

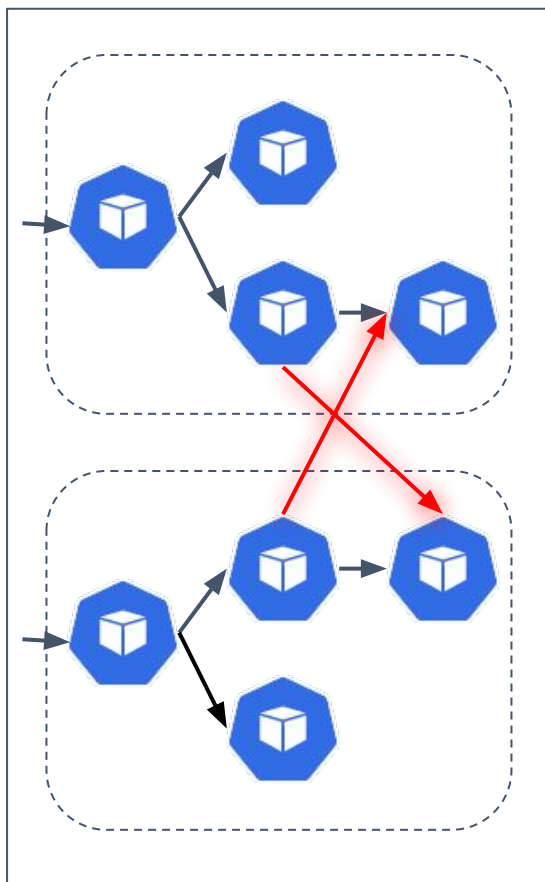
Monitoring Kubernetes with eBPF and Prometheus

KubeCon North America 2018

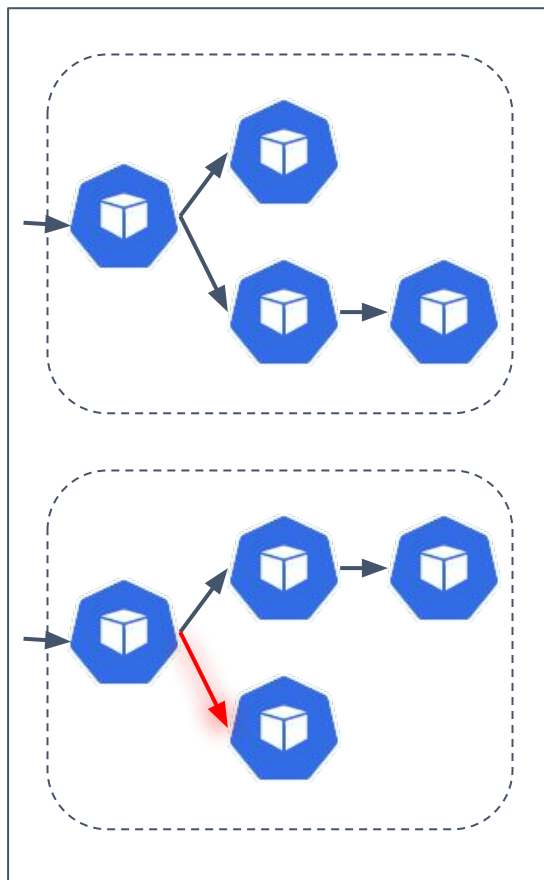
■ Agenda

- Flow monitoring: benefits
- Getting flow data
- Technology: eBPF
- Tour of our staging cluster
- Productizing: Challenges

