# Logs Told Us It Was DNS
# It Felt Like DNS
# It Had To Be DNS
## *It Wasn't DNS*

Elijah Andrews, Hemanth Malla
2023-03-21, SREcon23 Americas, Santa Clara

DATADOG

# Who are we?

**Elijah Andrews**

Software Engineer

*Datadog*

🐦 **@elijahca**

**Hemanth Malla**

Software Engineer

*Datadog*

🐦 **@hemanthmalla**

DATADOG

2

# Datadog

Over 600 integrations
Over 5,000 employees
Over 23,000 customers
Runs on millions of hosts
Tens of trillions of events per day

Tens of thousands of nodes
Hundreds of thousands of pods
100s of k8s clusters with 100-4000 nodes
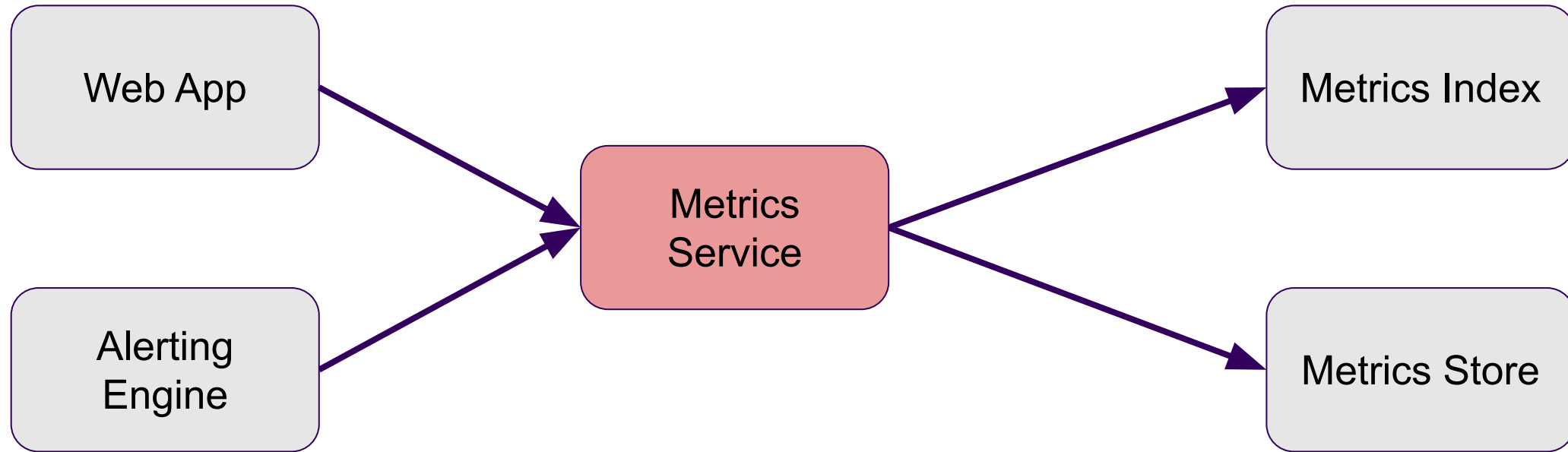Multi-cloud
Very fast growth



DATADOG

# Kubernetes primer

- *pods* are scheduled on *nodes*

- applications run in pods

- each pod has a unique IP address

- *Cilium* is our Container Network Interface (CNI)
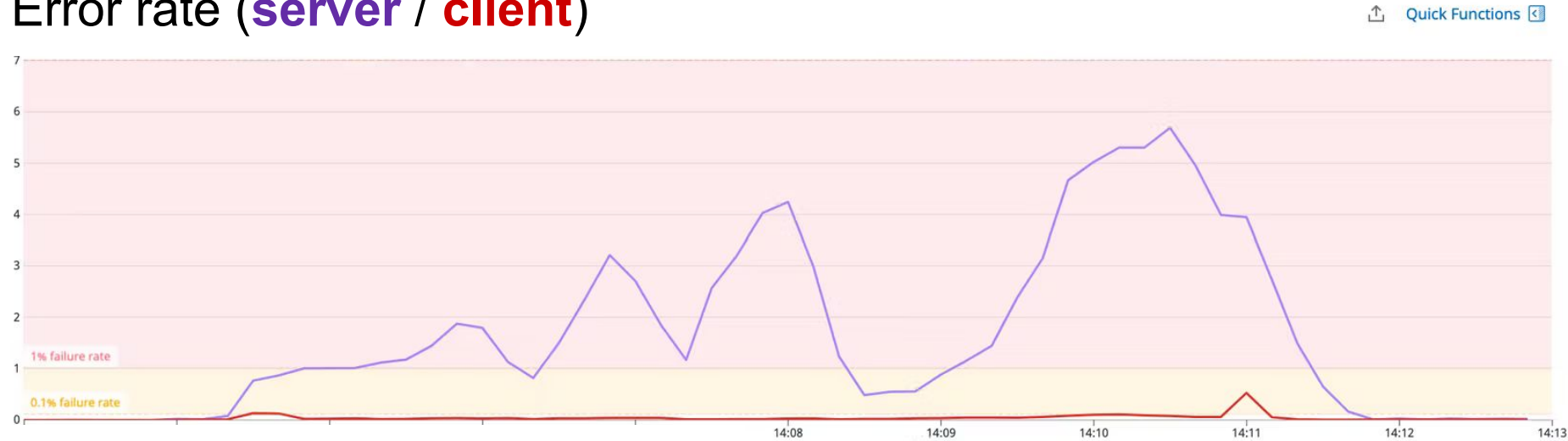
- Cilium glues node, pod, and AWS networking together
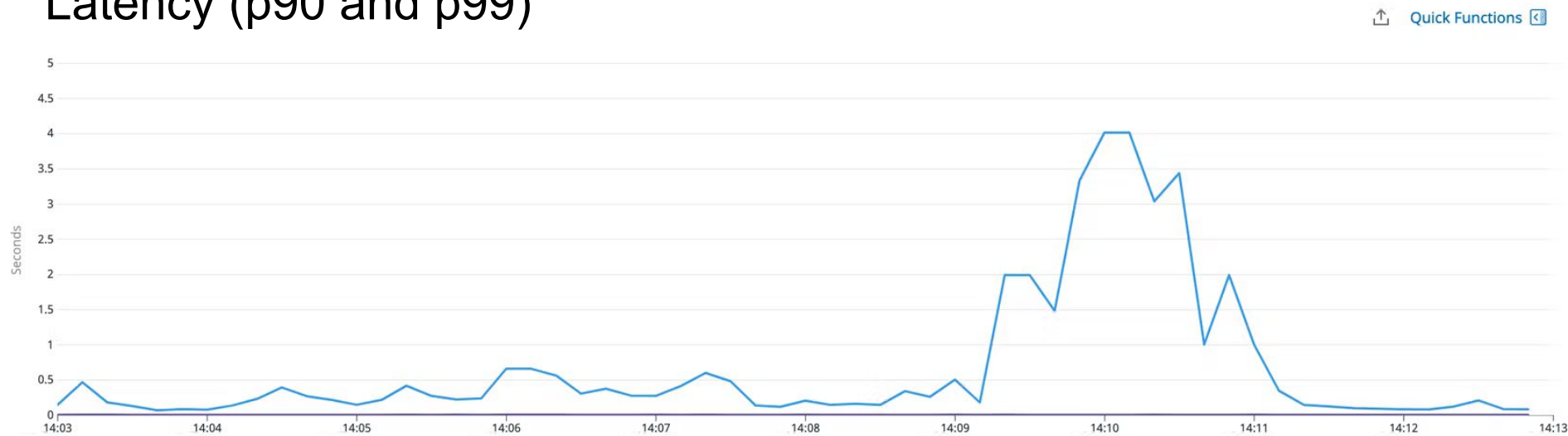
**DATADOG**

# How it all started

# Applications involved

```
[Web App] ──────┐
                ├──→ [Metrics Service] ──┬──→ [Metrics Index]
[Alerting       │                        │
 Engine] ───────┘                        └──→ [Metrics Store]
```

DATADOG

# Metrics service errors during rollouts

Error rate (**server** / **client**)

⬆ Quick Functions ◁



Latency (p90 and p99)

