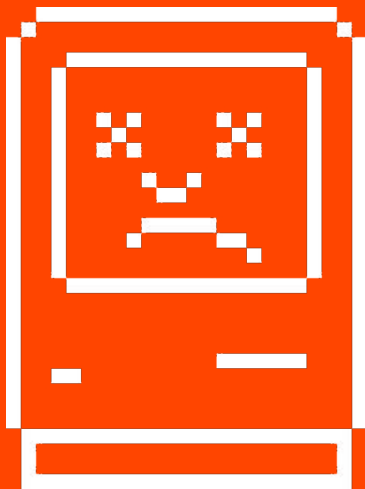
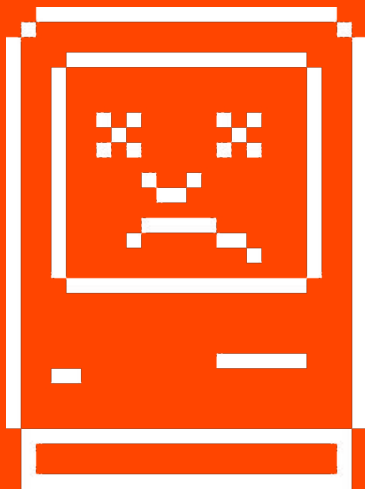


Noisy Neighbors



through networking

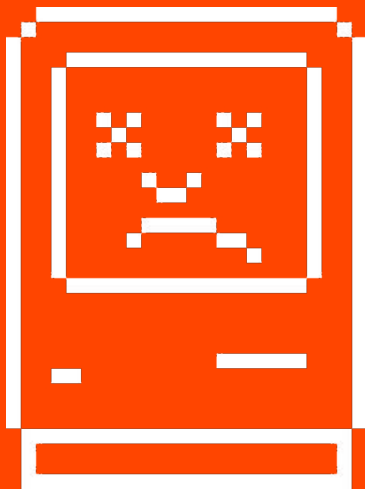
Noisy Neighbors



René Treffer
(he/him)



Noisy Neighbors

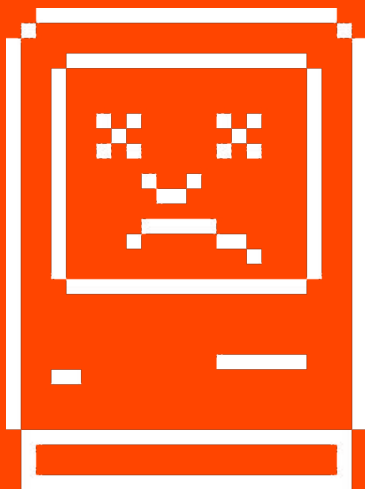


Ben Kochie (he/him)



Incidents

5

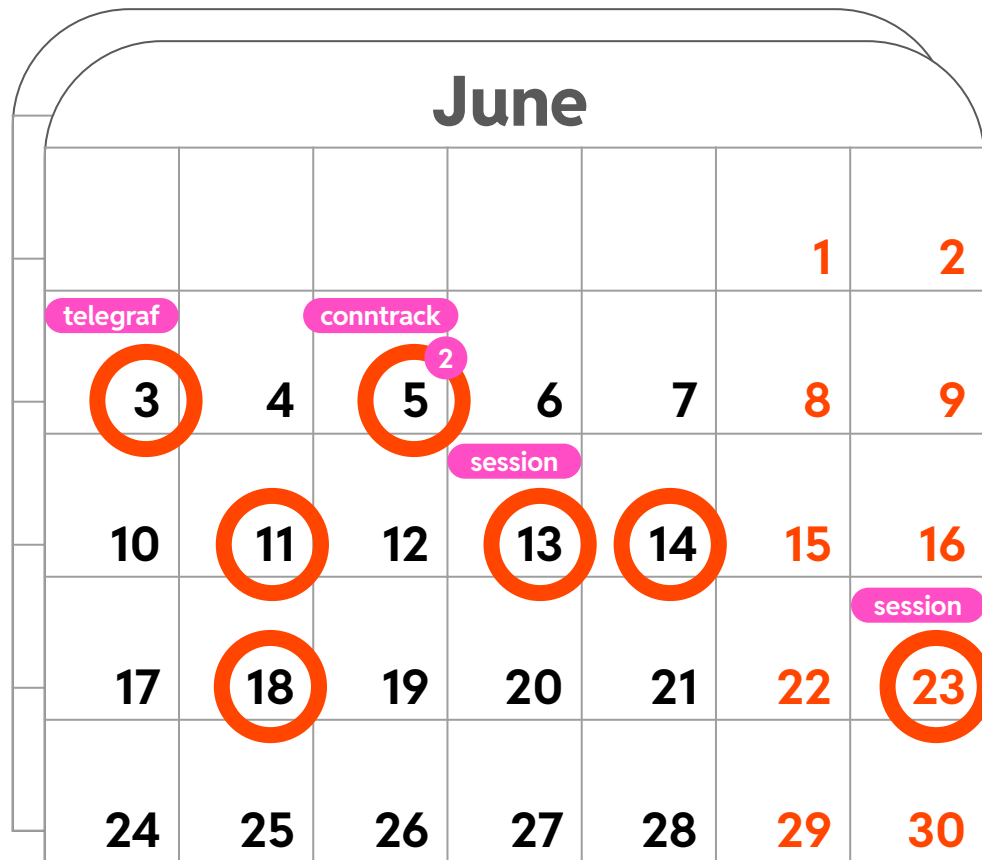
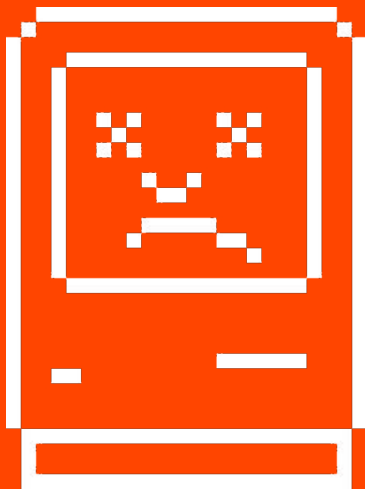


May						
		1	2	3	4	5
6	7	8	9	10	11	12
13	14	15	16	17 ²	18	19
20	21	22	23	24	25	26
			k8s			
27	28	29	30	31		



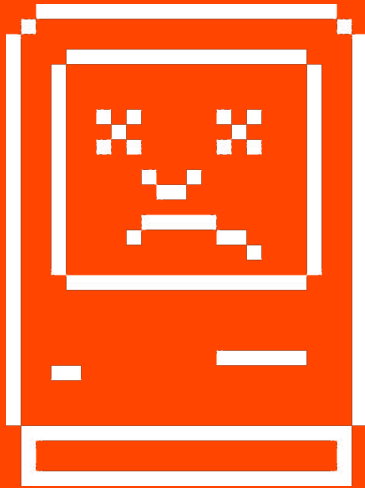
Incidents

5 + 7 = 12



Incidents

$$5 + 7 + 2 = 14$$

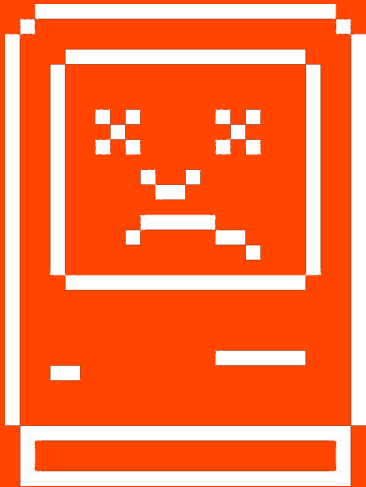


July

1	2	3	4	5	6	7
8	9	10	11	12	13	14
15	16	17	18	19	20	21
22	23	24	25	26	27	28
29	30	31				



Agenda



01

The telegraf case

02

The session case

03

Network = CPU noise

04

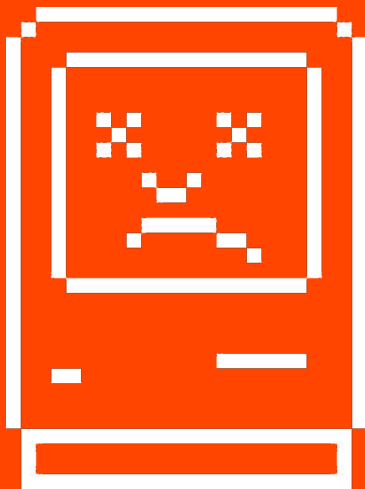
Conntrack

05

What now?



Reddit Shaped



What is “Reddit-Shaped”

- 100M Daily-Active Users
- Several “large” Kubernetes Clusters
 - >75k CPUs
 - >30k pods
 - 96 CPU / 384GB memory nodes



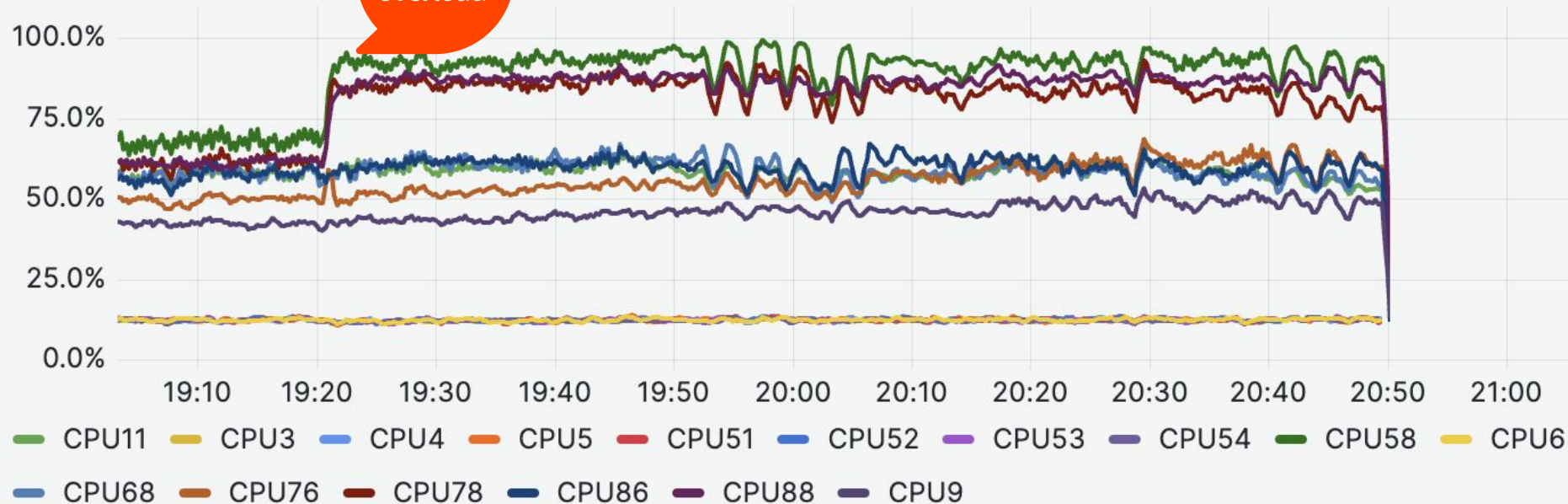
(01)

