

The logo for AWS re:Invent features the word "re:" in a smaller, gray sans-serif font positioned above the word "Invent". The word "Invent" is in a large, bold, white sans-serif font. The "i" in "Invent" has a vertical stroke that extends upwards, connecting it to the "e" in "re:". The background of the logo is a white rectangular area with a thin black border, set against a dark blue gradient background.

AWS
re:Invent

D A T 4 0 6

Netflix: Iterating on Stateful Services in the Cloud

Joey Lynch
Senior Software Engineer
Netflix

Speaker

Joey Lynch

Senior Software Engineer on Cloud Database
Engineering at Netflix

Distributed system addict and data wrangler



Agenda

Netflix cloud philosophy

What makes state difficult

Strategies for managing state

How to take it to 11

Netflix cloud philosophy

AMI is source of truth

- "Base" AMI
 - Operating system (such as Ubuntu)
 - Virtualization type
- Installed software
 - Linux kernel version
 - User software

```
context: !bunch.Bunch
ami: !bunch.Bunch
base_ami_name: xenialbase-x86_64-201805030044-ebs
creator: owner-team-email@netflix.com
ena_networking: true
enhanced_networking: true
suffix: 20181117192736-xenial-hvm-sriov
tags: !bunch.Bunch {}
vm_type: hvm

2018-11-17 19:28:56 [DEBUG] The following NEW packages will be installed:
cde-bootup{a} cde-cass-filebeat{a} cde-coredump-uploader{a}
cdecassostweak{a} cdecassutils{a} cdedatascrambler-spring{a} cphalo{a}
dstat{a} gandalf-agent{a} happycache{a} jvmquake{a} libgd3{a} libjbig0{a}
libjpeg-turbo8{a} libjpeg8{a} liblz4-tool{a} libtiff5{a} libvpx3{a}
libxpm4{a} libxslt1.1{a} linux-headers-4.13.0-19{a}
linux-headers-4.13.0-19-generic{a} linux-hwe-edge-tools-4.13.0-19{a}
linux-image-4.13.0-19-generic{a} linux-image-extra-4.13.0-19-generic{a}
linux-tools-4.13.0-19-generic{a} metatron-persistence{a} mlocate{a}
nf-spectator-agent{a} nflx-bolt{a} nflx-cde-kyber{a} nflx-cloudpassage{a}
nflx-ezconfig{a} nflx-python-3.5.5{a} nflx-python-3.6.4{a} nfpriam2{a}
nfpypy27-awscli{a} nfpypy27-base{a} nfpypy27-botocore{a} nfpypy27-certifi{a}
nfpypy27-chardet{a} nfpypy27-colorama{a}
nfpypy27-gevent{a} nfpypy27-greenlet{a} A "Cassandra" AMI
nfpypy27-pip{a} nfpypy27-pyasn1{a} nfpypy27-requests{a} nfpypy27-rsa{a} nfpypy27-s3transfer{a}
nfpypy27-setuptools{a} nfpypy27-six{a} nfpypy27-urllib3{a} nginx{a}
nginx-common{a} nginx-core{a} oss-cass3x{a} parallel{a} s3gof3r{a}
schedstat-collector{a}
```

Huzzah immutability!

- Build: bake an AMI
- Deployment: red-black ASGs with new AMIs
- Scaling + Failure: ASGs just take care of it

Huzzah immutability!



But what's in a typical stateful instance?

Service	Configuration	Tooling
• Cassandra	• Sysctls	• Perf tools
• Elasticsearch	• IO Scheduler	• State transfer
• Memcache	• TC qdisc	• Metrics
• Zookeeper	• Readahead tunings	• Monitoring
• More	• More	• More

And?

Mutable state

We made a state “closet”

SOAs pushed the challenges of state and state management down a layer into the cache/data tier

