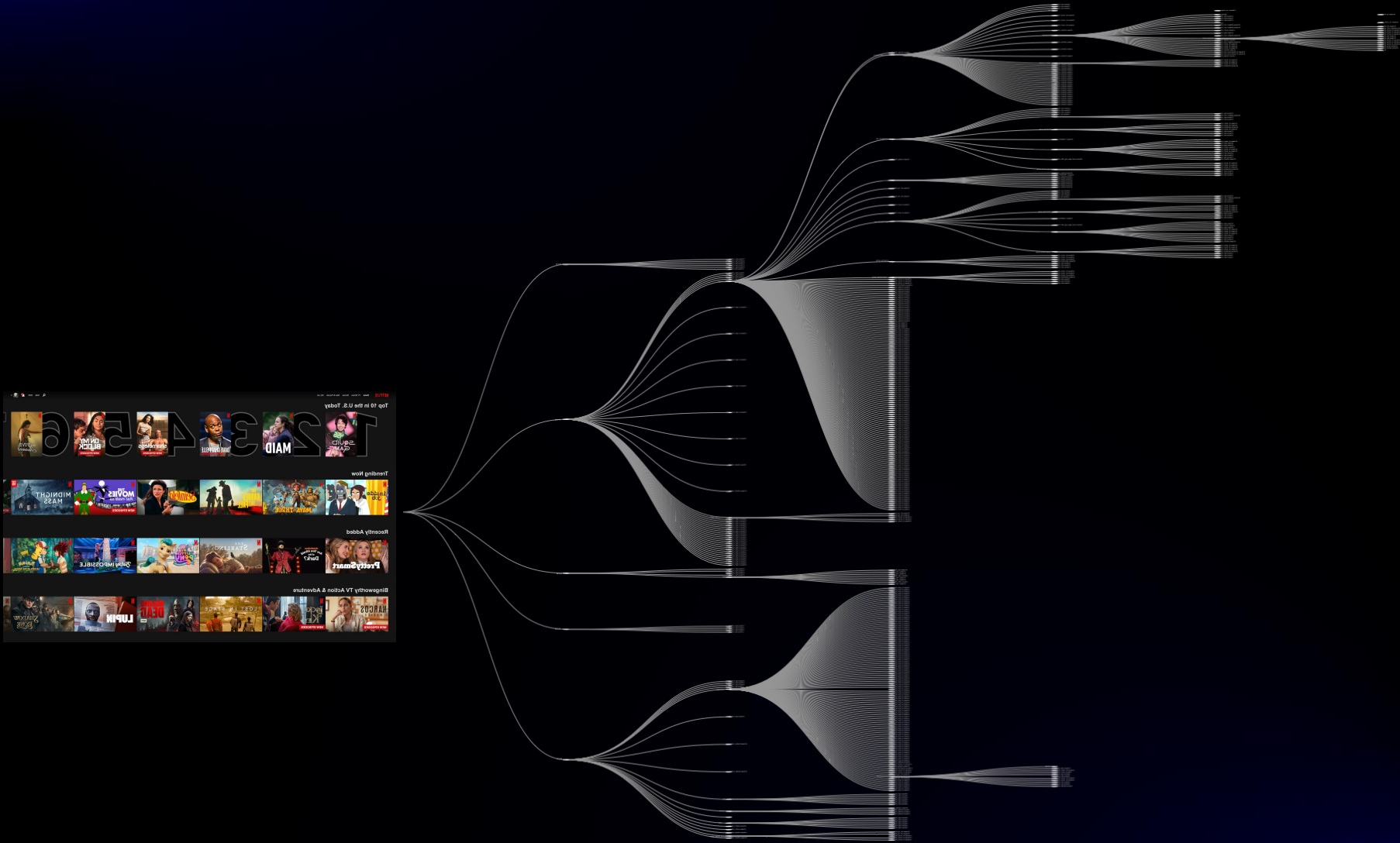