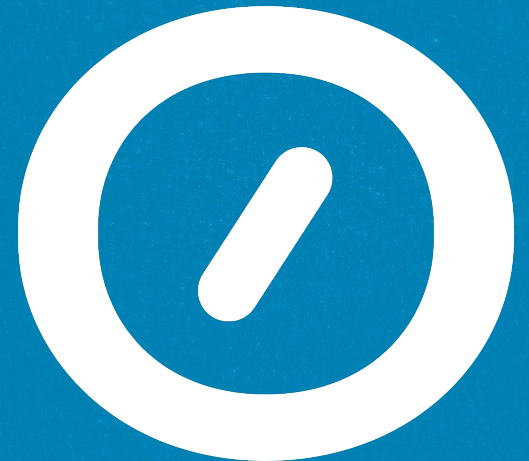
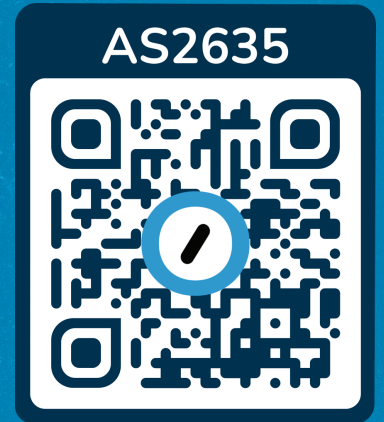


The Internet X-Ray: Diagnosing ECMP failures from the edge

Bela Toros, Chris Laffin, Tyler Leeds

AUTOMATTIC



AUTOMATTIC

 WordPress.com

 vip

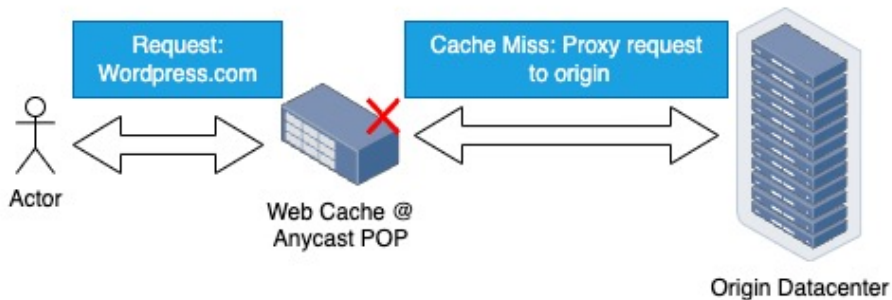
 Jetpack

 **WOO** COMMERCE

tumblr

Who we are:

Basics of an Anycast CDN



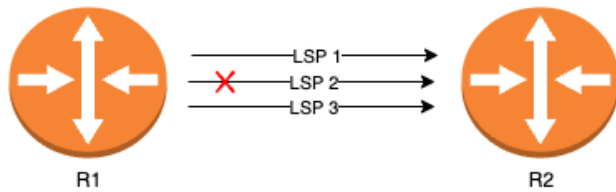
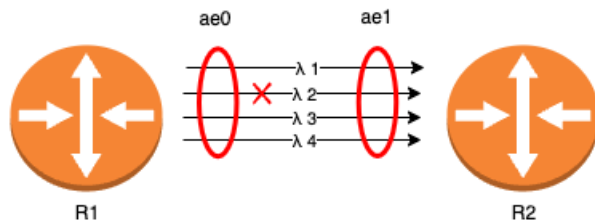
Anycast CDN Basics

- Same IP prefix announced from multiple places
- User requests are routed by BGP to nearest POP
- Media that is cached is served directly
- Media that isn't cached is proxied by the POP back to an origin datacenter

Automatic CDN Specifics

- Origin fetch happens over DFZ transit links
- Each POP and Origin connected to several Tier1 NSP
- No private links between POP and Origin

Basics of ECMP

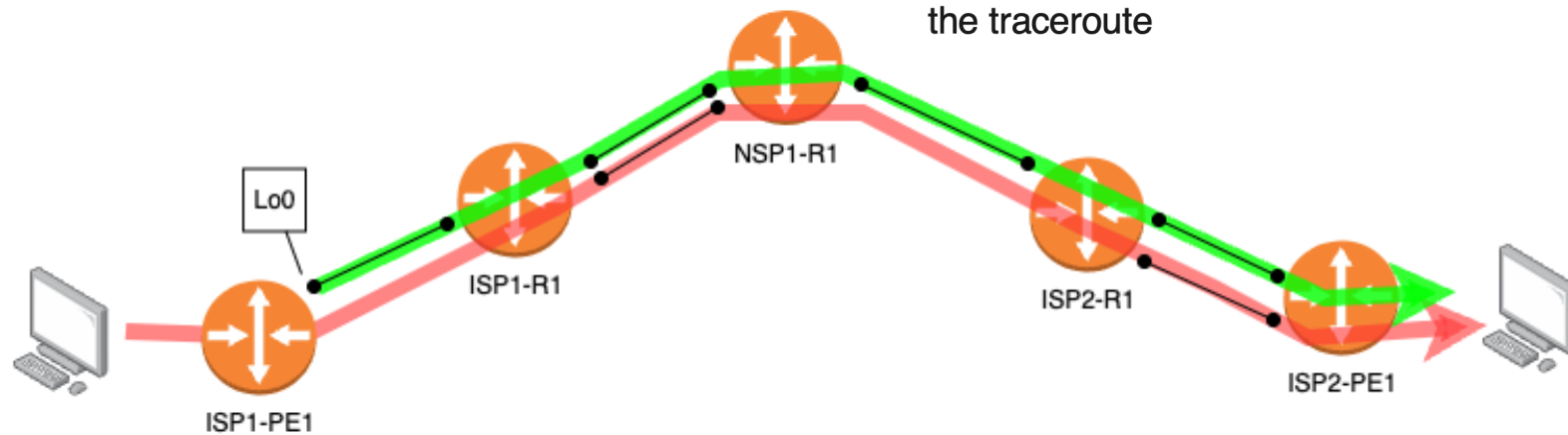


ECMP Basics

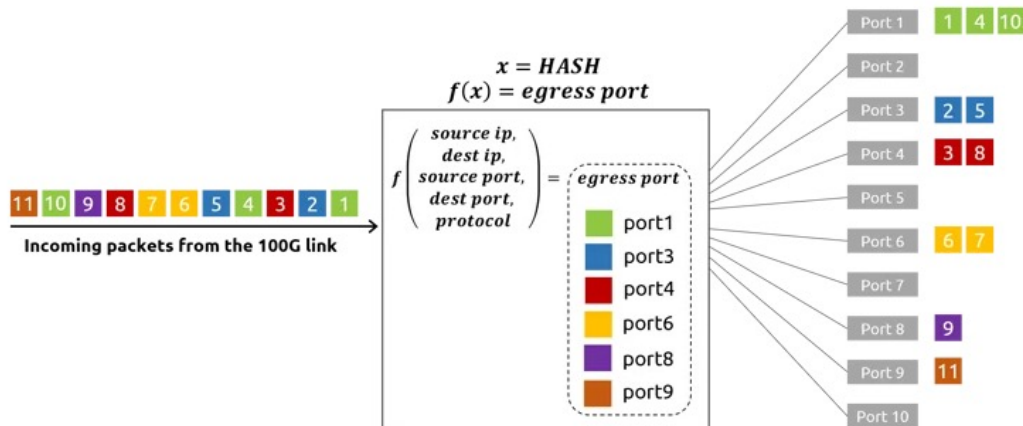
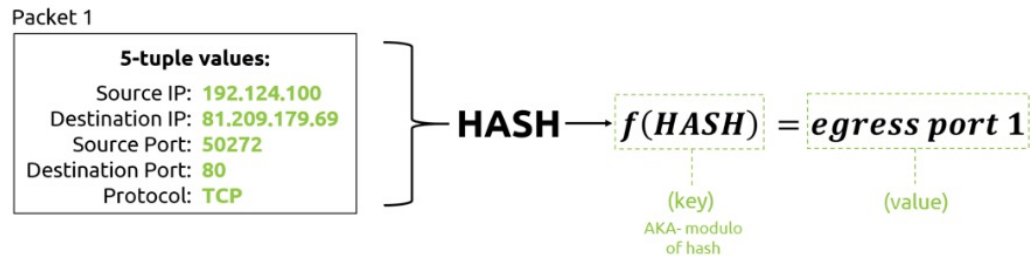
- Single hop represents multiple physical links
- Links might be aggregated into an LACP LAG
- Links might be several equal cost Layer 3 paths
- Links can congest or drop packets independently
- Layer 1 issues can drop packets without generating errors or discards on routers

The Internet is a Flow Switched Network

- Passing through the same routers does not mean you pass through the same physical links
- Pinging along the same path is insufficient to establish whether path is loss-free.
- Traceroute usually shows router Lo0 IP/hostname or the IP of the interface used to reach the source of the traceroute