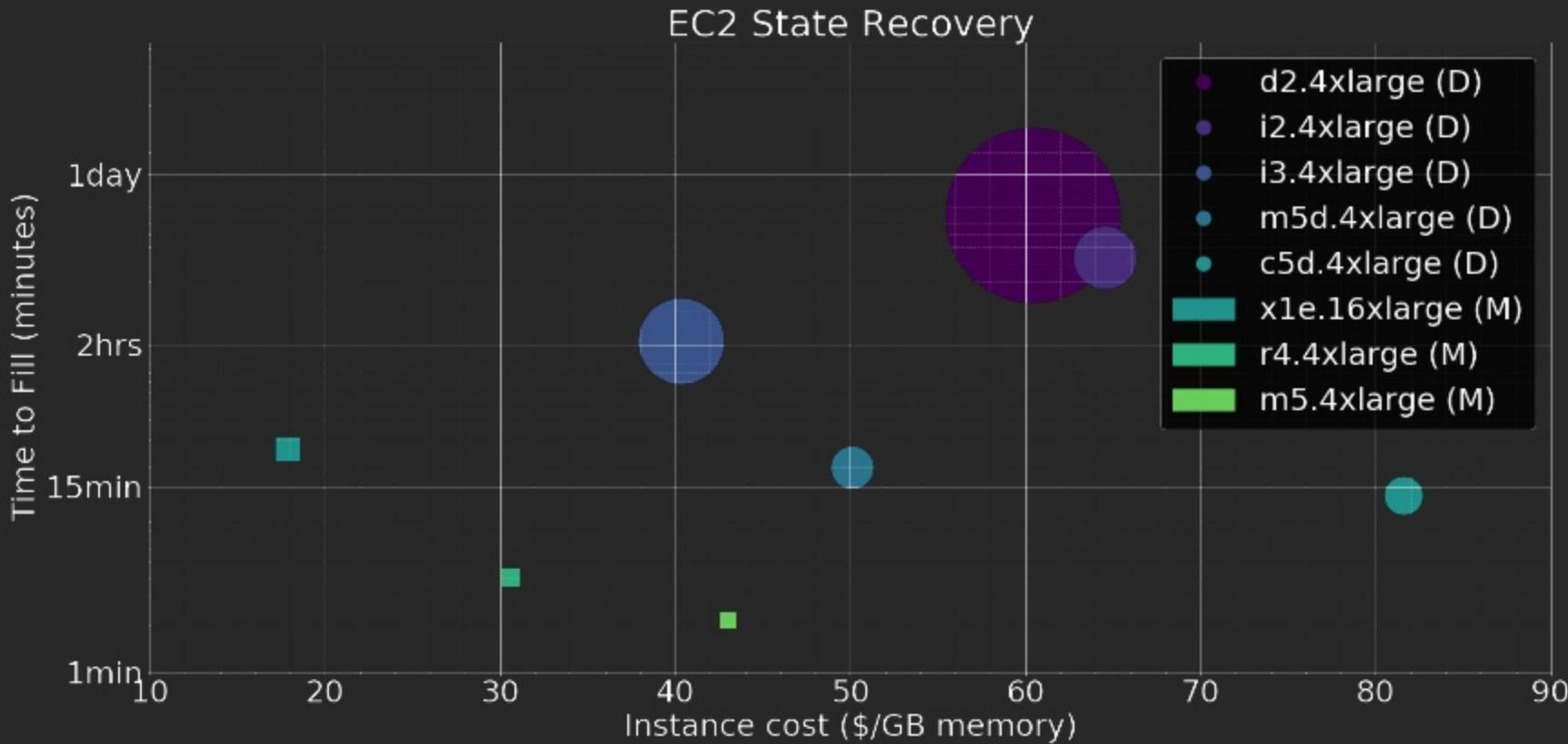
This didn't make the problem go away



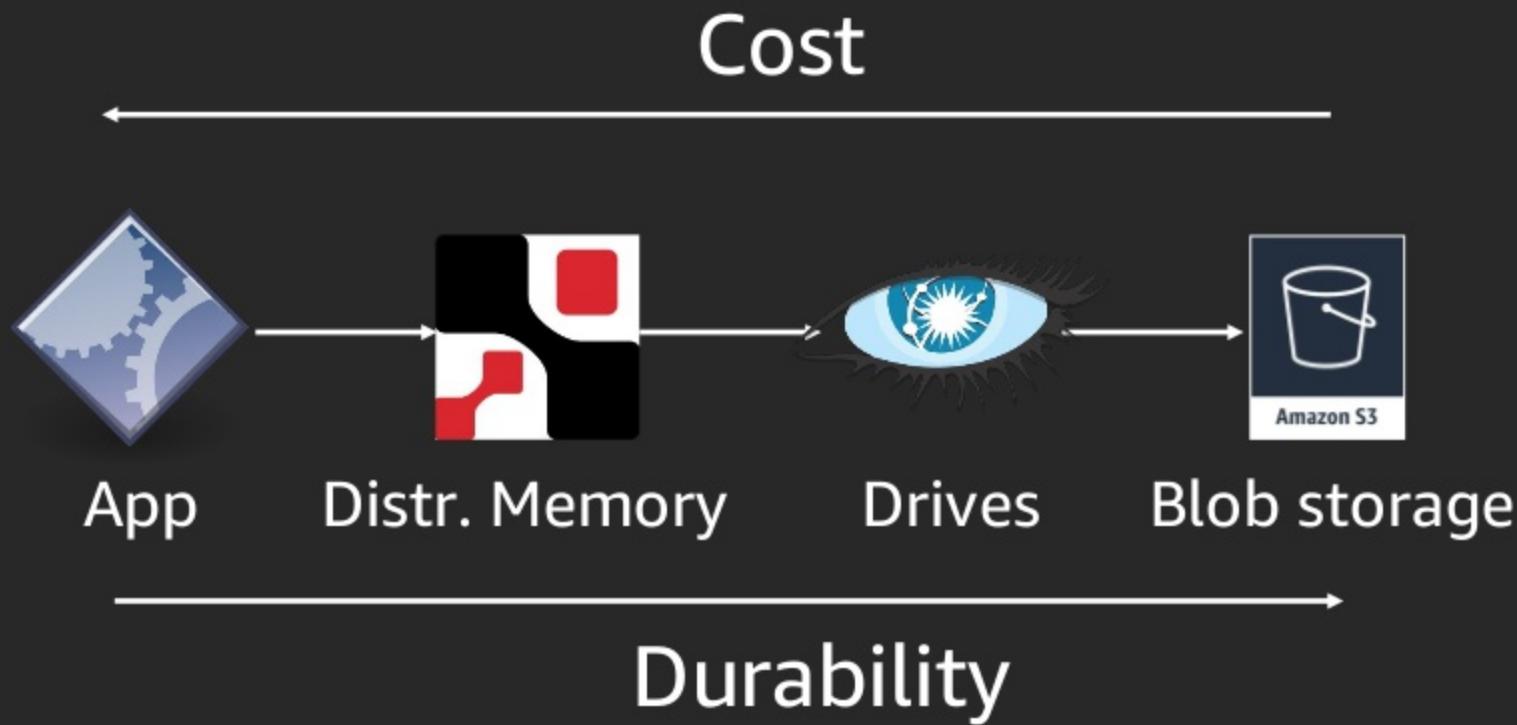
State is heavy



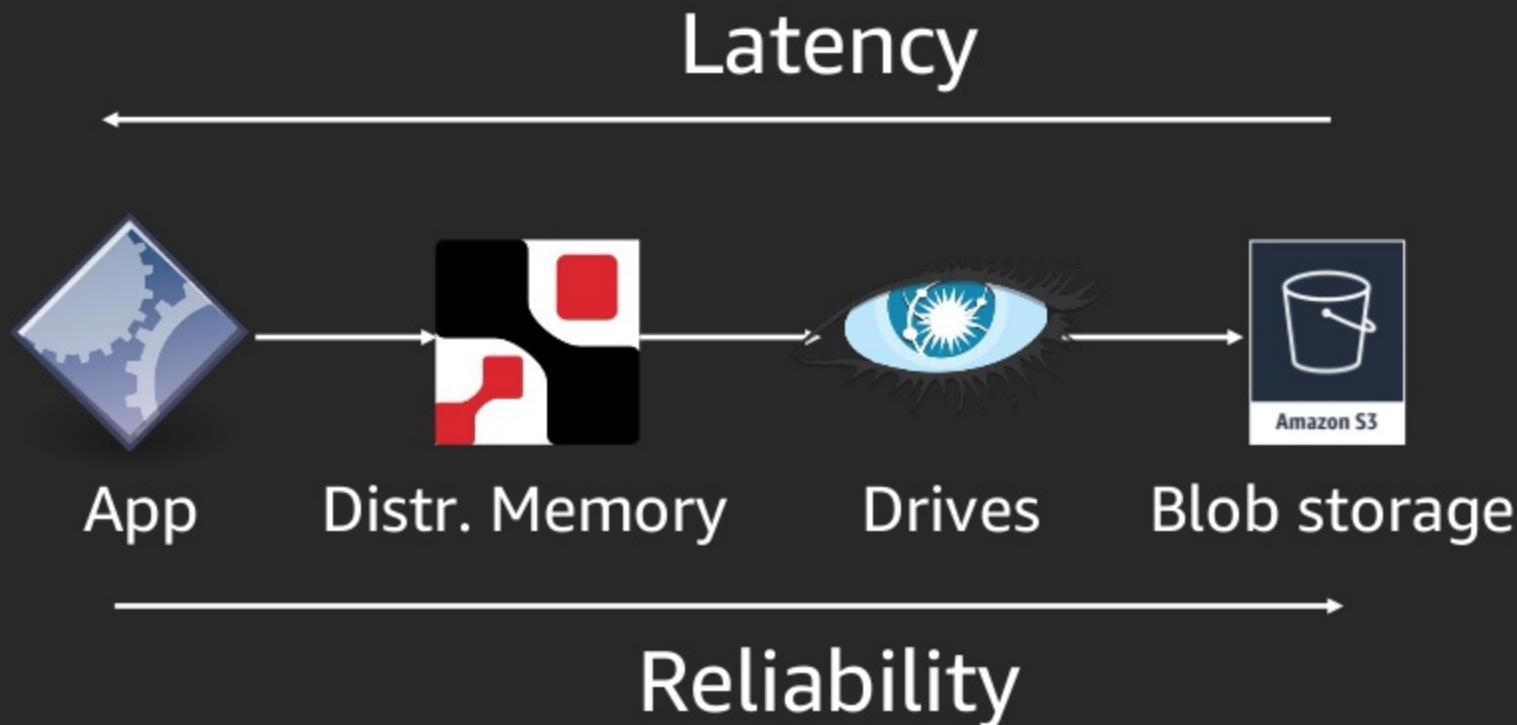
Instance fill time



It's caches all the way down



It's caches all the way down



What makes stateful services different?

1) Stateful services have, well ... *state*

Means you can't just throw the instance away

2) Stateful services have many components

Datastore, sidecars, cron jobs, monitoring, metrics

3) Require different iteration rates

Strategies for managing state

Strategy 0: Don't manage state

Use another closet
(Amazon Elastic Block Store
[Amazon EBS] gp2 and io1)

We can push the state even
further into the network layer

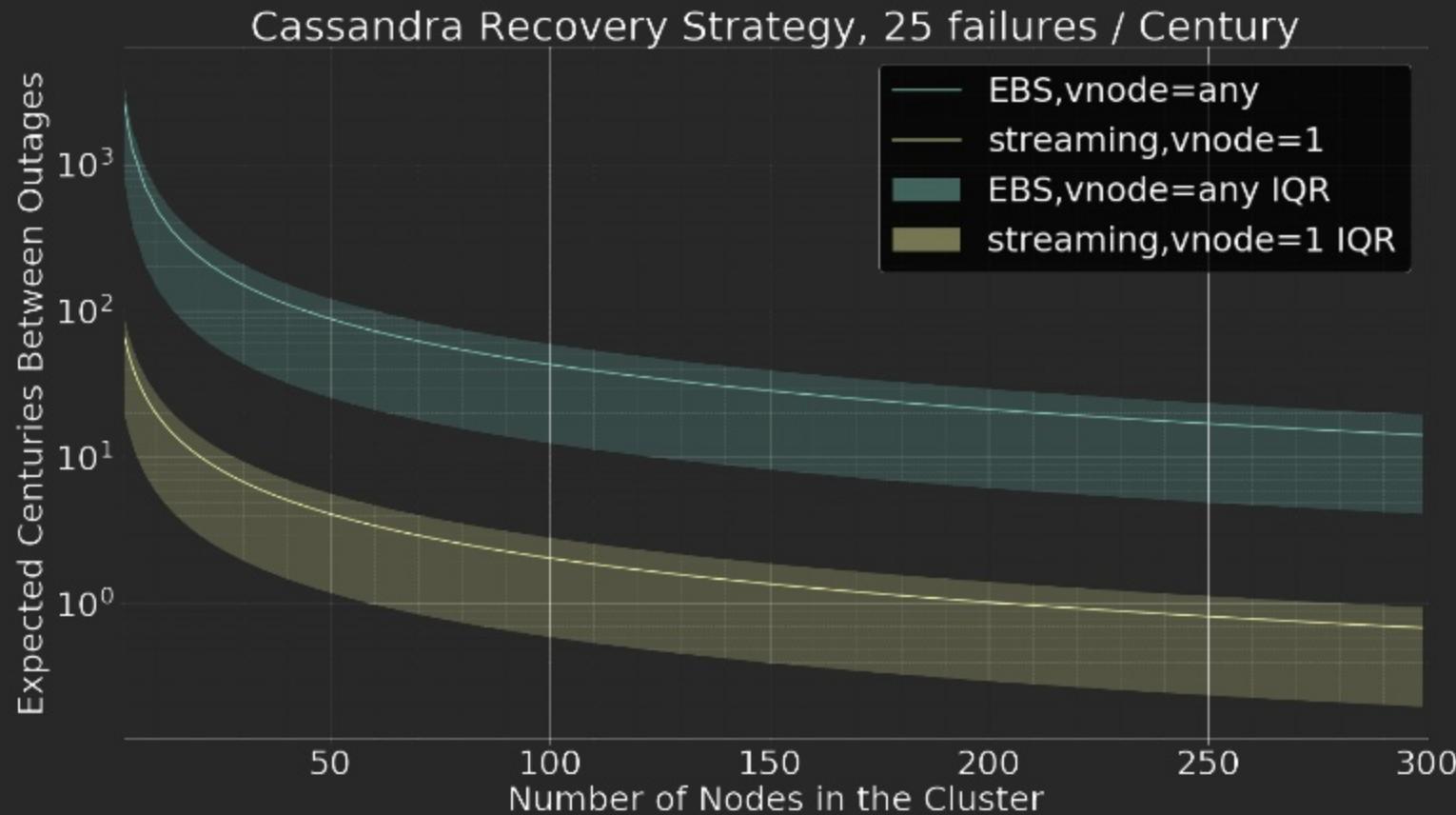
The database is now “stateless.”



AWS instance state is very fast

But ... Amazon EBS gp2 is comparatively slow

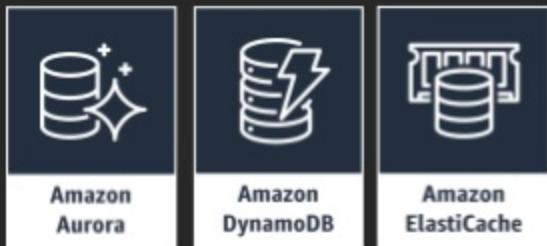
Amazon EBS GP2 means more availability



Use a bigger closet
(Amazon Aurora, Amazon
DynamoDB ...)

We can push the state into
datastore as a service options

The state is now someone else's
problem!



Should I run my own datastores?

- 1) Do you need specific datastore APIs or services?
- 2) Do you need to be under full control of the SLAs of that system?
- 3) Do you need to meet specific regulatory requirements?
- 4) Do you have a limited budget and need maximum efficiency?
- 5) Do you enjoy debugging difficult distributed systems problems?