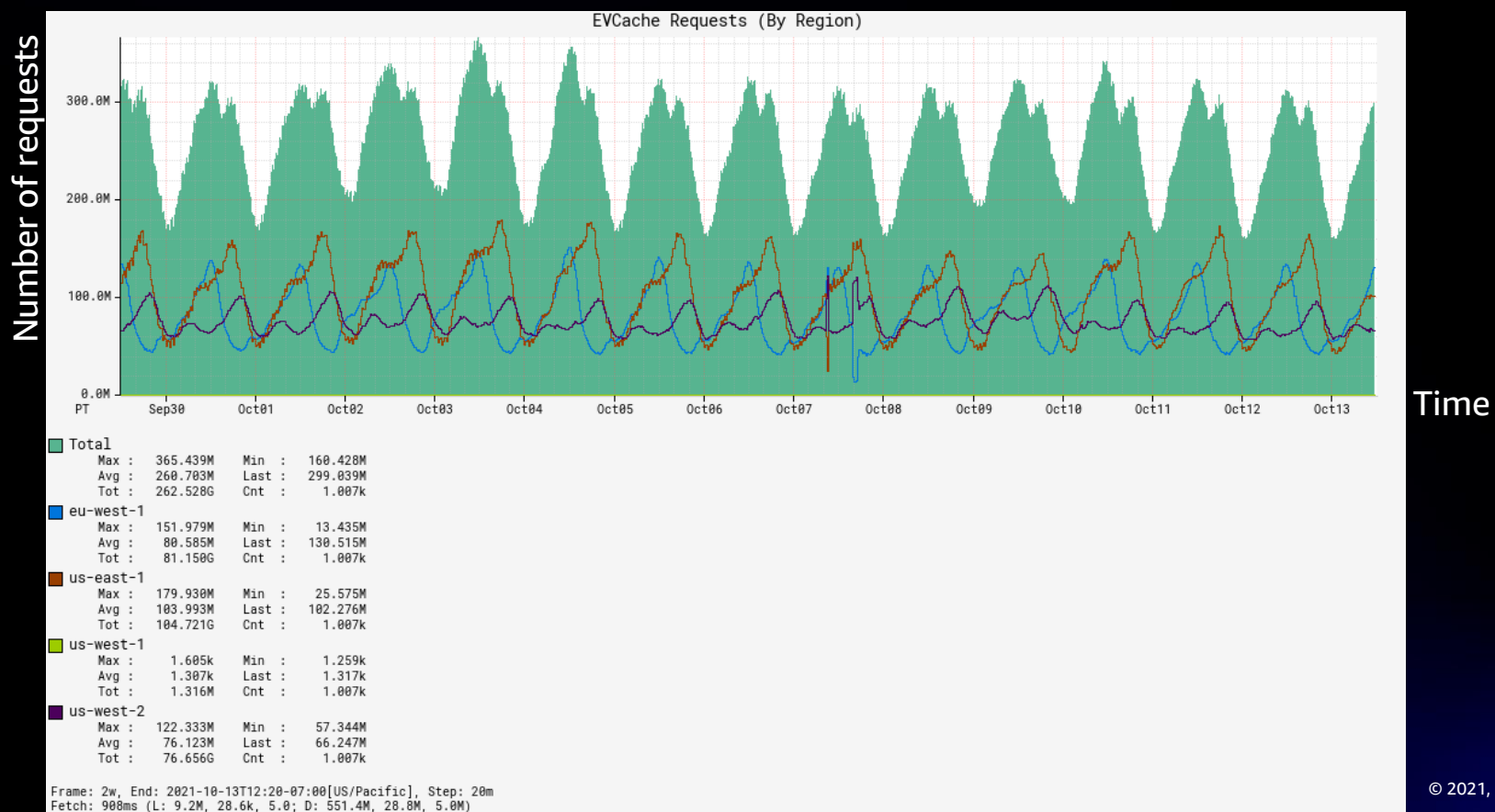
A day in the life of EVCache