Basics of Link Load Balancing



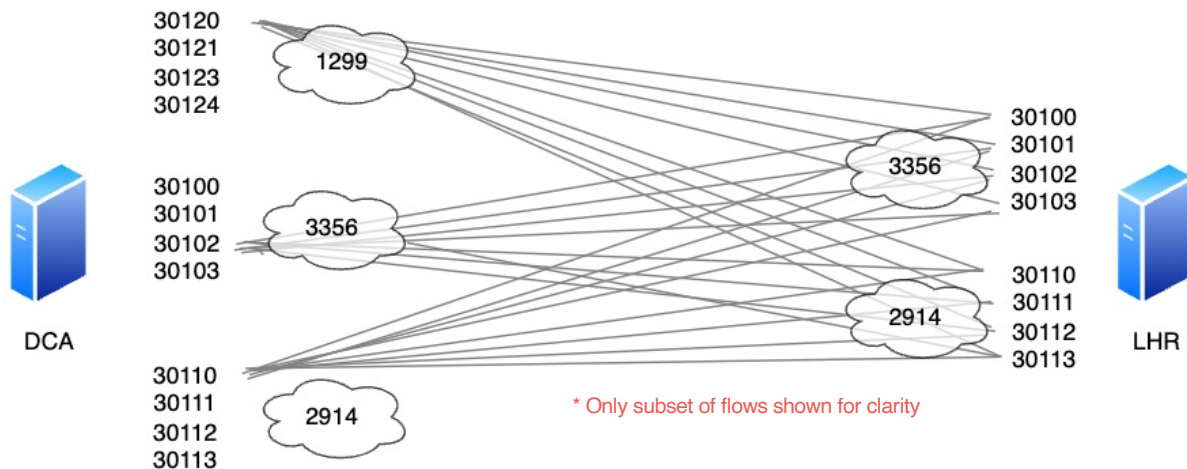
Load Balancing Basics

- Routers DO NOT use per-packet load balancing
- Routers use their own formula based on the 5tuple or 3tuple to choose the egress port
- Egress port decision is stateless (but stable)
- Hash formula is applied to every packet
- Specific flows (5tuples) will always use same physical ports
- Port hashes are usually recalculated on topology change

The Problem

Our Solution: PINGO

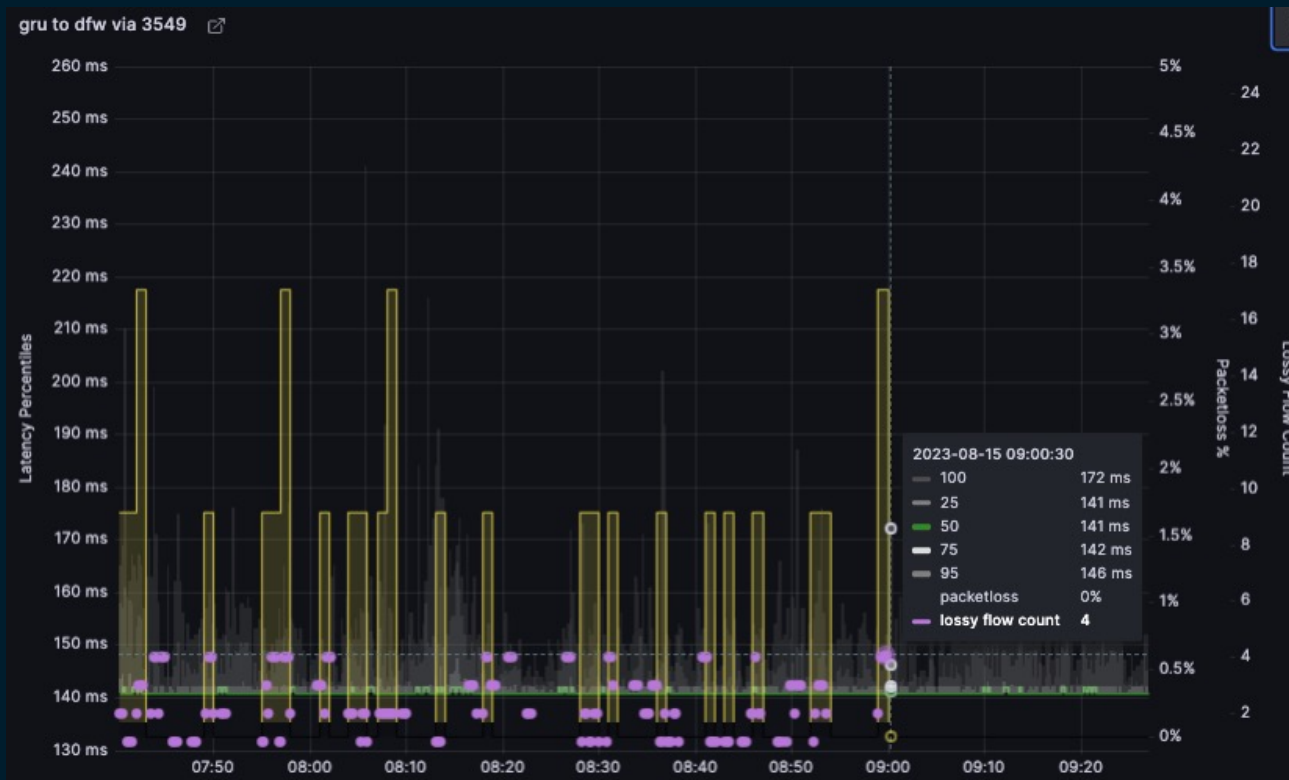
Pingo is a custom GoLang daemon that creates a matrix of probes between source-ports at each site. By varying the 5tuple, probes are load-balanced by the internet to different physical links and issues with a single link can be quantified and located. Pingo will be open sourced on Github shortly.



Specifications

- Sends timecode based probe every 500ms.
- Source-ports are assigned to each NSP
- Probe initiated from Origin DC's
- Each source port pings all ports on all NSPs at target site
- Calculates latency for each probe
- Loss is declared after 2000ms
- Stats are gathered in TSDB (Prometheus)
- Can determine unidirectional loss and provide data to pinpoint where
- Routinely sensitive to < 0.05% loss.
- Accurate to ICMP based latency calculation to within 2ms
- Currently running a 6x6 matrix for every dc/nsp combo. (4x4 shown)