Before starting: Isolate the state

Instance layout matters

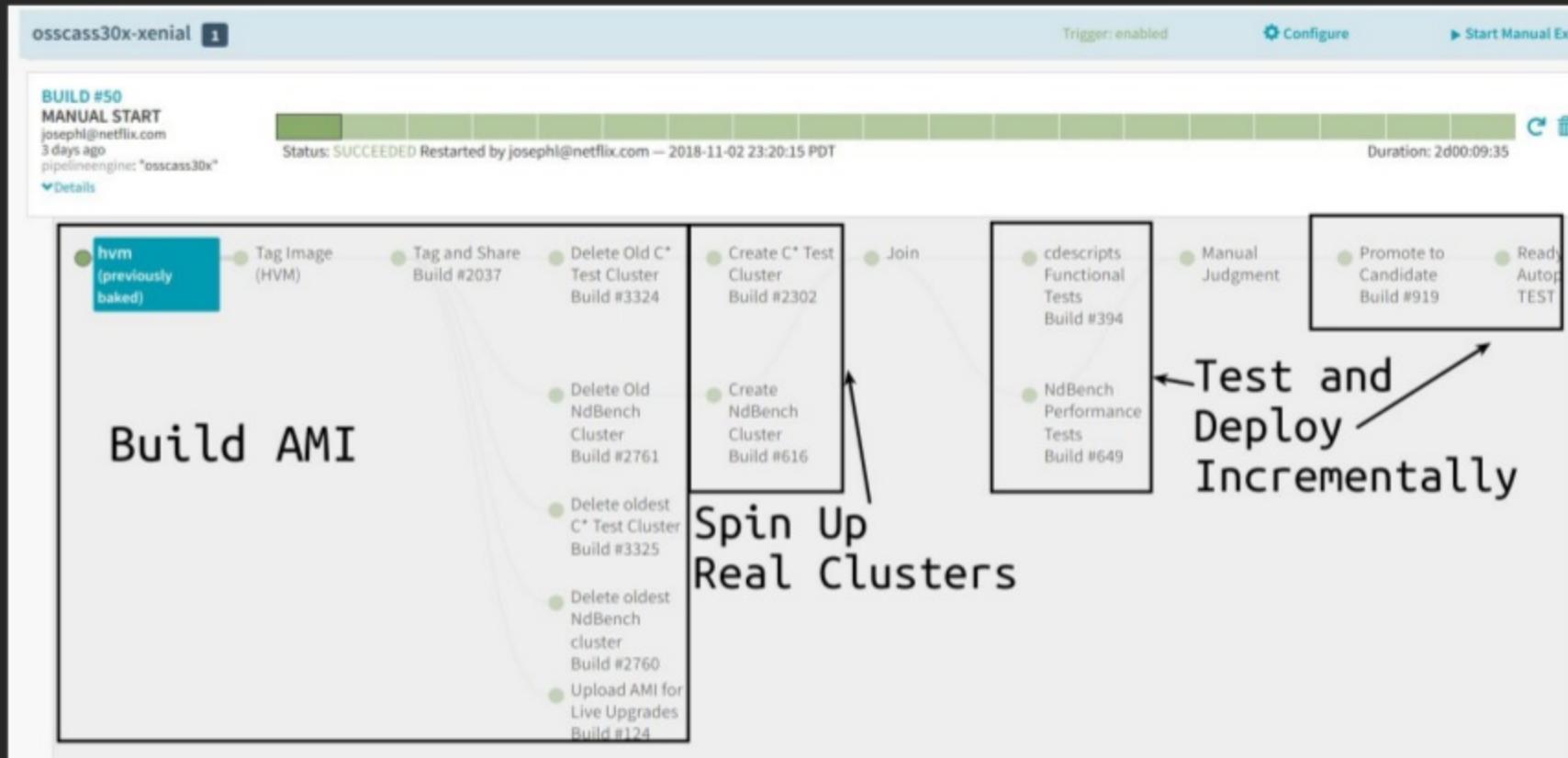
Stateful components

- Data: /mnt/data
- Logs: /var/log

Stateless components

- Primary service
- OS image
- System packages
- Background services
- Metrics, discovery
- More

Test your stateful systems E2E



Strategy 1: Mutate instances in place

Make mutation safe

- Configuration management works
 - Puppet and chef are mature products
- Package upgrades works
 - Linux OS packaging is mature and works
- Container upgrades work
 - Less mature but safer than packages

Why mutate: Speed!

	Interval	Mutation Time	Mutation Time (cluster)
Datastore Upgrade	~months	~1 hour	~days
Configuration Change	~days	~1 second	~minutes
OS Upgrade	~months	~15 minutes	~days
Background Process	~days	~30 seconds	~days
Hardware Change	~weeks	N/A*	N/A*

Mutate instances

Pros

- Super easy
- Fast, as you don't have to move data
- Tried-and-true technique
- Containers + Agents = Less foot shooting

Cons

- Totally violates AMI as source of truth
- Completely untested
 - Transitions
 - States
- Can easily have fleet drift

Which strategies does Netflix use?

Don't manage

Amazon EBS st1, gp2
Aurora
Amazon Relational
Database Service
(Amazon RDS)
DynamoDB
Amazon Simple Storage
Service (Amazon S3)