The telegraf case



(01) The telegraf case

node softirq ⓘ



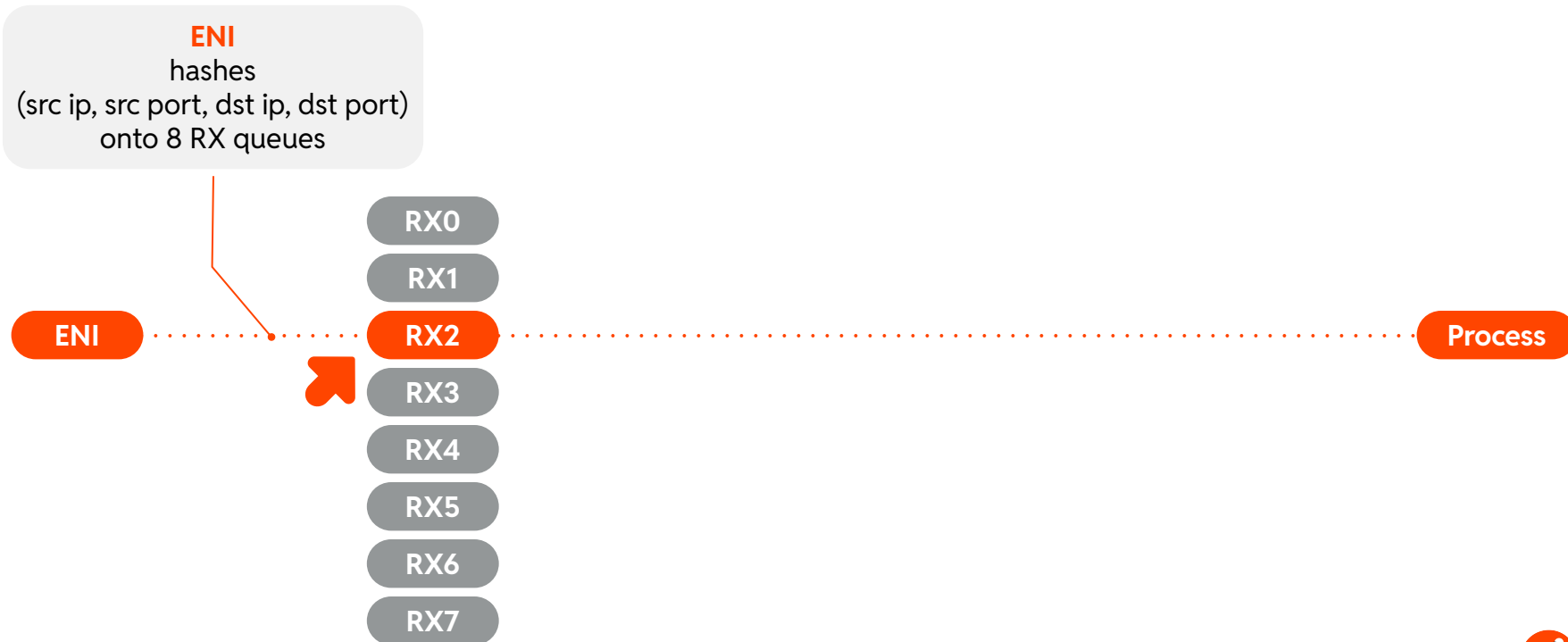
Receiving packets

AWS & Linux



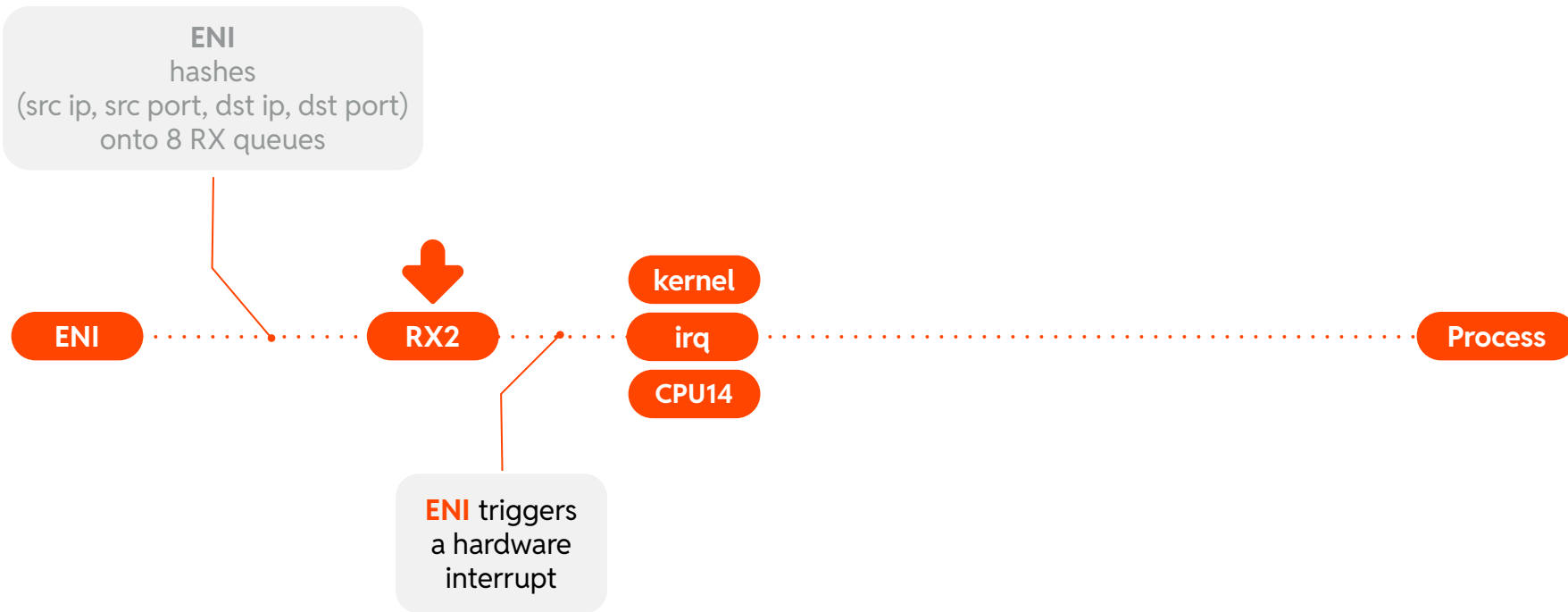
Receiving packets

AWS & Linux



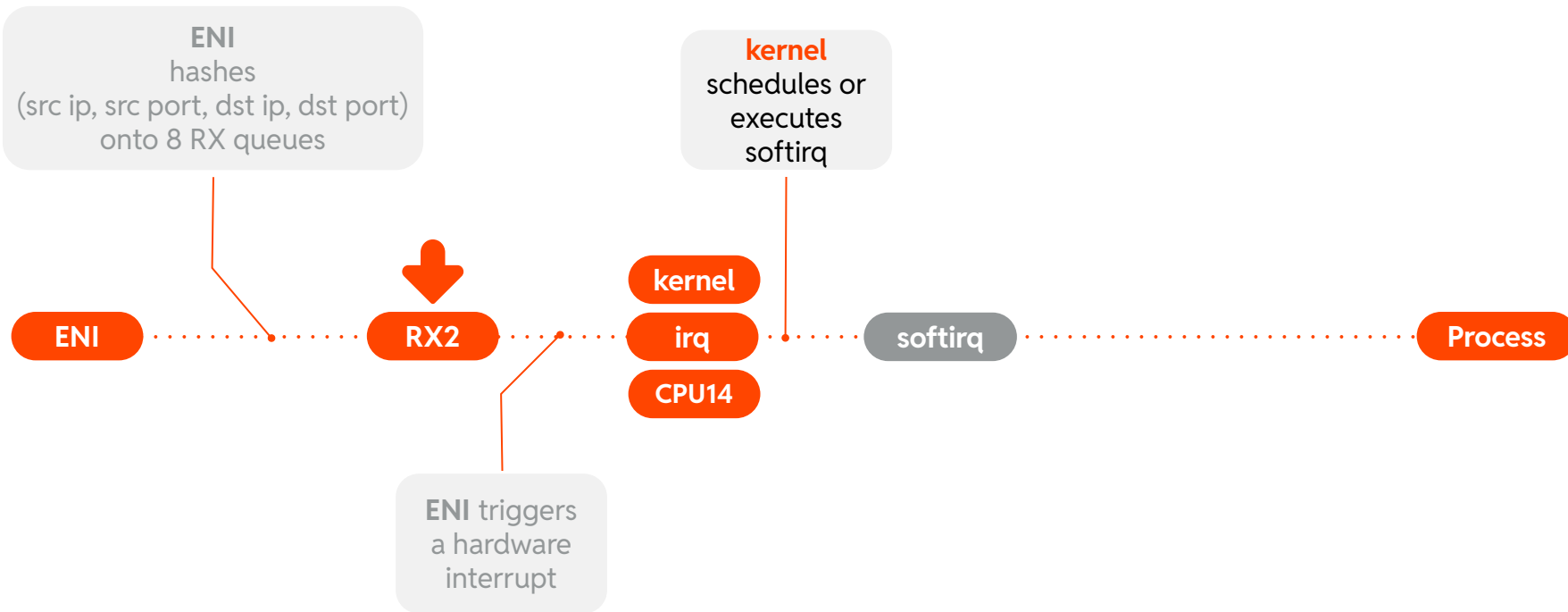
Receiving packets

AWS & Linux



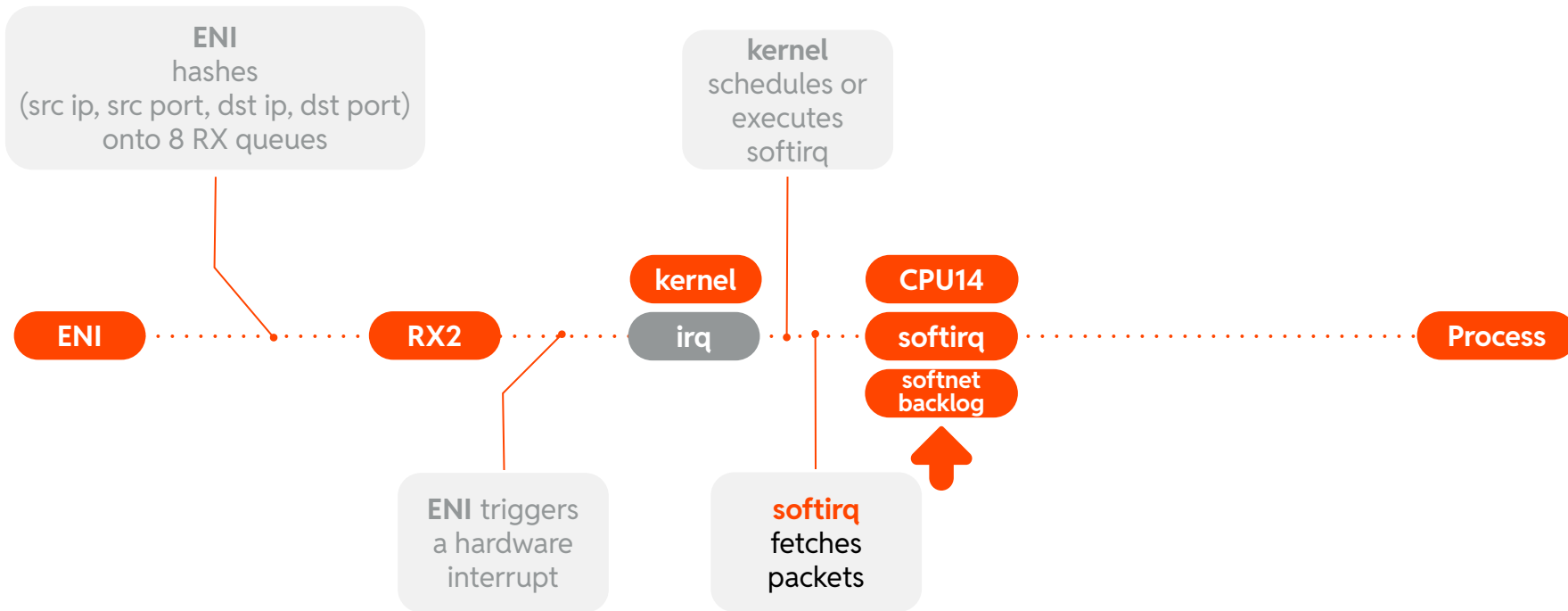
Receiving packets

AWS & Linux



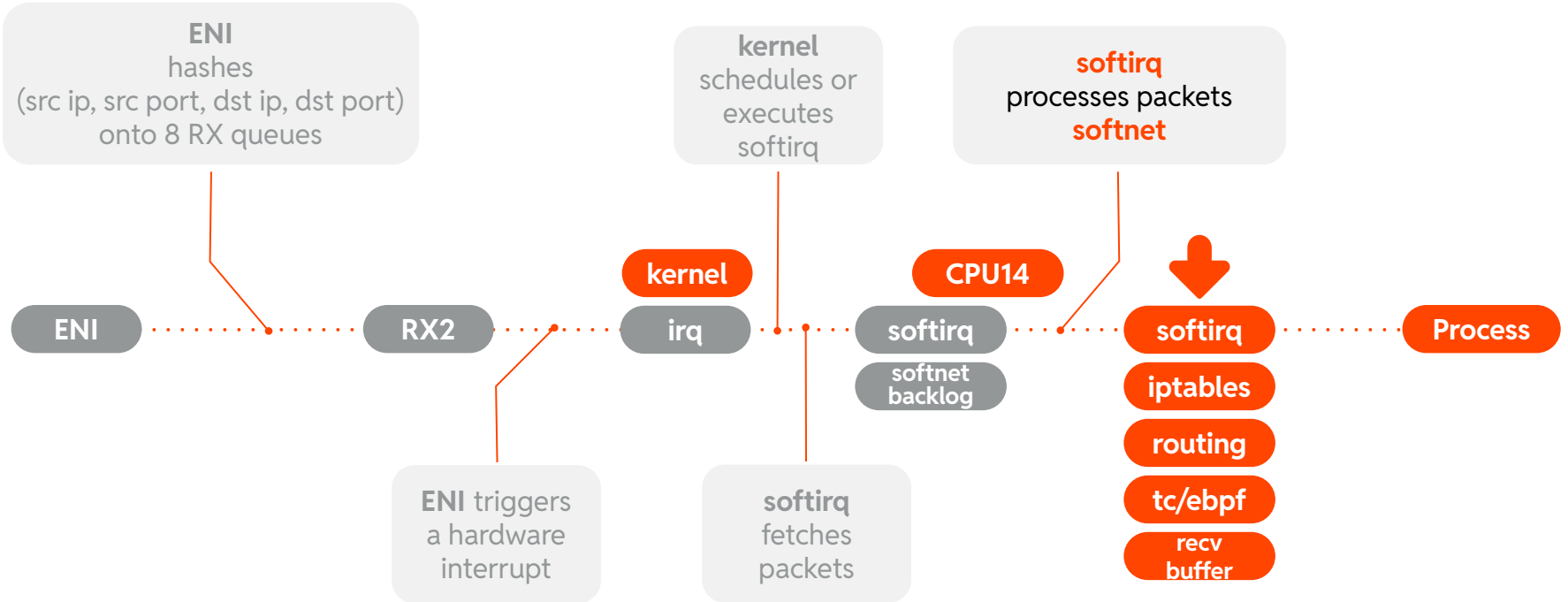
Receiving packets

AWS & Linux



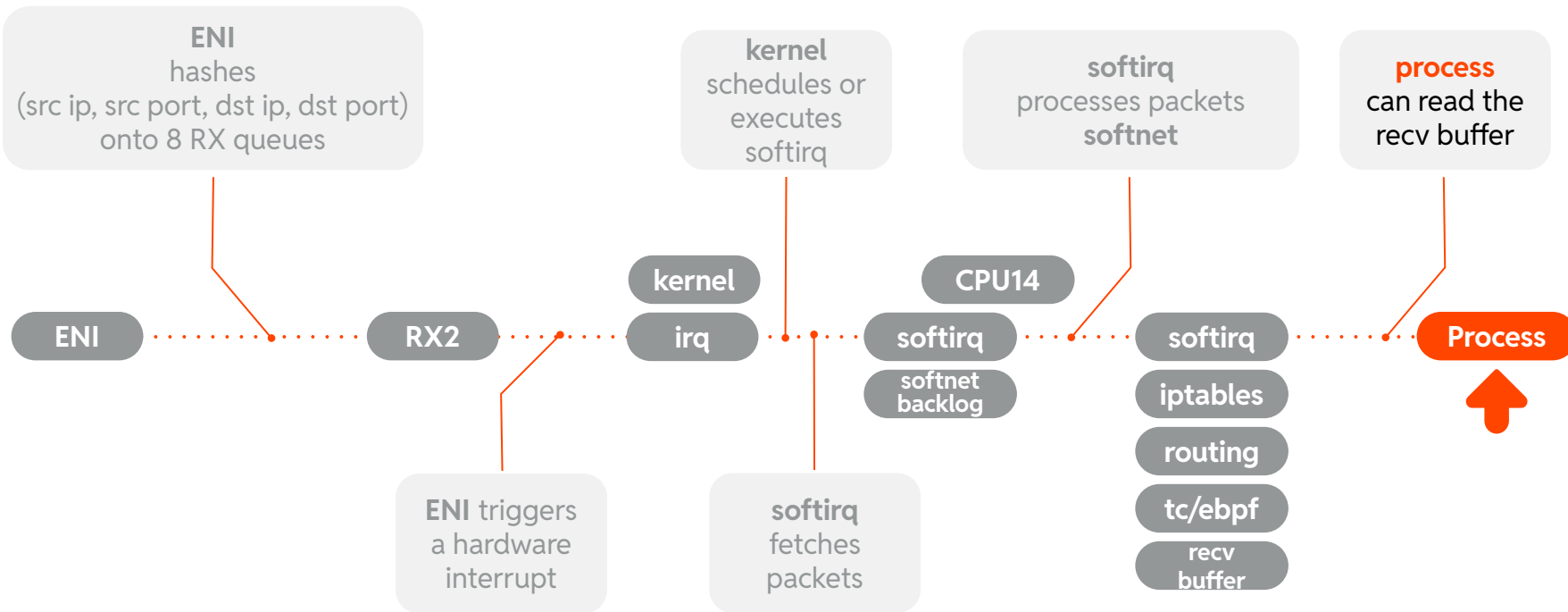
Receiving packets

AWS & Linux



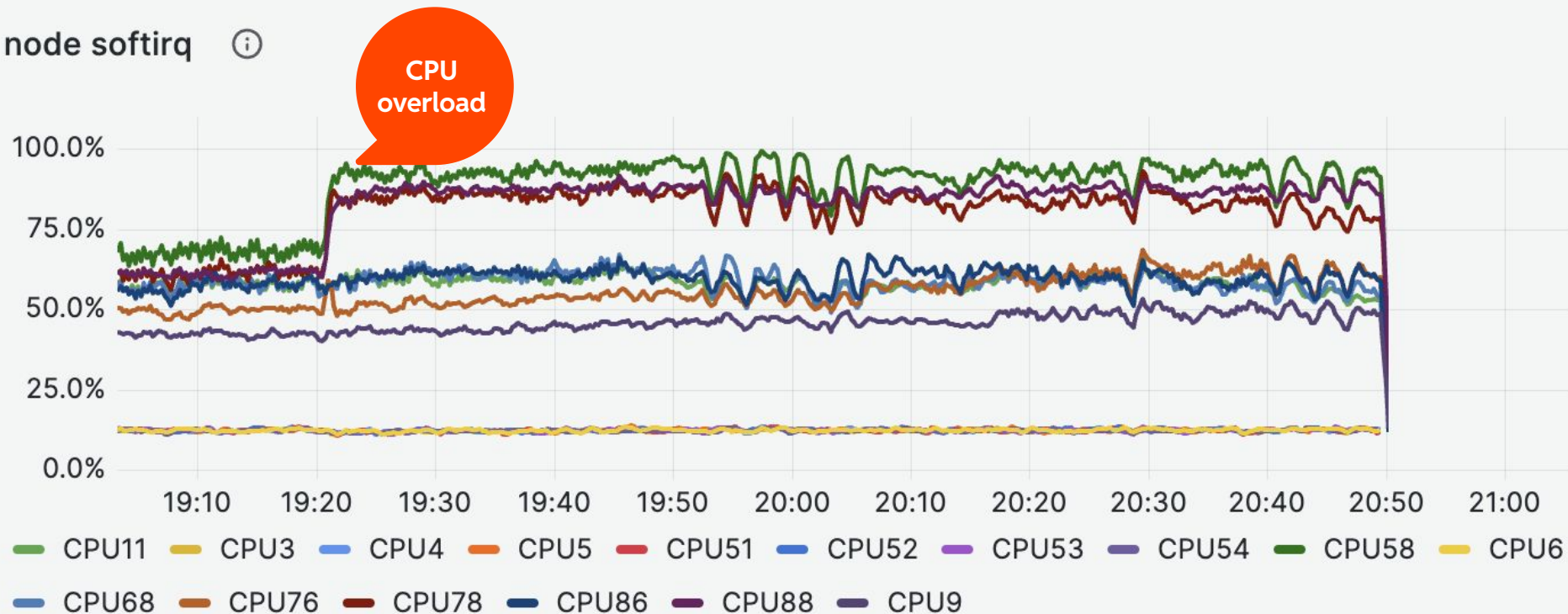
Receiving packets

AWS & Linux

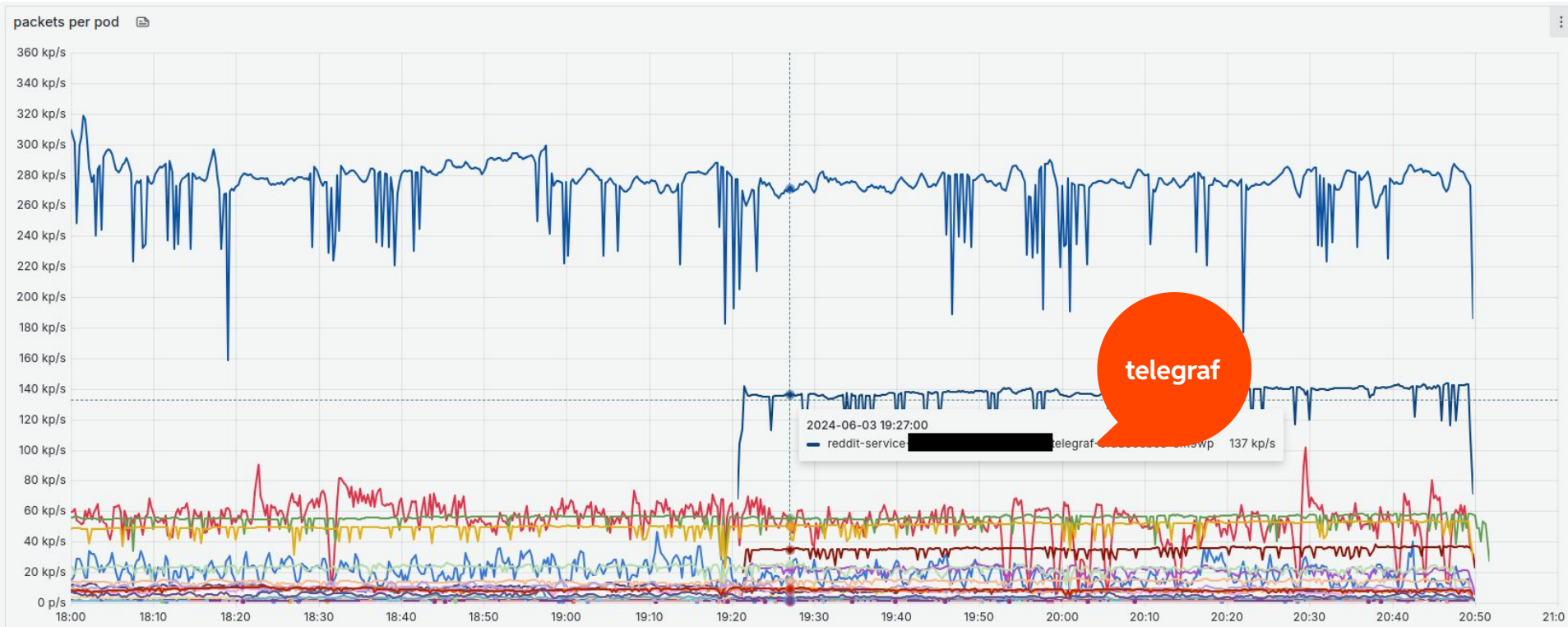


(01) The telegraf case

node softirq ⓘ



(01) The telegraf case



(01) The telegraf case - HOW?