⬆ Quick Functions ◁

# It's always DNS

Logs told us it was DNS

DATADOG

# It's always DNS

Logs told us it was DNS

It looked like DNS

DATADOG

# It's always DNS

Logs told us it was DNS

It looked like DNS

**It had to be DNS**
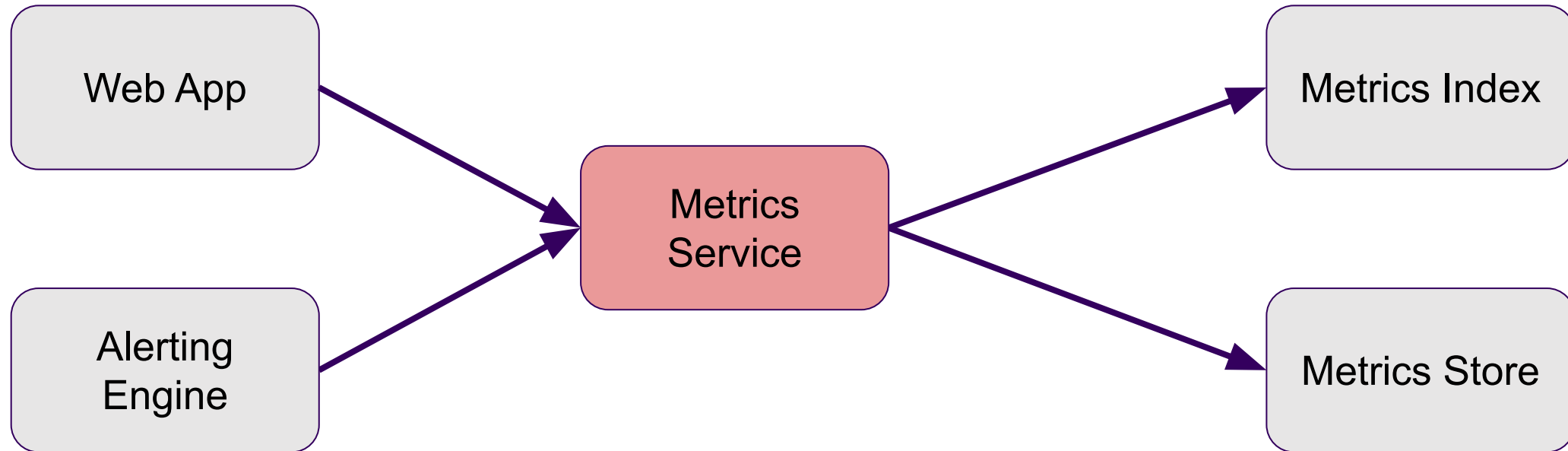
DATADOG

# It's always DNS

Logs told us it was DNS

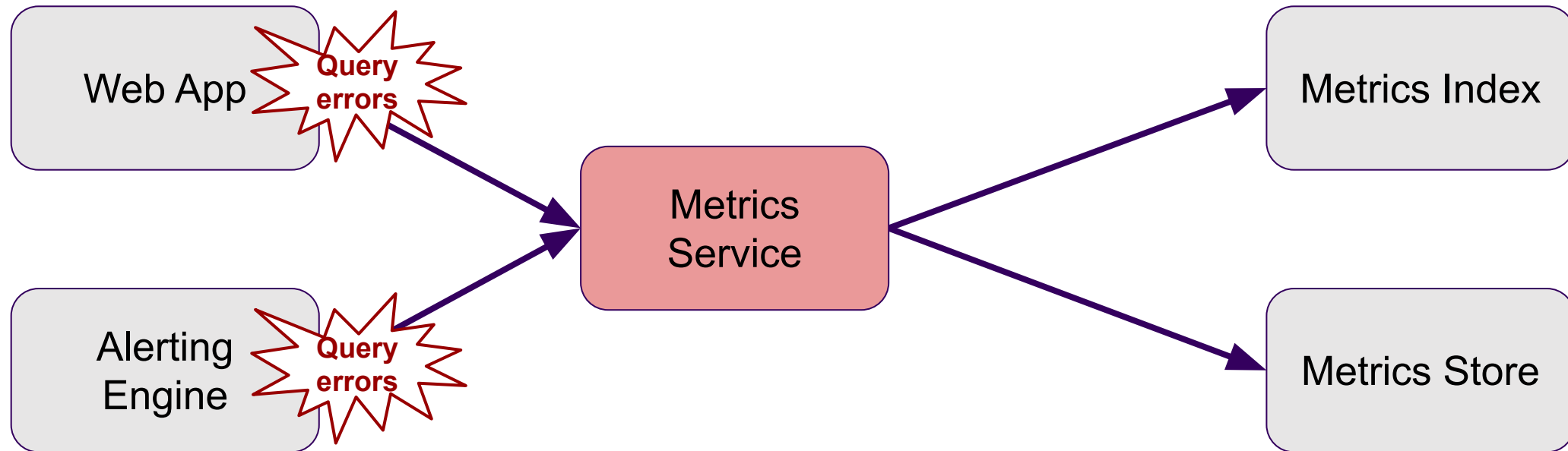It looked like DNS

**It had to be DNS**
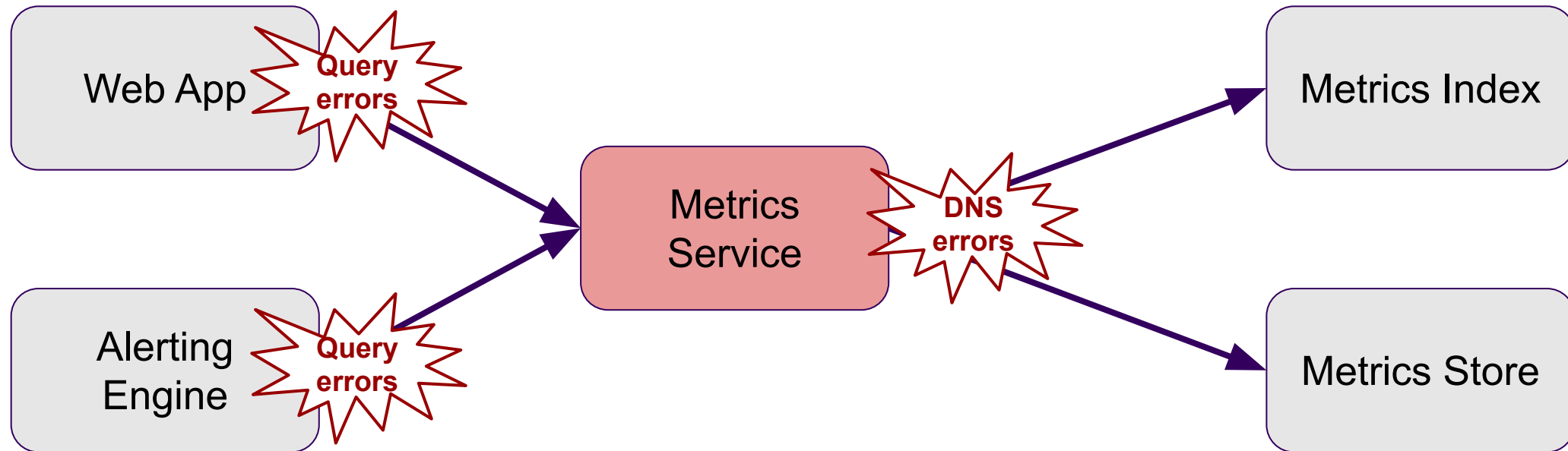