Mutate

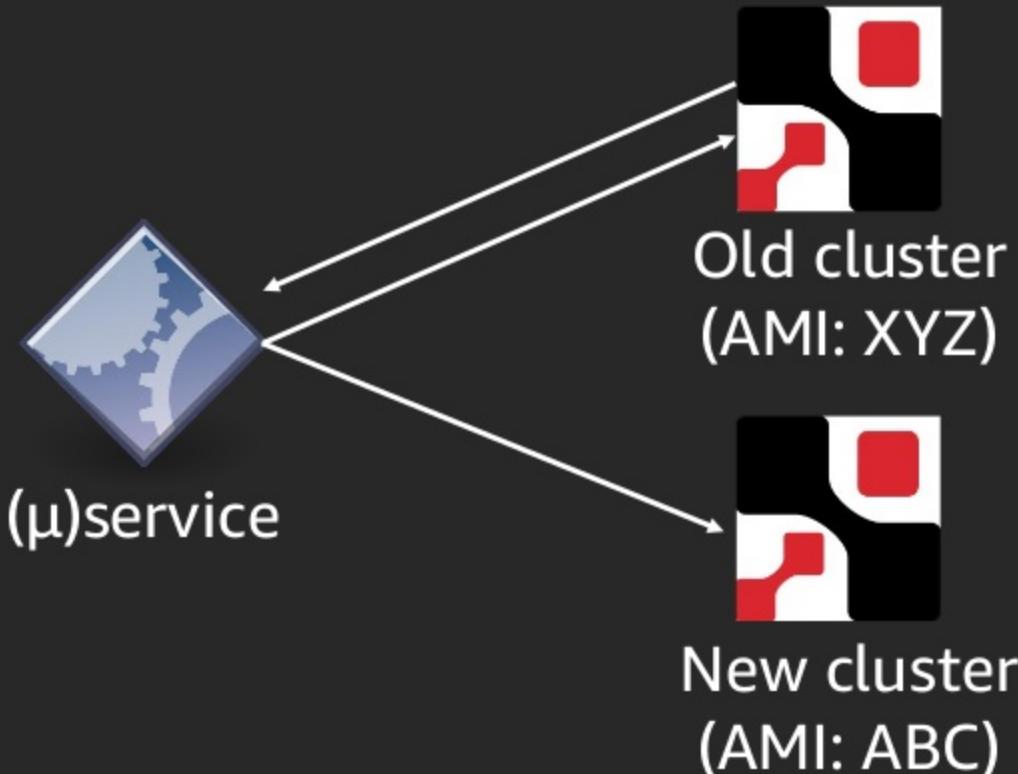
All of them

Strategy 2: Migrate state at the client

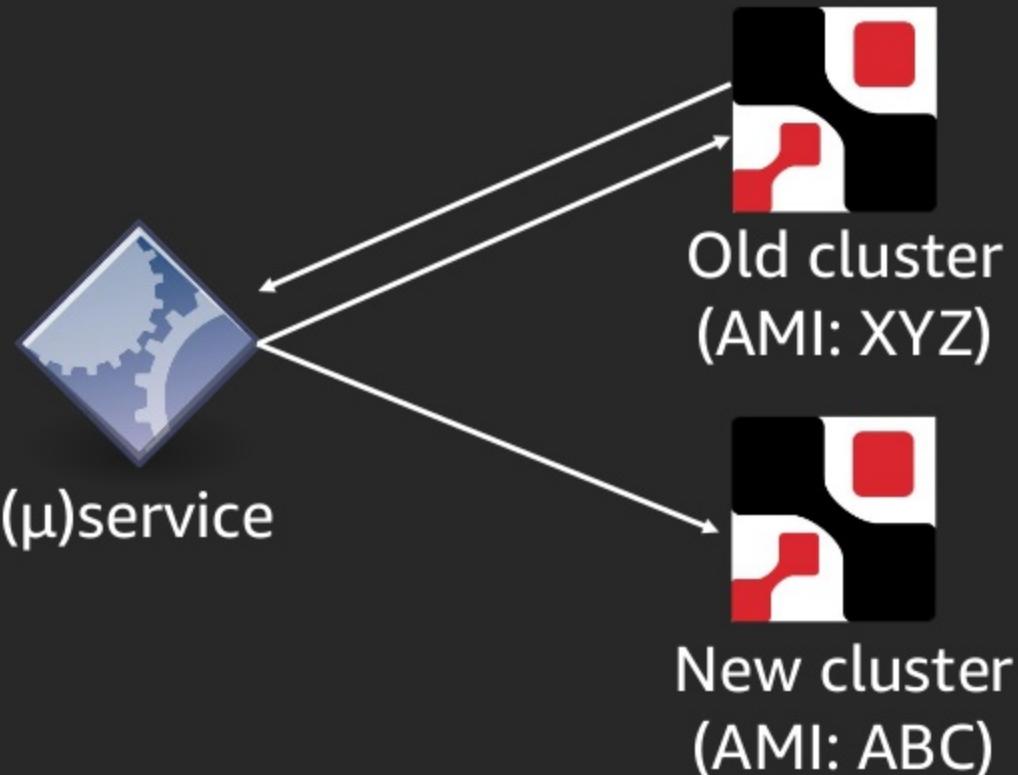
Mirroring state



Step 1: Mirror writes

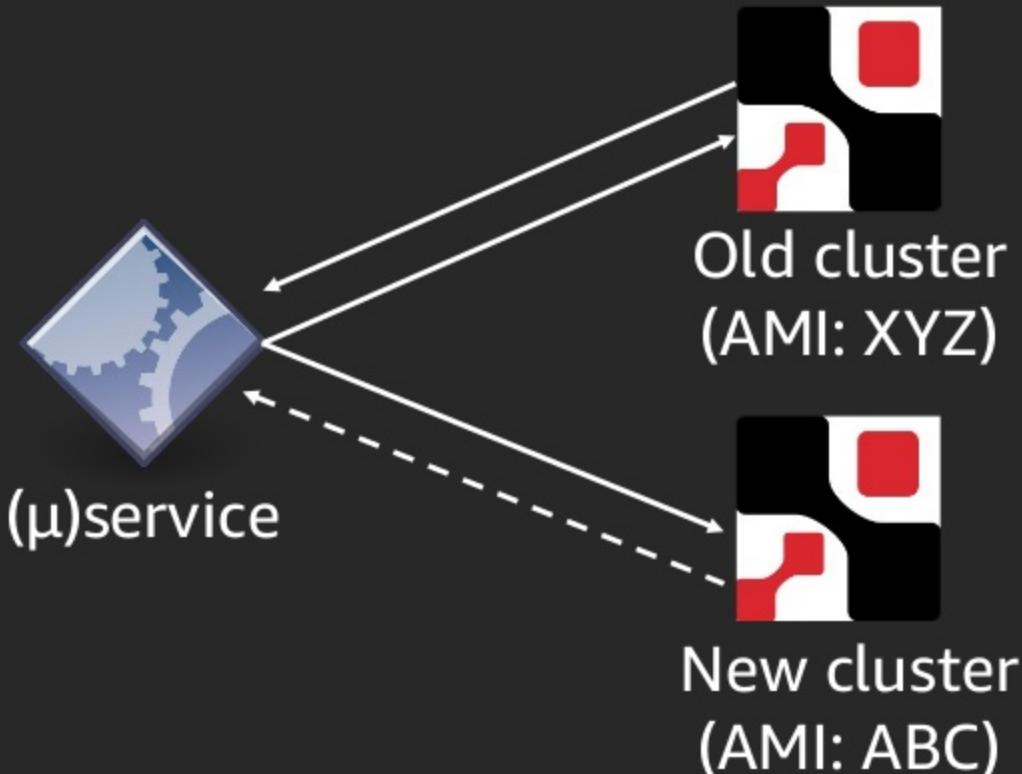


Step 2: Wait

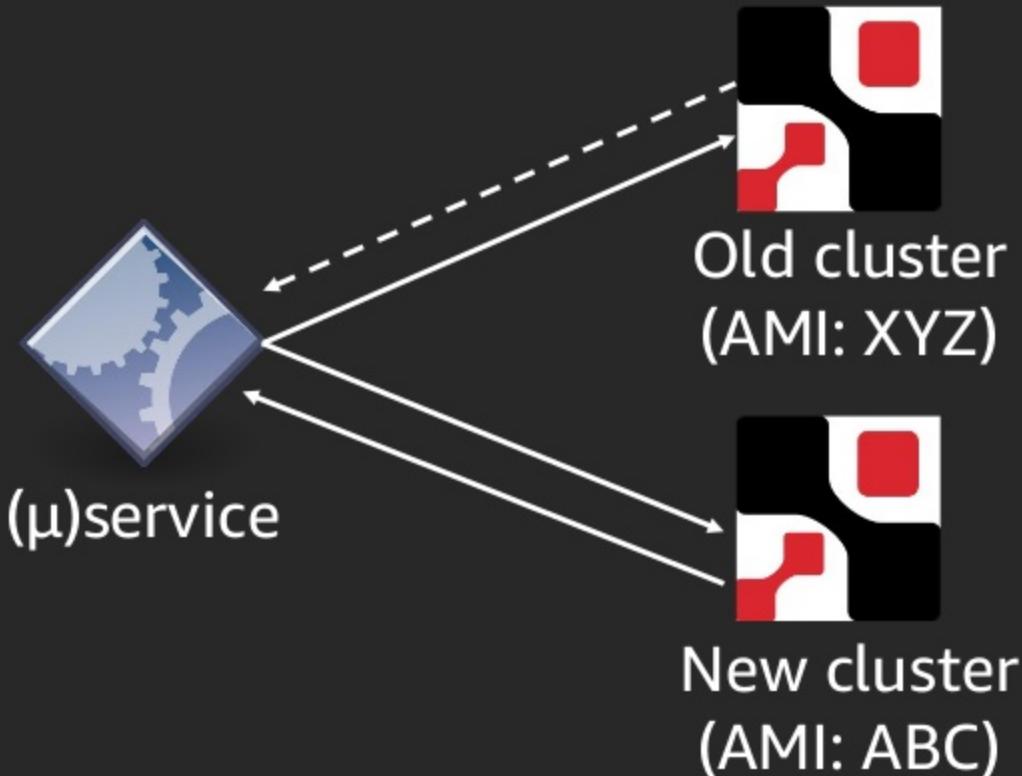


Wait for TTL

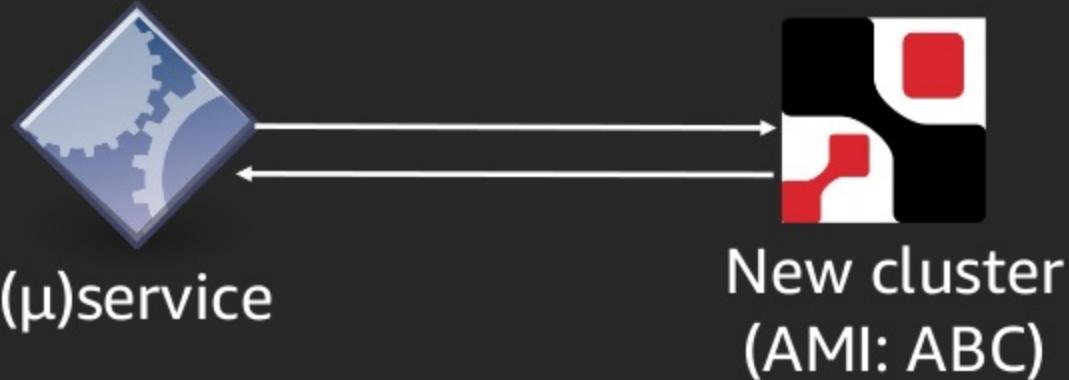
Step 3: Concurrent reads*



Step 4: Live launch reads



Step 5: Retire the old cluster



Migrate state at client

Pros

- Relatively easy to do
- Works well for ephemeral data
 - Caches
 - Search engines
 - Logs

Cons

- Not 100% consistent
- Hard to use with non ephemeral data
 - Persistent caches (yes)
 - Sources of truth
- Harder than mutating instances

Which strategies does Netflix use?