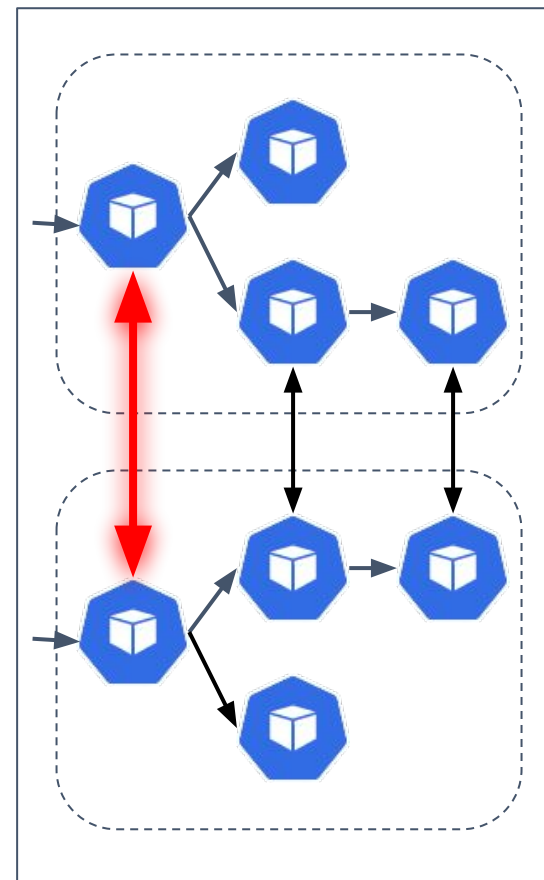
Flow Monitoring: benefits



Architecture, HA,
Env isolation

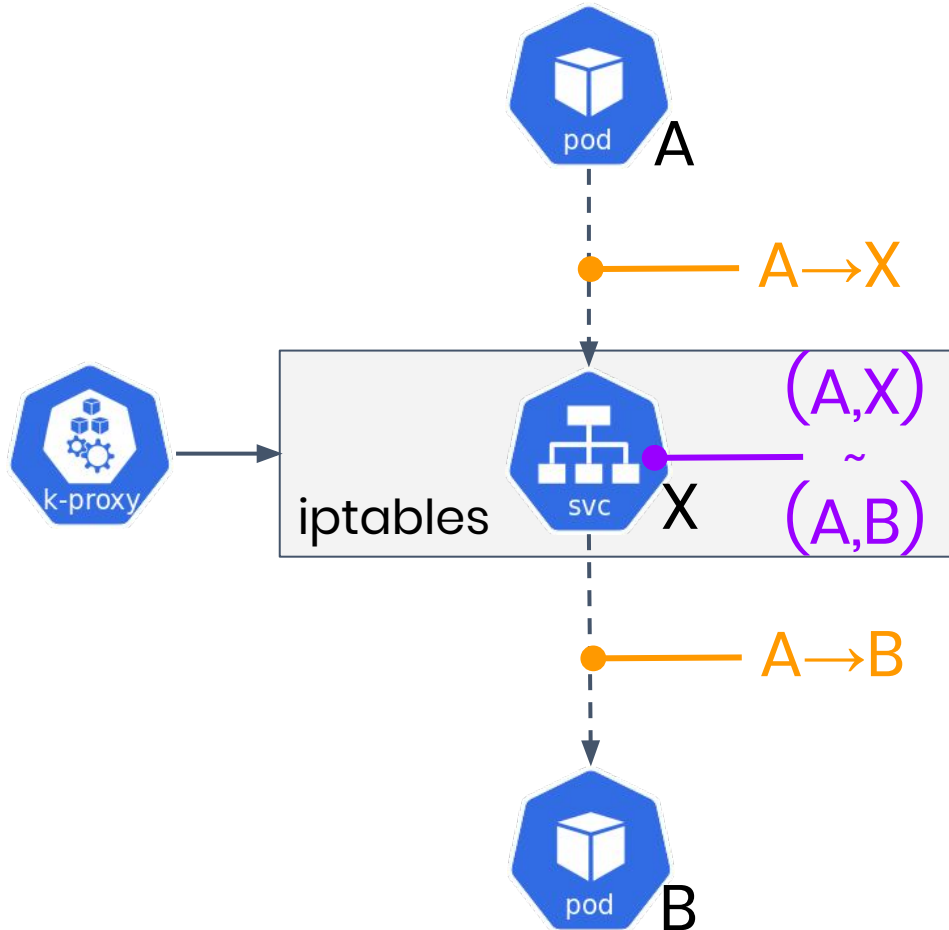


Health



Cost

Getting Flow Data



```
$ kubectl describe pod $POD
```

```
Name:          A
Namespace:     staging
...
Status:        Running
IP:            100.101.198.137
Controlled By: ReplicaSet/A
```

```
# PID=`docker inspect -f '{{.State.Pid}}' $CONTAINER` \
  nsenter -t $PID -n ss -ti
ESTAB  0  0  100.101.198.137:34940  100.65.61.118:8000
        cubic wscale:9,9 rto:204 rtt:0.003/0 mss:1448 cwnd:19
ssthresh:19 bytes_acked:2525112 segs_out:15664
segs_in:15578 data_segs_out:15662 send 73365.3Mbps
lastsnd:384 lastrcv:10265960 lastack:384 rcv_space:29200
minrtt:0.002
```

```
# conntrack -L
```

```
X tcp      6 86399 ESTABLISHED src=100.101.198.137
B dst=100.65.61.118 sport=34940 dport=8000
src=100.101.198.147 dst=100.101.198.137 sport=8000
dport=34940 [ASSURED] mark=0 use=1
```