**Right?**

**DATADOG**

# Chapter 1: DNS

# Applications involved
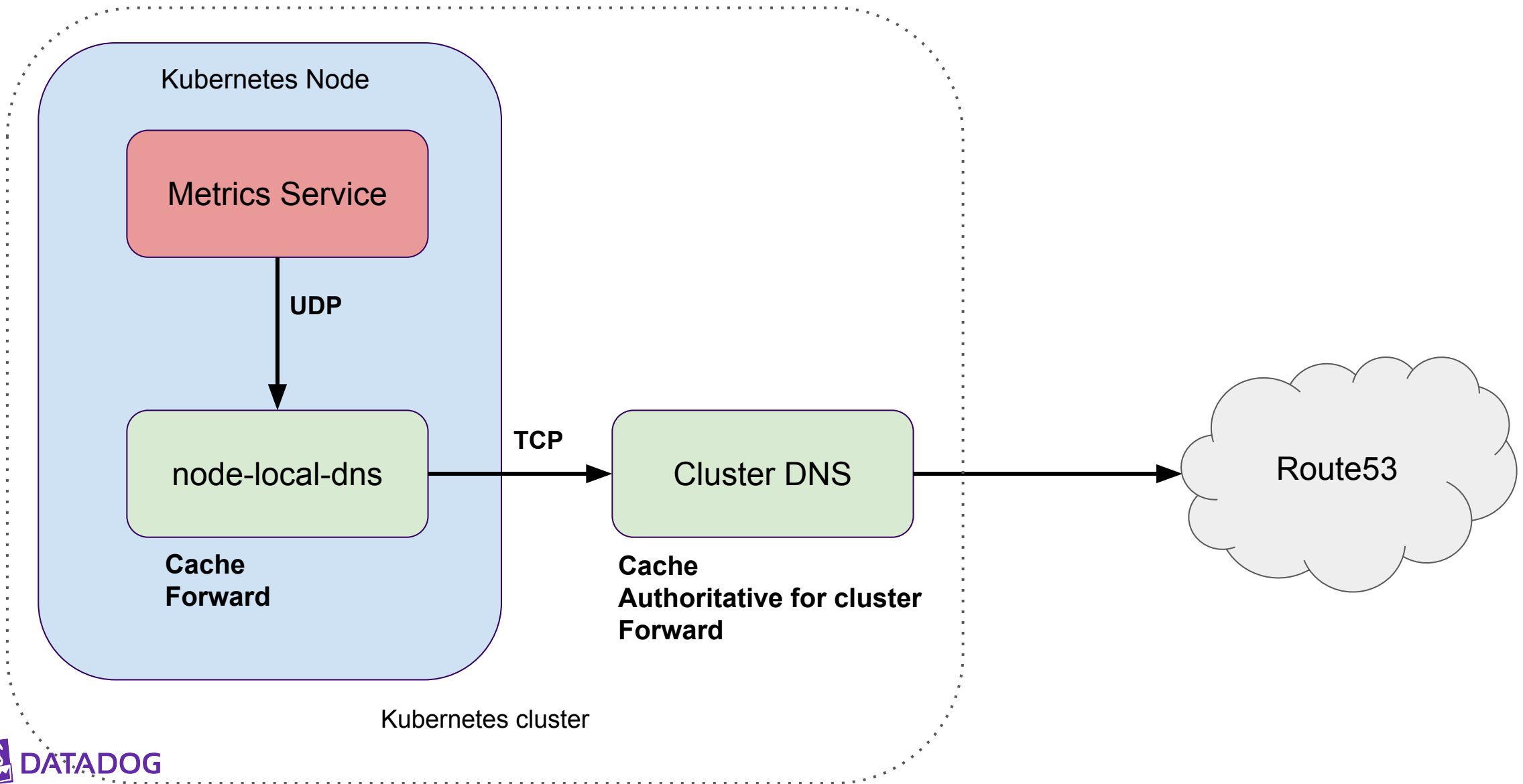
# Applications involved

# Applications involved

# DNS setup

Kubernetes Node

**Metrics Service**

**UDP**

node-local-dns

**Cache
Forward**

**TCP**

Cluster DNS

**Cache
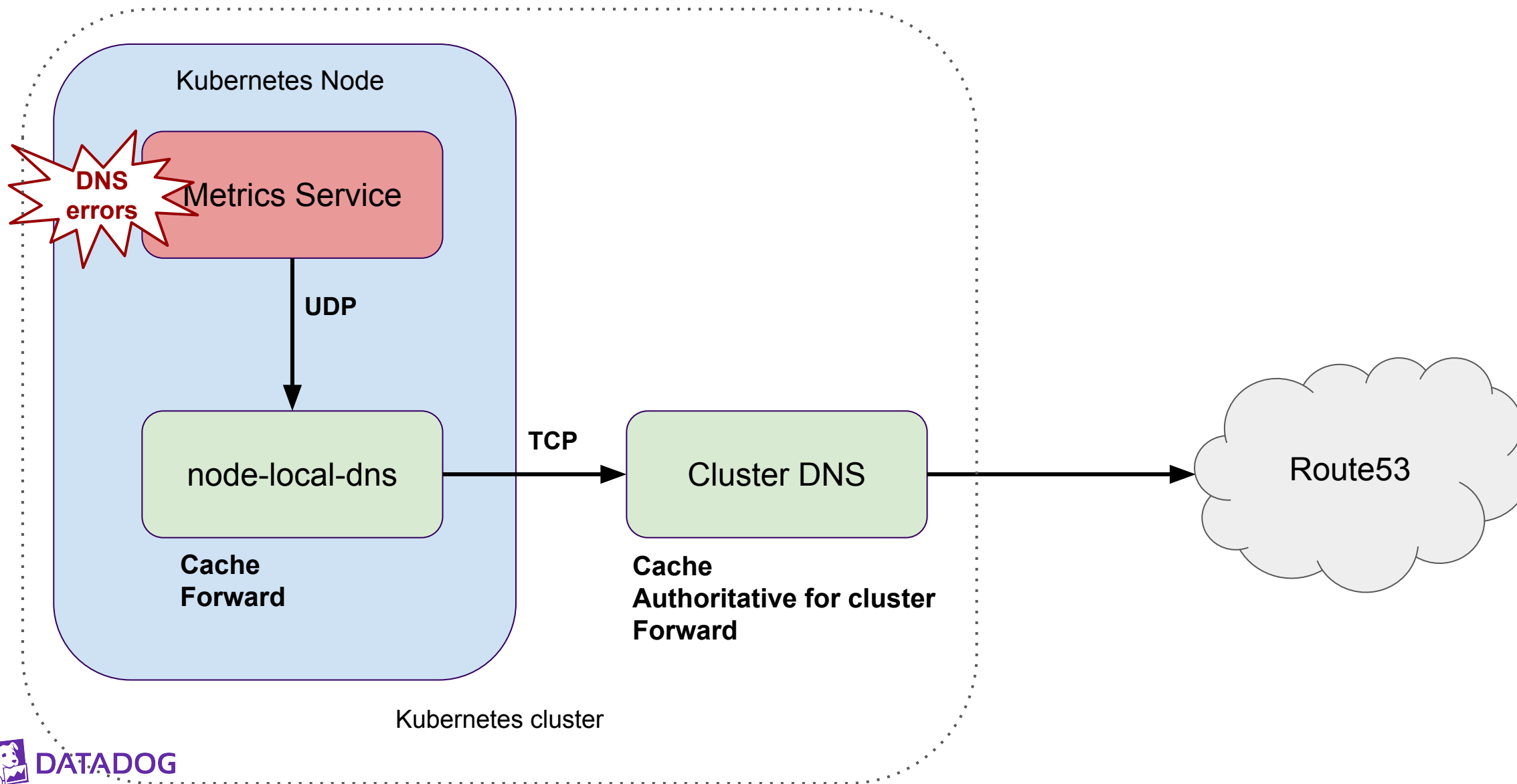Authoritative for cluster
Forward**

Route53

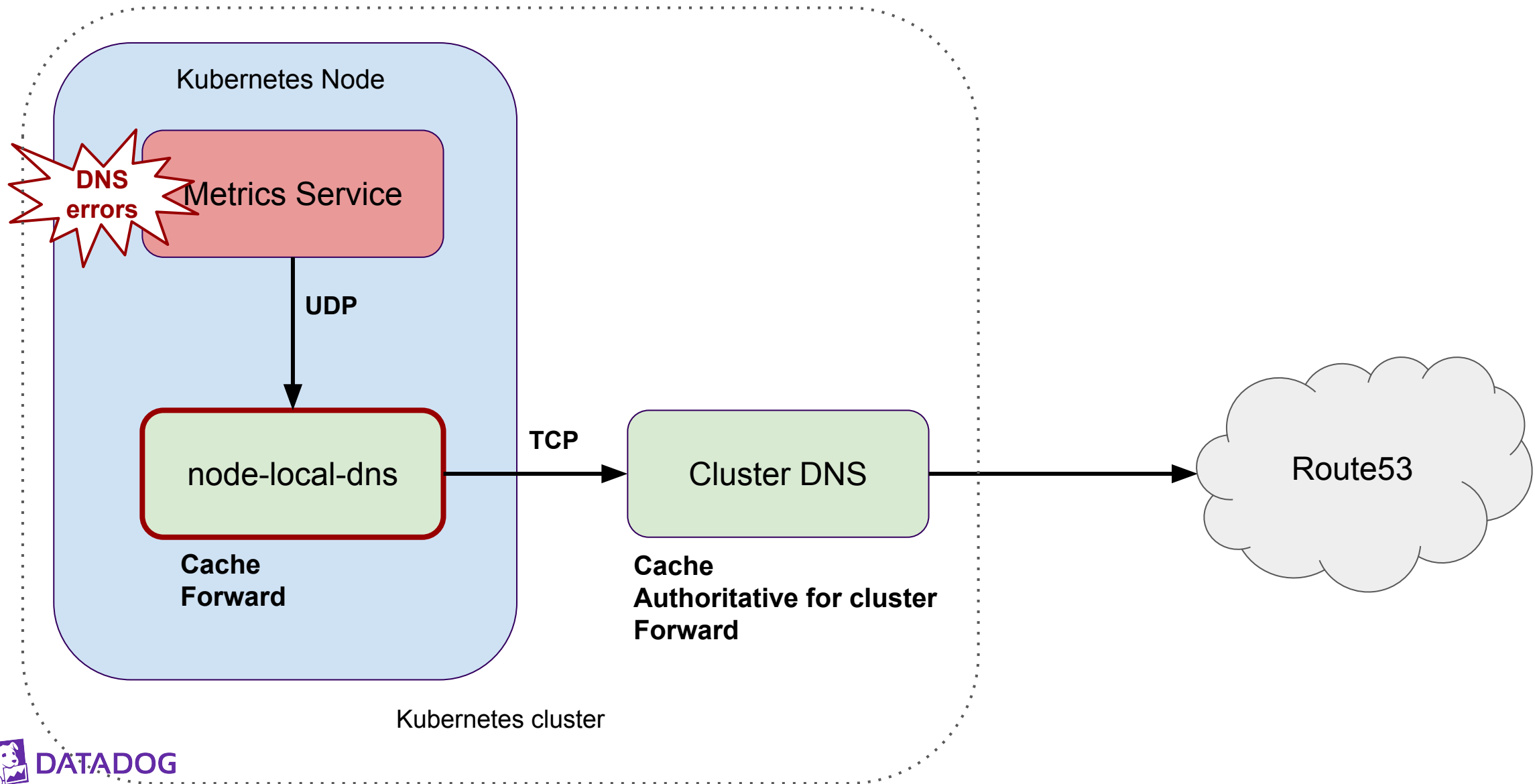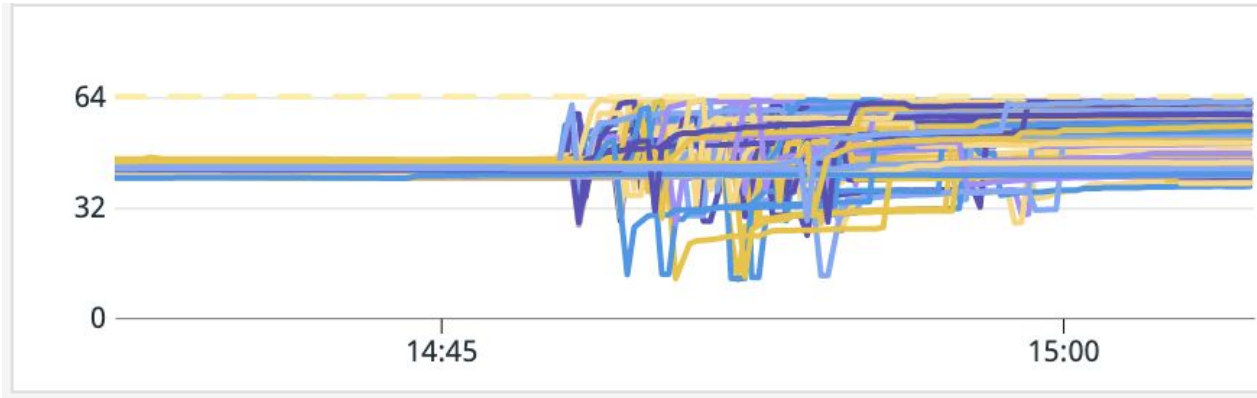Kubernetes cluster