(01) The telegraf case - HOW?

statsd metrics

many small udp packets



(01) The telegraf case - HOW?

statsd metrics

many small udp packets

4 sender, 1 receiver

loadtest like setup



(01) The telegraf case - HOW?

statsd metrics

many small udp packets

4 sender, 1 receiver

loadtest like setup

4 streams, 8 queues

load imbalance



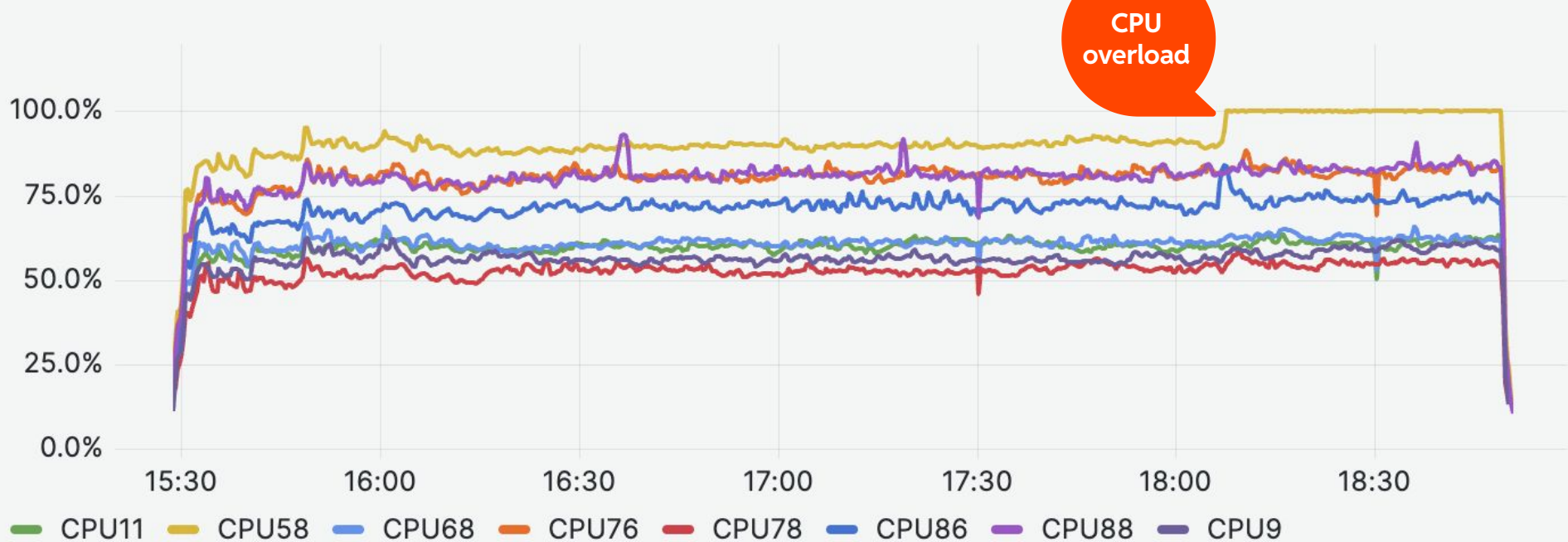
(02)

The session case



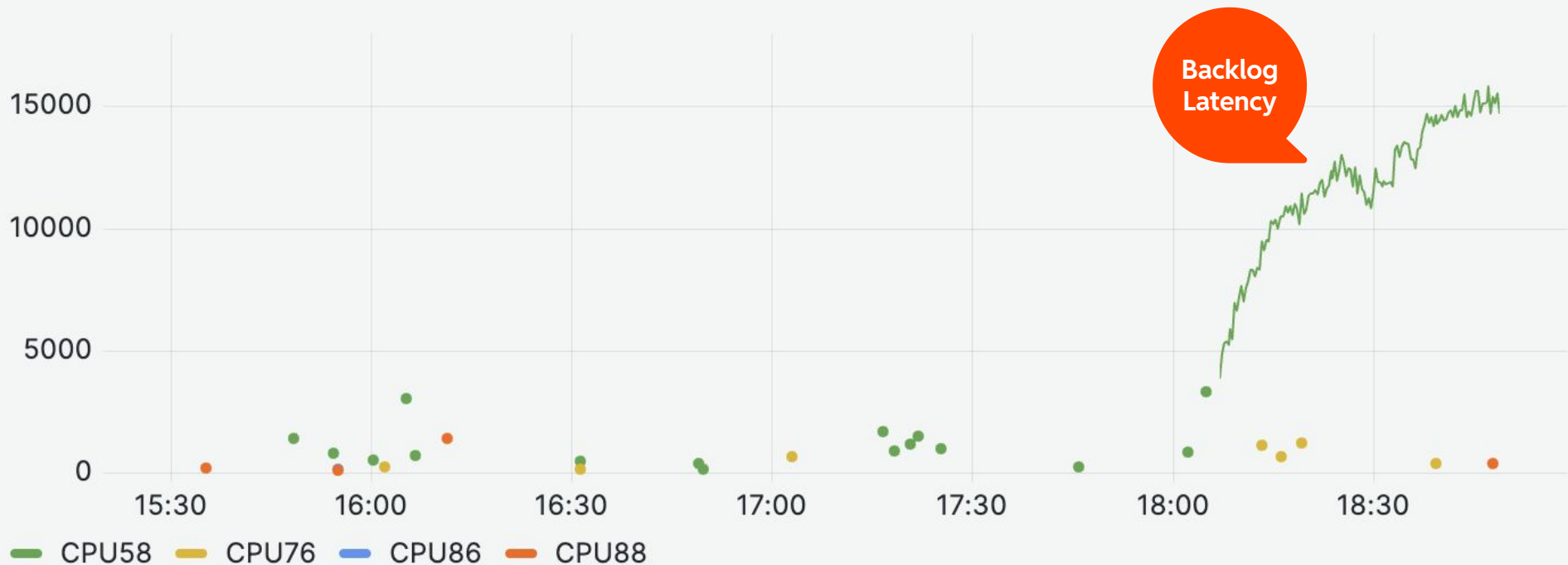
(02) The session case

node softirq ⓘ



(02) The session case

softnet backlog ⓘ



(02) The session case - HOW?



(02) The session case - HOW?

16 session pods

scale-up event



(02) The session case - HOW?

16 session pods

scale-up event

All cores pinned

throughput problem



(02) The session case - HOW?

16 session pods

scale-up event

All cores pinned

throughput problem

netdev_max_backlog

>>1000



(02) The session case - core pinning

16 session pods

5 pinned cpus

CORE0 | CPU0
CPU48

CORE1 | CPU1
CPU49

CORE2 | CPU2
CPU50

CORE3 | CPU3
CPU51



(02) The session case - core pinning

16 session pods

5 pinned cpus

CORE0 | CPU0
CPU48

CORE1 | CPU1
CPU49

CORE2 | CPU2
CPU50

CORE3 | CPU3
CPU51



(02) The session case - core pinning

16 session pods

5 pinned cpus

CORE0 | CPU0
CPU48

CORE1 | CPU1
CPU49

CORE2 | CPU2
CPU50

CORE3 | CPU3
CPU51

$$16 \times 2.5 = 40$$

$$16 \times 3 = 48$$



(02) The session case - core pinning

16 session pods

5 pinned cpus

CORE0

CPU0
CPU48

CORE1

CPU1
CPU49

CORE2

CPU2
CPU50

CORE3

CPU3
CPU51

$$16 \times 2.5 = 40$$

$$16 \times 3 = 48$$

8 cores for network?



(02) The session case - core pinning

16 session pods

5 pinned cpus

CORE0

CPU0
CPU48

CORE1

CPU1
CPU49

CORE2

CPU2
CPU50

CORE3

CPU3
CPU51

$$16 \times 2.5 = 40$$

$$16 \times 3 = 48$$

8 cores for network?

hyper threading != isolation



(03)

Network = CPU

noise

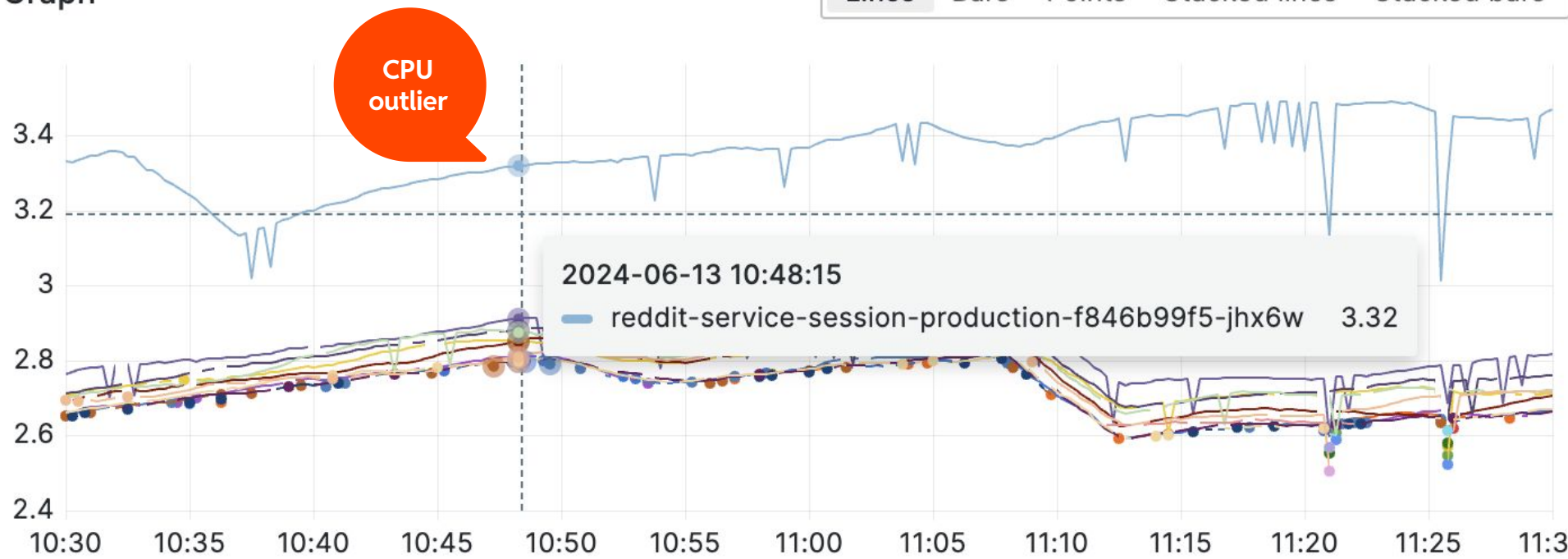
(session again)