Incidents and takeaways
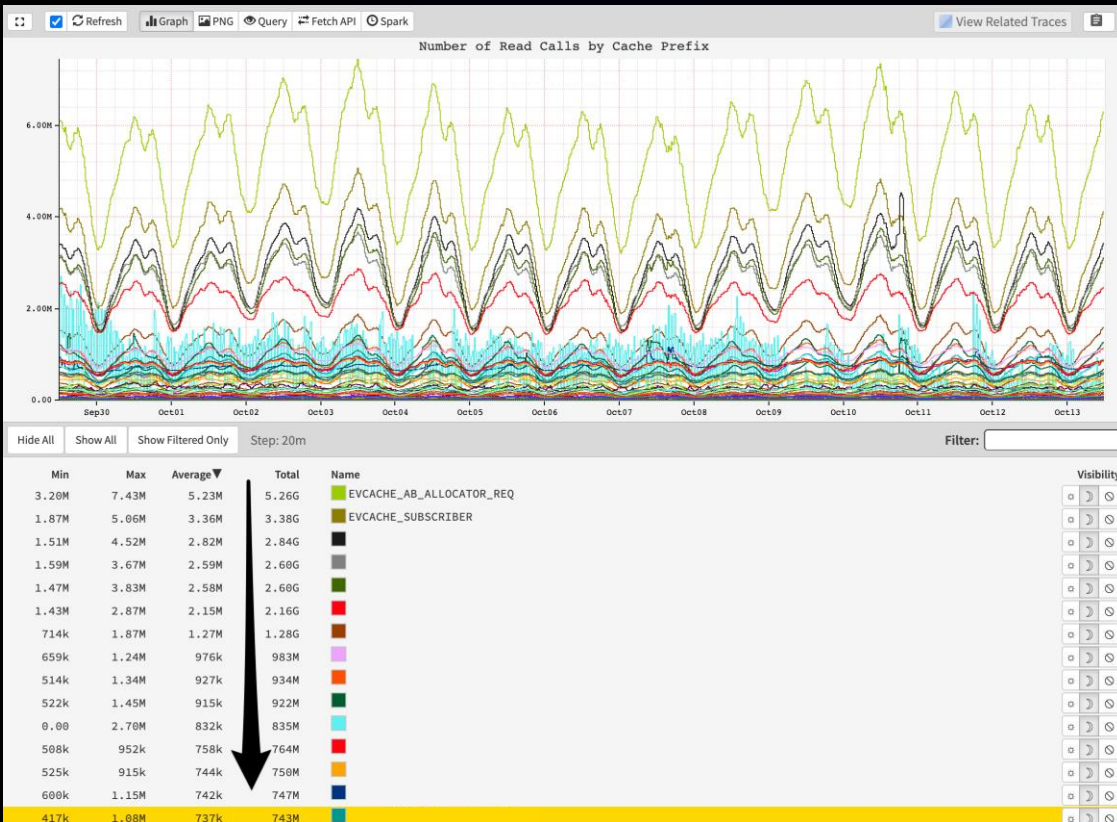
# Memcached workload characteristics at Netflix