PINGO Dashboard Primer



- 50th Percentile RTT
- Qty of lossy flows
- Loss percentage

PINGO Dashboard Primer - How to interpret

Current active path

Outbound on 2914

Outbound on 1299

Outbound on 3356

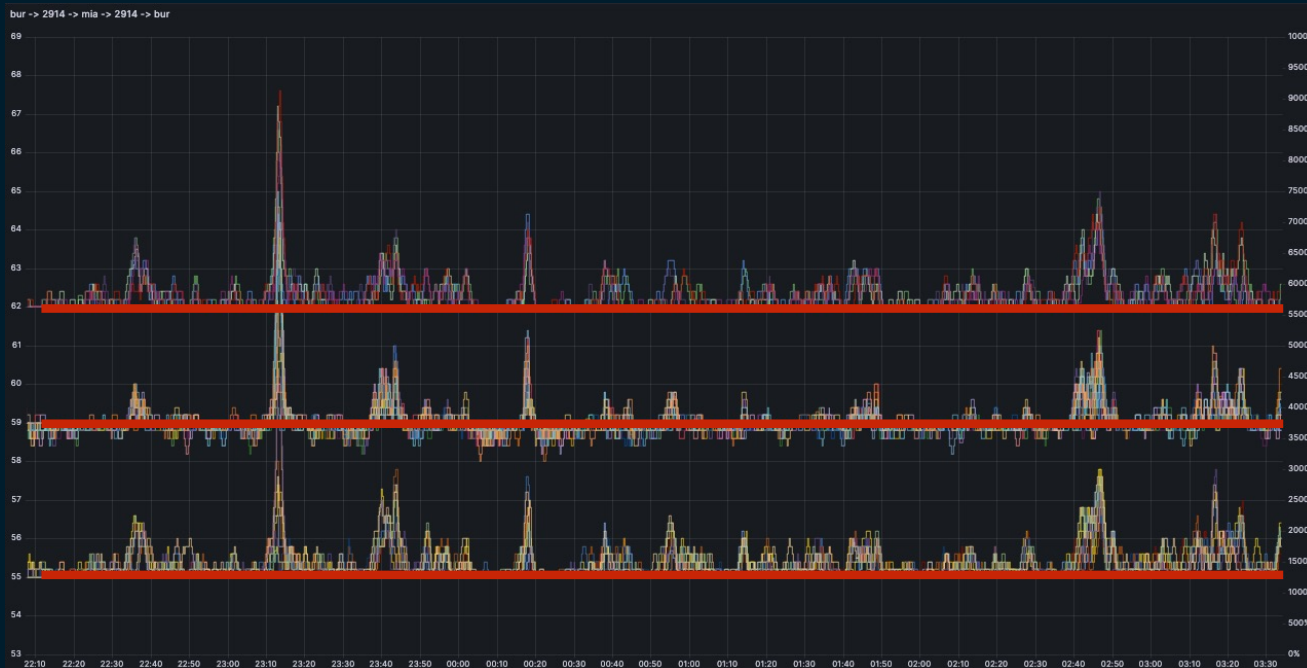


Return on 2914

Return on 3356

Return on 7922

PINGO Dashboard Primer - Drill down



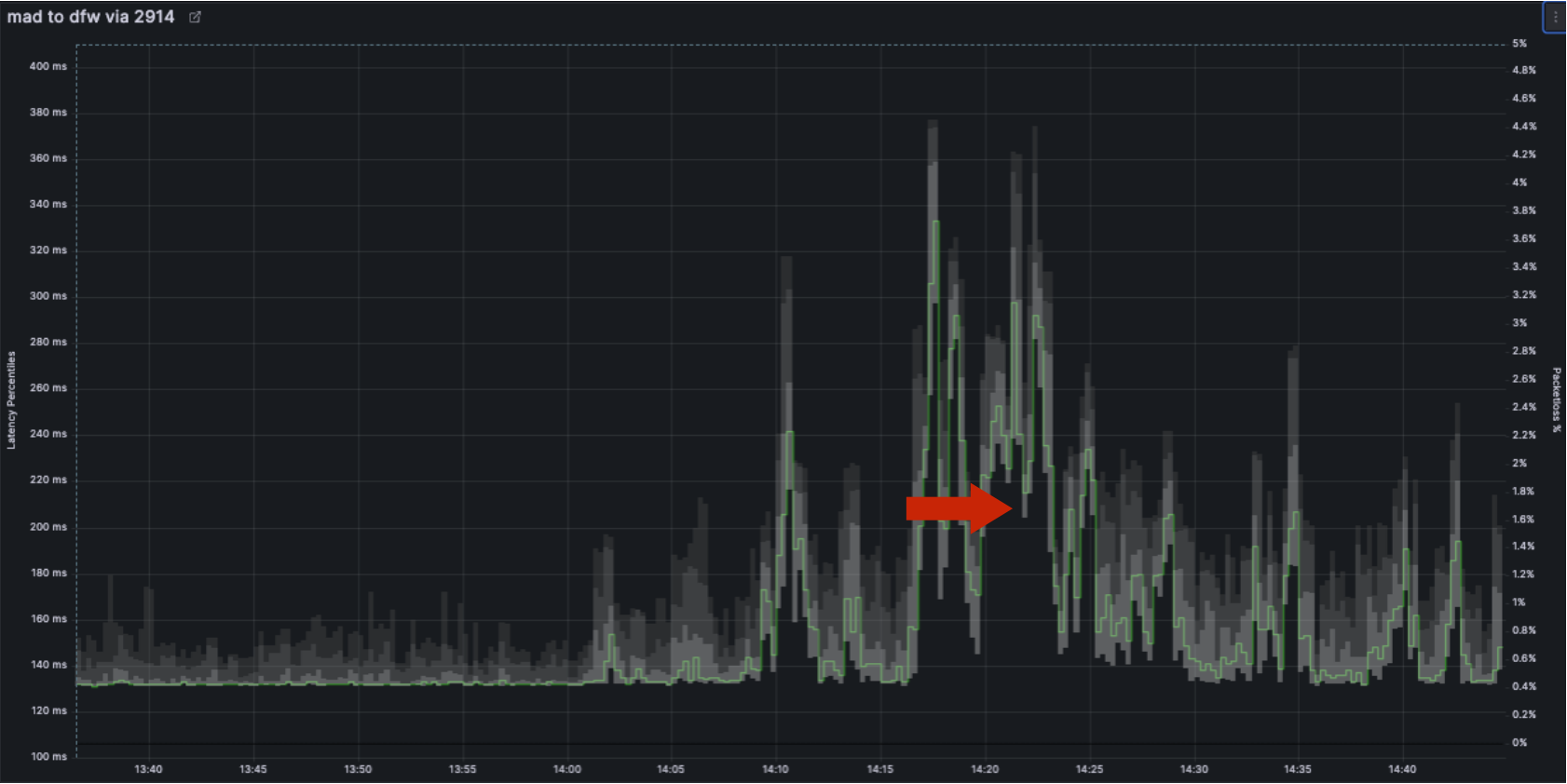
Latency Strata Plot

- Flows self organize into obvious latency strata.
- Each strata represents a distinct path across the internet
- Number of strata is unfixed
- Strata will frequently disappear or appear as paths are changed by NSP's