(03) Network = CPU noise

Graph

Lines Bars Points Stacked lines Stacked bars



(03) Network = CPU noise

CORE9 | CPU9
CPU57

CORE10 | CPU10
CPU58

CORE11 | CPU11
CPU59



(03) Network = CPU noise

CORE9

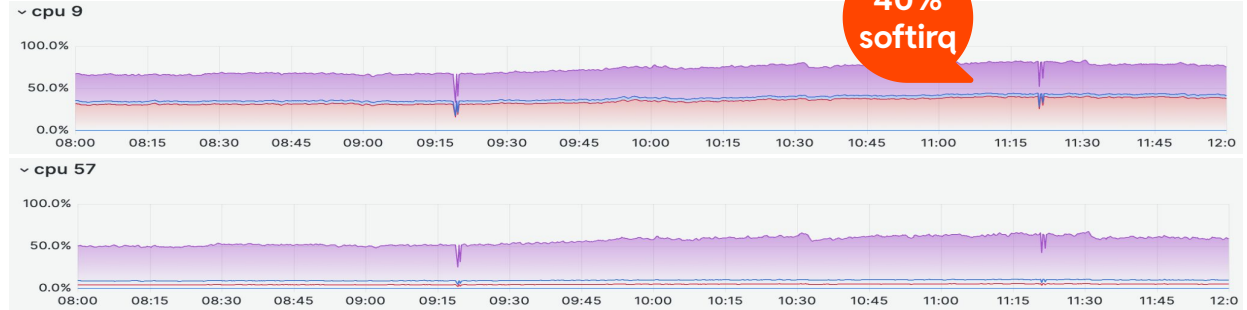
CPU9
CPU57

CORE10

CPU10
CPU58

CORE11

CPU11
CPU59



(03) Network = CPU noise

CORE9 | CPU9
CPU57

CORE10 | CPU10
CPU58

CORE11 | CPU11
CPU59



40%
softirq



(03) Network = CPU noise

CORE9

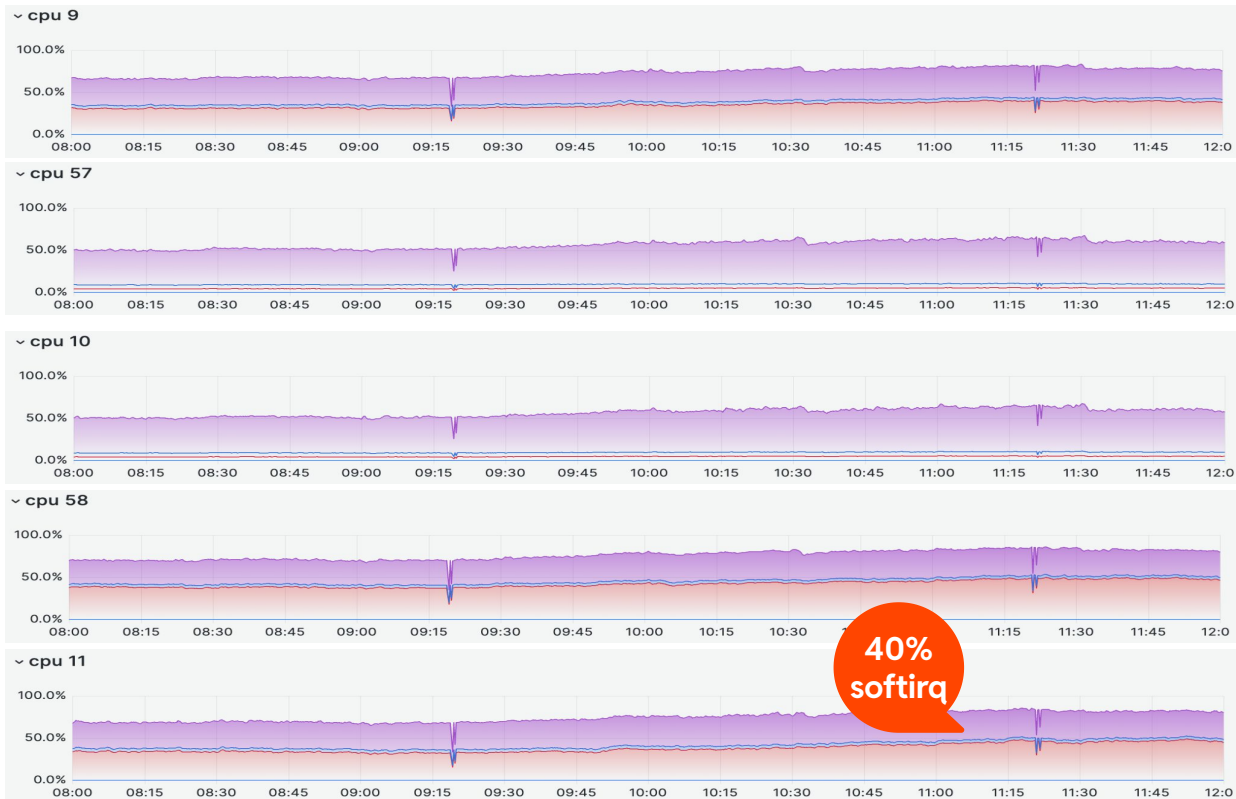
CPU9
CPU57

CORE10

CPU10
CPU58

CORE11

CPU11
CPU59

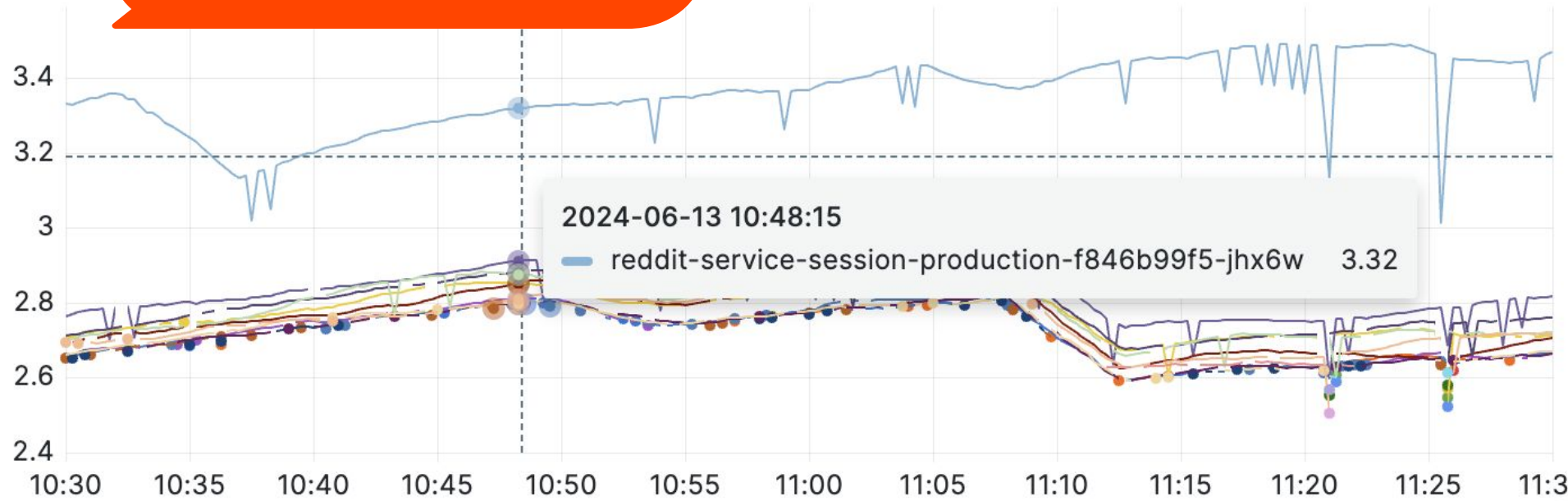


(03) Network = CPU noise

Graph

$$5 - 3 \times 0.4 = 3.8$$

Lines Bars Points Stacked lines Stacked bars



(03) Network = CPU noise

CORE9

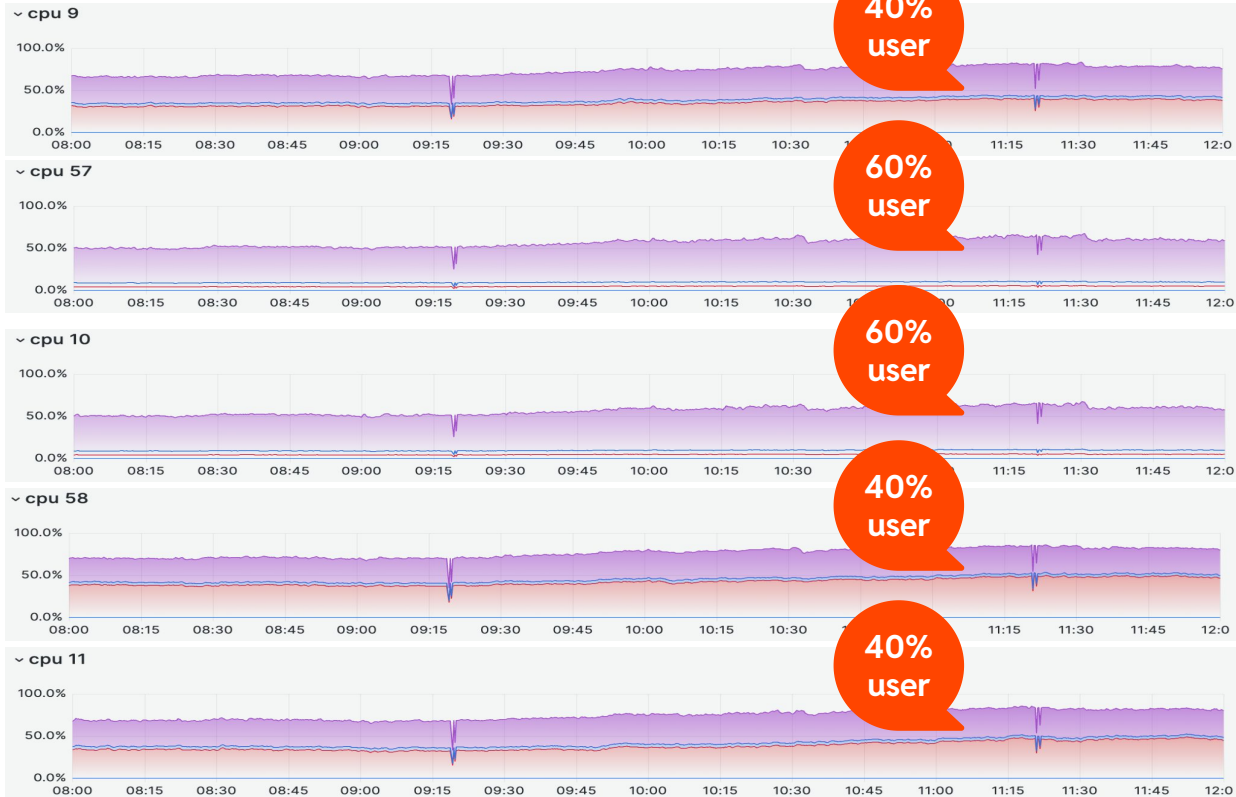
CPU9
CPU57

CORE10

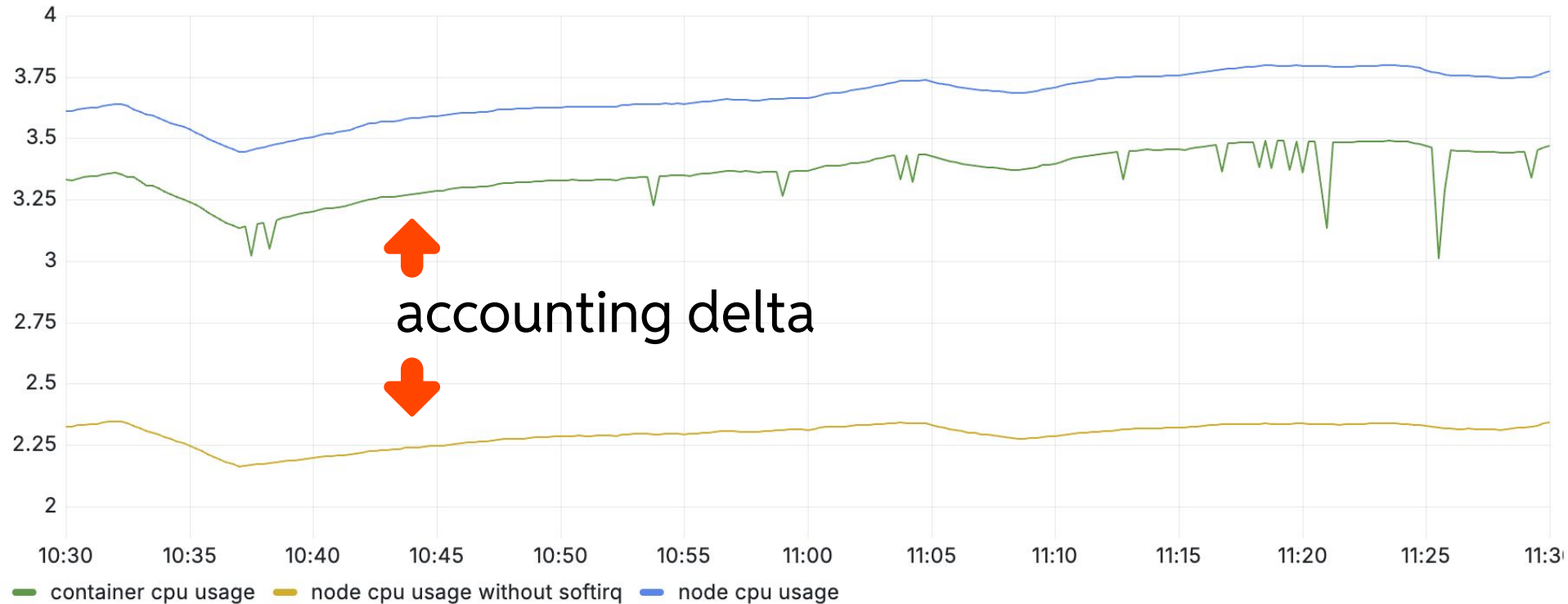
CPU10
CPU58

CORE11

CPU11
CPU59



(03) Network = CPU noise



(03) Network = CPU noise - HOW?



(03) Network = CPU noise - HOW?

pinned cpus

kernel won't move you



(03) Network = CPU noise - HOW?

pinned cpus

kernel won't move you

high softirq node

reduced cpu available



(03) Network = CPU noise - HOW?

pinned cpus

kernel won't move you

high softirq node

reduced cpu available

unlucky core pick

probabilities...



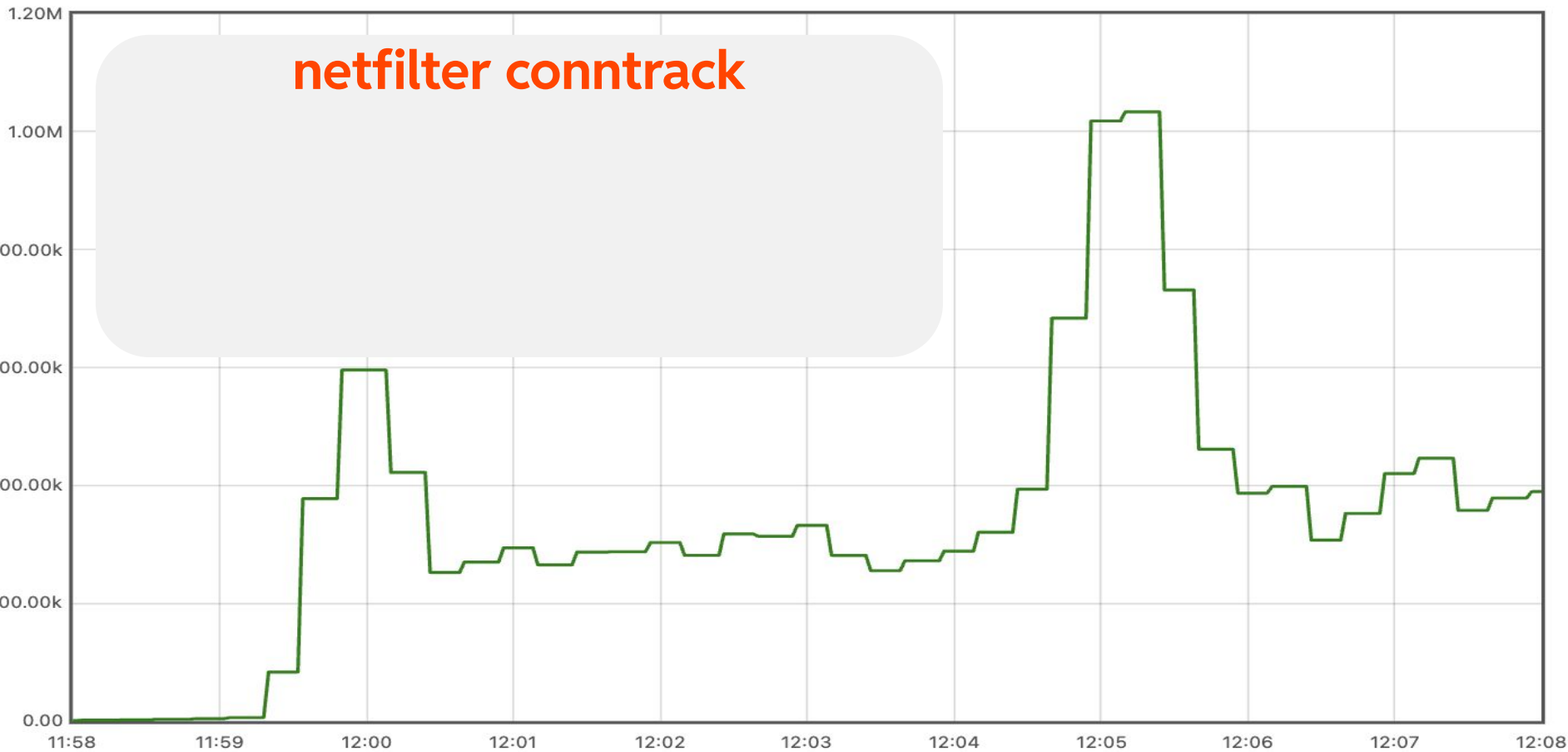
(04)

Conntrack



(04) Conntrack

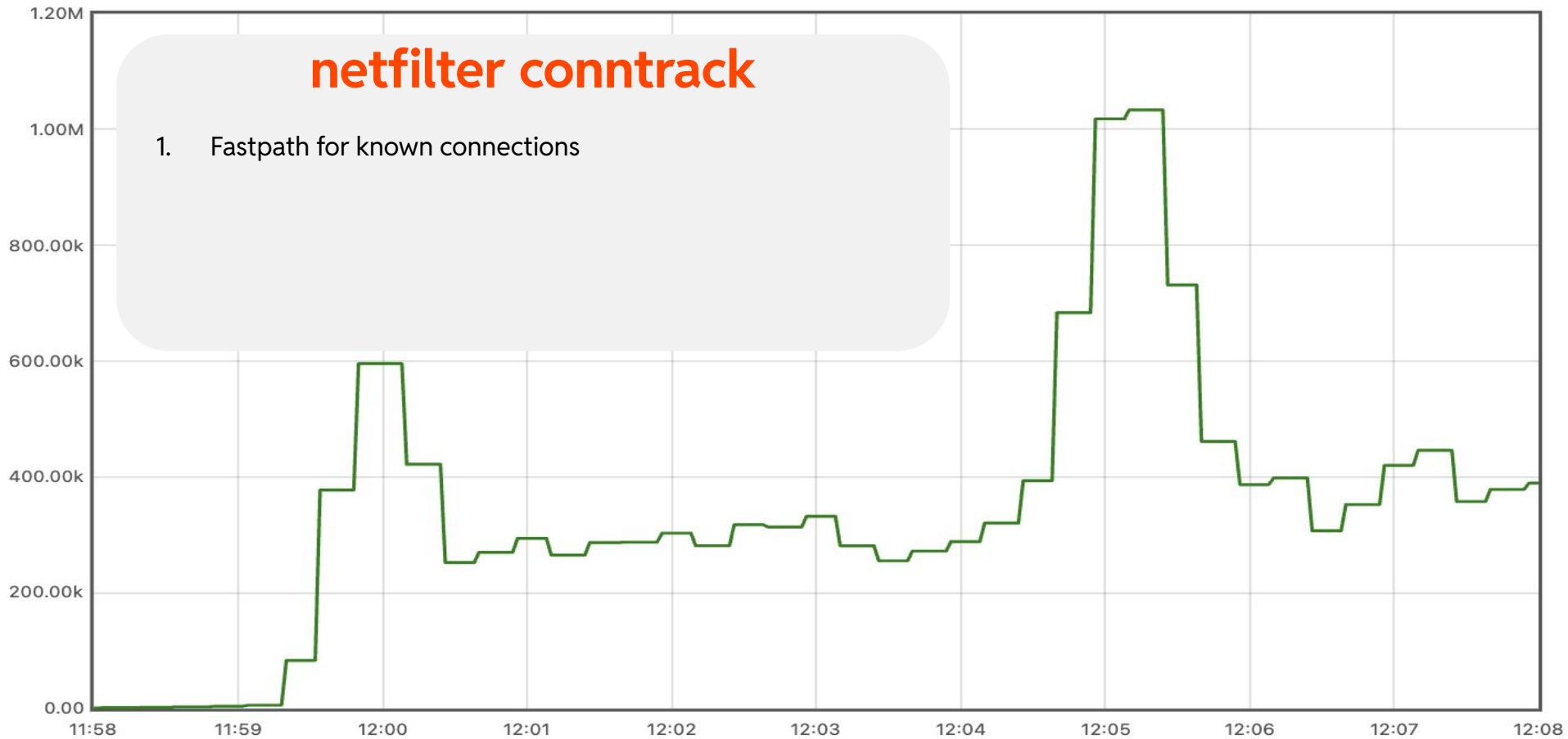
netfilter conntrack



(04) Conntrack

netfilter conntrack

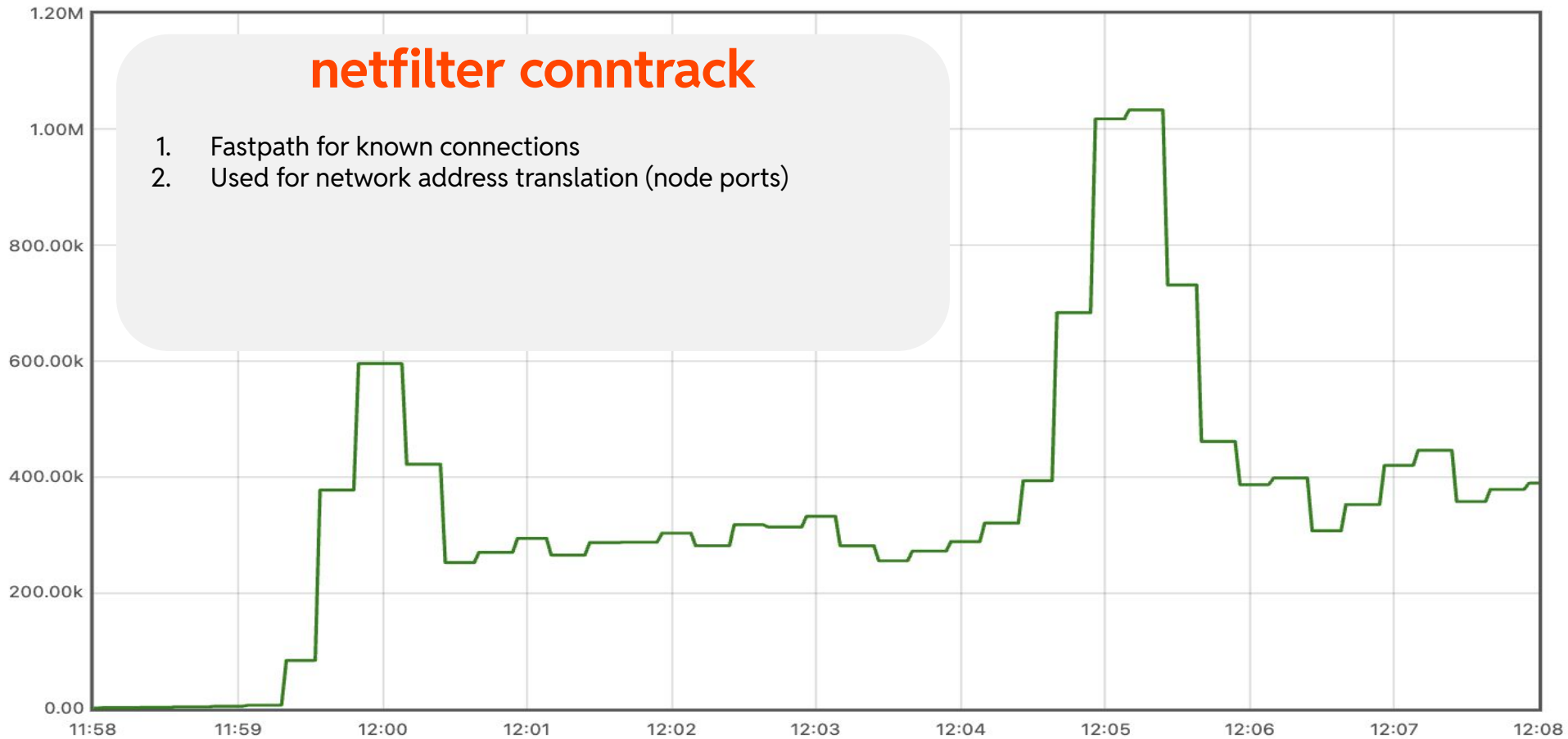
1. Fastpath for known connections



(04) Conntrack

netfilter conntrack

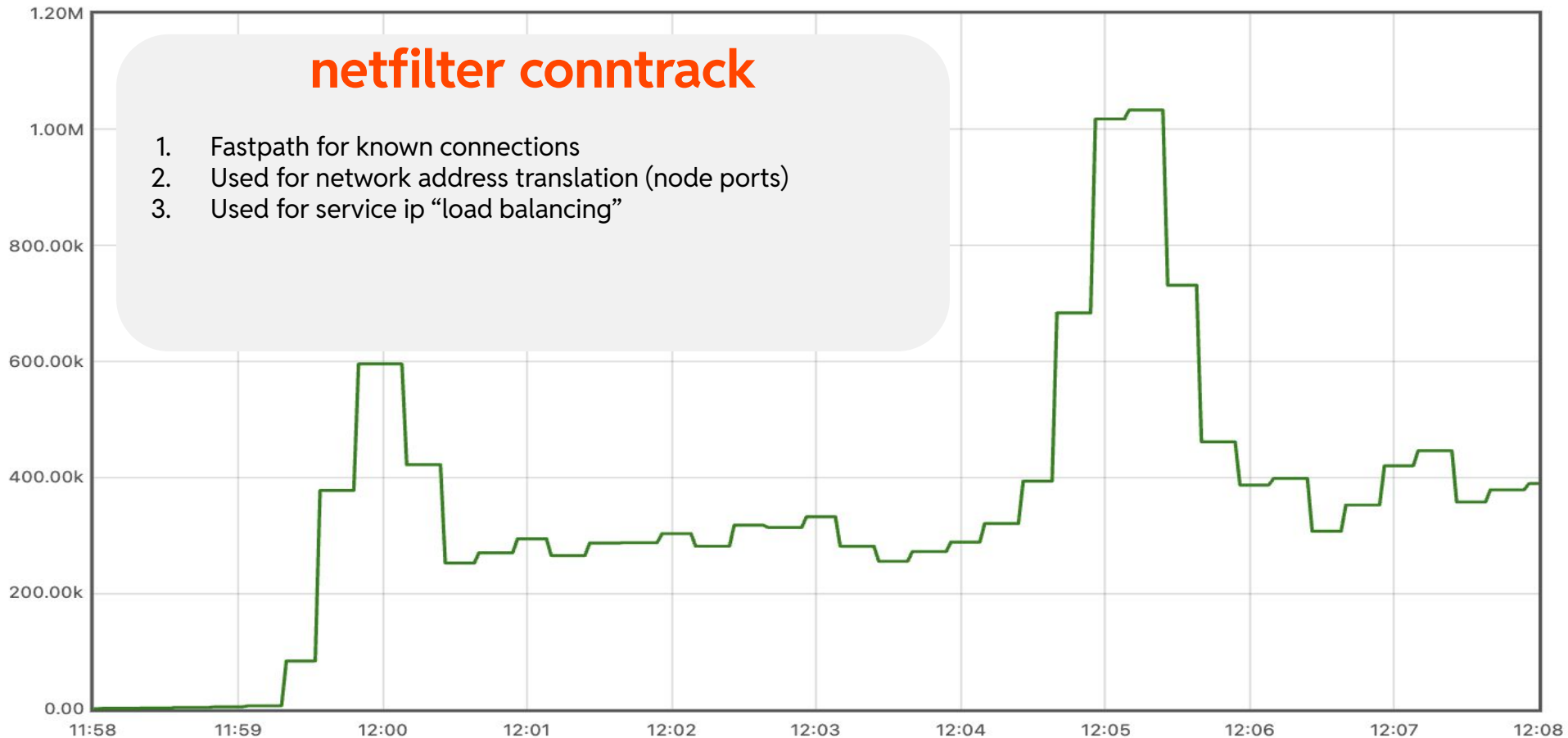
1. Fastpath for known connections
2. Used for network address translation (node ports)