DATADOG

# DNS setup

# DNS setup

# Node-Local-DNS (NLD)

NLD Memory per pod on Metrics Service hosts (and limit)



ran out of memory (OOM killed)

during rollouts

Should *never* happen

# Node-Local-DNS (NLD)

NLD Memory per pod on Metrics Service hosts (and limit)





max_concurrent is working

Sizing is wrong

# Node-local-dns, 64MB => 256MB

NLD Memory per pod on Metrics Service hosts (and limit)



Metrics Service Error rate (**server** / **client**)



Node-local-dns max concurrent rejects



No more OOM-kills

*But* not any better for Metrics Service

# Too many queries at startup?

Node-local-dns queries

Max_concurrent: **1000**

Upstream queries: ~5ms

=> NLD should do **> 200k rps**

=> with **<400 rps** we hit max_concurrent

What's happening?

# Too many queries at startup?

## Node-local-dns queries



## Node-local-dns forward healthcheck failures



Upstream marked unhealthy

Upstream is TCP

Connections are reused

  but expire=10s

*NLD can't create connections?*

# Why we hit **max_concurrent**

- NLD can't establish connections to upstreams

- The Forward plugin has a 5s timeout by default

- Incoming queries occupy a query slot for 5s

=> We hit max_concurrent=1000 with only 200rps

# Networking issues?

Throughput (Gb/s) on Metrics Service nodes (+: Received / -: Sent)



m5.4xlarge

    Max: 10Gb/s

      Sustained: 5Gb/s

=> looks ok

# Networking issues?

Throughput (Gb/s) on Metrics Service nodes (+: Received / -: Sent)



m5.4xlarge

    Max: 10Gb/s

    Sustained: 5Gb/s

=> looks ok

TCP retransmits on Metrics Service nodes



But we are dropping packets

Microbursts?

=> Elastic Network Adapter (ENA) metrics

DATADOG

# Status

- DNS errors in Metrics Service on rollouts

- Node-local-DNS can't establish connections


=> *Network issue?*

DATADOG

# Chapter 2: AWS Networking

# Are we bursting over the instance limits?

Metrics Service Error rate (**server** / **client**)



Average DNS response time by pod



ENA: Bandwidth exceeded (**+:in** / **-:out**)



We are saturating the interface

But no correlation with errors

DATADOG

# Are we bursting over the instance limits?

Metrics Service Error rate (**server** / **client**)

ENA: Conntrack exceeded

Average DNS response time by pod

**conntrack allowance exceeded?**

# aws.ec2.conntrack_allowance_exceeded

*The number of packets dropped because connection tracking exceeded the maximum for the instance and new connections could not be established. This can result in packet loss for traffic to or from the instance*

Connection tracking is required for security groups (stateful)

# Let's test with network optimized instances

ENA: Bandwidth exceeded



**blue**/**yellow** => m5.4xlarge

**purple**/**grey** (~0) => m5n.4xlarge

Promising!

# Let's test with network optimized instances

## ENA Limits - conntrack exceeded

| | | | |
|---|---|---|---|
| 2k | | | |
| 1k | | | |
| 0k | 17:45 | 18:00 | 18:15 |

## Average DNS response time by pod

| | | | |
|---|---|---|---|
| 100 | | | |
| 50 | | | |
| 0 | 17:45 | 18:00 | 18:15 |

■ m5.4xlarge

■ m5n.4xlarge

No impact on

- Conntrack

- Metrics Service errors / latency

# What about bigger instances?

Average DNS response time by pod



■ m5.4xlarge

■ m5.8xlarge

Much better!

ENA Limits - conntrack exceeded



TCP retransmits on Metrics Service nodes



DATADOG

# Conntrack limits?

**From AWS**

- Hypervisor conntrack can track hundreds of thousands of flows

- m5.8xlarge : can track 2x the flows compared to m5.4xlarge

- m5n.4xlarge : same as m5.4xlarge

=> Makes sense based on our tests

# How can we saturate this conntrack?

Conntrack count for 4xls (blue) and 8xls (red)

80k

60k

40k

20k

0k
13:30                    14:00                    14:30

Stable state: **~13k** connections

Rollouts: **~60k**

Pretty high but **60k vs X00k ????**

# VPC Flow Logs

- Capture IP flow information on Elastic Network Interfaces (ENI)

- **Flow level: 5-tuple, 2 flows per TCP connection**

- Flow record: 5 tuple, bytes, packets, TCP flags...

- Aggregated every 1mn and delivered to S3

- **Not always complete**

- Huge amount for large VPCs (we filtered with Athena)

**DATADOG**

# Flows from a Metrics Service node

Egress flows by source

— Node IP  — Old IP  — New IP



Old pod IP disappears after ~60s

Spike in flows at pod deletion

50k flows in 1mn feels very high

# What about ingress flows?

Egress flows by source

— Node IP  — Old IP  — New IP



Ingress flows by destination

— Node IP  — Old IP  — New IP



Ingress flows should ~match Egress

Very weird second spike

**What are these flows?**

39

# Zoom on ingress flows to old IP

Ingress Flows by TCP flag

— None — FIN — SYN



**None**: already established

**FIN**: terminating

**SYN**: reconnect attempts: **130k over 90s!**

# What about egress?

**Ingress Flows by TCP flag**

— None  — FIN  — SYN

RST match first SYN spike

What about this second spike?

**Egress Flows by TCP flag**

— None  — FIN  — RST

# Why do we get RST for a few seconds only?

- Metrics Service performs a grpc.GracefulStop with 10s timeout

  - Server stops accepting new connections

  - Server waits for existing RPC to finish

  - Server tells clients to disconnect (HTTP2 GoAway)

- During these 10s, incoming connection attempts get an RST

- After these 10s, the pod is deleted and its IP is not bound by anything

DATADOG

# Where are these attempts coming from?

Only a few IPs => Alerting Engine

# Where are these attempts coming from?

Web App

Alerting Engine

Metrics Service

Metrics Index

Metrics Store

#conntrack entries on Alerting Engine nodes

AME conntrack entries by host

200k

100k

0k
17:00    17:05    17:10    17:15

**Seems to confirm!**

DATADOG

# Status

- DNS errors in Metrics Service on rollouts

- Node-local-DNS can't establish connections

- AWS conntrack for instance is saturated

- Alerting Engine is SYN-Flooding Metrics Service on rollouts

=> *Why don't we see these connections on Metric Service Nodes?*

# Chapter 3: Node Networking

# Cilium Architecture



*Source : Cilium documentation*

# Routing on nodes

# Routing on nodes

Metrics Service Node

Cilium Operator

ENI IPs:
- 10.x.y.z
- ....

ens6
(pods)

vethX

Metrics Service
IP: 10.x.y.z

Cilium agent

ens5
(node)

**main route table**
0.0.0.0/0 => ens5
**10.x.y.z => vethX**

# Routing on nodes

Metrics Service Node

Cilium Operator

ENI IPs:
- 10.x.y.z
- ....

ens6
(pods)

vethX

Metrics Service
IP: 10.x.y.z

Cilium agent

**main route table**
0.0.0.0/0 => ens5
**10.x.y.z => vethX**

ens5
(node)

**ens6 route table**
**0.0.0.0/0 => ens6**

**ip rules**
to other pods => main route table
**from 10.x.y.z => ens6 route table**

# Stable state

**Metrics Service Node**

**Alerting node conntrack**
10.a.b.c:portN => 10.x.y.z:6415 SYN

ENI IPs:
- 10.x.y.z
- ….

ens6
(pods)

SYN

vethX

**Metrics Service
IP: 10.x.y.z**

**Alerting pod
IP: 10.a.b.c**

**Hypervisor conntrack**
10.a.b.c:portN => 10.x.y.z:6415 SYN

**Node conntrack**
10.a.b.c:portN => 10.x.y.z:6415 SYN

**main route table**
0.0.0.0/0 => ens5
**10.x.y.z => vethX**

**ens6 route table**
**0.0.0.0/0 => ens6**

**ip rules**
to other pods => main route table
**from 10.x.y.z => ens6 route table**

Connection tracked with **SYN-RECV** state

# Stable state

**Metrics Service Node**

**Alerting node conntrack**
10.a.b.c:portN => 10.x.y.z:6415 EST

ENI IPs:
- 10.x.y.z
- ....

ens6
(pods)

SYN

vethX

Metrics Service
IP: 10.x.y.z

Alerting pod
IP: 10.a.b.c

SYN/ACK

**Hypervisor conntrack**
10.a.b.c:portN => 10.x.y.z:6415 EST

**Node conntrack**
10.a.b.c:portN => 10.x.y.z:6415 EST

**main route table**
0.0.0.0/0 => ens5
**10.x.y.z => vethX**

**ens6 route table**
**0.0.0.0/0 => ens6**

**ip rules**
to other pods => main route table
**from 10.x.y.z => ens6 route table**

Connection transitions to **ESTABLISHED** state

# What happens on pod deletion?

Metrics Service Node

ENI IPs:
- 10.x.y.z
- ....

ens6
(pods)

ens5
(node)

**main route table**
0.0.0.0/0 => ens5
~~10.x.y.z => vethX~~

**ens6 route table**
**0.0.0.0/0 => ens6**

**ip rules**
to other pods => main route table
~~from 10.x.y.z => ens6 route table~~

DATADOG

# What about traffic to old IP?

Metrics Service Node

ens6
(pods)

ENI IPs:
- 10.x.y.z
- ….