We have 100s of clusters operating in production, out of which 45% are SSD based and the other 55% use RAM-based Memcached only

EVCache Requests (By Region)

Number of requests

300.0M

200.0M

100.0M

0.0M
PT    Sep30    Oct01    Oct02    Oct03    Oct04    Oct05    Oct06    Oct07    Oct08    Oct09    Oct10    Oct11    Oct12    Oct13

Time

Total
    Max :    365.439M    Min :    160.428M
    Avg :    260.703M    Last:    299.039M
    Tot :    262.528G    Cnt :    1.007k

eu-west-1
    Max :    151.979M    Min :    13.435M
    Avg :    80.585M    Last:    130.515M
    Tot :    81.150G    Cnt :    1.007k

us-east-1
    Max :    179.930M    Min :    25.575M
    Avg :    103.993M    Last:    102.276M
    Tot :    104.721G    Cnt :    1.007k

us-west-1
    Max :    1.605k    Min :    1.259k
    Avg :    1.307k    Last:    1.317k
    Tot :    1.316M    Cnt :    1.007k

us-west-2
    Max :    122.333M    Min :    57.344M
    Avg :    76.123M    Last:    66.247M
    Tot :    76.656G    Cnt :    1.007k

Frame: 2w, End: 2021-10-13T12:20-07:00[US/Pacific], Step: 20m
Fetch: 908ms (L: 9.2M, 28.6k, 5.0; D: 551.4M, 28.8M, 5.0M)

aws

# Memcached workload characteristics at Netflix
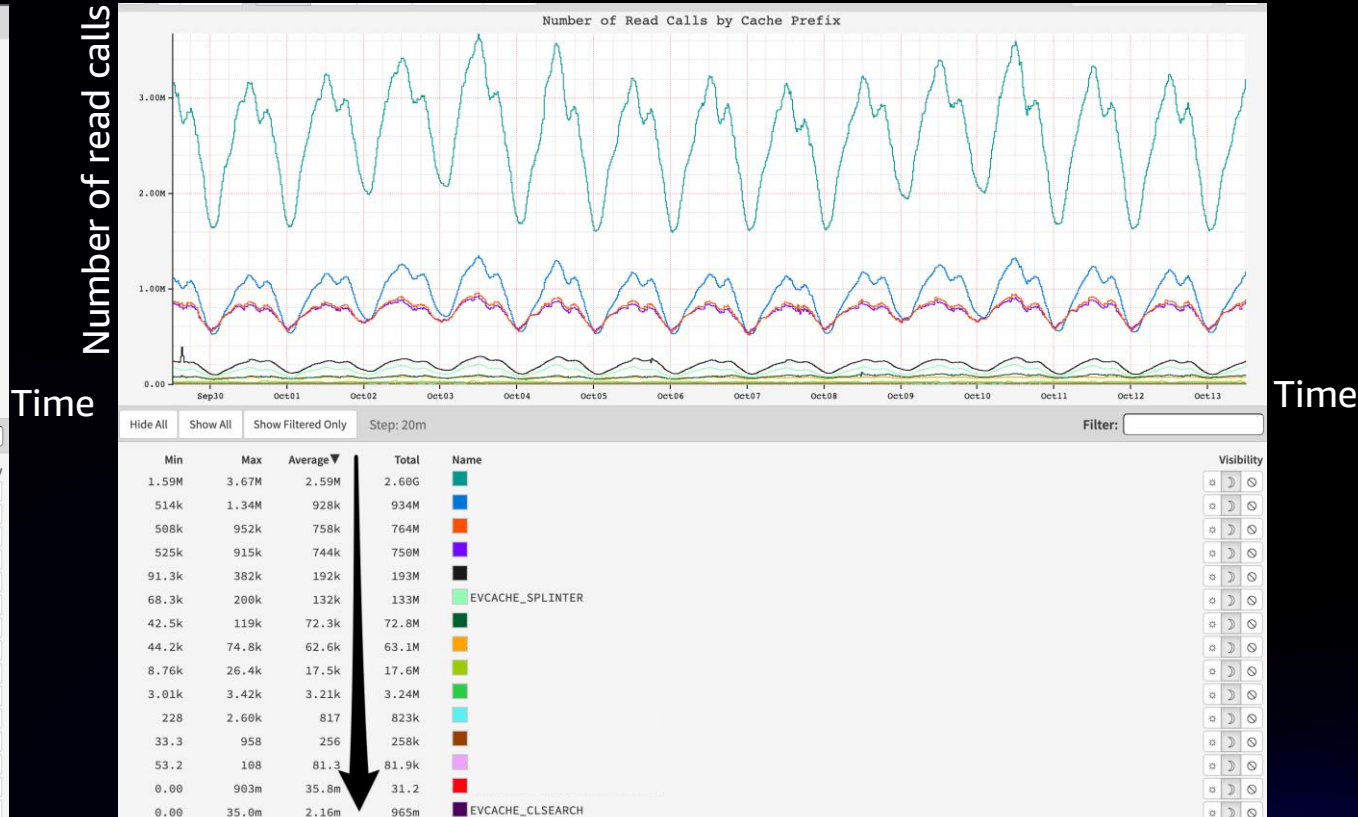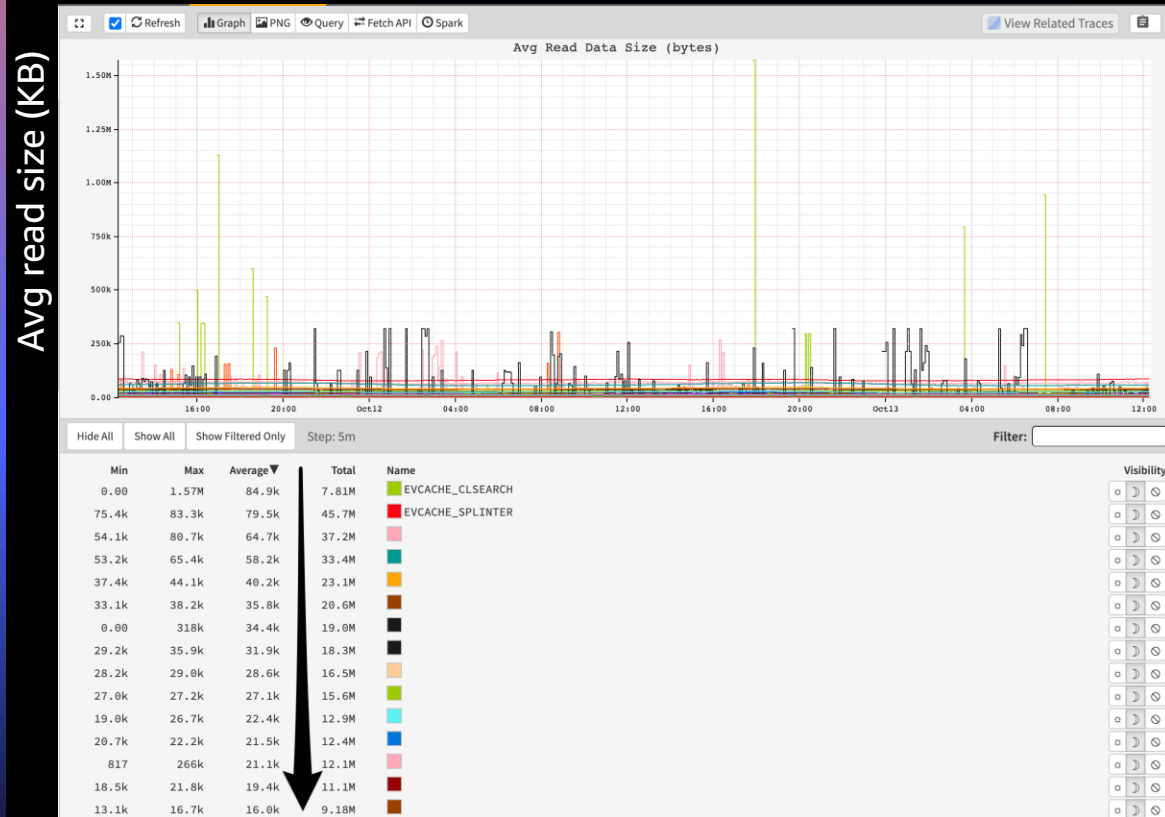
Clusters with high RPS and their average payload sizes