(04) Conntrack

netfilter conntrack

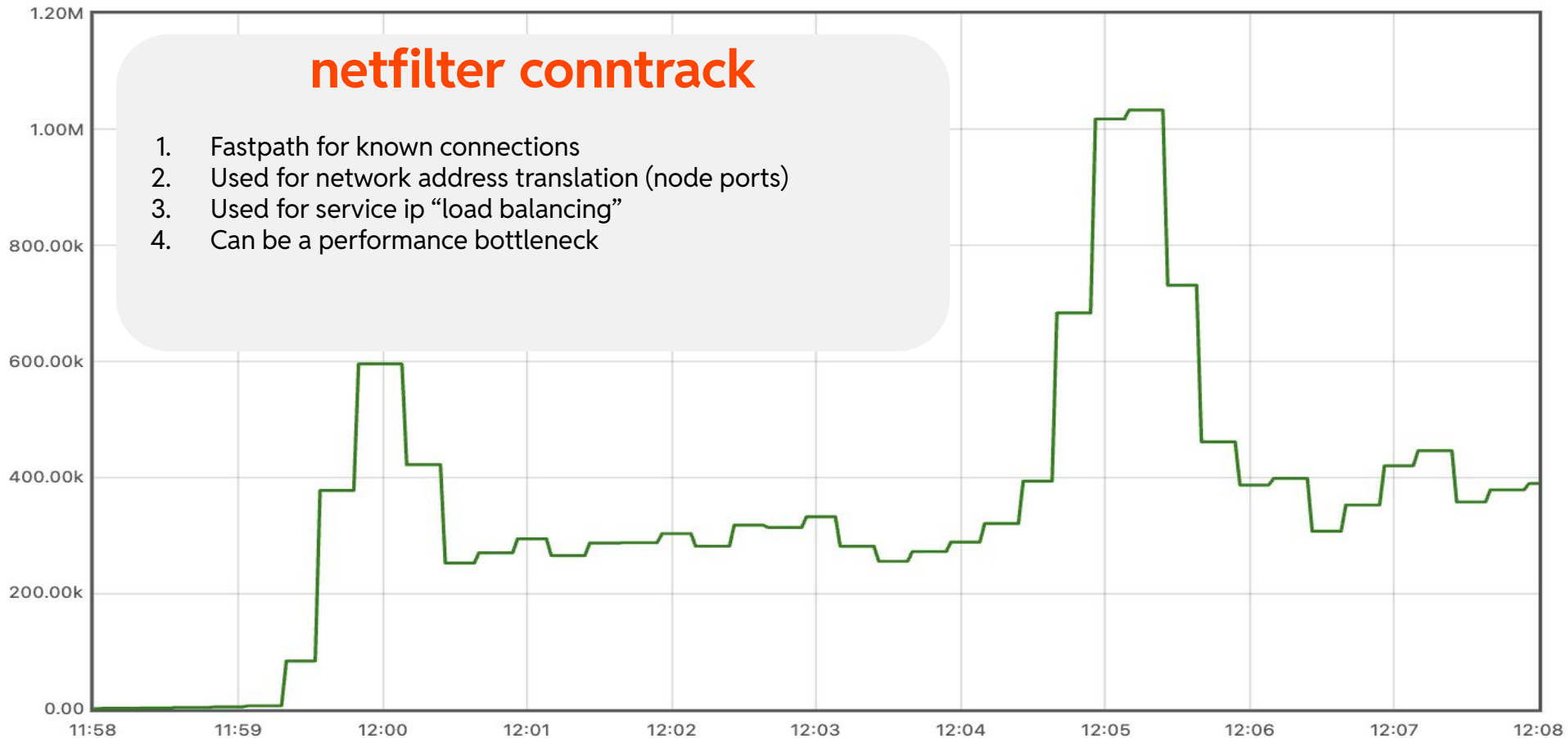
1. Fastpath for known connections
2. Used for network address translation (node ports)
3. Used for service ip “load balancing”



(04) Conntrack

netfilter conntrack

1. Fastpath for known connections
2. Used for network address translation (node ports)
3. Used for service ip “load balancing”
4. Can be a performance bottleneck

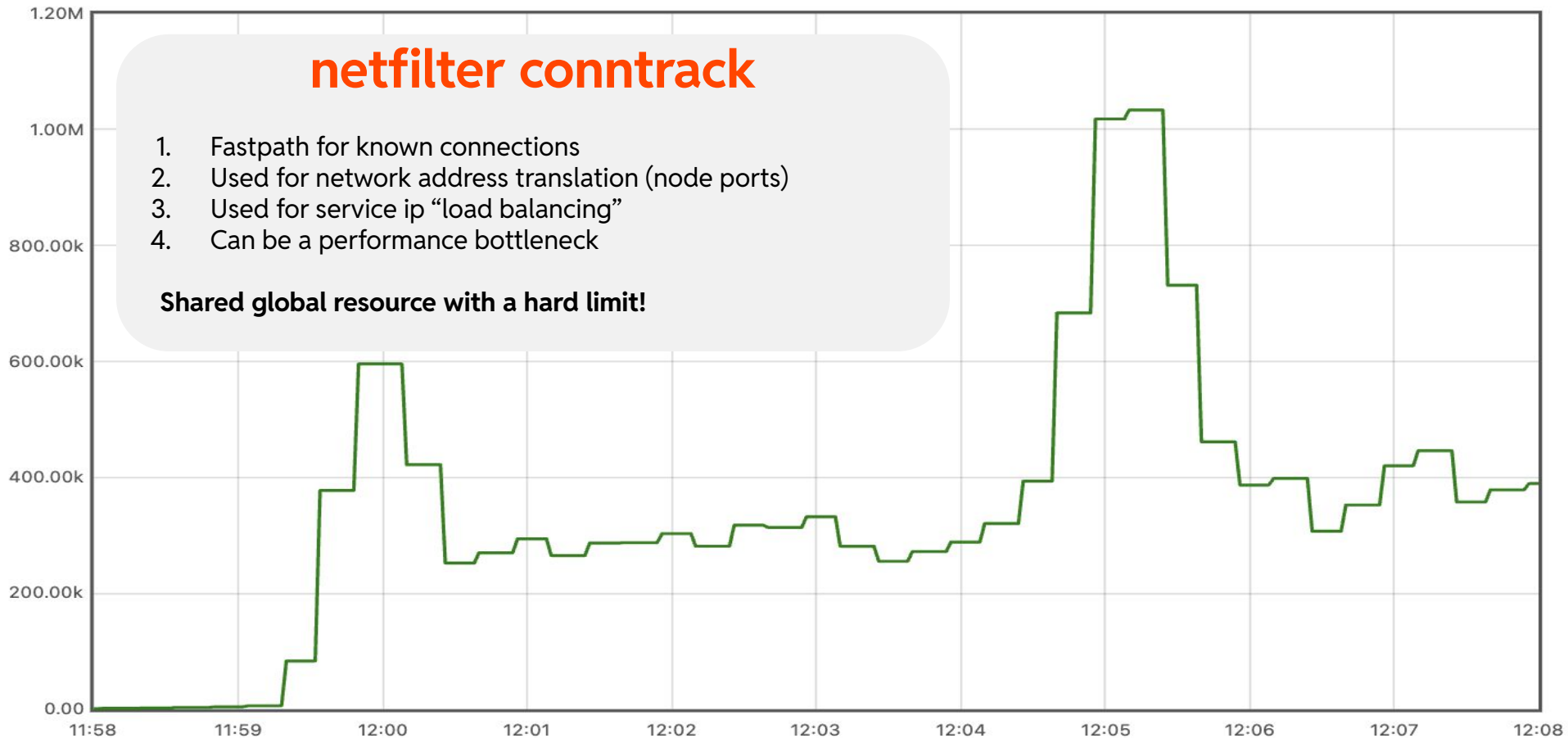


(04) Conntrack

netfilter conntrack

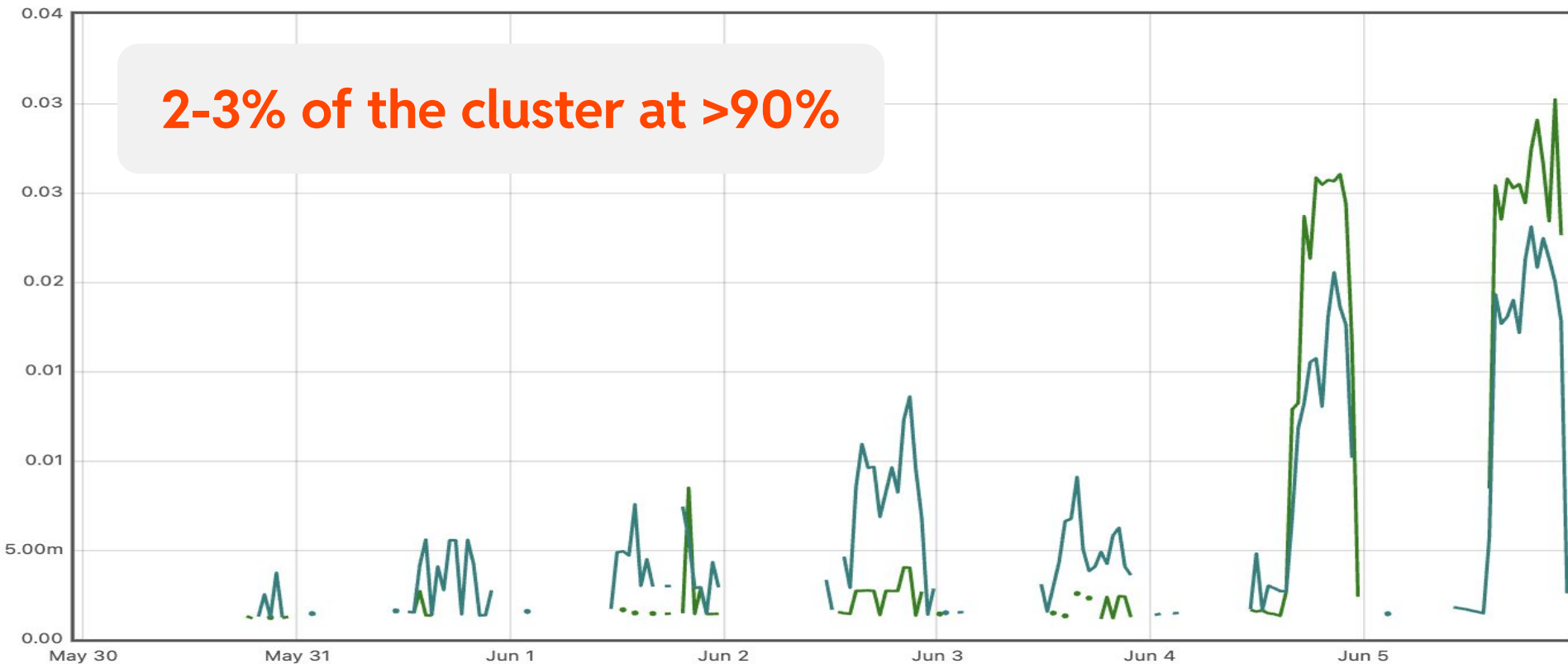
1. Fastpath for known connections
2. Used for network address translation (node ports)
3. Used for service ip “load balancing”
4. Can be a performance bottleneck

Shared global resource with a hard limit!



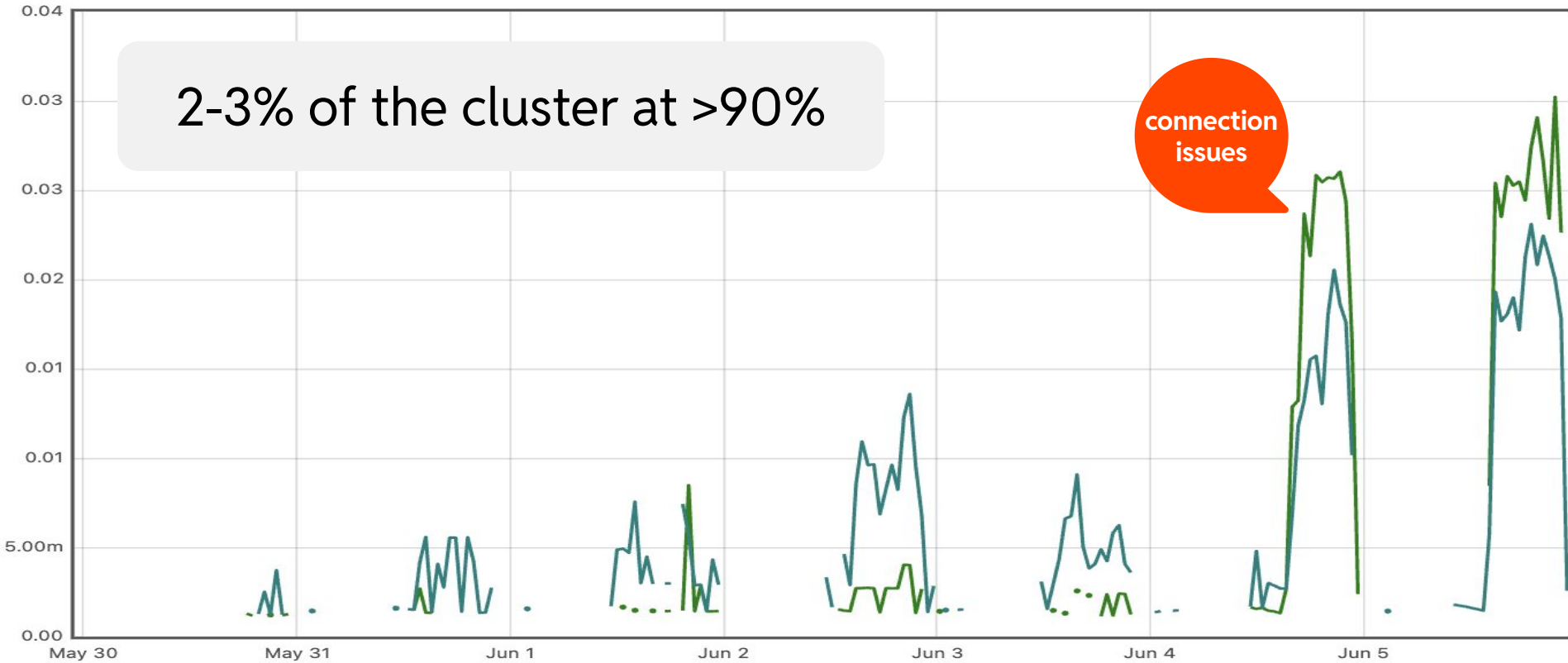
(04) Contrack

2-3% of the cluster at >90%



■ {cluster="prod-3d"}
■ {cluster="prod-3e"}

(04) Contrack



■ {cluster="prod-3d"}

■ {cluster="prod-3e"}

(04) Conntrack - HOW?



(04) Conntrack - HOW?

k8s AMI update

undetected config change



(04) Contrack - HOW?

k8s AMI update

undetected config change

1M contrack limit

way lower



(04) Contrack - HOW?

k8s AMI update

undetected config change

1M contrack limit

way lower

k8s upgrade

node rollover



(04) Contrack - HOW?

k8s AMI update

undetected config change

1M contrack limit

way lower

k8s upgrade

node rollover

1 week buildup

kube-proxy bugfix



05

What
now?



(05) What now? **conntrack**



(05) What now?