Case Study 1: Single Path Congestion

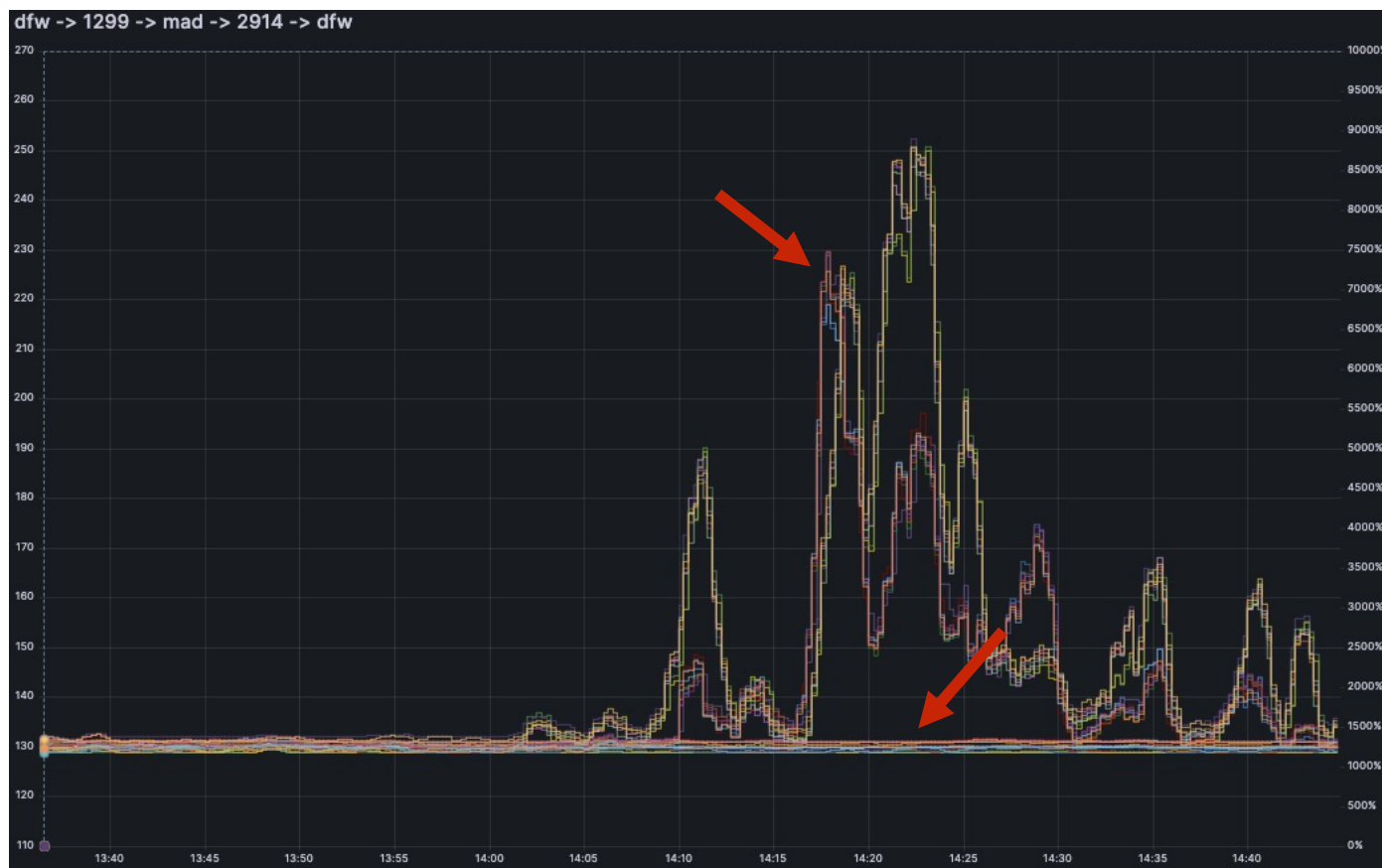


Case Study 1: Single Path Congestion



* Only subset of flows shown for clarity

Case Study 1: Single Path Congestion



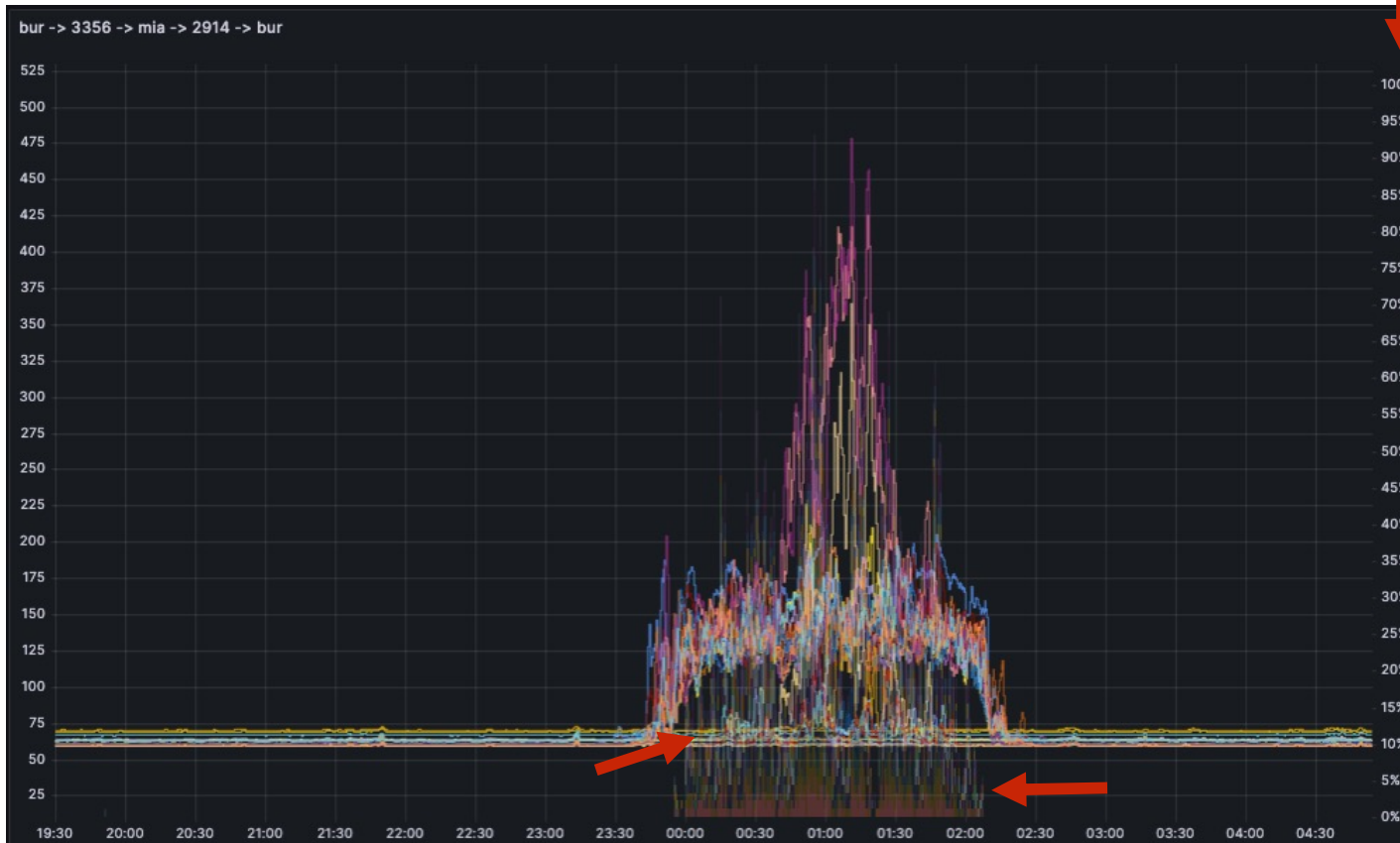
What's happening

- Showing all 36 flows
- Shows several flows are experiencing latency while others are fine
- Root cause was a customer sending lots of traffic via a single flow. The links in that flow congested while most were fine.
- This affected 50% of flows. Normal monitoring would miss this 50% of the time.

Case Study 2: Single Path Congestion + Link dropping



Case Study 2: Single Path Congestion (Tail dropping)



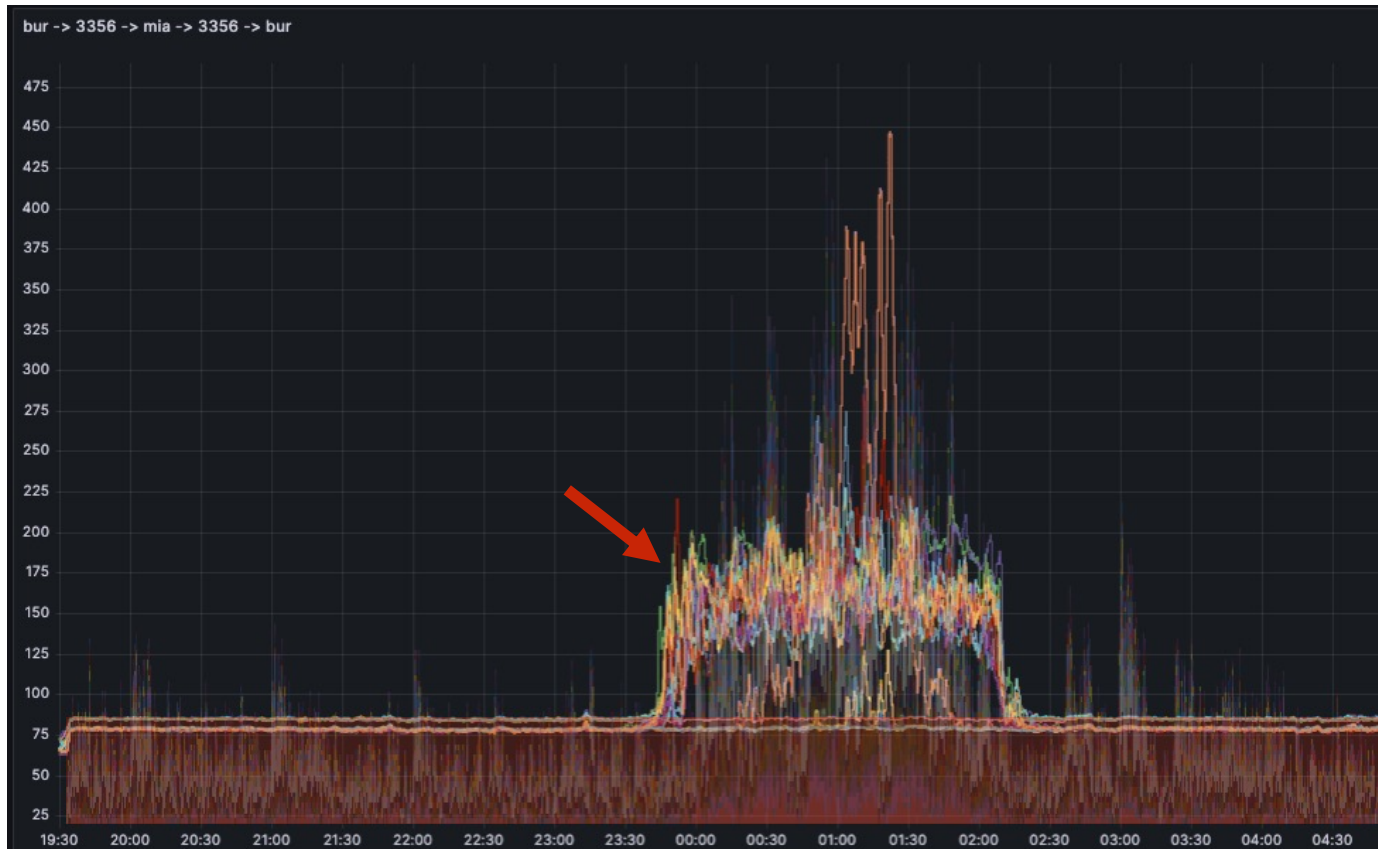
What's happening

- Showing all 36 flows
- Shows several flows are experiencing latency while others are fine
- Congestion on high-latency circuit is sufficient to induce loss.
- Low latency path is loss-free