# Memcached workload characteristics at Netflix

## Clusters with large payload sizes and their RPS

# Agenda

Introduction
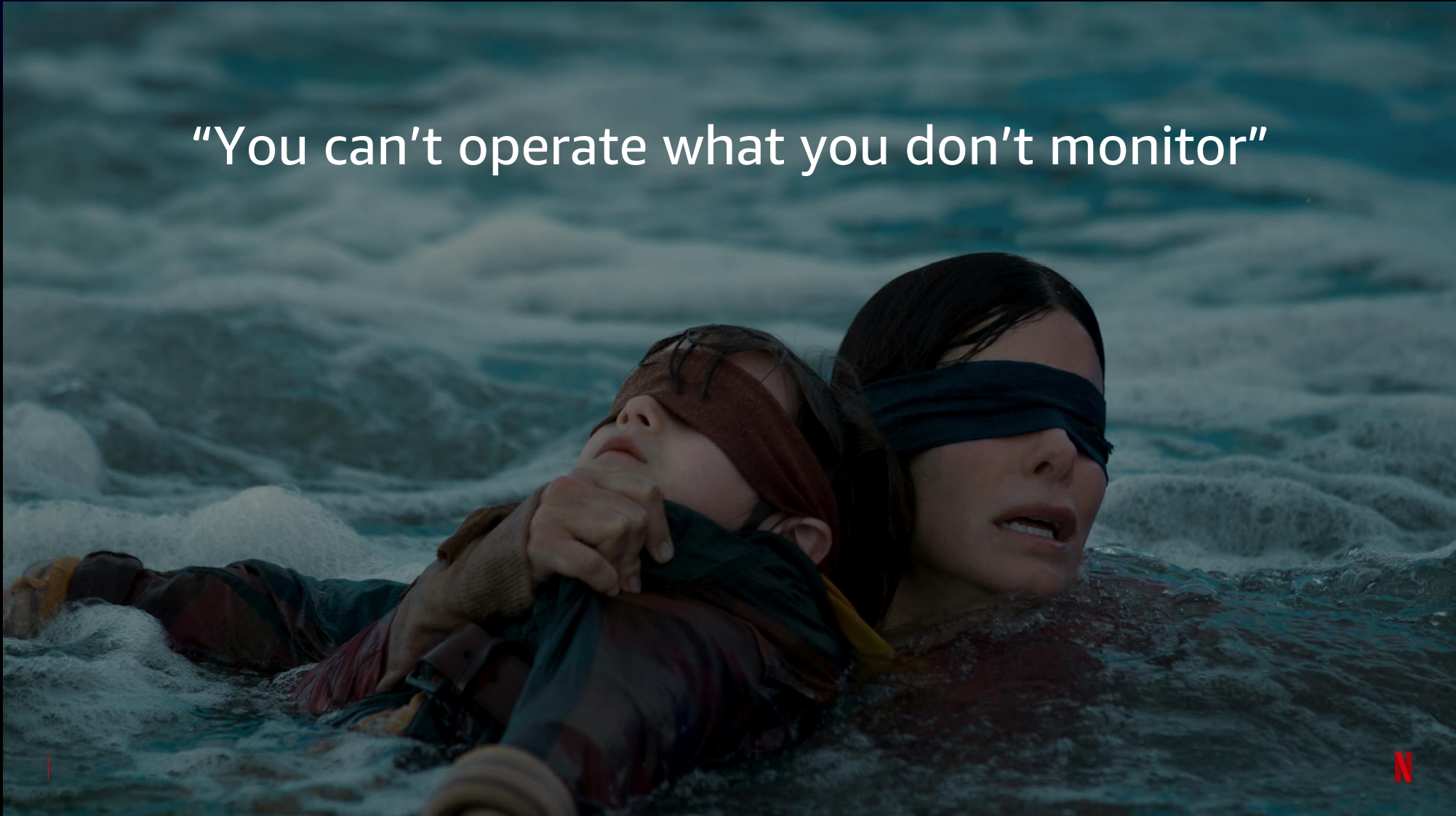
High-level architecture

Workload characterization

Observability

A day in the life of EVCache

Incidents and takeaways

# Observability

"You can't operate what you don't monitor"

# Observability



Client dashboard

Server dashboard

Alerts

# Observability

## Read latency of clients

# Observability

## Server network usage

# Observability

## Memcached page pools

# How do we provision the caching servers

**3 dimensions to monitor**

- Compute

- Network

- Memory

# How do we provision the caching servers

Let's pick an example and see how we provision memory for Memcached

```
/usr/bin/memcached  -m 27045 -o hashpower=27 <...>
```

- The above instance can hold up to $2^{27}$ items = 128 million items with average item size of 220 bytes

- Total memory used by Memcached

- Overcommit ratio is disabled

General rule of the thumb = aim for 90% memory provisioned for Memcached

Total RAM

Memcached

**Memcached (item memory)**

Memcached (hash table)

Memcached (connection objects, buffers, etc.)