Don't manage

Amazon EBS
st1, gp2

Aurora

Amazon RDS

DynamoDB

Amazon S3

Mutate

All of them

Client

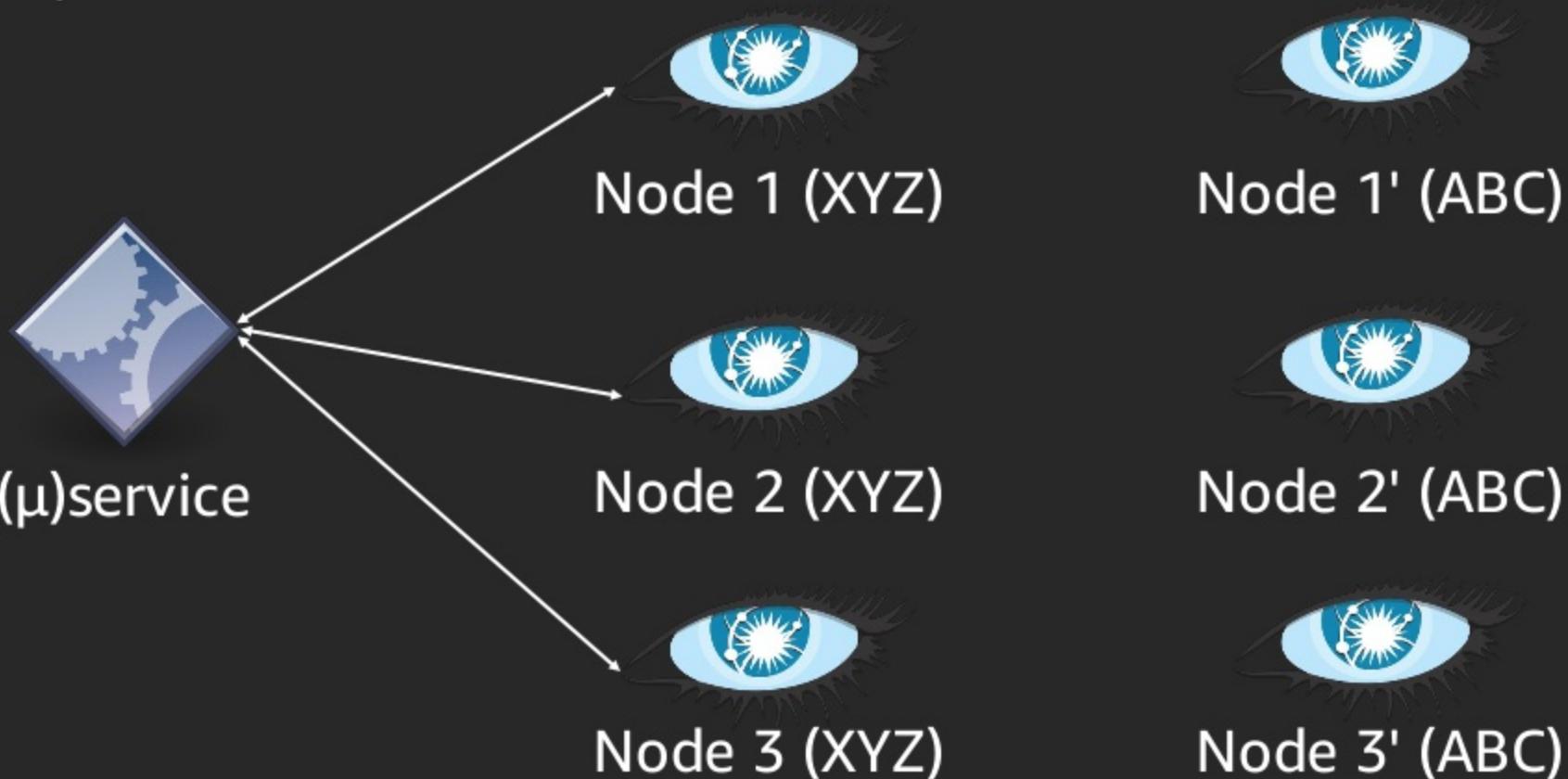
Caches
Search engines
Logs

Strategy 3: Migrate state at the server

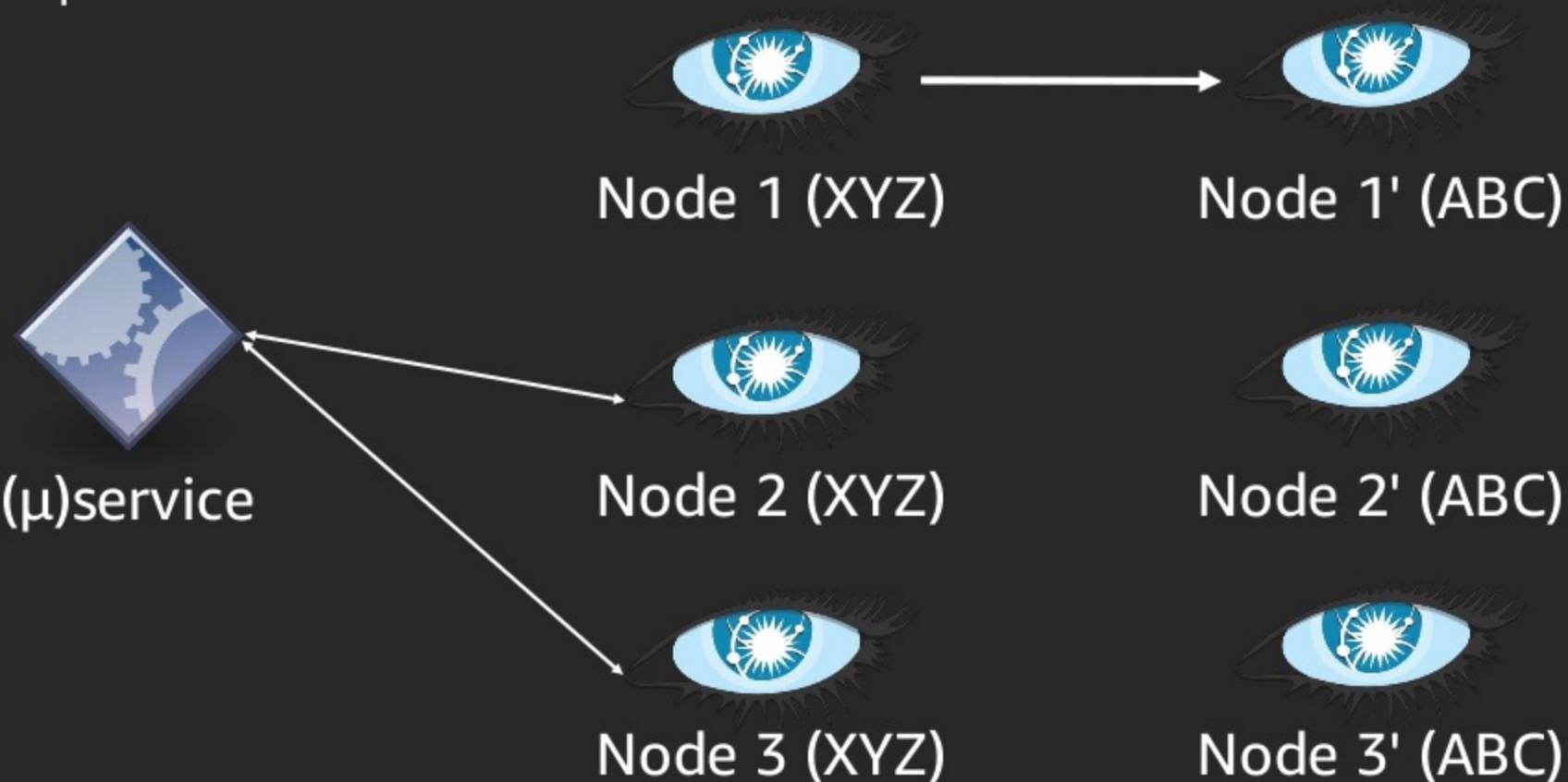
Migrate state



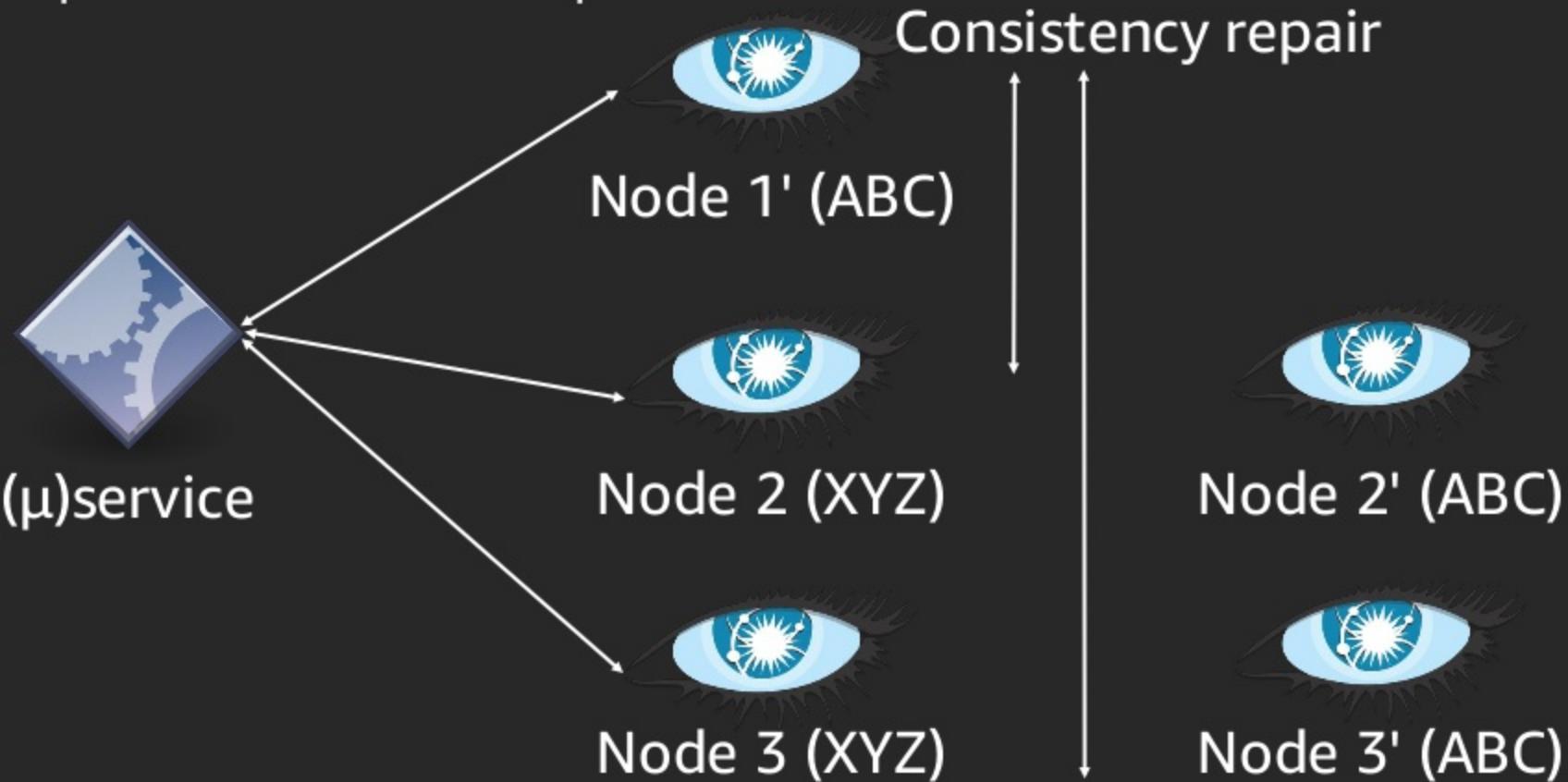
Step 1: Launch new nodes



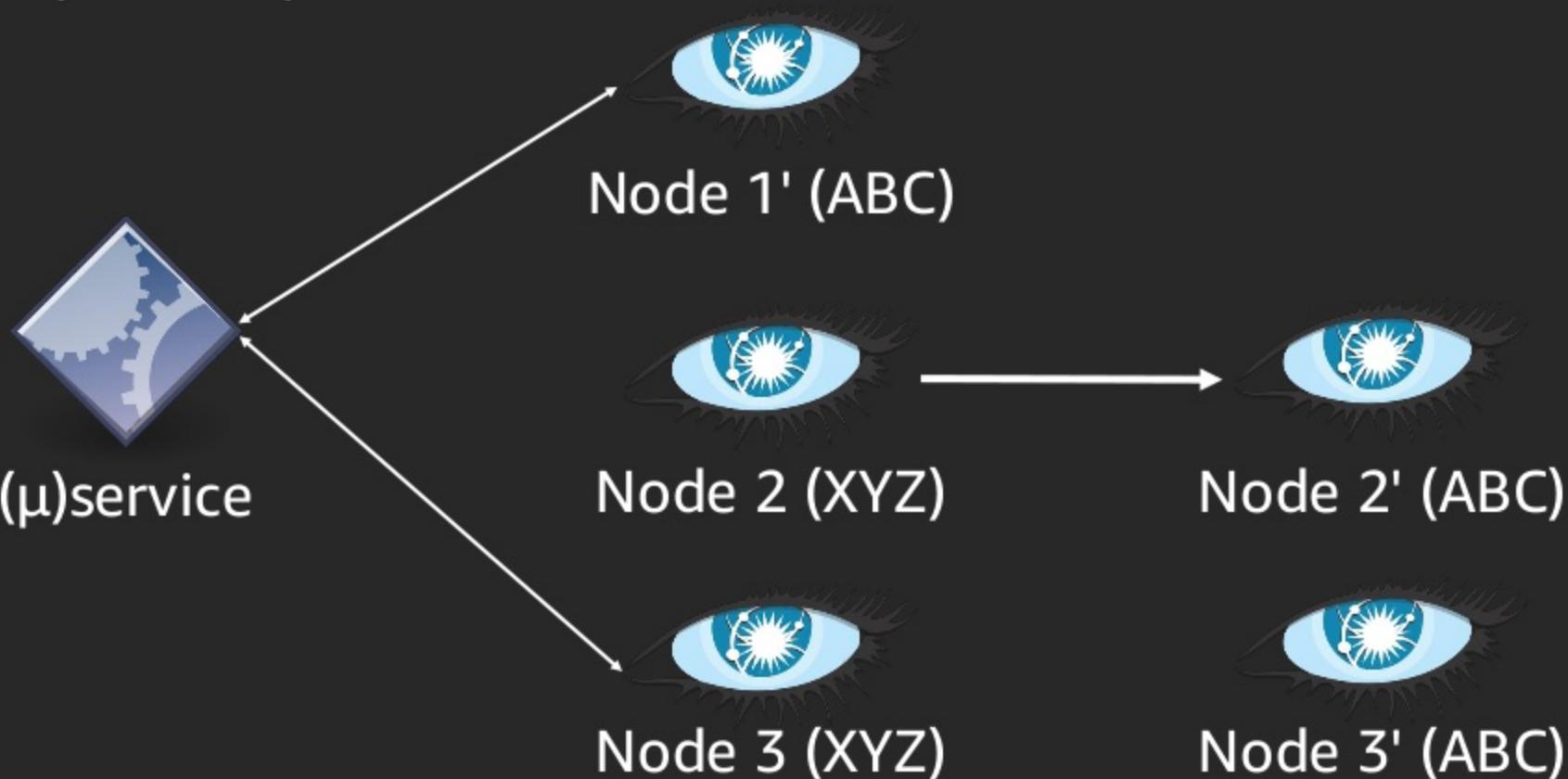
Step 2: Partition and transfer state



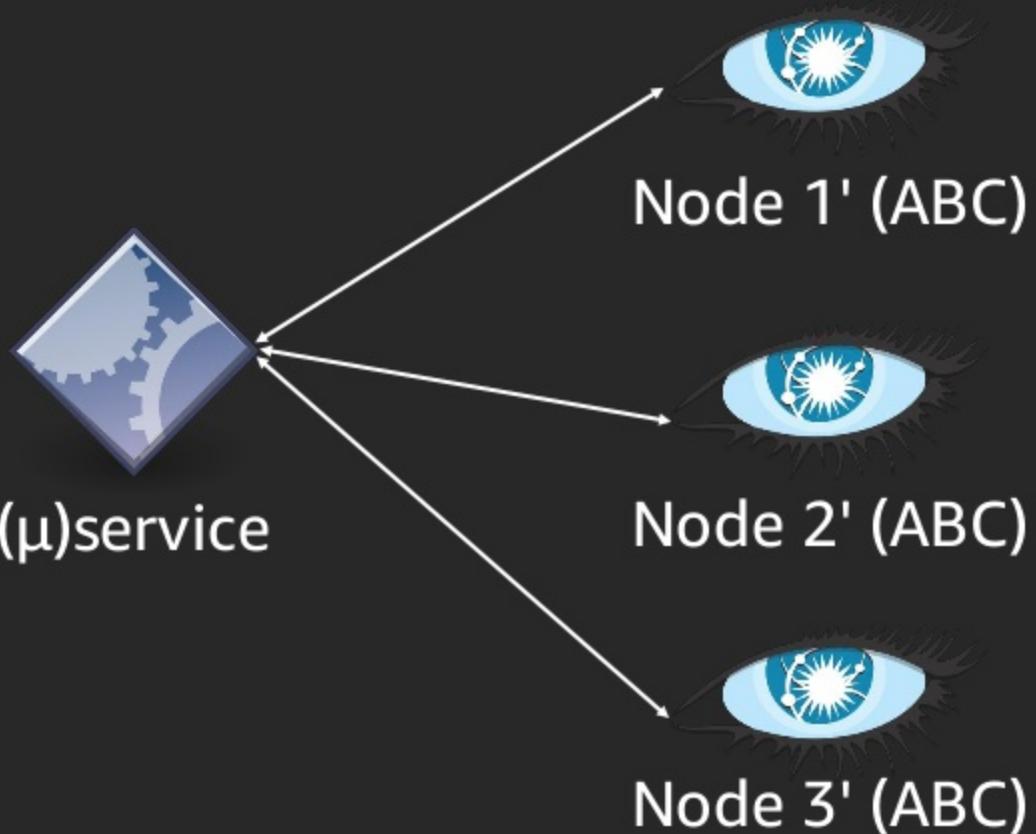
Step 3: Heal and repair



Step 4: Repeat



Step 5: Heal inconsistencies



Data movement

- Consistency repair
- Exploit "single node" failures
- Must ensure data fidelity

Migrate state at server

Pros

- Works for almost any distributed database
 - "A single node can fail"
 - Exercise chaos on your clients
 - "But we can't restart our database!!??"

Cons

- Pretty complicated
- Moving state is slow
 - Faster with Amazon EBS
- Presents performance challenges
- Requires API backwards compatibility

Which strategies does Netflix use?