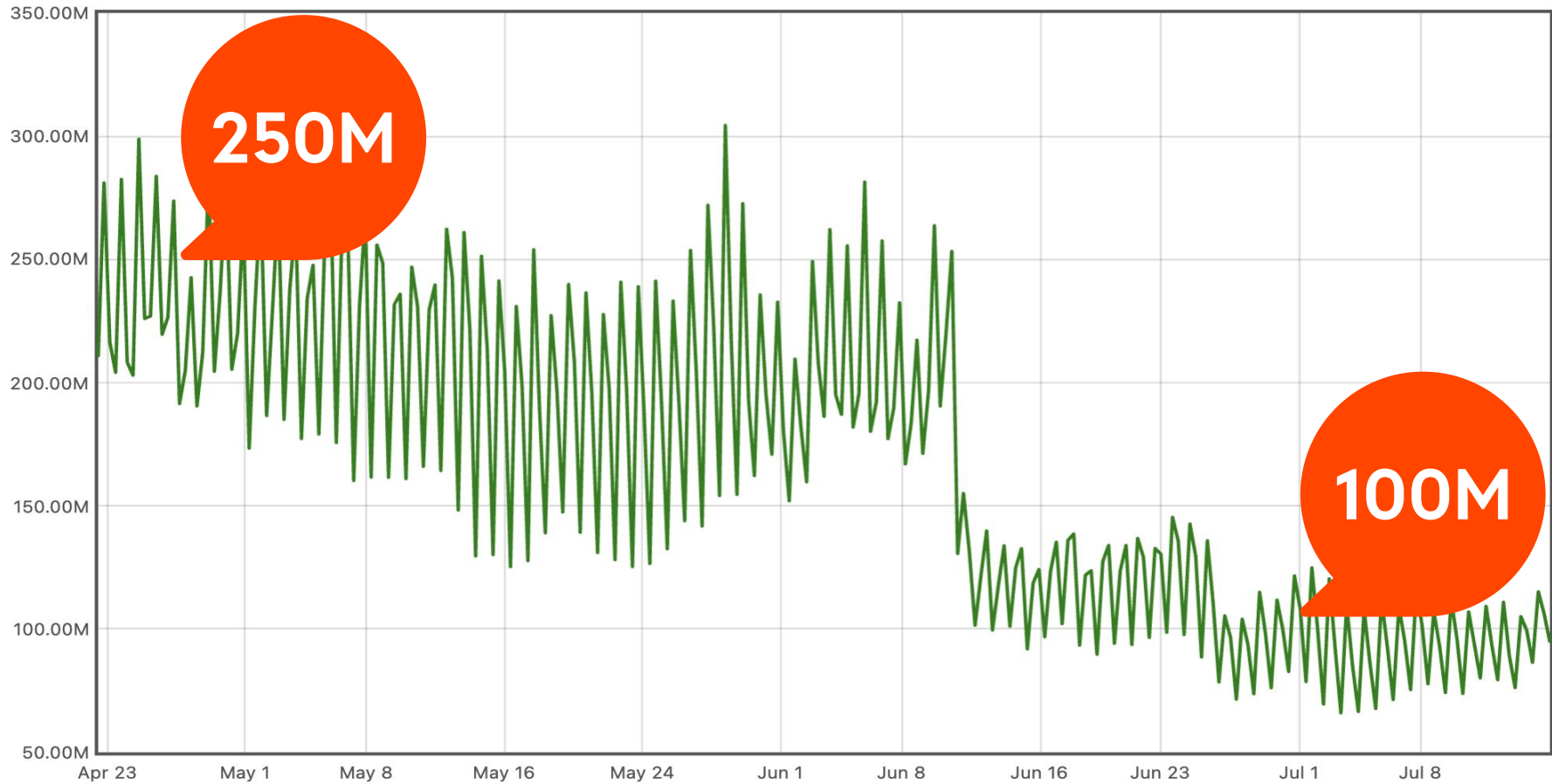
conntrack / optimize

wg-conntrack-conn-numb

[@anton.kuklin](#) created this channel on June 10th.

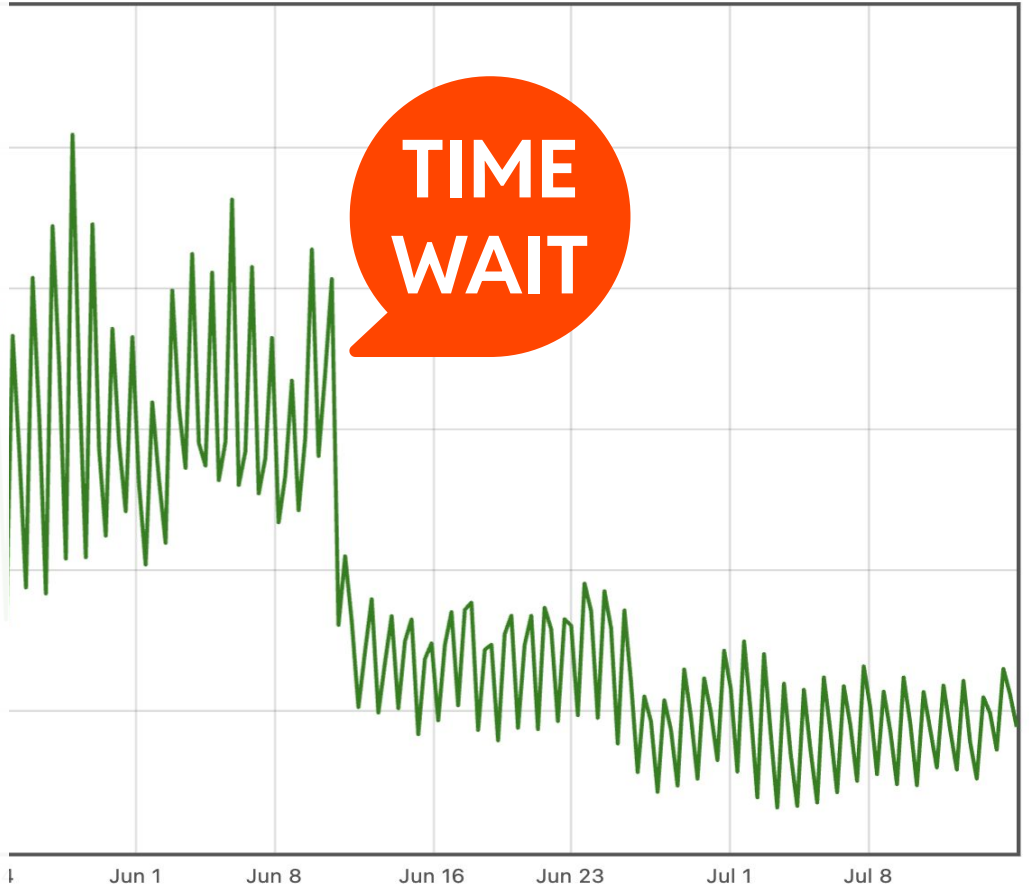
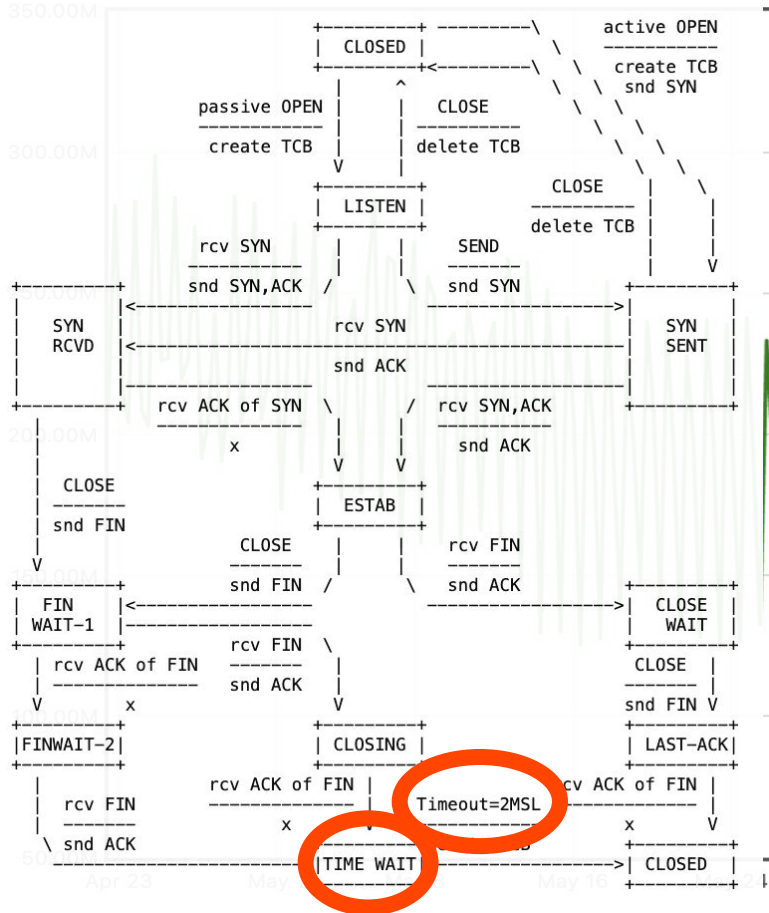


(05) What now? **contrack / optimize**



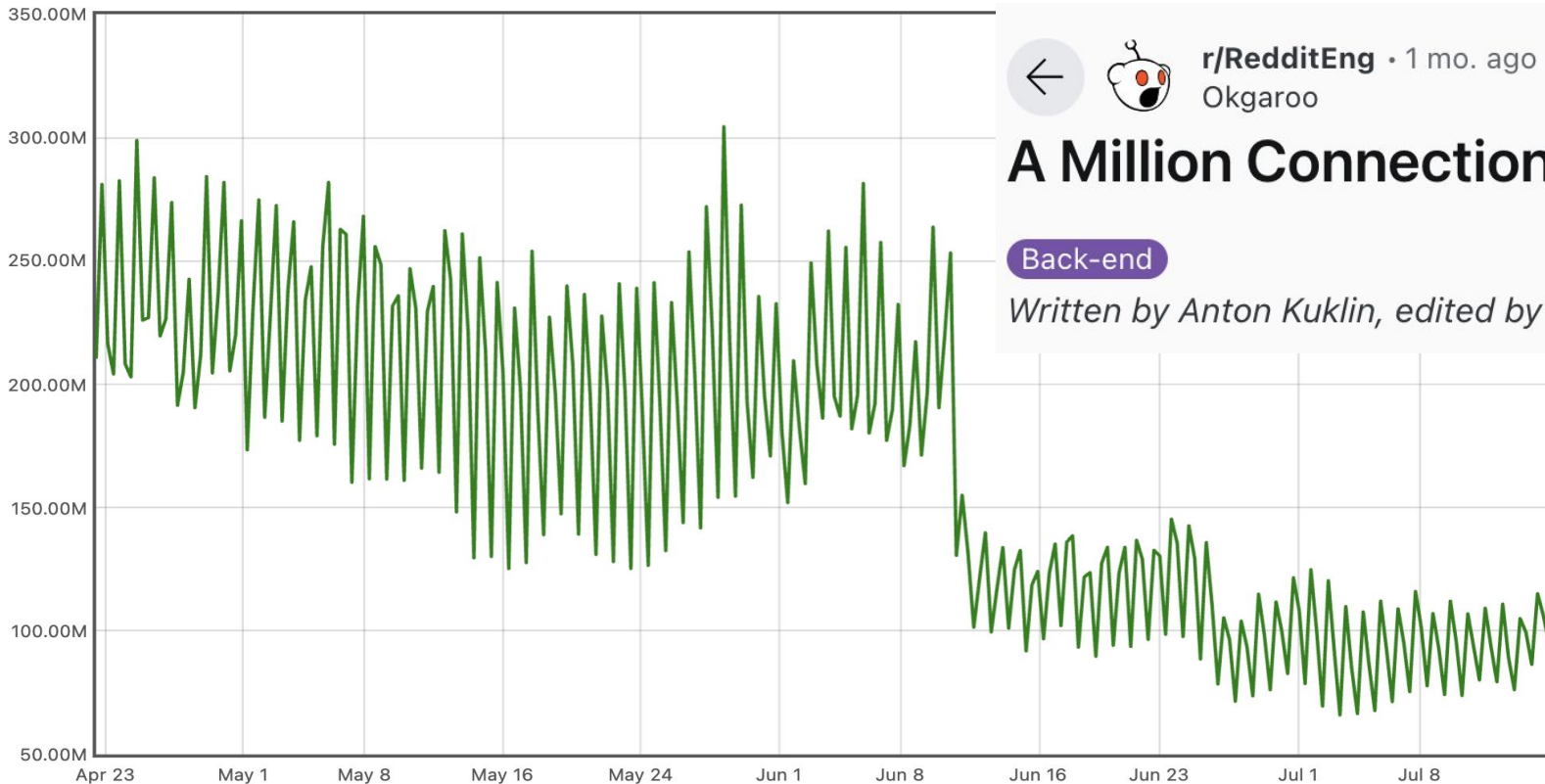
(05) What now?

conntrack / optimize



(05) What now?

conntrack / optimize



r/RedditEng · 1 mo. ago

Okgaroo

A Million Connection Problem

Back-end

Written by Anton Kuklin, edited by René Treffer

https://www.reddit.com/r/RedditEng/comments/1fnt8af/a_million_connection_problem/

(05) What now? **conntrack / alerting**

alert: NodesNumbWithTooManyConntrackRowsTooHigh

annotations:

summary: "prod-3d: Too many nodes exceeding conntrack rows threshold"

description: "prod-3d: Too many nodes exceeding conntrack rows threshold, which is $\leq 15\%$ of nodes with $\geq 400k$ rows"

dashboard: <https://grafana.kubernetes.ue1.sno.westeurope.net/d/00cf4fdb-46f5-4063-9eb8-f22a76caf614/node-network-issues?orgId=1>

expr: |

```
count(
  node_nf_conntrack_entries{} > 400000
```

```
)
```

```
/
```

```
count(
  node_nf_conntrack_entries{
```

```
) > 0.15
```