Case Study 2: Single Path Congestion + Link dropping



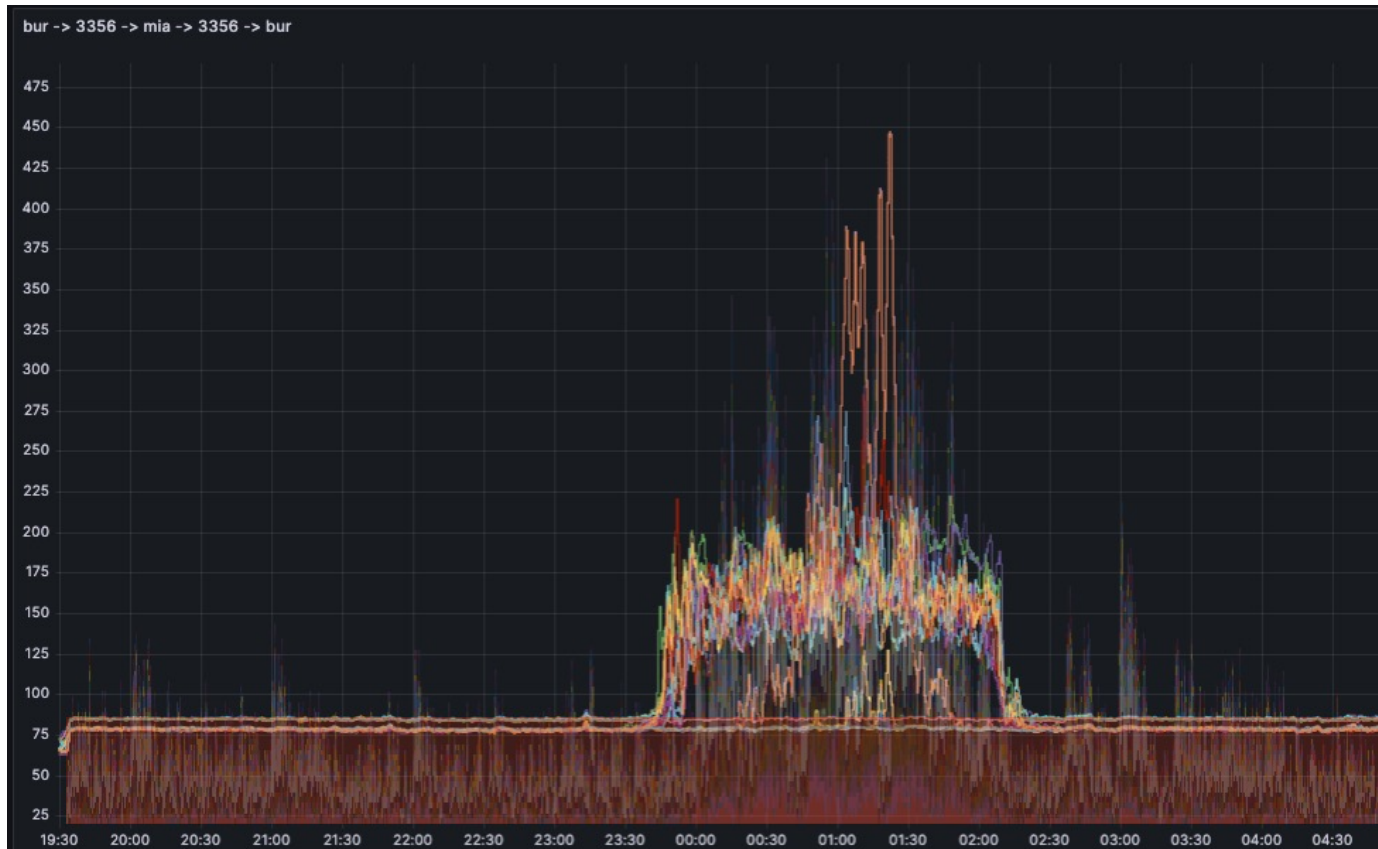
Case Study 2: Link loss



What's happening

- Not congestion
- Showing all 36 flows
- All paths experiencing loss
- Still shows obvious latency "strata" indicating multipathing
- No jump in latency at all.
- Caused by single "bottleneck" link taking errors.

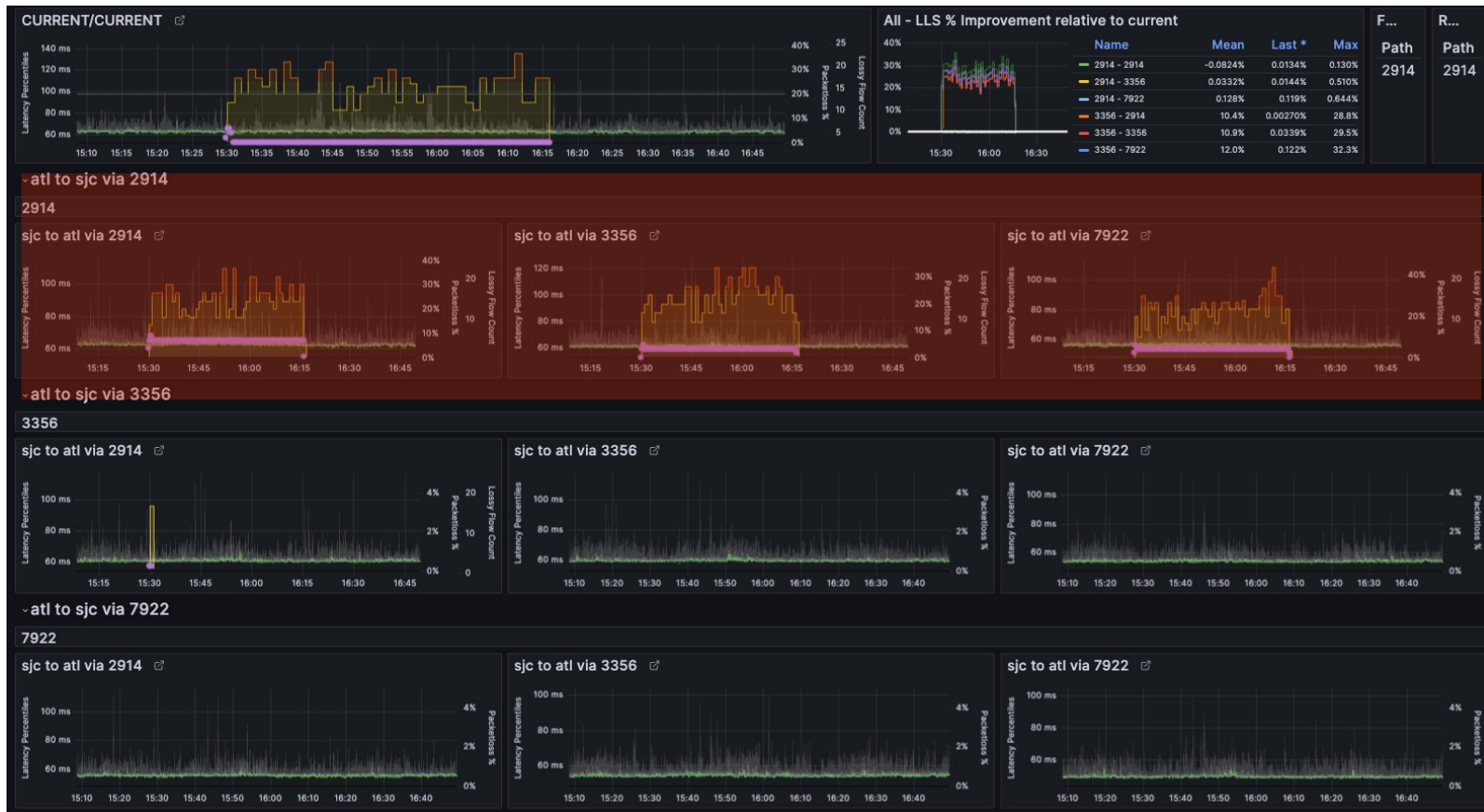
Case Study 2: Link loss



What's happening

- Not congestion
- Showing all 36 flows
- All paths experiencing loss
- Still shows obvious latency "strata" indicating multipathing
- No jump in latency at all.
- Caused by single "bottleneck" link taking errors.