Don't manage

Amazon EBS
st1, gp2

Aurora

Amazon RDS

DynamoDB

Amazon S3

Mutate

All of them

Client

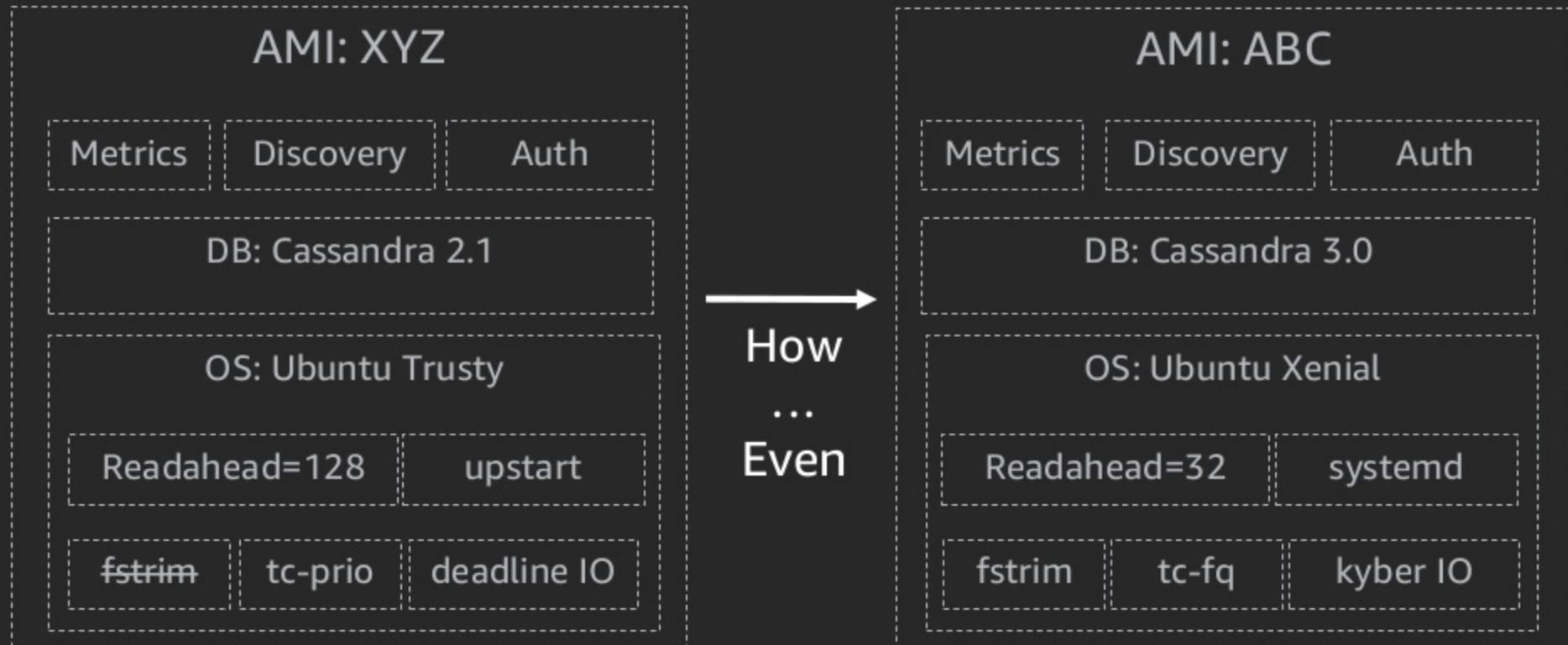
Caches
Search engines
Logs

Server

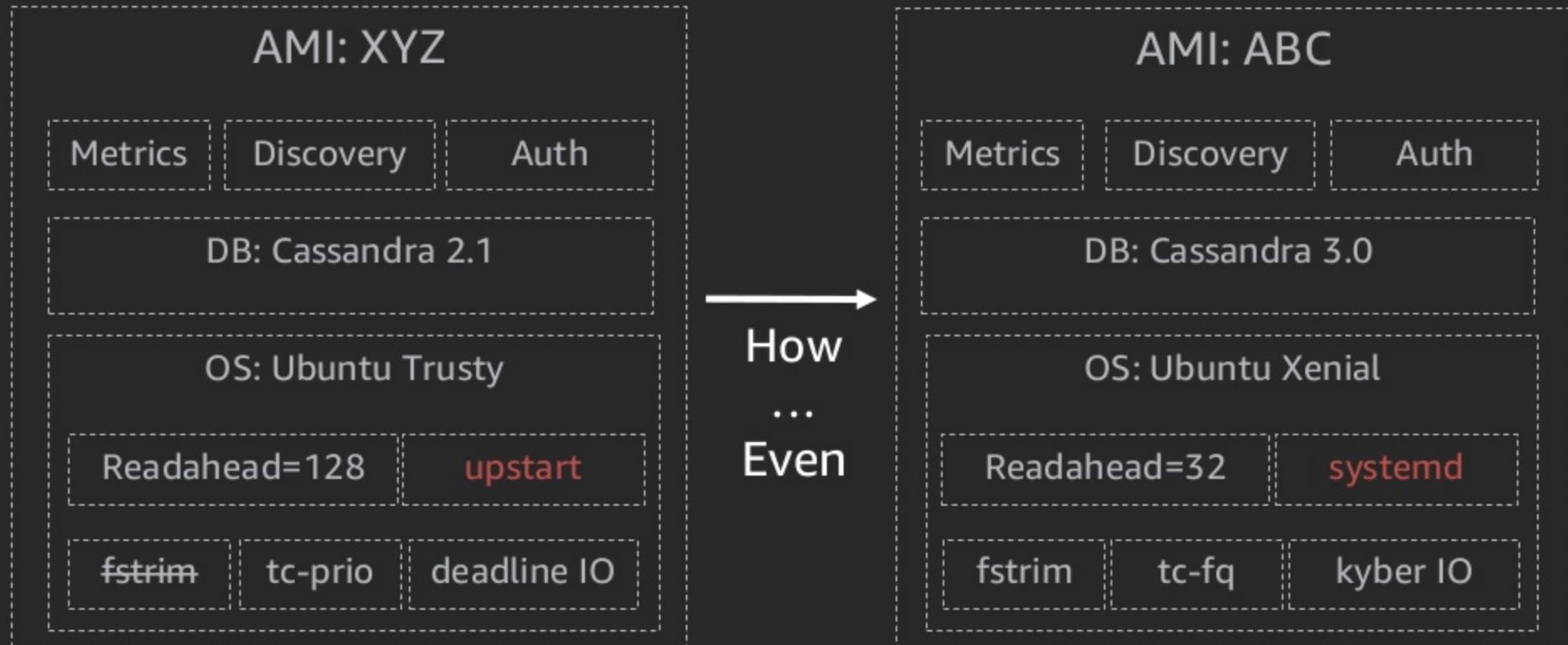
Sources of truth

Take it to 11

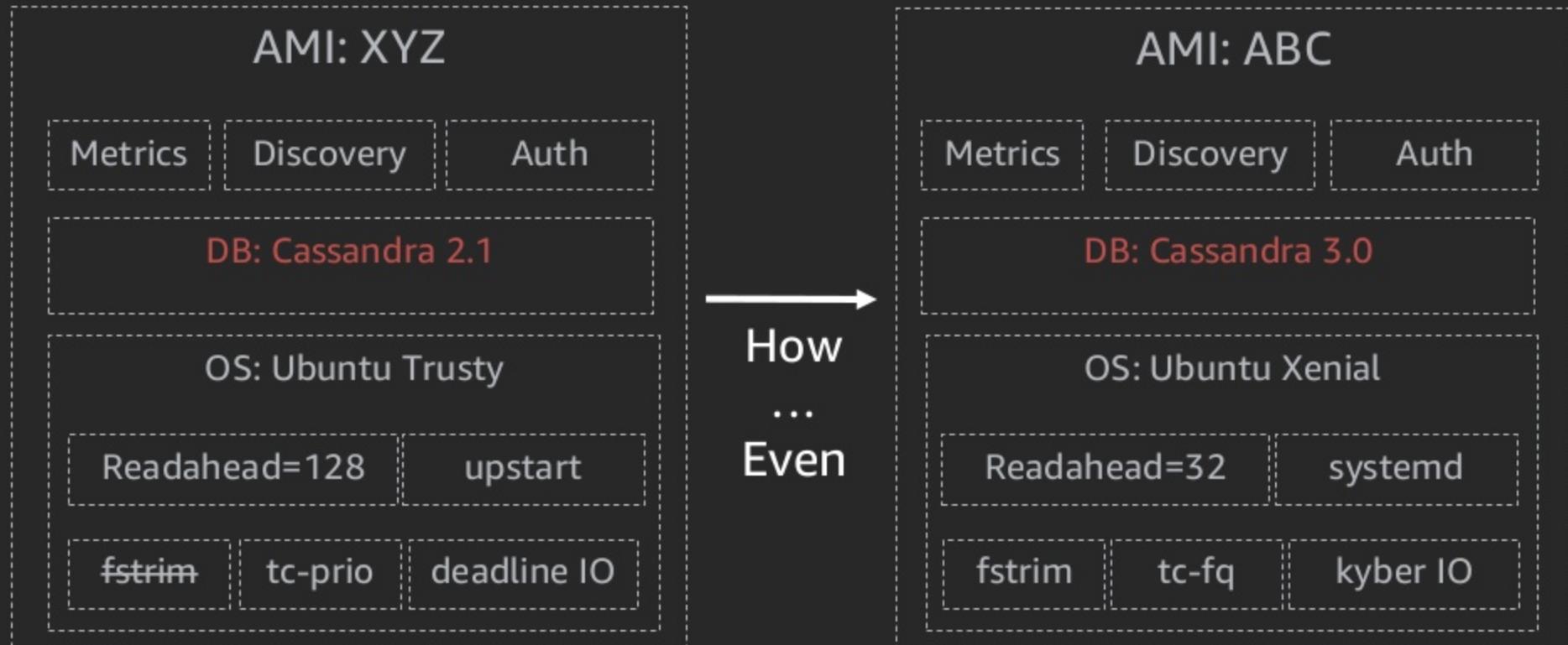
Problem with mutation: Risk



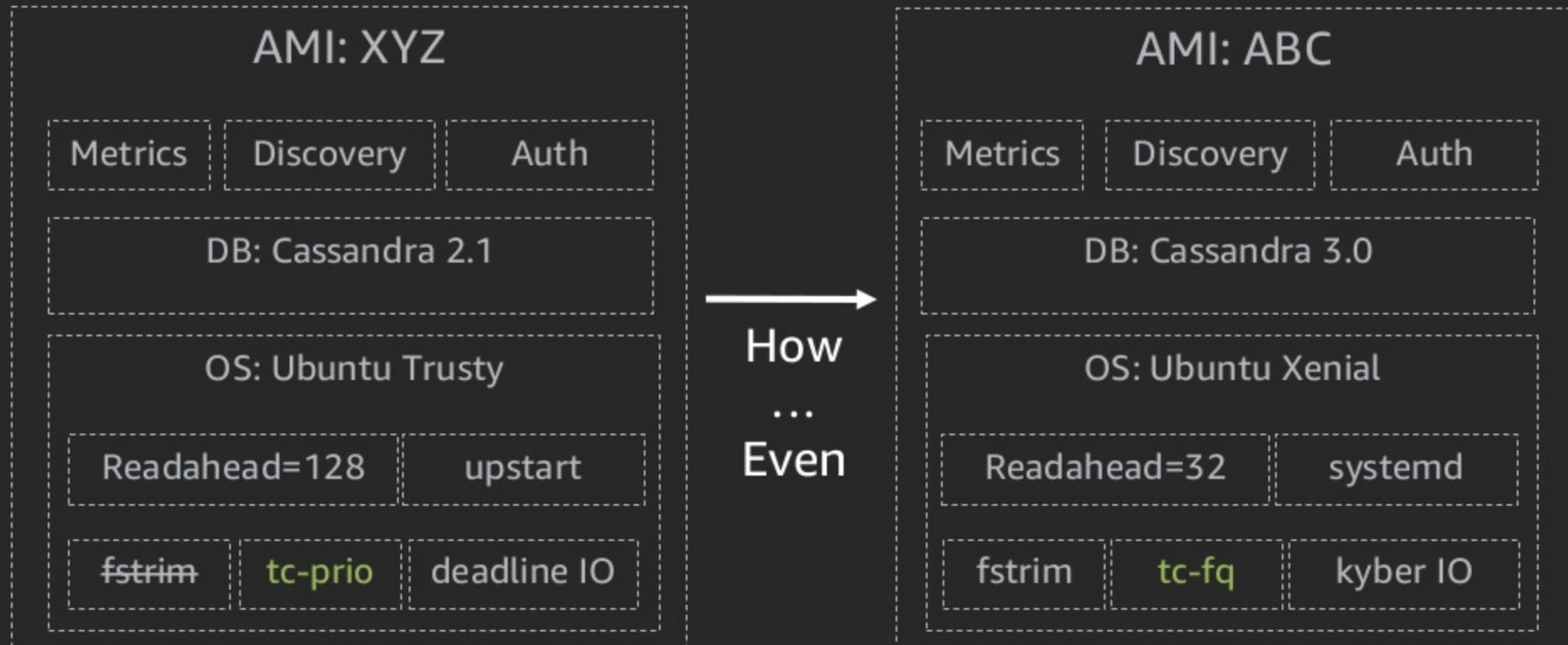
Problem with mutation: Risk



Problem with mutation: Risk



Problem with mutation: Risk



Solution: Image new AMIs directly

1. Prepare

Running: XYZ

Fetch AMI Image ID,
Write to Bootloader

→
reboot

2: Image AMI

Running: Imager

Download ABC Image,
Write to /dev/xvda



3. Bootup

Running: ABC

Load /mnt/data

→
reboot

AMI live upgrade

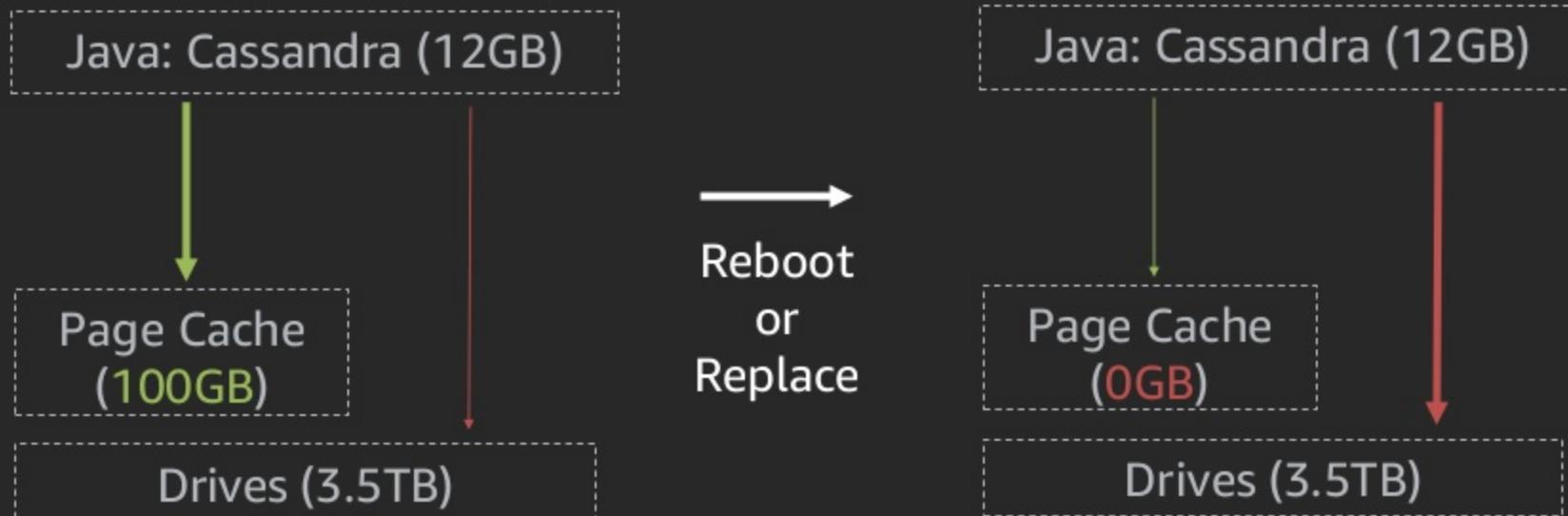
Pros

- Only ever run tested AMIs
- No data transfer or IP change
 - $P(\text{data corruption}) = \text{low}$
- Pretty fast (~10 minutes per node)

Cons

- Probably better ways
 - Containers
 - Active/Passive flashing
- Doesn't work for instance type upgrades
- Confuses your security team

Problem with state transfer: Page cache



Solution: Warm the caches

1. Dump

Running: XYZ

happycache dump

2: Load

Running: ABC

happycache load

3: Startup

Running: ABC

Cassandra



Reboot
or
Replace



Start
Datastore

<https://github.com/hashbrowncipher/happycache>

Problem with state transfer: Scale



Instance type: i3.4xlarge
Dataset: 400GB
Time to replace:
~20 minutes per node

= 4 hours for cluster



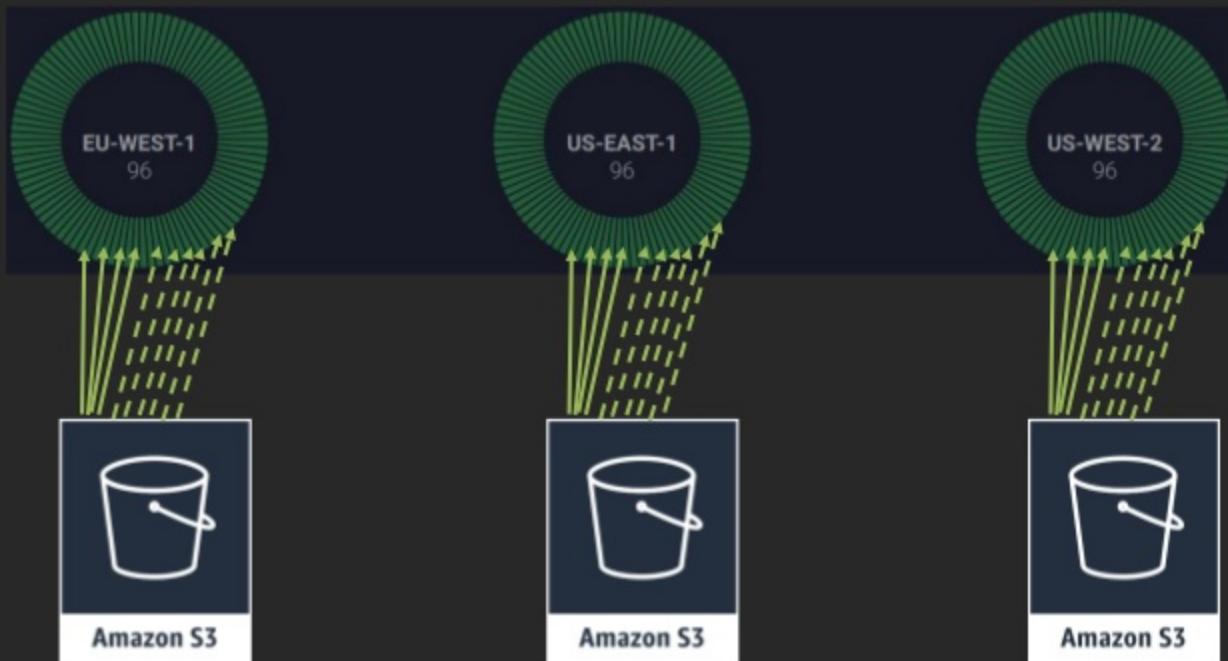
Instance type: i3.4xlarge
Dataset: 400GB
Time to replace:
~20 minutes per node

= 4 days for cluster

Solution: Massive parallelism

Instance type: i3.4xlarge
Dataset: 400GB
4 nodes per AZ at once

Time to replace:
~7 minutes per node + 20
= 3.2 hours for cluster!!



With entire AZ, 50 minutes

Problem with state transfer: Validation is slow

```
$ du -shc .
885G .
885 total

$ find . -type f | xargs -IX du --block-size=1G X | cut -f 1 | histogram.py -l
 1.0000 -      1.7263 [    691]: #####
 1.7263 -      3.1789 [     8]: 
 3.1789 -      6.0841 [     2]: 
 6.0841 -     11.8944 [     1]: 
11.8944 -     46.7566 [     1]: 
46.7566 -   372.1369 [     1]: 
372.1369 -   744.0000 [     1]: 

$ time find . -type f | xargs -IX -P 16 sha256sum X > /dev/null
real    1h22m27.694s
user    1h36m44.131s
sys     4m40.120s
```

Problem with state transfer: Validation is slow

```
1 [||||| 100.0%] 5 [||||| 100.0%] 9 [||||| 100.0%] 13 [||||| 100.0%]
2 [||||| 100.0%] 6 [||||| 100.0%] 10 [||||| 100.0%] 14 [||||| 100.0%]