400k soft limit

15% of nodes

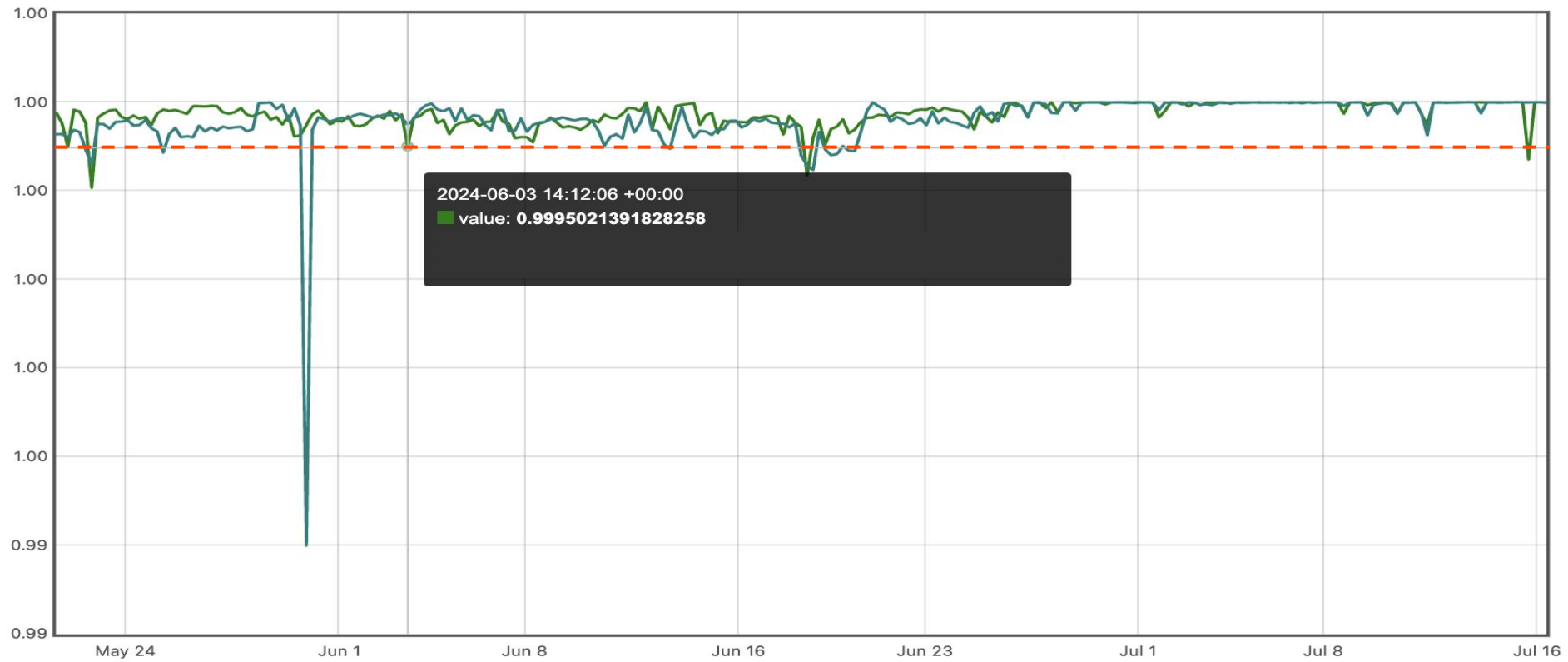


(05) What now? **Session**



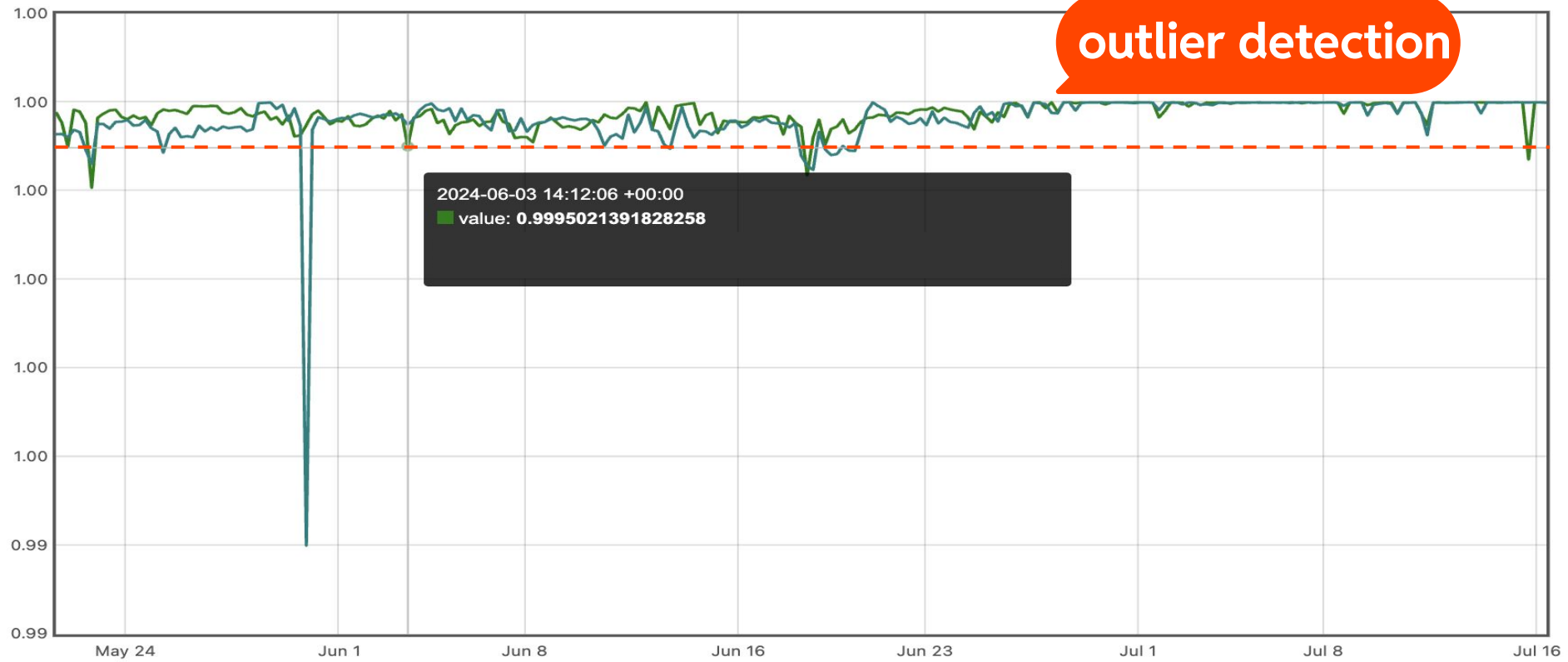
(05) What now? **Session**

Only raw data

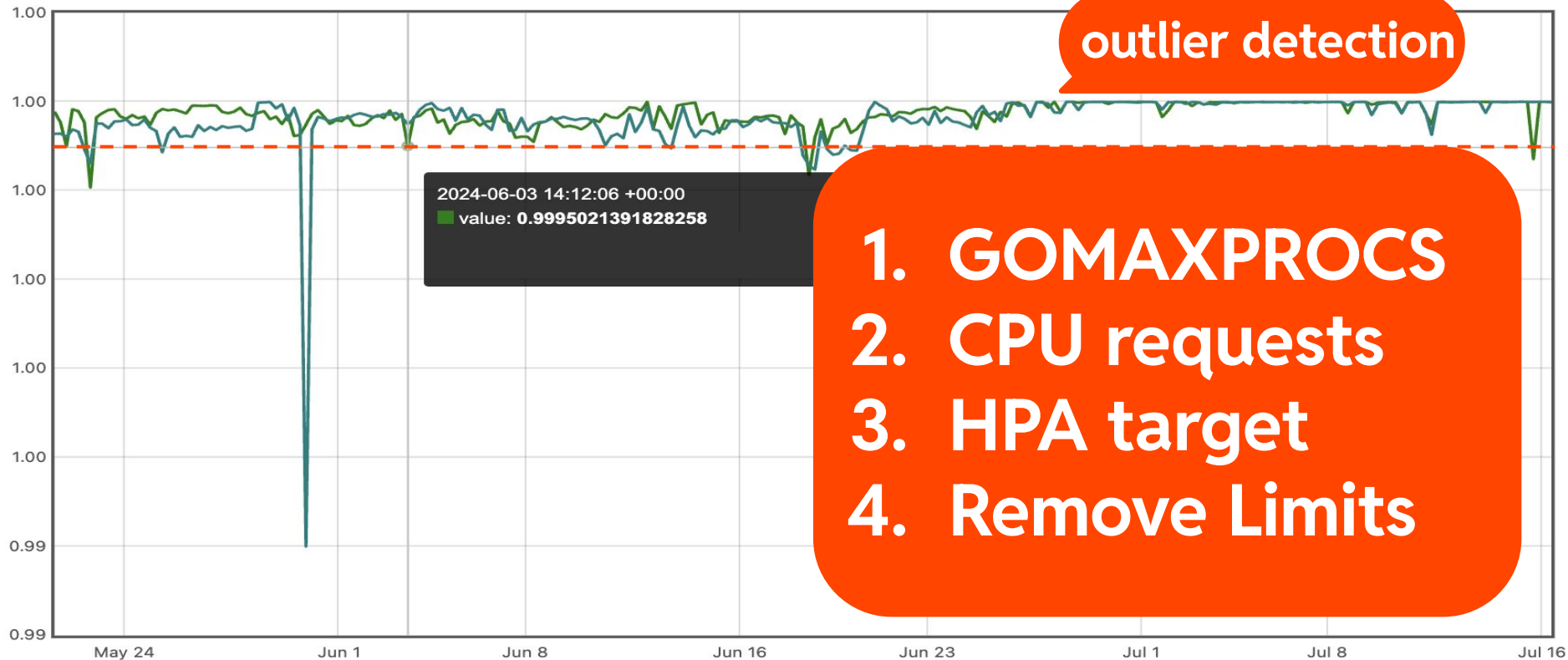


(05) What now? Session

Only raw data



(05) What now? Session



(05) What now? **Telegraf**



(05) What now? **Telegraf**

```
sum by (container) (container:cpu_usage:5m{prometheus="monitoring/monitoring",cluster_group="[REDACTED]",namespace="reddit-service-[REDACTED]"}))
```

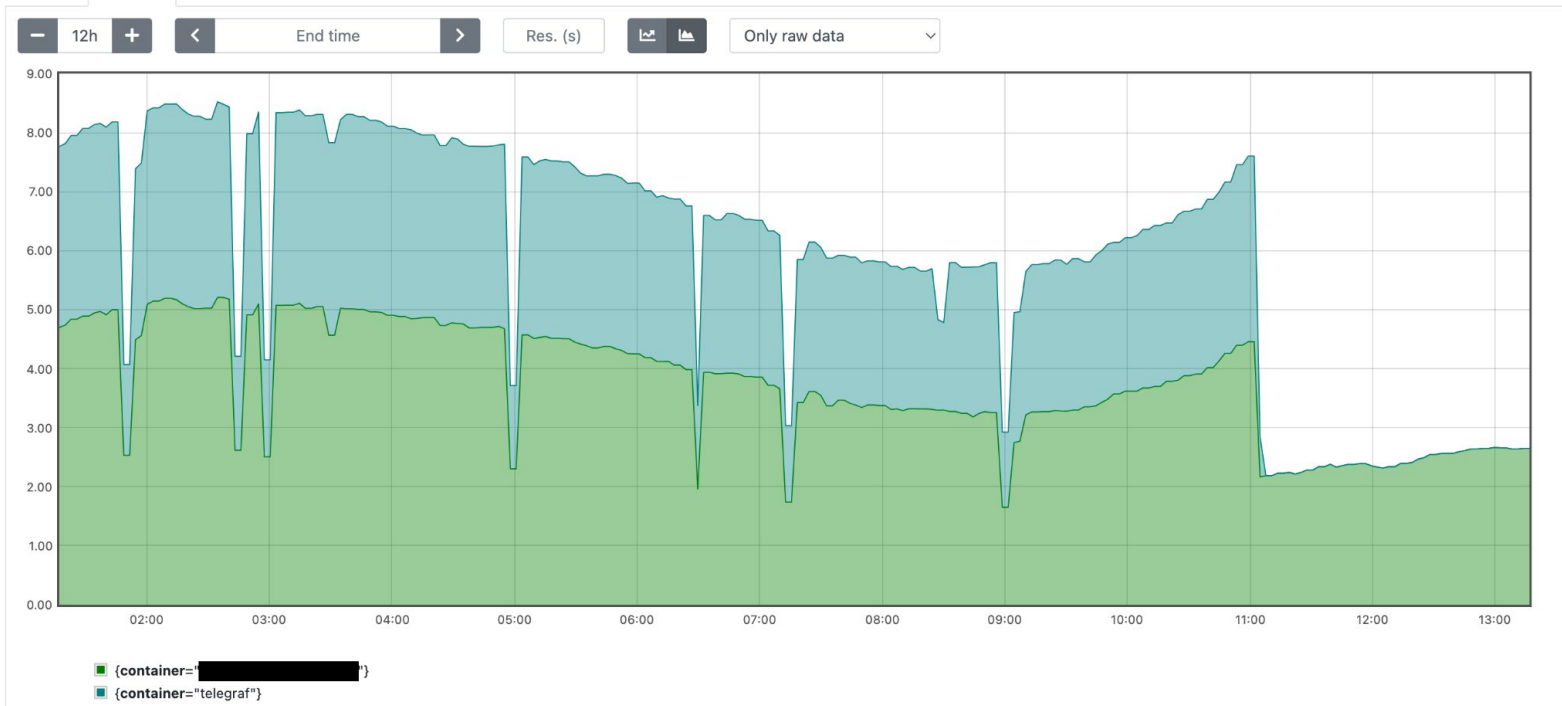
Execute Explain

Use Deduplication Use Partial Response Force Tracing Engine Prometheus

Analyze

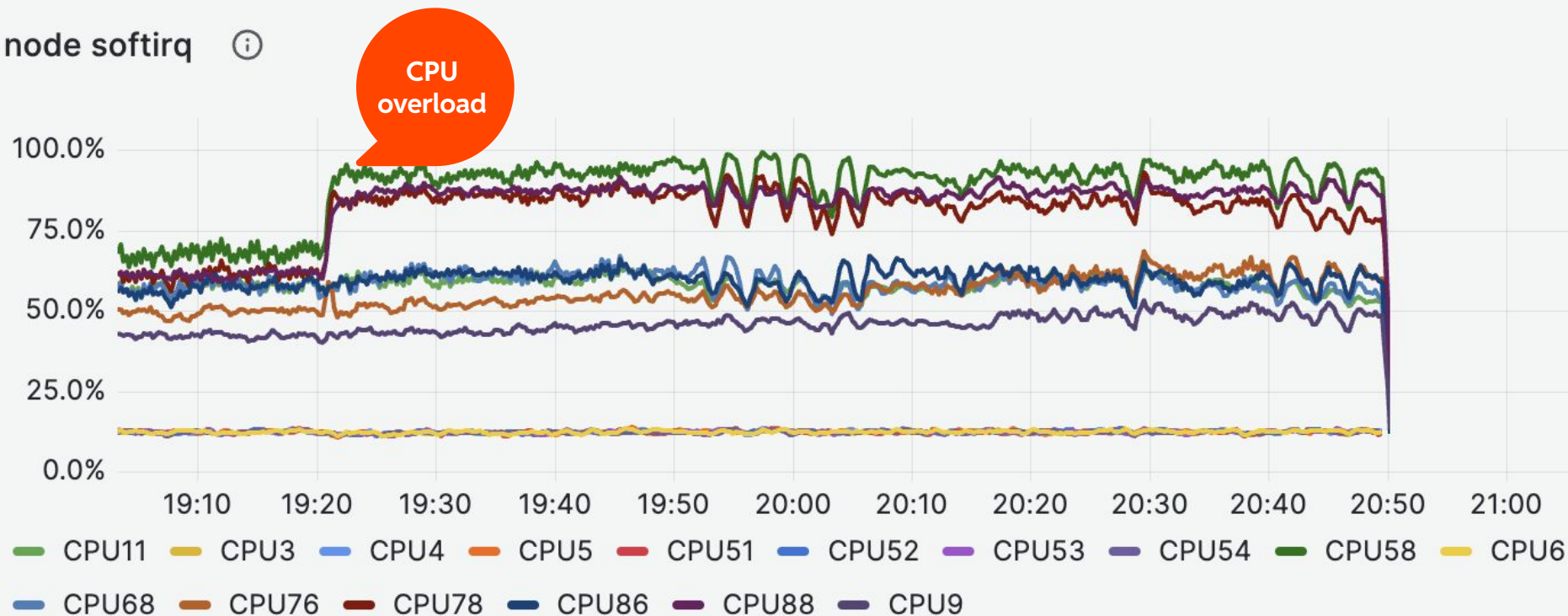
Load time: 611ms Resolution: 172s Result series: 2 Trace ID: 690d84ba9fe7dc0991c33f1c275447ca

Table Graph



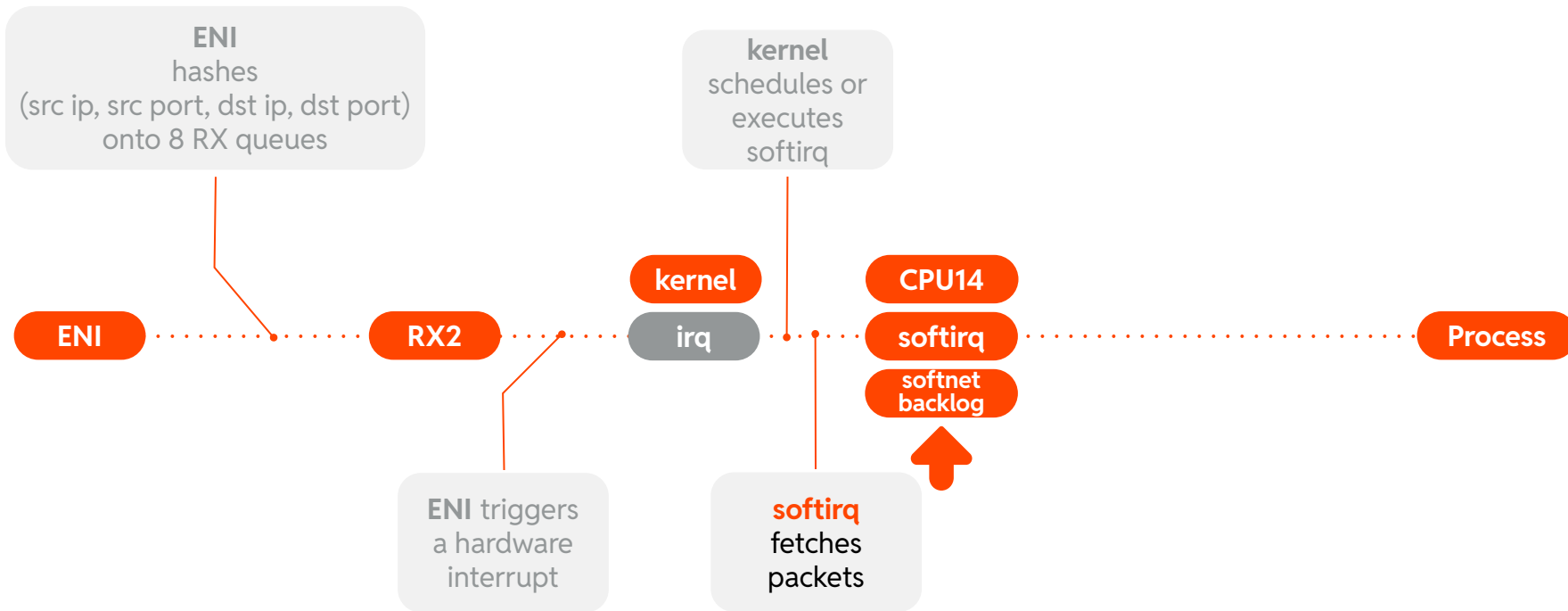
(05) What now? CPU overload...