**?**

SYN

ens5
(node)

**main route table**
0.0.0.0/0 => ens5
~~10.x.y.z => vethX~~

**ens6 route table**
**0.0.0.0/0 => ens6**

**ip rules**
to other pods => main route table
~~from 10.x.y.z => ens6 route table~~

DATADOG

54

# Let's simulate

Delete pod with IP 10.x.y.z on nodeB and attempt to connect from nodeA

Connection attempt

```
nodeA:~$ nc -vz 10.x.y.z 12345
```

# Let's simulate

Delete pod with IP 10.x.y.z on nodeB and attempt to connect from nodeA

Connection attempt

```
nodeA:~$ nc -vz 10.x.y.z 12345
```

On nodeB => SYN without an answer

```
nodeB:~$ sudo tcpdump -pni ens6 "port 12345"
listening on ens6, link-type EN10MB (Ethernet), capture size 262144 bytes
08:28:52.086251 IP 10.a.b.c.51718 > 10.x.y.z.12345: Flags [S], seq 4126537246, win 26883, options [mss
8961,sackOK,TS val 2002199904 ecr 0,nop,wscale 9], length 0
```

# Let's simulate

Delete pod with IP 10.x.y.z on nodeB and attempt to connect from nodeA

Connection attempt

```
nodeA:~$ nc -vz 10.x.y.z 12345
```

On nodeB => SYN without an answer

```
nodeB:~$ sudo tcpdump -pni ens6 "port 12345"
listening on ens6, link-type EN10MB (Ethernet), capture size 262144 bytes
08:28:52.086251 IP 10.a.b.c.51718 > 10.x.y.z.12345: Flags [S], seq 4126537246, win 26883, options [mss
8961,sackOK,TS val 2002199904 ecr 0,nop,wscale 9], length 0
```

Where would the SYN be routed to? => Reverse Path filter!

```
$ ip route get 10.x.y.z from 10.a.b.c iif ens6
RTNETLINK answers: Invalid cross-device link
```

DATADOG

# Let's simulate

Delete pod with IP 10.x.y.z on nodeB and attempt to connect from nodeA

Connection attempt

```
nodeA:~$ nc -vz 10.x.y.z 12345
```

On nodeB => SYN without an answer

```
nodeB:~$ sudo tcpdump -pni ens6 "port 12345"
listening on ens6, link-type EN10MB (Ethernet), capture size 262144 bytes
08:28:52.086251 IP 10.a.b.c.51718 > 10.x.y.z.12345: Flags [S], seq 4126537246, win 26883, options [mss
8961,sackOK,TS val 2002199904 ecr 0,nop,wscale 9], length 0
```

Where would the SYN be routed to? => Reverse Path filter!

```
$ ip route get 10.x.y.z from 10.a.b.c iif ens6
RTNETLINK answers: Invalid cross-device link
```

Sure enough, martian packet warning in kernel logs

```
Oct 28 08:25:54 nodeB kernel: IPv4: martian source 10.x.y.z from 10.a.b.c, on dev ens6
```

# Reverse Path filtering

- Security feature from the kernel to prevent IP spoofing

  - If return path uses incoming interface accept the packet

  - Otherwise drop it

- Log these events : "Martian Packets"

- Loose mode: only drop if there is no return route

# Back to our node

**Metrics Service Node**

**Reverse path filter**
Return path? => ens5
=> Martian packet
=> Drop

**ens6**
**(pods)**

ENI IPs:
- 10.x.y.z
- ....

**Drop**

**main route table**
0.0.0.0/0 => ens5
~~10.x.y.z => vethX~~

**ens5**
**(node)**

**ens6 route table**
**0.0.0.0/0 => ens6**

**ip rules**
to other pods => main route table
~~from 10.x.y.z => ens6 route table~~

```
Oct 28 08:25:54 kernel: IPv4: martian source
10.x.y.z from 10.a.b.c, on dev ens6
```

# What about conntracks?

**Metrics Service Node**

**Reverse path filter
Return path? => ens5
=> Martian packet
=> Drop**

**Drop**

ens6
(pods)

ENI IPs:
- 10.x.y.z
- ….

SYN

**Node conntrack**

**main route table**
0.0.0.0/0 => ens5

**ens6 route table**
**0.0.0.0/0 => ens6**

**ip rules**
to other pods => main route table

**Alerting node conntrack**
**10.a.b.c:portN** => **10.x.y.z:6415** SYN
…
**10.a.b.c:portZ** => **10.x.y.z:6415** SYN

Alerting pod
IP: 10.a.b.c

**Hypervisor conntrack**
**10.a.b.c:portN** => **10.x.y.z:6415** SYN
…
**10.a.b.c:portZ** => **10.x.y.z:6415** SYN

*Repeated for all Alerting clients*

Connections in **SYN-RECV** state expire after **60s**

# But, we use "loose" mode

```
$ ip route get 10.x.y.z from 10.a.b.c iif ens6
RTNETLINK answers: Invalid cross-device link

$ sysctl net.ipv4.conf.ens6.rp_filter
net.ipv4.conf.ens6.rp_filter = 2
```

- rp_filter = 2 => loose mode
- Loose + default route (ens5) => we should not drop
- What's happening?

# Let's have a look

https://github.com/torvalds/linux/blob/master/net/ipv4/fib_frontend.c#L344

```
336    /* Given (packet source, input interface) and optional (dst, oif, tos):
337     * - (main) check, that source is valid i.e. not broadcast or our local
338     *   address.
339     * - figure out what "logical" interface this packet arrived
340     *   and calculate "specific destination" address.
341     * - check, that packet arrived from expected physical interface.
342     * called with rcu_read_lock()
343     */
344    static int __fib_validate_source(struct sk_buff *skb, __be32 src, __be32 dst,
345                                     u8 tos, int oif, struct net_device *dev,
346                                     int rpf, struct in_device *idev, u32 *itag)
347    {
```

# Let's have a look

```
416    e_rpf:
417            return -EXDEV;
```

```
22    #define EXDEV          18     /* Cross-device link */
```

DATADOG

# Let's have a look

```
416    e_rpf:
417            return -EXDEV;


408    last_resort:
409            if (rpf)
410                    goto e_rpf;
```

# Let's have a look

```
416    e_rpf:
417            return -EXDEV;


408    last_resort:
409            if (rpf)
410                    goto e_rpf;



395            if (no_addr)
396                    goto last_resort;
```

# Let's have a look

```
416   e_rpf:
417            return -EXDEV;



408   last_resort:
409          if (rpf)
410                    goto e_rpf;



395          if (no_addr)
396                    goto last_resort;          Interface IP check is made after evaluating
                                                   loose mode


367     no_addr = idev->ifa_list == NULL;          ifa_list => List of IPs associated with device
```

# _But_ pod interfaces don't have IPs assigned

Let's test

```
$ ip route get 10.x.y.z from 10.a.b.c iif ens6
RTNETLINK answers: Invalid cross-device link
```

Expected, Let's now give ens6 a random IP unrelated to our network

```
$ ip addr add 192.168.1.1/32 dev ens6

$ ip route get 10.x.y.z from 10.a.b.c iif ens6
10.x.y.z from 10.a.b.c via 10.m.n.1 dev ens5
    cache iif ens6
```

We are hitting reverse path filtering because the pod interface has no IP...
- Recent versions of Cilium give it an IP
- If it has an IP, SYN are still dropped but conntrack sizes are consistent (and no martian packet warnings)
- We contributed a PR to make old IPs unreachable and send ICMP errors to clients
  https://github.com/cilium/cilium/pull/18505