Case Study 3: Locating a lossy link



Case Study 3: Locating a lossy link

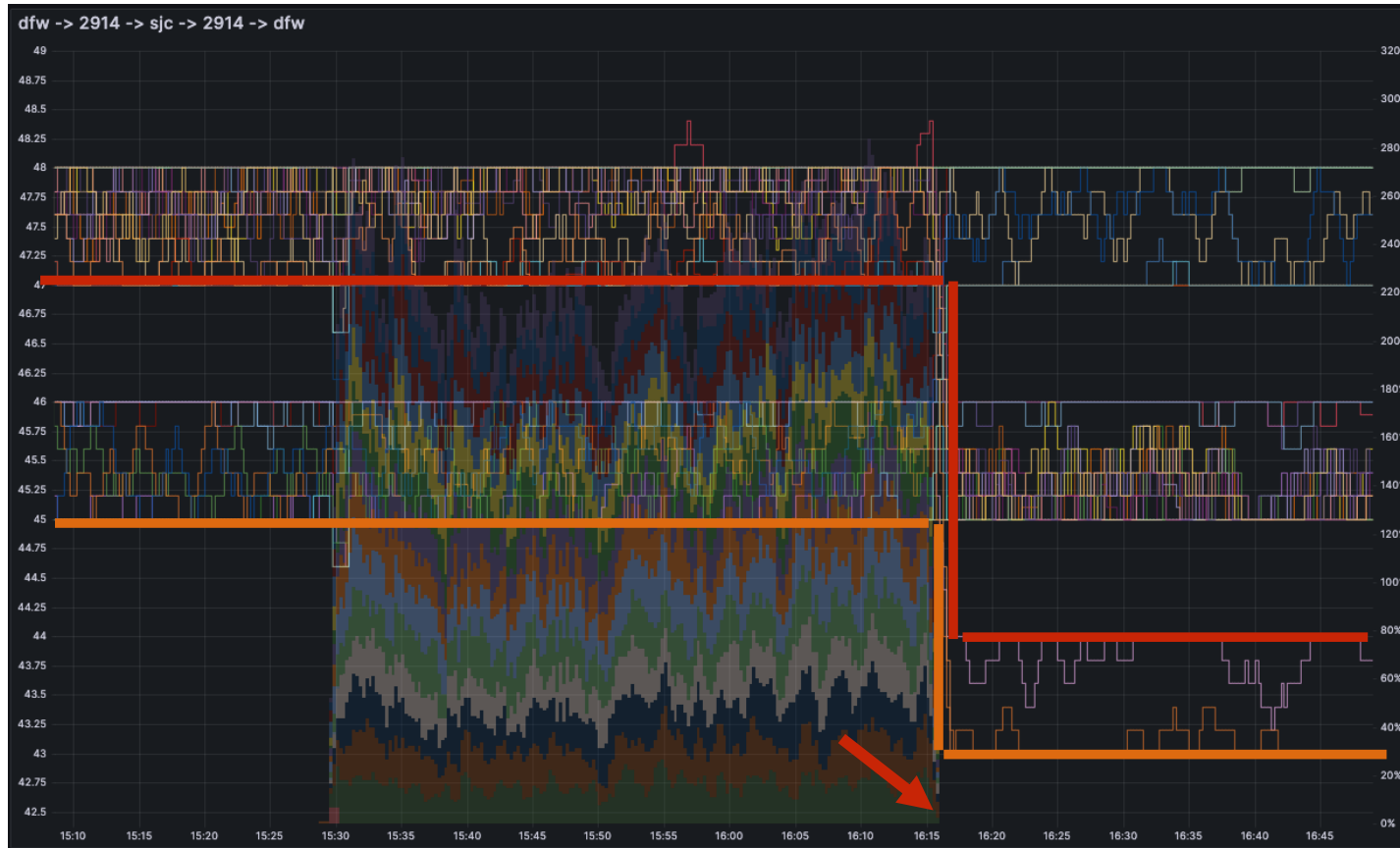
```
root@net1.atl.wordpress.com:~# mtr -z -b -u -4 -B253 -L50002 -P50000 192.0.70.252 -z -i.3 -w -c 50 -o LSDNA
Start: 2023-09-21T15:54:24+0000
HOST: net1.atl.wordpress.com
```

		Loss%	Snt	Drop	Last	Avg
1.	AS2635 atl-core-a01.atl.automattic.net (198.181.118.1)	0.0%	50	0	0.5	1.3
2.	AS2914 xe-2-4-2-2.a03.atlnga05.us.bb.gin.ntt.net (129.250.206.181)	0.0%	50	0	0.5	3.1
3.	AS2914 ae-2.r25.atlnga05.us.bb.gin.ntt.net (129.250.2.120)	0.0%	50	0	32.1	5.5
4.	AS2914 ae-6.r21.dllstx14.us.bb.gin.ntt.net (129.250.4.116)	0.0%	50	0	17.2	17.6
5.	AS2914 ae-2.r25.snjsca04.us.bb.gin.ntt.net (129.250.4.154)	22.0%	50	11	59.1	62.7
6.	AS2914 ae-4.r25.sttlwa01.us.bb.gin.ntt.net (129.250.3.125)	14.0%	50	7	84.7	84.1
7.	AS2914 ae-1.a03.sttlwa01.us.bb.gin.ntt.net (129.250.2.207)	16.0%	50	8	80.8	81.6
8.	AS2914 ae-0.automattic.sttlwa01.us.bb.gin.ntt.net (129.250.202.66)	14.0%	50	7	70.2	71.0
9.	AS2635 192.0.70.252	18.0%	50	9	70.3	70.7