OS memory

# Agenda

Introduction

High-level architecture

Workload characterization

Observability

A day in the life of EVCache

Incidents and takeaways

# A day in the life of EVCache

- Instance crashes/disappears, AWS health check
- Memcached gets terminated
- Disks go bad
- Data get wiped out

"Instance warmer refills data"

# Instance warming



us-east-1c (replica 1)

us-east-1d (replica 2)

Healthy cluster

# Instance warming

Instance 1 | Instance 2 | Instance 3

us-east-1c (replica 1)

Instance 1 | Instance 2 | Instance 3

us-east-1d (replica 2)

Unhealthy cluster in replica 2

aws

# Instance warming

Instance 1 | Instance 2 | Instance 3

us-east-1c (replica 1)

Instance 1 | Instance 2 | Instance 3 (new instance with no data)

us-east-1d (replica 2)

New instance coming up in replica 2

# Instance warming

Trigger cache dump for new instance

Cache warmer controller

Instance 1    Instance 2    Instance 3    us-east-1c (replica 1)

Instance 1    Instance 2    Instance 3 (new instance with no data)    us-east-1d (replica 2)

Cache warmer controller triggers cache dump for the new instance based on the new IP address

# Instance warming



us-east-1c (replica 1)

Upload relevant data chunks to
Amazon S3/Amazon EBS

**Amazon EBS**    **Amazon S3**

Instance 1    Instance 2    Instance 3 (new instance with no data)

us-east-1d (replica 2)

Cache warmer controller triggers cache dump for
the new instance based on the new IP address

aws

# Instance warming



us-east-1c (replica 1)

Amazon EBS    Amazon S3

Download data

Instance 1    Instance 2    Instance 3 (with all the data)

us-east-1d (replica 2)

Warm-up cache

Cache populator

Cache populator downloads the data from Amazon EBS/ Amazon S3 and writes to the destination IP(s)

# A day in the life of EVCache

- Workload changes, scale out
- Data set sizes increase from TBs to PBs

"Cache warmer to move data, faster with Amazon EBS Multi-Attach"

# Scaling up fast

- New deployments are quite common
  - Increase/decrease capacity
  - Memcached upgrade (fixes/new features)
  - Base Amazon Machine Image (AMI) upgrades
  - Instance type changes

- New deployments cannot take read traffic right away
  - Wait until item TTL duration
  - TTL can range from a few hours to a few weeks

# Scaling out fast



More context: https://netflixtechblog.com/cache-warming-agility-for-a-stateful-service-2d3b1da82642

# Scaling up fast(er)

More context: https://netflixtechblog.medium.com/cache-warming-leveraging-ebs-for-moving-petabytes-of-data-adcf7a4a78c3

# A day in the life of EVCache

Data corruptions by apps

"Snapshots and restores"


EDEN DUNCAN-SMITH    DANTÉ CRICHLOW

FROM PRODUCER **SPIKE LEE** AND DIRECTOR **STEFON BRISTOL**

**SEE YOU YESTERDAY**

A NETFLIX FILM

GOING BACK IS THE ONLY WAY FORWARD.

MAY 17 | NETFLIX

# A day in the life of EVCache

- Data center/Availability Zone outages
- Co-related network failures
- Region failures



You're very brave.

"Serve traffic from another Availability Zone or
failover to another Region"

# A day in the life of EVCache

## Serving traffic from another data center/Availability Zone

# A day in the life of EVCache

Serving traffic from another Region

# Agenda

Introduction

High-level architecture

Workload characterization

Observability

A day in the life of EVCache

Incidents and takeaways

# Incidents and takeaways

P95 latency issues

Heterogeneous client environments

Memory issues
- Application memory leaks
- Incorrect budgeting of resources
- Memory fragmentation

aws

# Triaging P95 latencies

- SLA with P95 <= ~5 ms*
- Disk latencies can come into play

# Triaging P95 latencies

- Is online traffic contending with any offline traffic?

- Is the workload spiky?

- eBPF tools – biolatency, biosnoop, etc.