■ Technology: eBPF

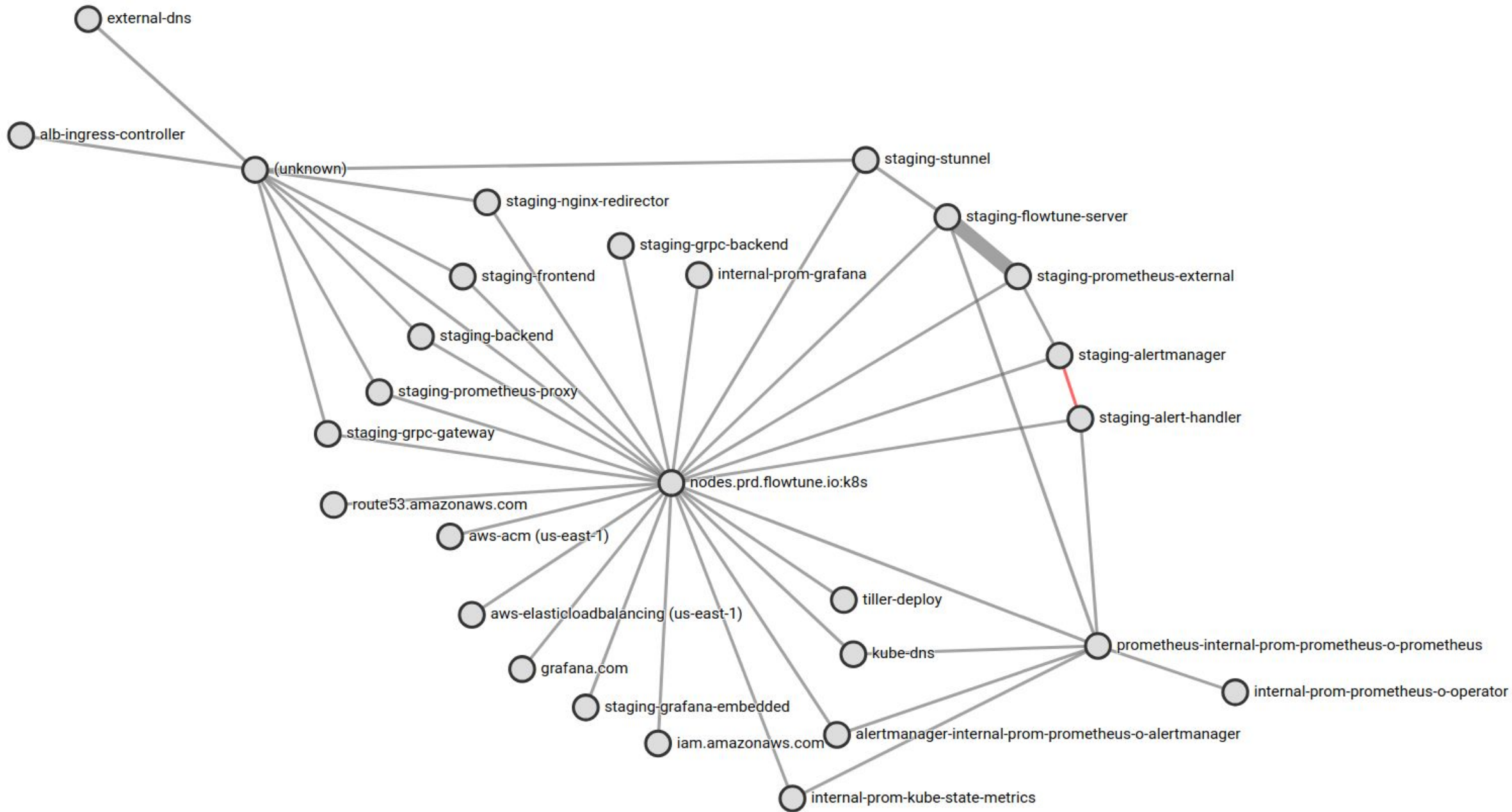
- Linux `bpf()` system call since 3.18
- Run code on kernel events
- Only changes, more data