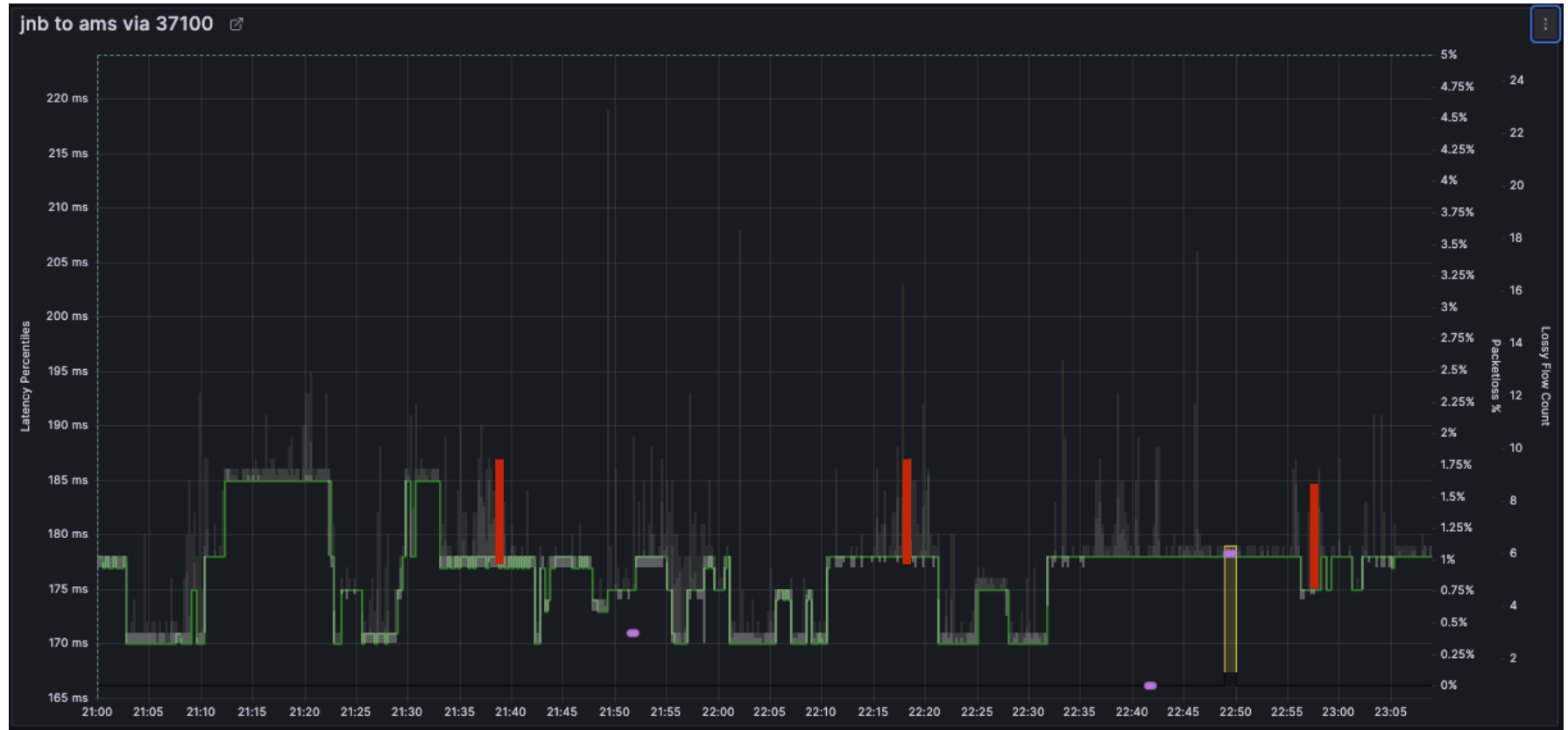
Case Study 3: Locating a lossy link



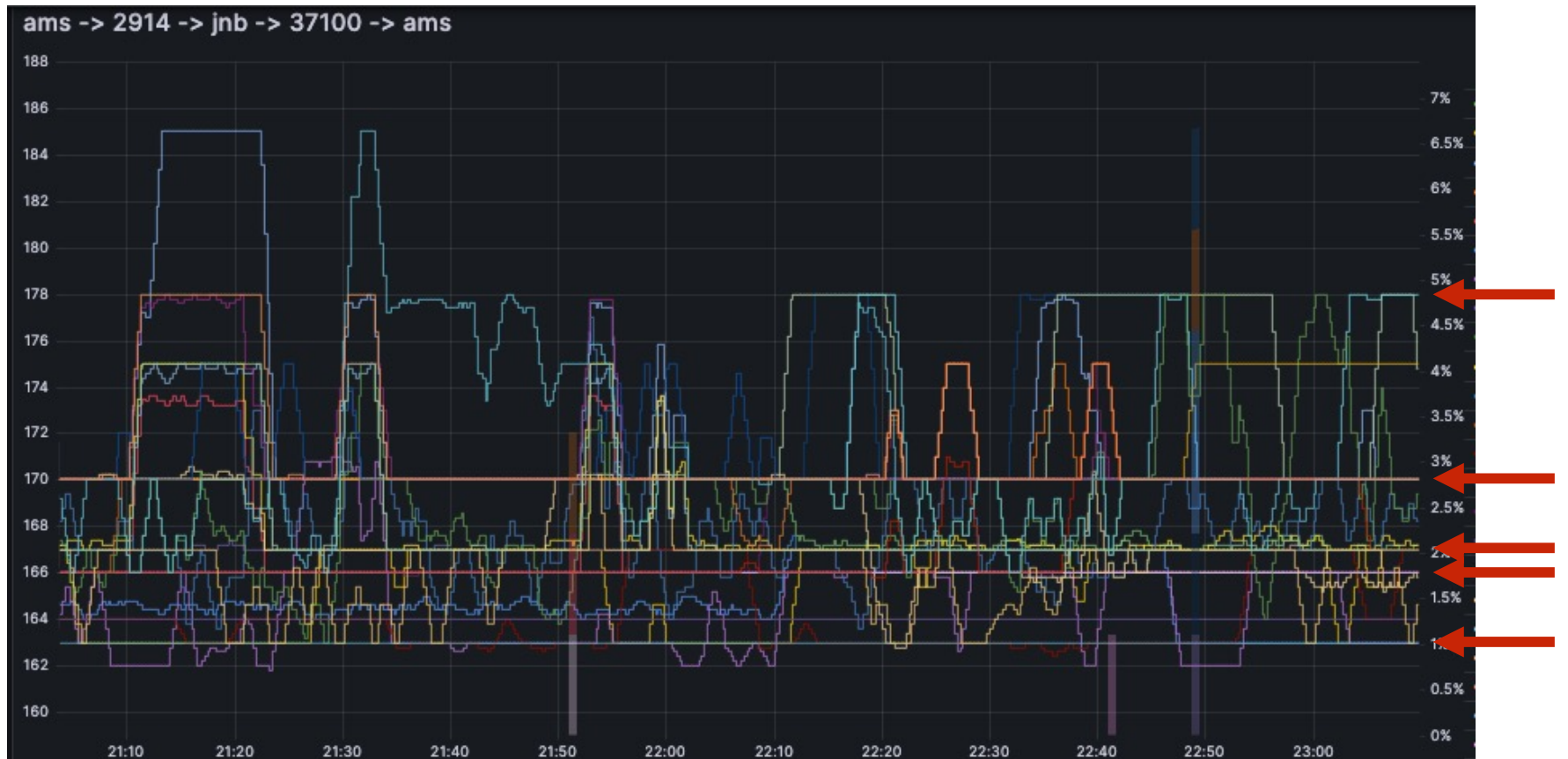
Case Study 3: Locating a lossy link



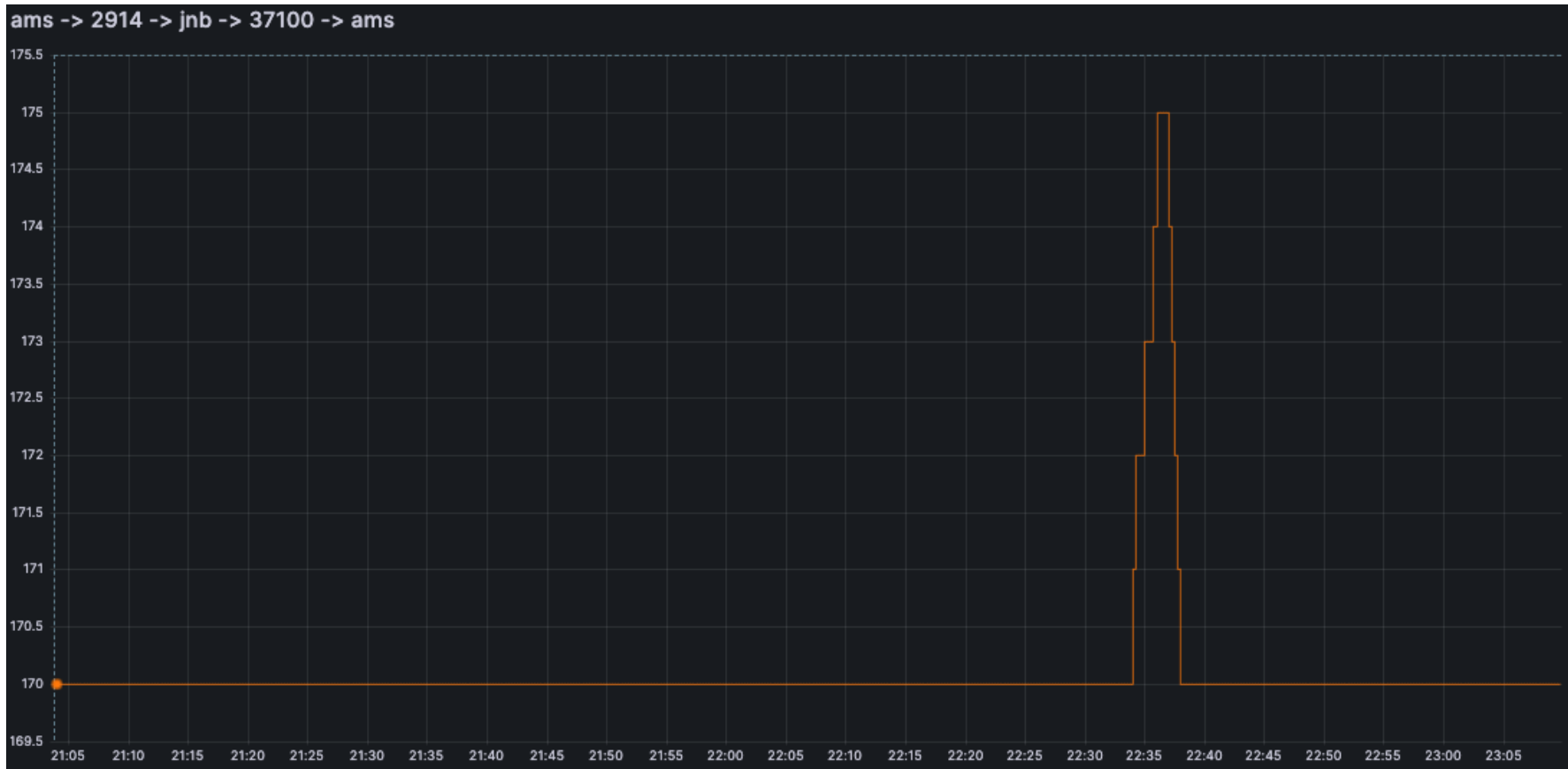
Case Study 4: WTH?



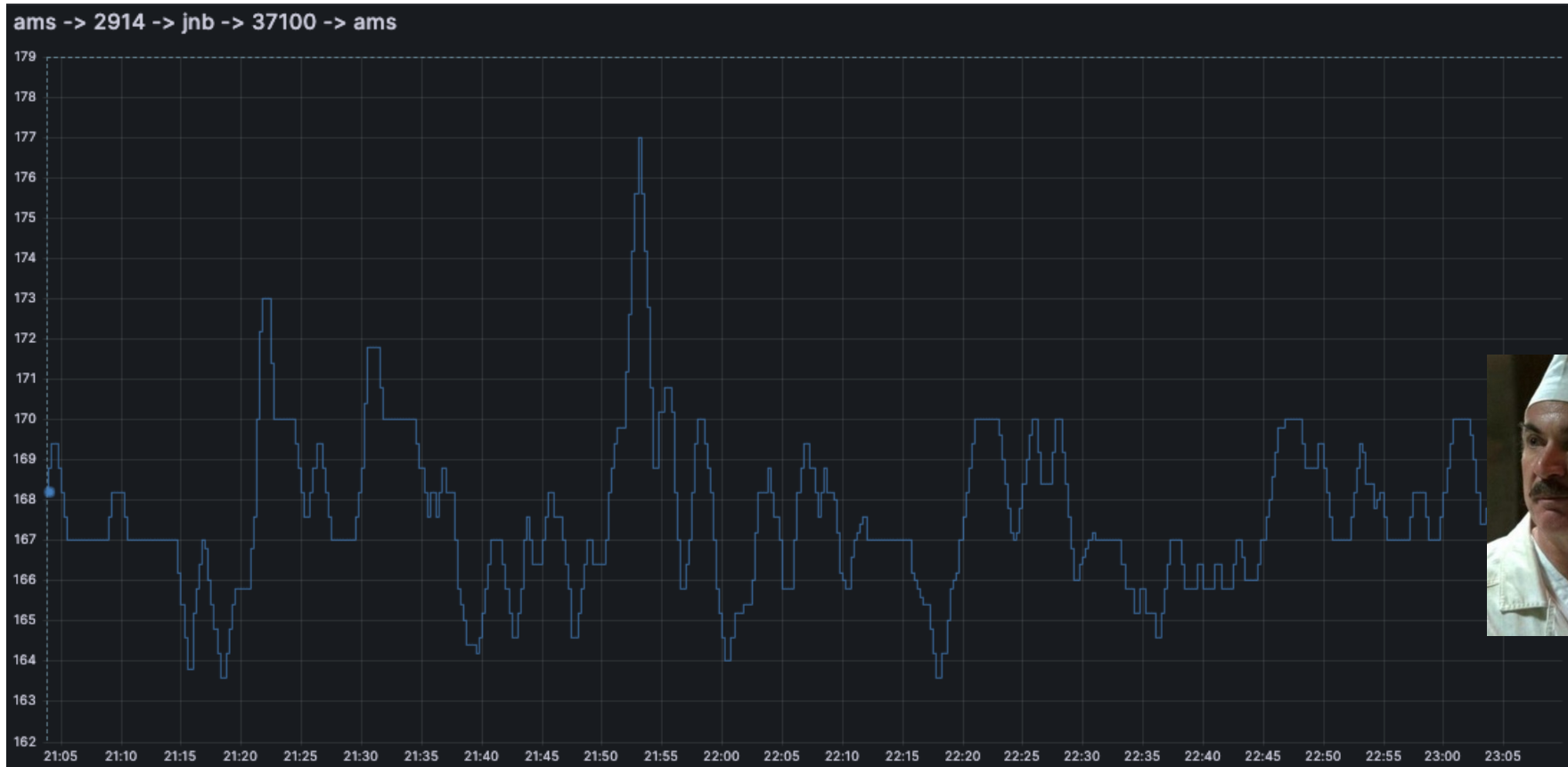
Case Study 4: WTH?