node softirq ⓘ



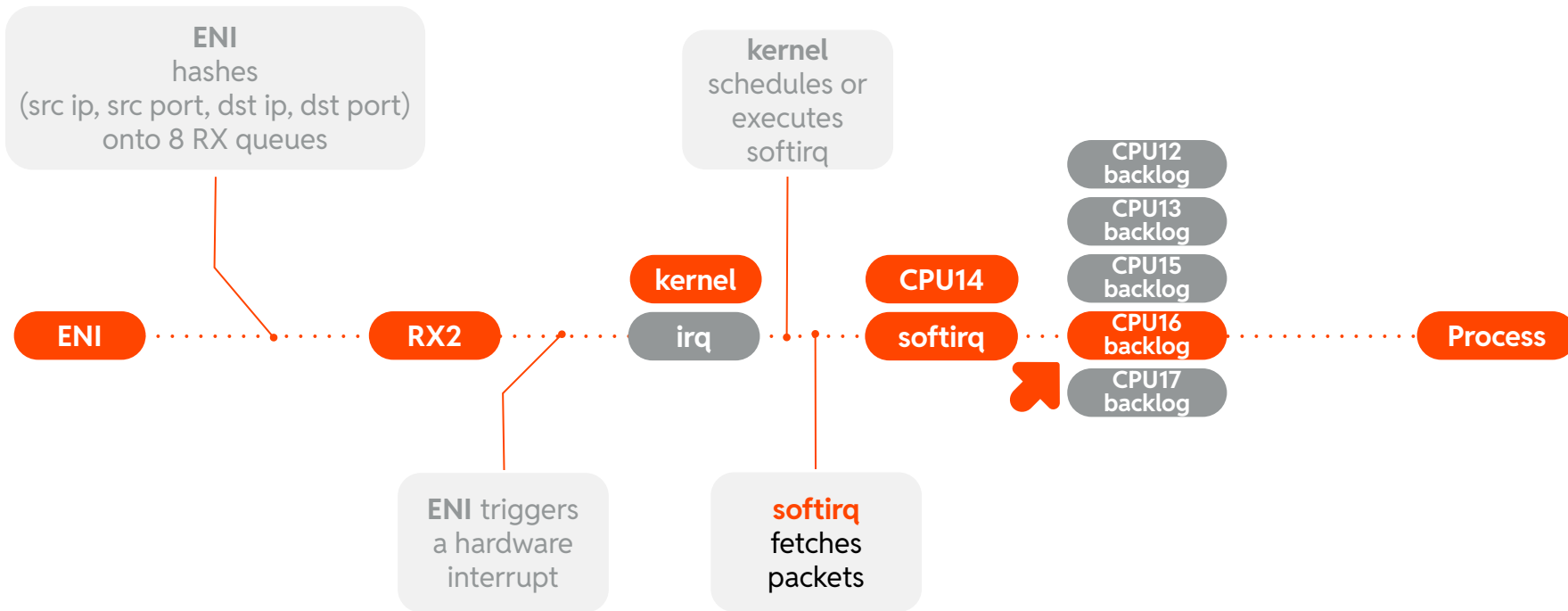
Receiving packets

AWS & Linux



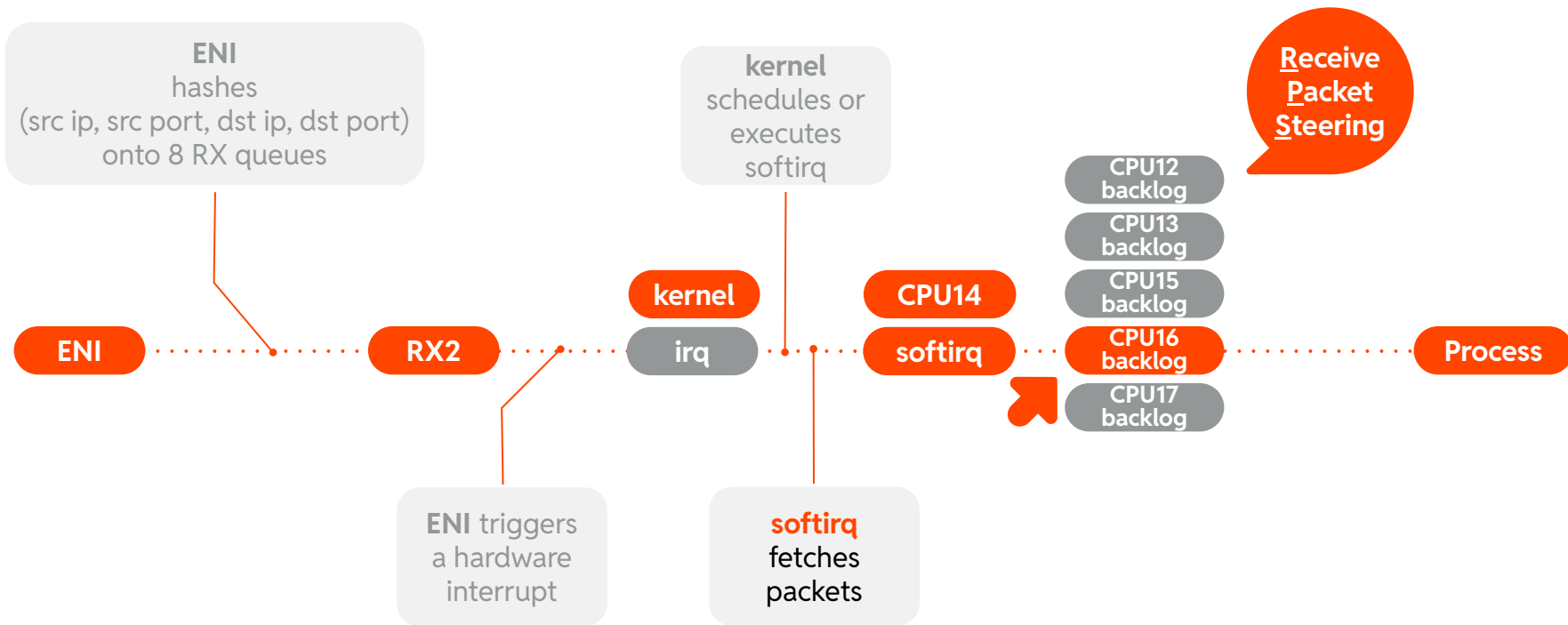
Receiving packets

AWS & Linux

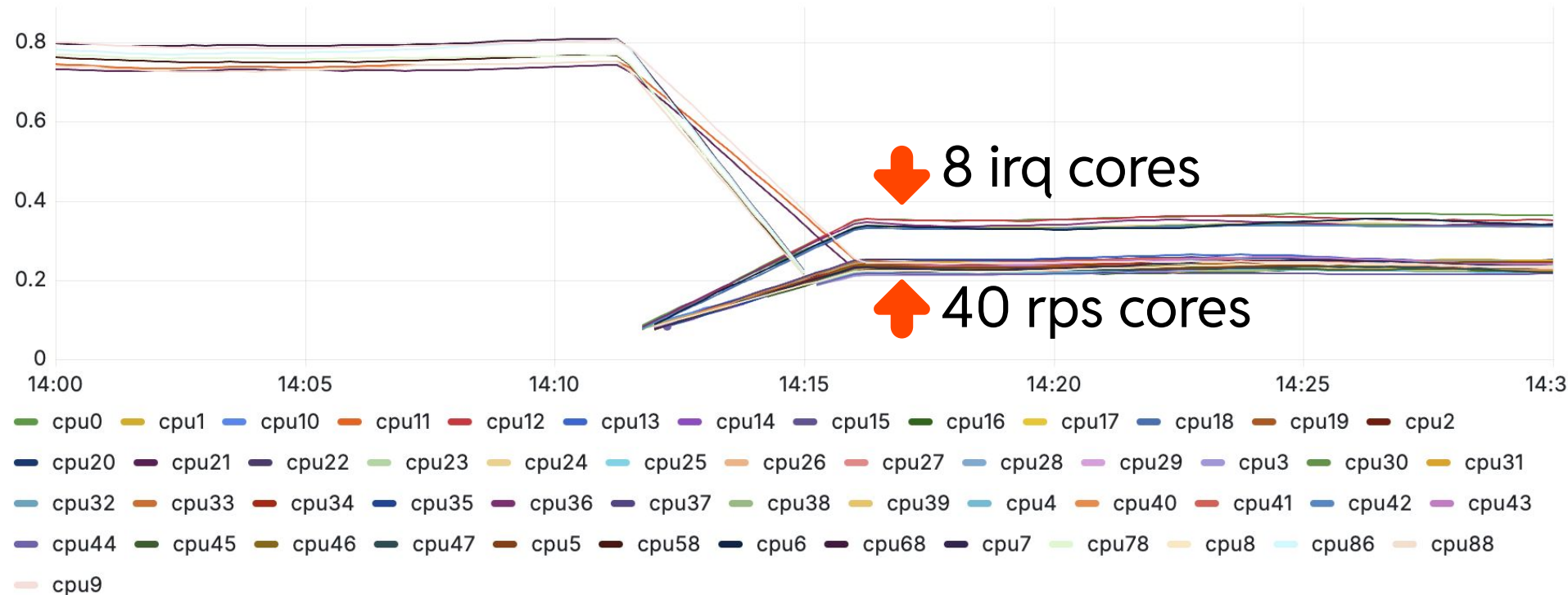


Receiving packets

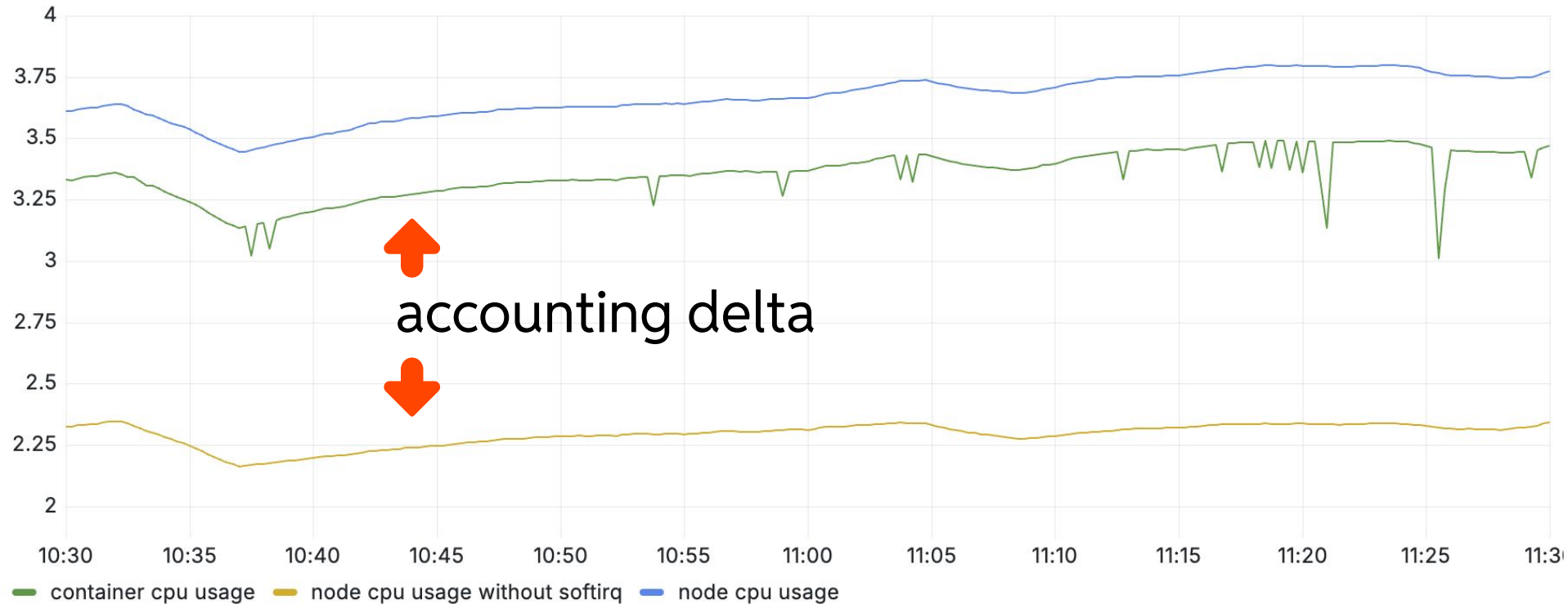
AWS & Linux



Receiving packets - The Fix

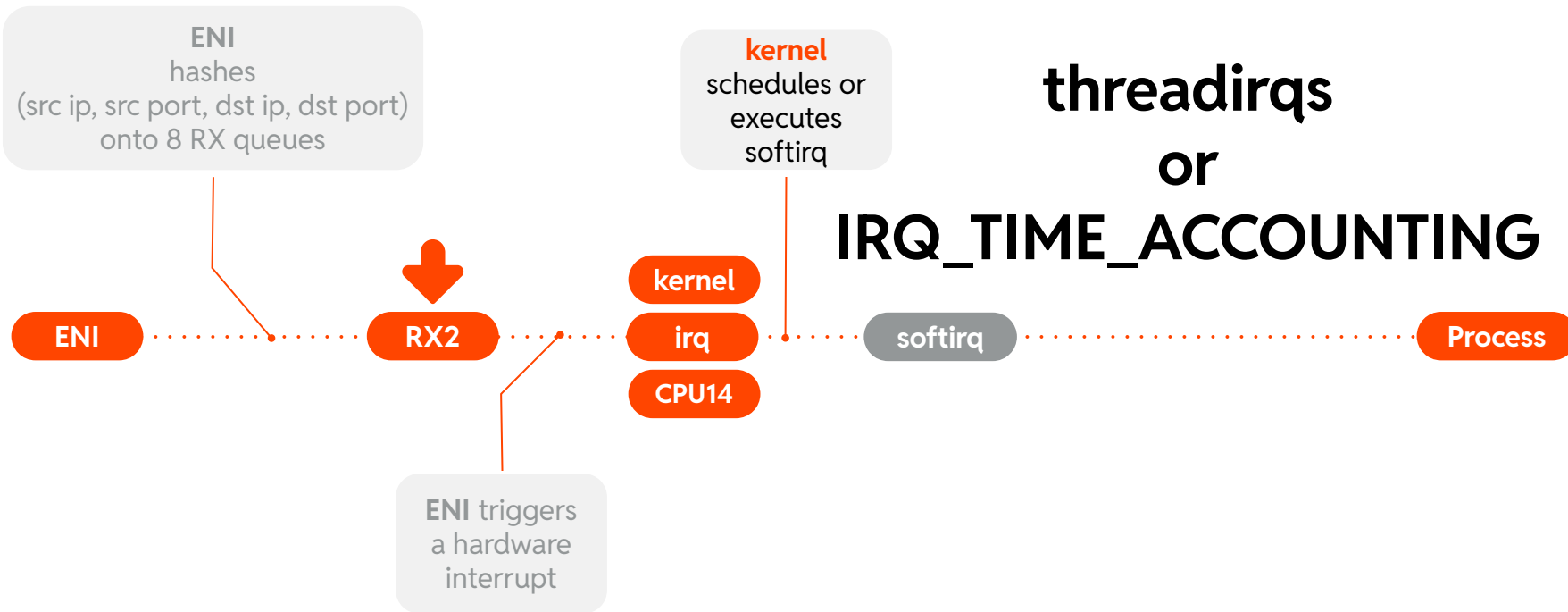


(05) What now? Accounting

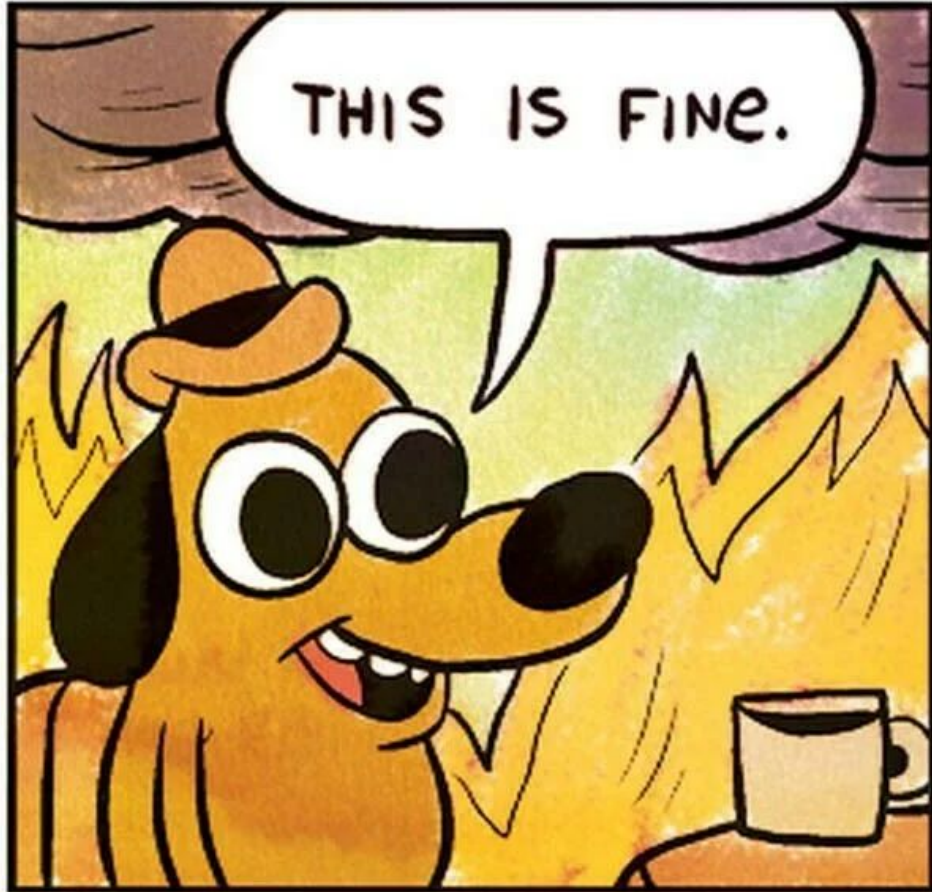


Receiving packets

AWS & Linux



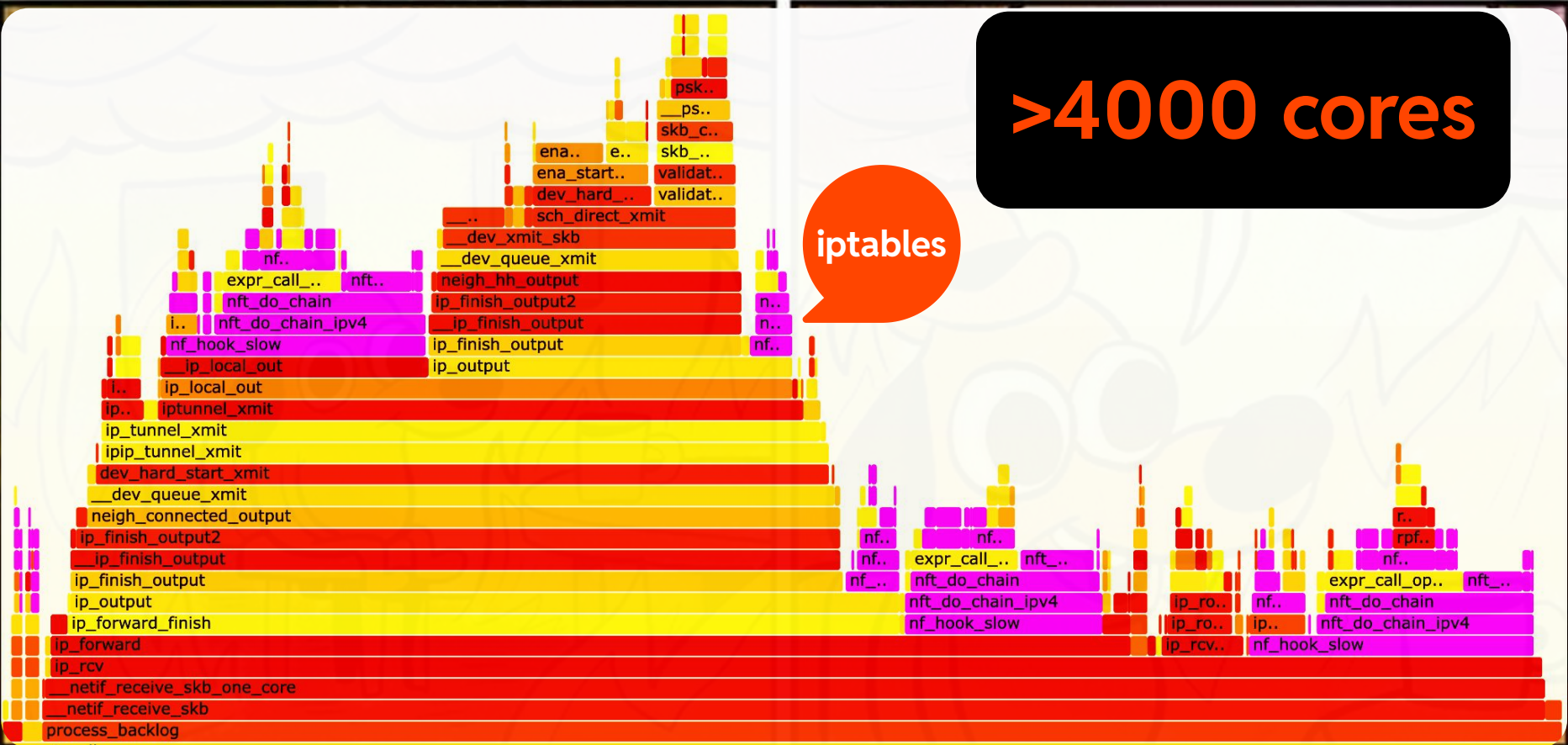
(05) What now? CPU cost (unsolved)



(05) What now? CPU cost (unsolved)

>4000 cores

iptables



Special thanks

To the teams
involved





Thank you