DATADOG

# Status

- DNS errors in Metrics Service on rollouts

- Node-local-DNS can't establish connections

- AWS conntrack for instance is saturated

- Alerting Engine is SYN-Flooding Metrics Service on rollouts

- Conntracks are not consistent because Reverse Path Filtering drops SYNs

- We hit Reverse Path filtering because of an edge case in the kernel

=> *Why do we have so many SYNs?*
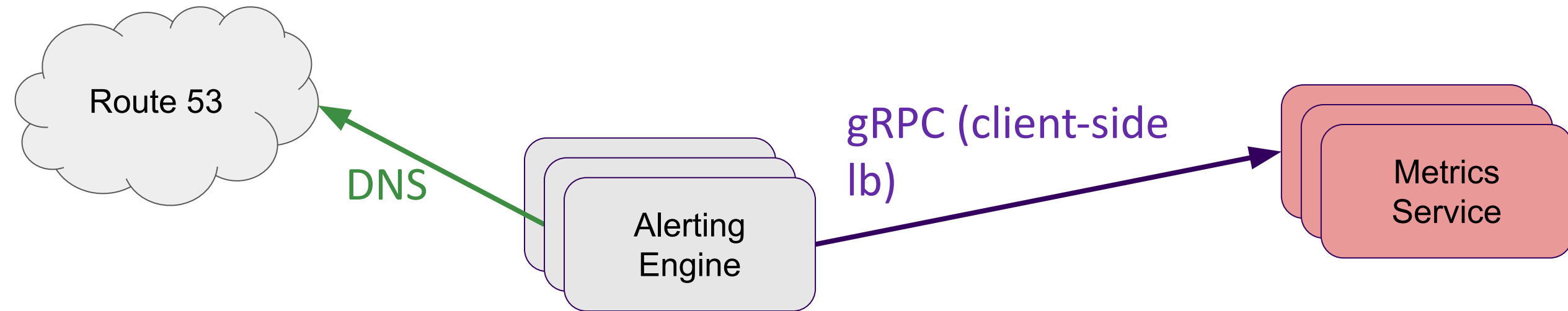
**DATADOG**

# Chapter 4: gRPC client configuration

# 2 questions

1. Why were clients sending SYN requests for so long?

2. Why were clients sending SYN requests so frequently?

**DATADOG**
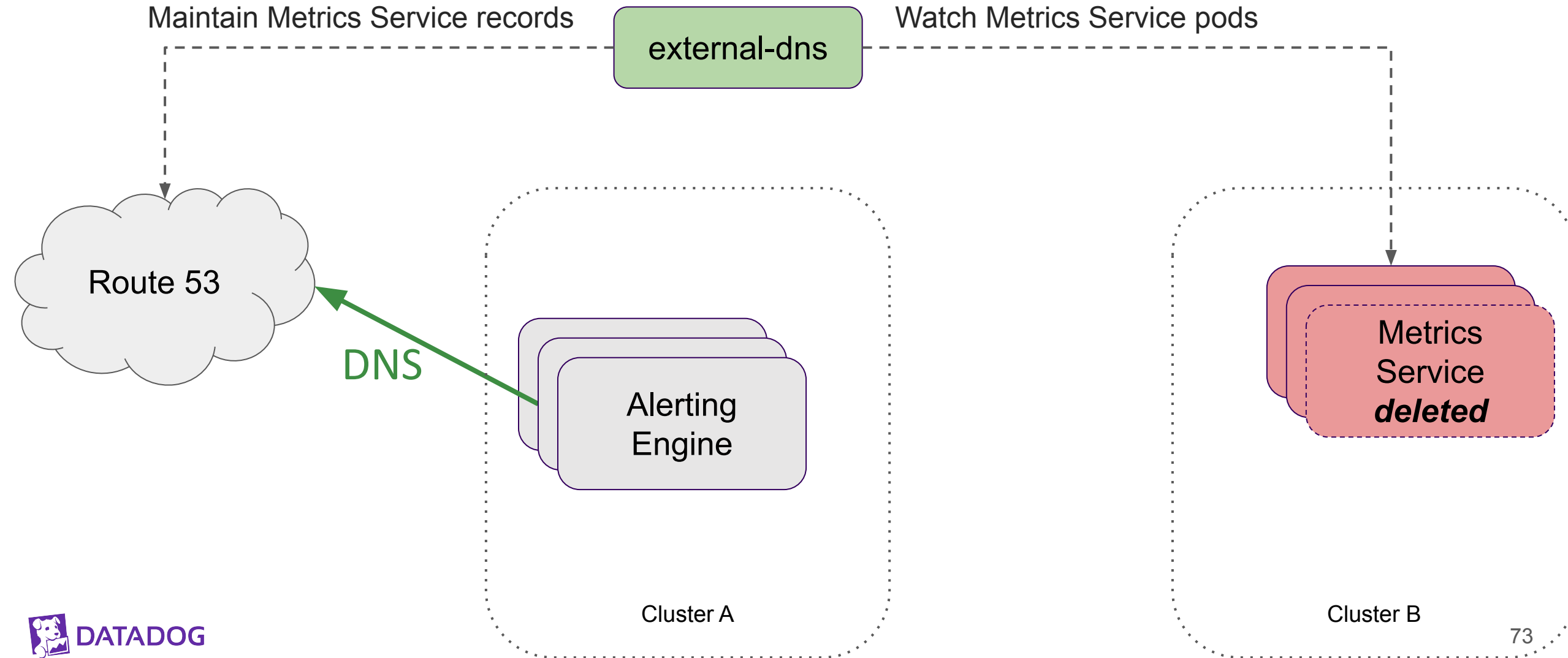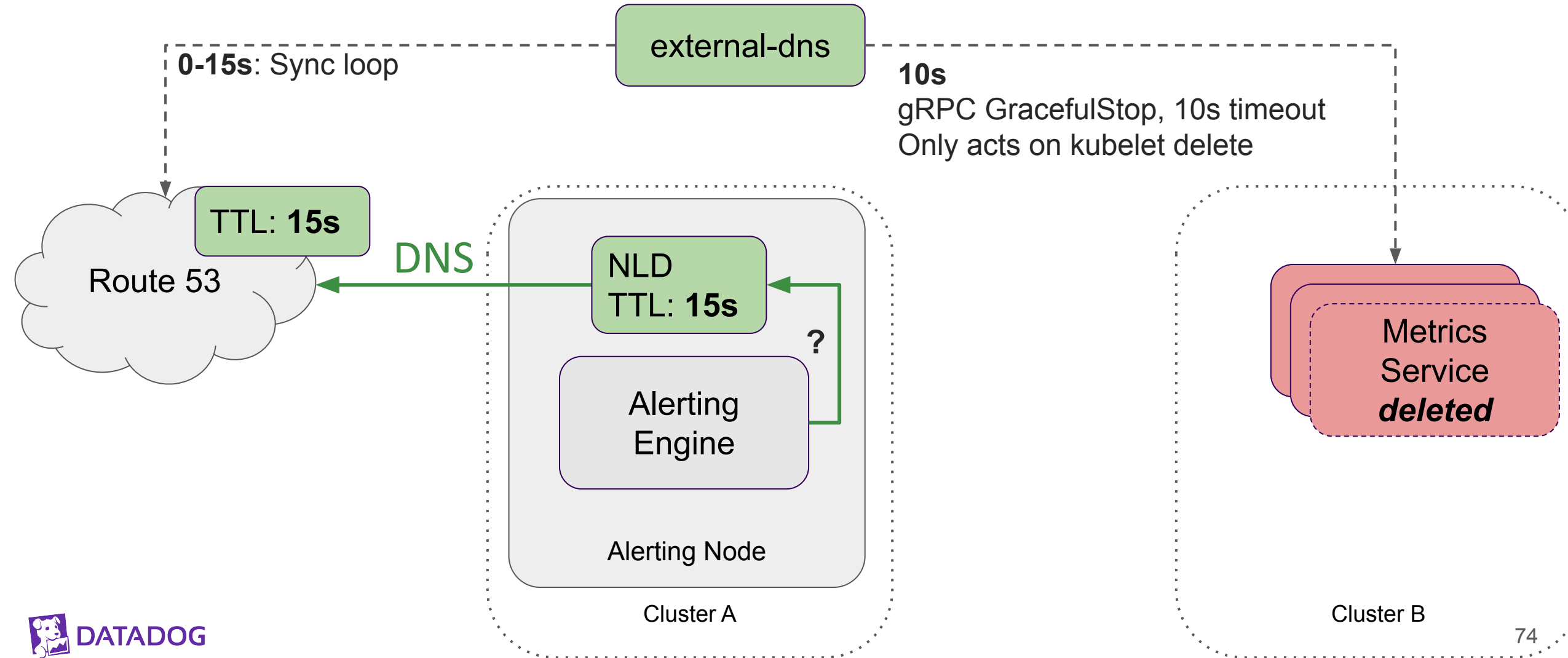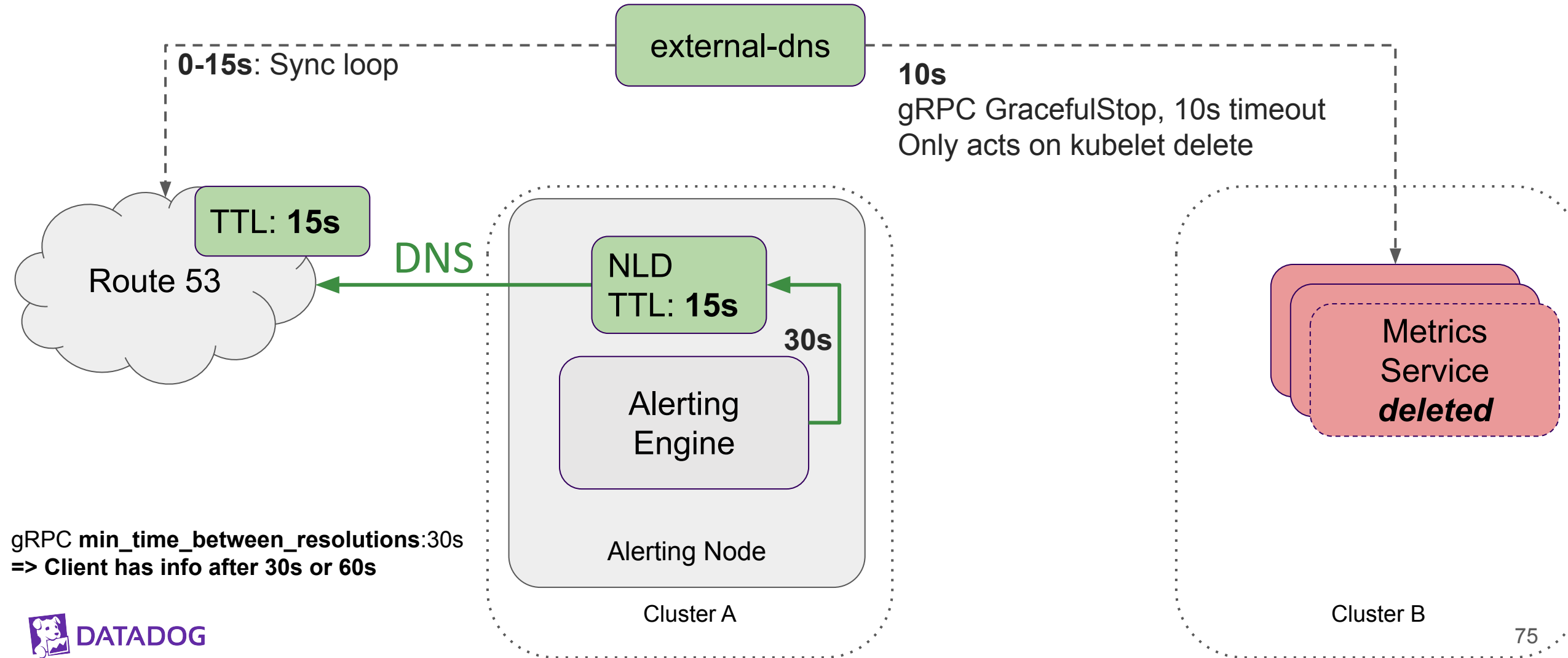
# RPC setup