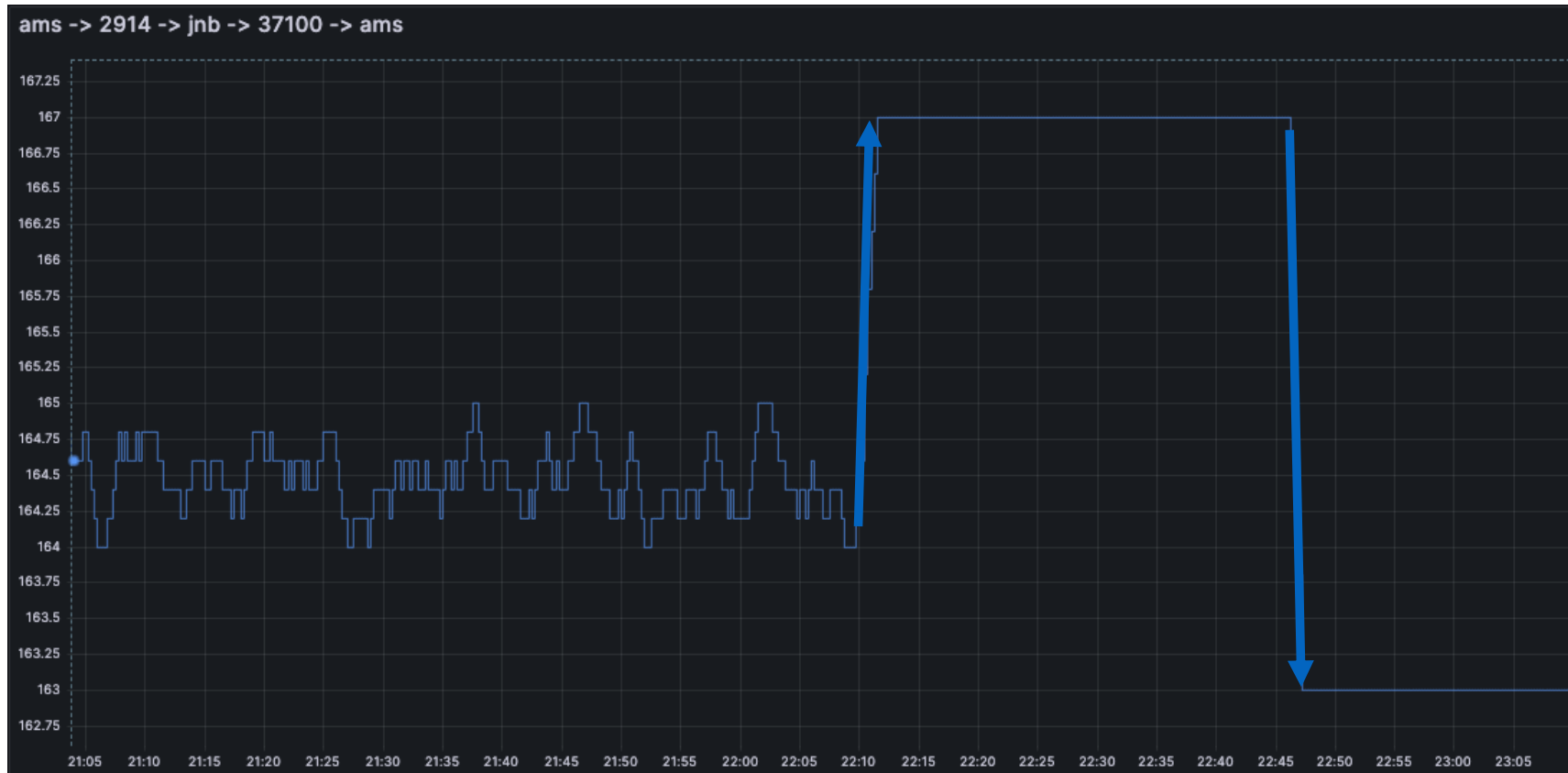
Case Study 4: WTH?



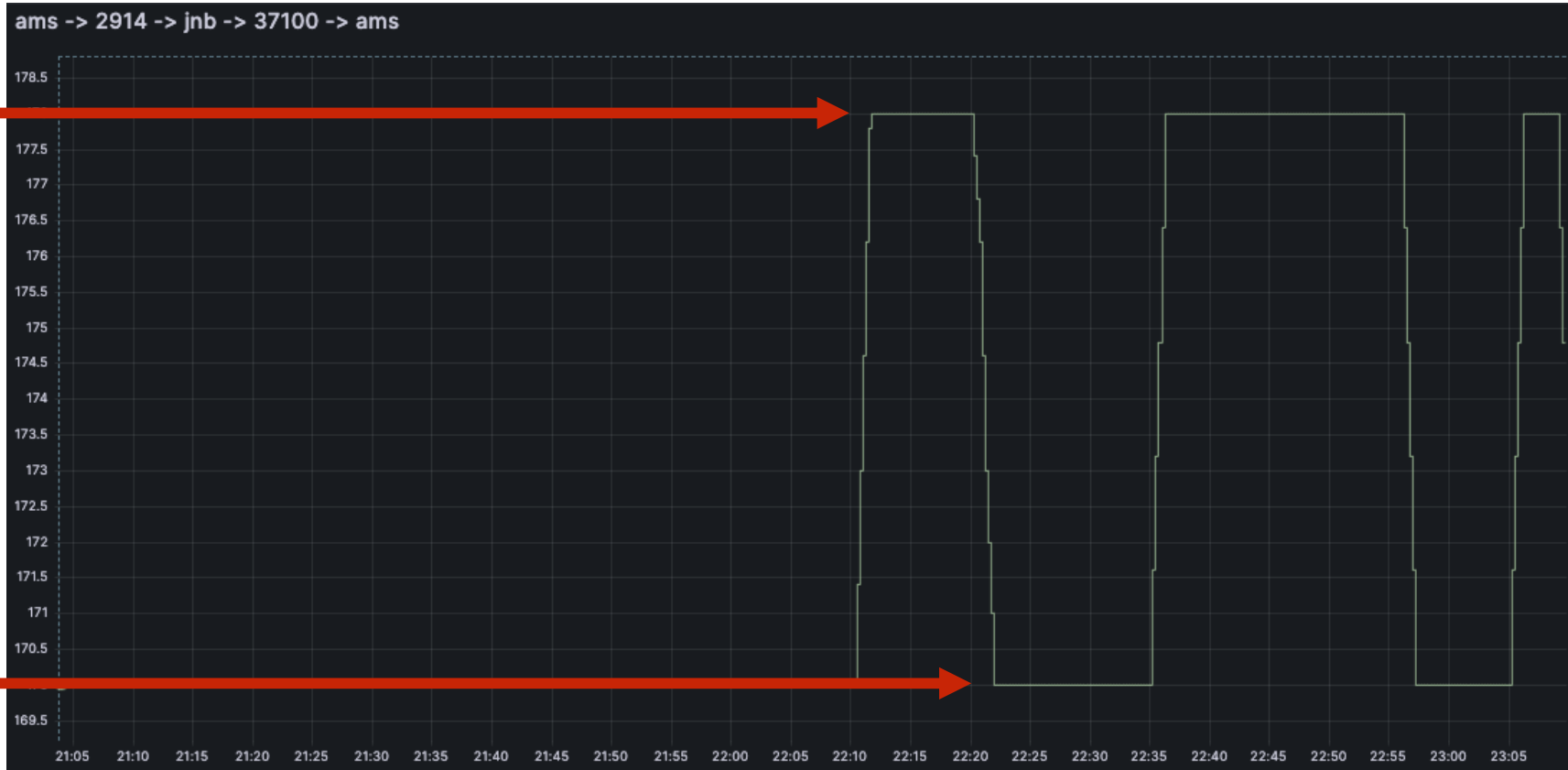
Case Study 4: WTH?



Case Study 4: WTH?



Case Study 4: WTH?



Case Study 4: WTH?



Thank you.



Work with us!

AUTOMATTIC