3 [||||| 100.0%] 7 [||||| 100.0%] 11 [||||| 100.0%] 15 [||||| 100.0%]
4 [||||| 100.0%] 8 [||||| 100.0%] 12 [||||| 100.0%] 16 [||||| 100.0%]
Mem[||||| 16.5G/120G] Tasks: 114, 1222 thr; 19 running
Swp[ 0K/0K] Load average: 6.15 3.13 1.32
Uptime: 4 days, 11:31:04
```

PID	USER	PRI	NI	VIRT	RES	SHR	S	CPU%	MEM%	TIME+	Command
119908	[REDACTED]	20	0	6032	668	596	R	99.6	0.0	0:02.84	sha256sum
119933	[REDACTED]	20	0	6032	676	604	R	99.6	0.0	0:02.42	sha256sum
119890	[REDACTED]	20	0	6032	728	656	R	99.6	0.0	0:02.91	sha256sum
119910	[REDACTED]	20	0	6032	668	600	R	99.6	0.0	0:02.83	sha256sum
119846	[REDACTED]	20	0	6032	836	764	R	99.0	0.0	0:02.93	sha256sum
119895	[REDACTED]	20	0	6032	768	700	R	99.0	0.0	0:02.89	sha256sum
119924	[REDACTED]	20	0	6032	684	612	R	99.0	0.0	0:02.61	sha256sum
119930	[REDACTED]	20	0	6032	832	756	R	99.0	0.0	0:02.59	sha256sum
119938	[REDACTED]	20	0	6032	724	656	R	98.3	0.0	0:02.10	sha256sum
119847	[REDACTED]	20	0	6032	684	612	R	98.3	0.0	0:02.78	sha256sum
119874	[REDACTED]	20	0	6032	680	608	R	96.4	0.0	0:02.88	sha256sum
119876	[REDACTED]	20	0	6032	696	624	R	96.4	0.0	0:02.88	sha256sum

Problem with state transfer: Validation is slow

Device:	tps	MB_read/s	MB_wrtn/s	MB_read	MB_wrtn
xvda	2.00	0.00	0.01	0	0
nvme1n1	4623.00	144.00	0.00	287	0
nvme0n1	4638.00	144.35	410 MBps ... very bad	0	0
md0	9312.50	409.17	0.00	818	0

Device:	tps	MB_read/s	MB_wrtn/s	MB_read	MB_wrtn
xvda	2.00	0.00	0.01	0	0
nvme1n1	2508.50	78.37	0.04	156	0
nvme0n1	2511.00	78.34	166 MBps single threaded ... Sad	0	0
md0	5019.50	166.48	0.11	332	0

Solution: Use a better hash (xxHash)

```
$ instance-type  
i3.4xlarge  
  
$ du -shc .  
885G .  
885 total  
  
$ time find . -type f | xargs -IX -P 16 ~/xxHash/xxh64sum --quiet X > /dev/null  
  
real    12m10.694s  
user    1m44.548s  
sys     7m17.611s
```

<https://github.com/Cyan4973/xxHash>

Solution: xxHash

Device:	tps	MB_read/s	MB_wrtn/s	MB_read	MB_wrtn
xvda	4.50	0.01	0.03	0	0
nvme1n1	30558.00	1908.87	0.05	3817	0
nvme0n1	30567.50	1908.78	4 GBps is more like it	3817	0
md0	61126.00	3817.69	0.07	7635	0

Device:	tps	MB_read/s	MB_wrtn/s	MB_read	MB_wrtn
xvda	2.00	0.00	0.01	0	0
nvme1n1	8937.50	558.53	0.00	1117	0
nvme0n1	8938.00	558.55	1.1 GBps single threaded	1117	0
md0	17875.50	1117.09	0.00	2234	0

Manage state Netflix style

Don't manage

This is a very
viable option

Often the right
choice

Mutate

Make it safe
Avoid drift

Client

Careful with
consistency
Coordination is
key

Server

Safest option
Can be fast as
well!

Thank you!

Joey Lynch
josephl@netflix.com



Please complete the session
survey in the mobile app.