1. **Service Discovery**

**2. Query**



Route 53

DNS

Alerting Engine

gRPC (client-side lb)

Metrics Service

# DNS propagation time during Rollouts



Maintain Metrics Service records

Watch Metrics Service pods

external-dns

Route 53

DNS

Alerting Engine

Cluster A

Metrics Service *deleted*

Cluster B

73

# DNS propagation time during Rollouts

# DNS propagation time during Rollouts

SRE CON _ AMERICAS

external-dns

**0-15s**: Sync loop

**10s**
gRPC GracefulStop, 10s timeout
Only acts on kubelet delete

TTL: **15s**

Route 53

DNS

NLD
TTL: **15s**

**30s**

Alerting
Engine

Alerting Node

Cluster A

gRPC **min_time_between_resolutions**:30s
**=> Client has info after 30s or 60s**

Metrics
Service
*deleted*

Cluster B
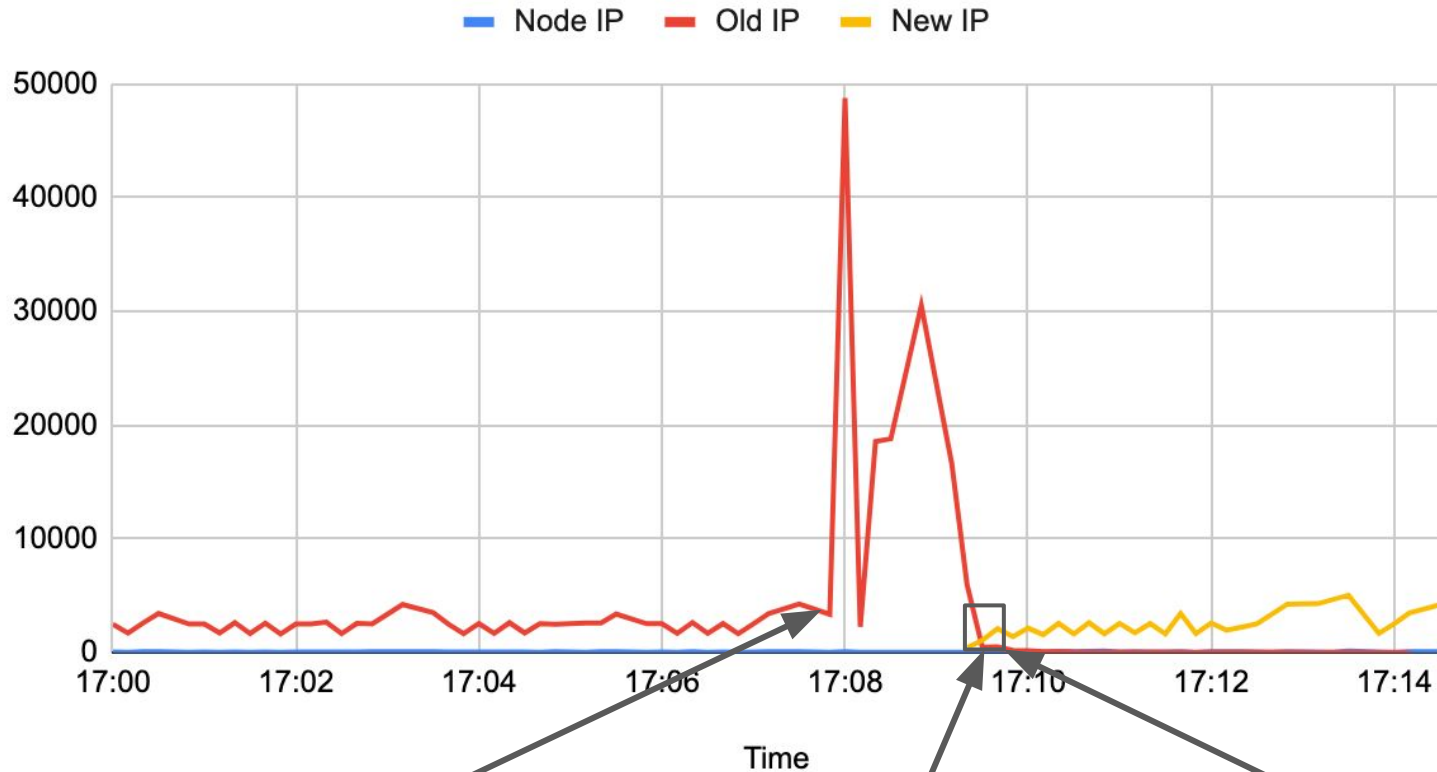
DATADOG

# DNS propagation time during Rollouts



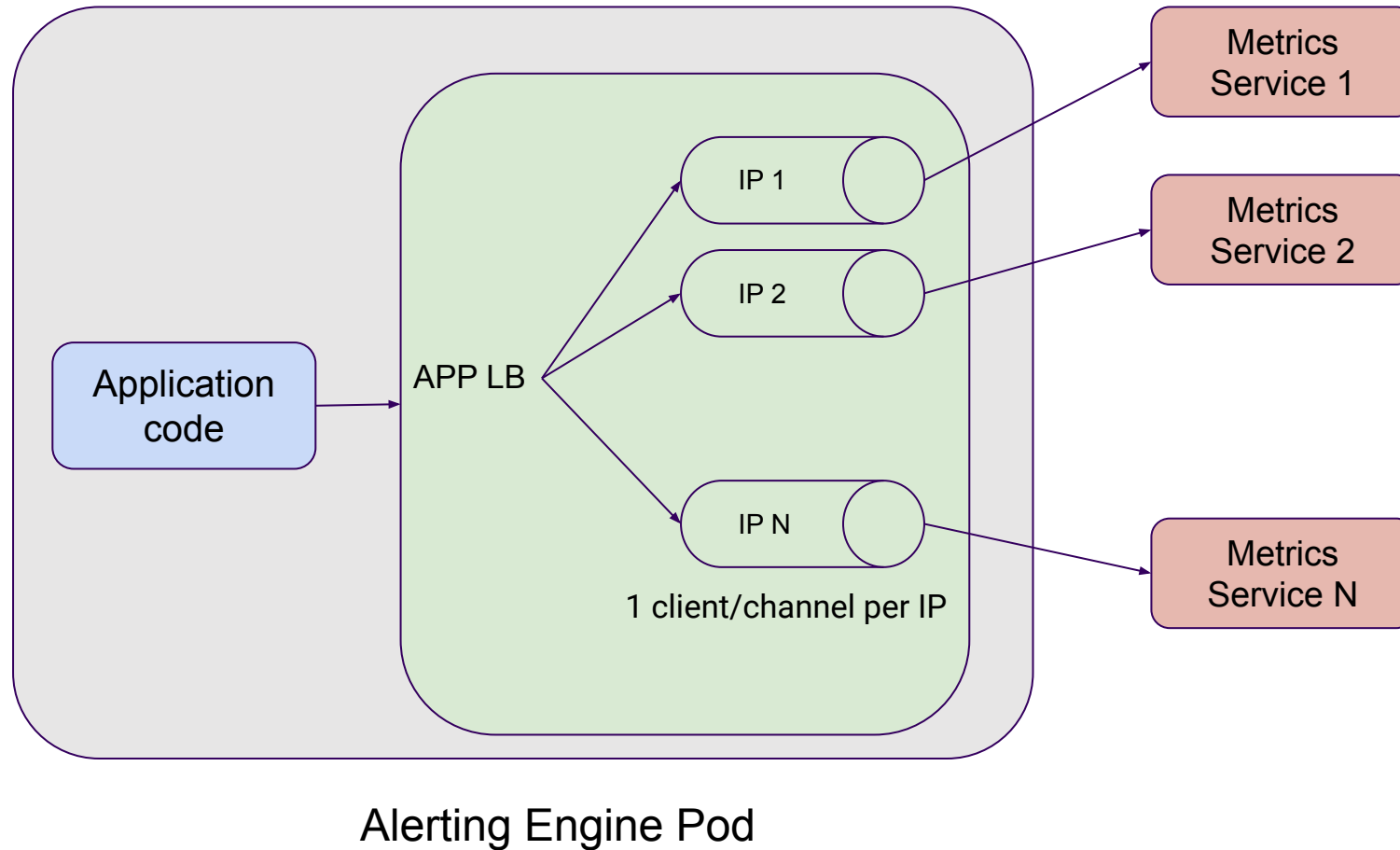Ingress flows by destination

Deletion starts

Clients progressively
start using new IP

No client is using the Old IP

# gRPC history at Datadog

- Originally, clients optimized for complex logic

  - DNS resolution in application code

  - One channel per backend IP

  - **pick_first** gRPC load balancing

- We changed the default to gRPC "standards"