# Heterogenous client environments

- Unable to connect from Mesos agents due to port exhaustion
- Standard tools/metrics don't help

Identify network namespaces

```
$ ip netns list
720 (id: 6)
16190 (id: 0)
```

# Heterogenous client environments

## Locate the cgroup and the container

```
$ systemctl status `pidof java`

daemontools.service - Daemontools service supervision
   Loaded: loaded (/lib/systemd/system/daemontools.service; enabled; vendor preset: enabled)
   Active: active (running) since Mon 2021-04-12 01:33:36 UTC; 4 weeks 1 days ago
 Main PID: 3771 (svscanboot)
    Tasks: 1365
   CGroup: /system.slice/daemontools.service
           ├─ 720 /apps/mesos-1.7.2/libexec/mesos/mesos-containerizer launch
```

# Heterogenous client environments

- Then we grouped connections by a server IP and observed the connection build-up

```
$ sudo ss -N 720 | grep "<server.ip.address.100>" | awk '{print $5}' | cut -d ':' -f 2 | wc -l
1020
```

- Clients were reconnecting and leaving old connections in the CLOSE_WAIT state
- Aggressive connection cleanup without leaving JVM to close them

# Triaging memory issues

```
[Thu Oct 28 05:28:26 2021] nvme nvme2: pci function 0000:00:1e.0
[Thu Oct 28 05:28:26 2021] nvme 0000:00:1e.0: enabling device (0000 -> 0002)
[Thu Oct 28 05:28:26 2021] nvme nvme2: 2/0/0 default/read/poll queues
[Thu Oct 28 05:28:26 2021] kworker/u24:2: page allocation failure: order:4,
mode:0x40dc0(GFP_KERNEL|__GFP_COMP|__GFP_ZERO), nodemask=(null),cpuset=/,mems_allowed=0
[Thu Oct 28 05:28:26 2021] CPU: 1 PID: 16990 Comm: kworker/u24:2 Tainted:
P           OE      5.3.0-1023-aws #25~18.04.1-Ubuntu
[Thu Oct 28 05:28:26 2021] Hardware name: Amazon EC2 i3en.3xlarge/, BIOS 1.0 10/16/2017
[Thu Oct 28 05:28:26 2021] Workqueue: nvme-wq nvme_scan_work
```

# Triaging memory issues

## First step is knowing state of memory

```
$ free -m
```

| | total | used | free | shared | buff/cache | availableMem: |
|---|---|---|---|---|---|---|
| | 31704 | 31348 | 242 | 8 | 113 | 37 |
| Swap: | 0 | 0 | 0 | | | |

^ Looking at this, we can deduce that the system is running low on memory

```
$ cat /proc/buddyinfo | ./analyze.sh
```

| | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| 15.54 MiB    Node 0, zone DMA | 1 | 0 | 0 | 1 | 2 | 1 | 1 | 0 | 1 | 1 | 3 |
| 134.90 MiB    Node 0, zone DMA32 | 1891 | 990 | 512 | 553 | 322 | 221 | 121 | 31 | 1 | 0 | 0 |
| 94.72 MiB    Node 0, zone Normal | 1646 | 2053 | 3238 | 693 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

^ and is also highly fragmented, as higher order pages are not available

# Triaging memory issues

Next is to look at top consumers of memory

```
$ ps -o pid,user,%mem,command ax | sort -b -k3 -r | head -n 2
  PID USER      %MEM COMMAND
2284 nfsuper   95.2 /apps/memcached/bin/memcached <masked the parameters>

$ pmap -x 2284 | tail -n 1
total kB          32168392 30923660 30922808
```

^ Since application in question is Memcached, we have a good understanding of how much
the Memcached process itself should take, but if for some reason the RSS memory keeps climbing,
then there is a high likelihood of memory leak in the application

# Triaging memory issues

- Amount of memory consumed by TCP or UDP buffers on the system can be profiled by checking the output of
  <span style="color:red">cat /proc/net/sockstats</span>

  - Sockstats shows the number of sockets opened on that machine and the memory occupied by the node

  - Whenever the application is slow or there is a spike in traffic, TCP memory usage goes up, leaving user space processes little to no memory to operate

aws

# Triaging memory issues

A few other things to look at

- What is the QPS or rate at which the requests are served on this cluster?

- What is the read/write network bandwidth on these nodes?

- Amount of memory consumed by kernel – monitor via slabtop

# Takeaways

- Architectural patterns used in operating data stores at scale on AWS
- Know your use cases and workload patterns
- Invest in observability from day one
- Continuously improve on reliability and scalability

# Thank you!

Tharanga Gamaethige

https://www.linkedin.com/in/tgamaethige/

Prudhviraj Karumanchi

https://www.linkedin.com/in/prudhviraj9/