  - Channels get a domain name and gRPC resolves

  - **round_robin** load balancing policy

  - This is when the issue started!

**DATADOG**

# Alerting still had one channel per backend



Alerting Engine Pod

# Reconnection differences

- pick_first and round_robin **have very different policies on connection failures**

  - pick_first: do not attempt to reconnect until the application asks for it

  - round_robin: automatically attempt to reconnect using reconnect options

- when using pick_first, we used **max_reconnect_backoff_ms=300 ms**

- ~reasonable for on-demand reconnects

DATADOG

# Does it add up?

Alerting Engine $* \mathbf{X000}$

$$\frac{}{\text{reconnect every 0.3 s}} = \mathbf{X0,000\ SYN\ /\ sec}$$

to each Metric Service Pod!

# The fix

Use default reconnect parameters

Browse files

⎇ lbernail/default-reconnect

👤 lbernail committed yesterday ✓ Verified

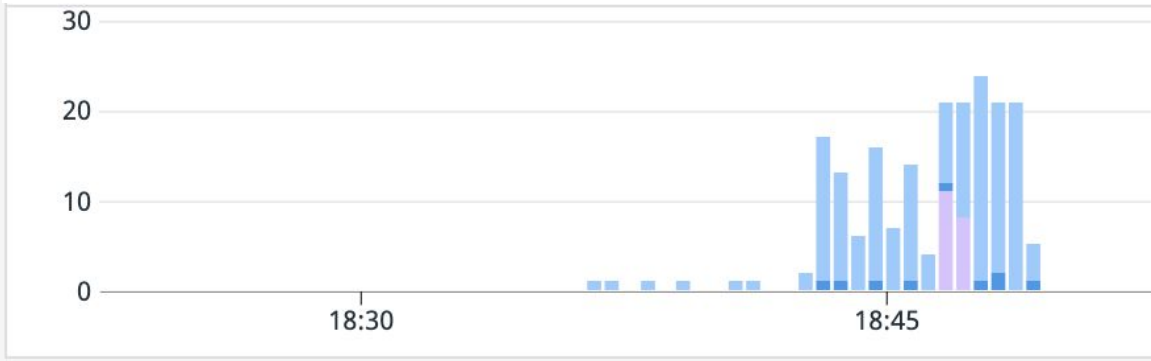± Showing **1 changed file** with **0 additions** and **4 deletions**.

Split | Unified

∨ ✥ 4 ■■■■□ grpc-reconnect.py ⧉ ⋯
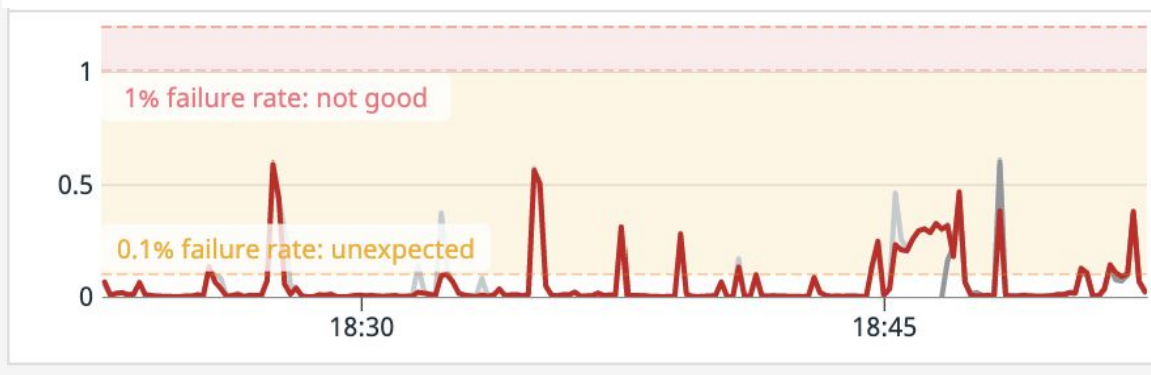
```
      @@ -981,10 +981,6 @@ def get_channel_for_service(host, dns_provider=None):
981              ("grpc.max_send_message_length", (16 << 20) - 1),          981              ("grpc.max_send_message_length", (16 << 20) - 1),
982              # receive max size is max uint, 2 GB                        982              # receive max size is max uint, 2 GB
983              ("grpc.max_receive_message_length", (1 << 31) - 1),        983              ("grpc.max_receive_message_length", (1 << 31) - 1),
984   -          # default is 20s, let's retry faster
985   -          ("grpc.min_reconnect_backoff_ms", 100),
986   -          ("grpc.initial_reconnect_backoff_ms", 200),
987   -          ("grpc.max_reconnect_backoff_ms", 300),
```
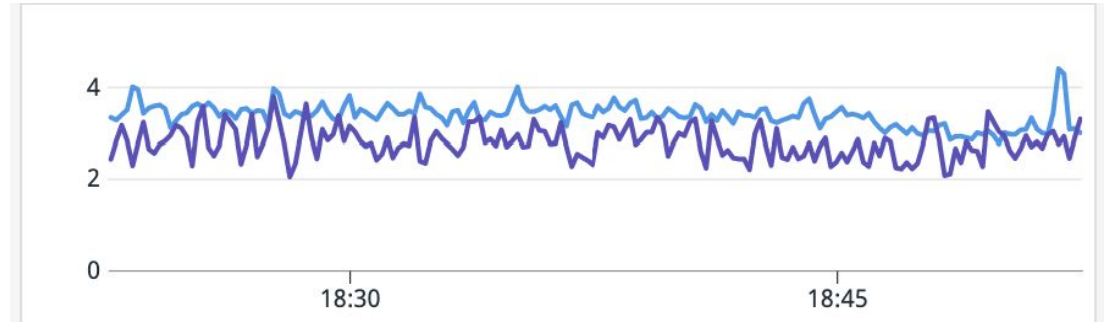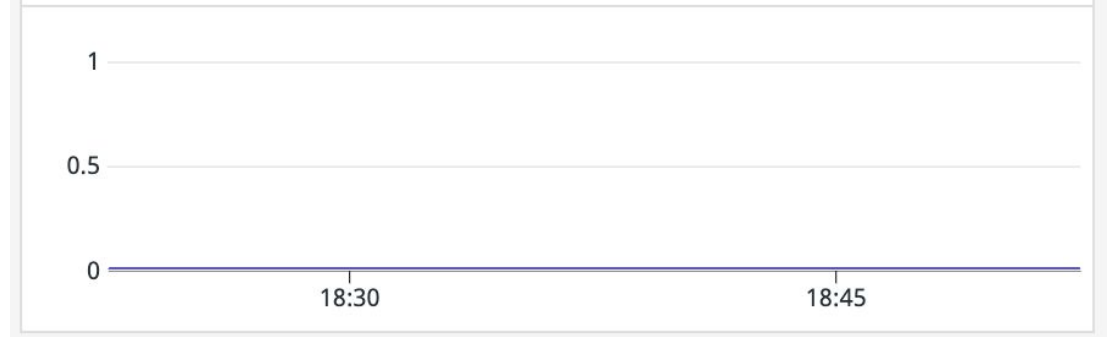
# Finally

average DNS response time by Metrics Service pod (ms)

Metrics Service pod replacements (rollout)



ENA Limits - conntrack exceeded



Metrics Service Error rate (**server** / **client**)



1% failure rate: not good

0.1% failure rate: unexpected

conntrack count for Alerting Engine



DATADOG

82

# Lessons Learned

# Lessons Learned

- Debugging this incident was long and painful but we learned a lot

- Sometimes it's not DNS

- Powerful abstractions leak in complex ways

- gRPC setup can be complex, making changes dangerous

- ENA metrics and VPC flow logs are extremely useful

- Required complex team efforts (thanks Laurent, Wendell, Matt, Nayef!)

**DATADOG**

# Thank you

See the blog post for more details: dtdg.co/not-always-dns

We're hiring: datadoghq.com/careers/

elijah@datadoghq.com
hemanth.malla@datadoghq.com

@elijahca
@hemanthmalla

DATADOG