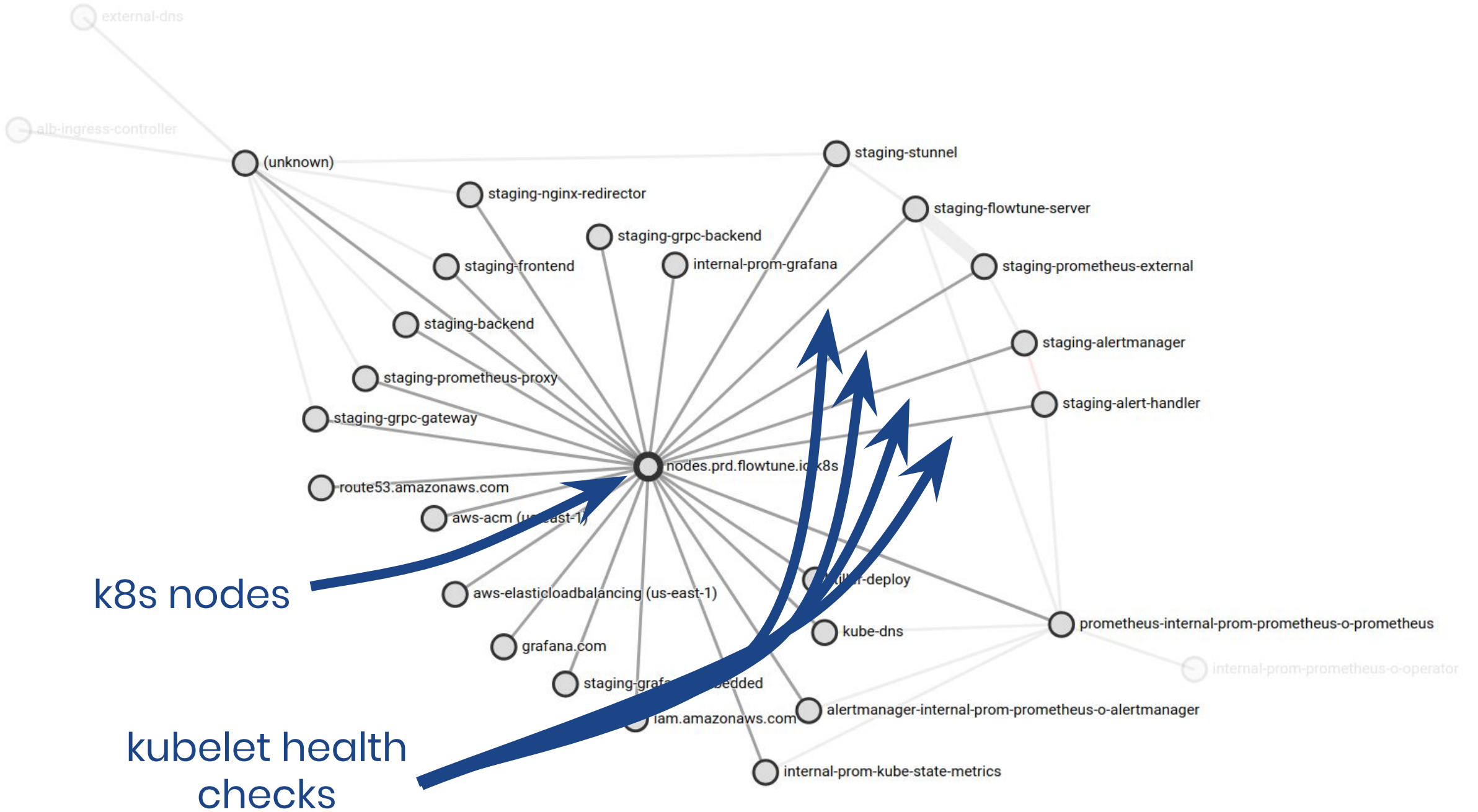
- Safe: In-kernel verifier, read-only
- Fast: JIT-compiled

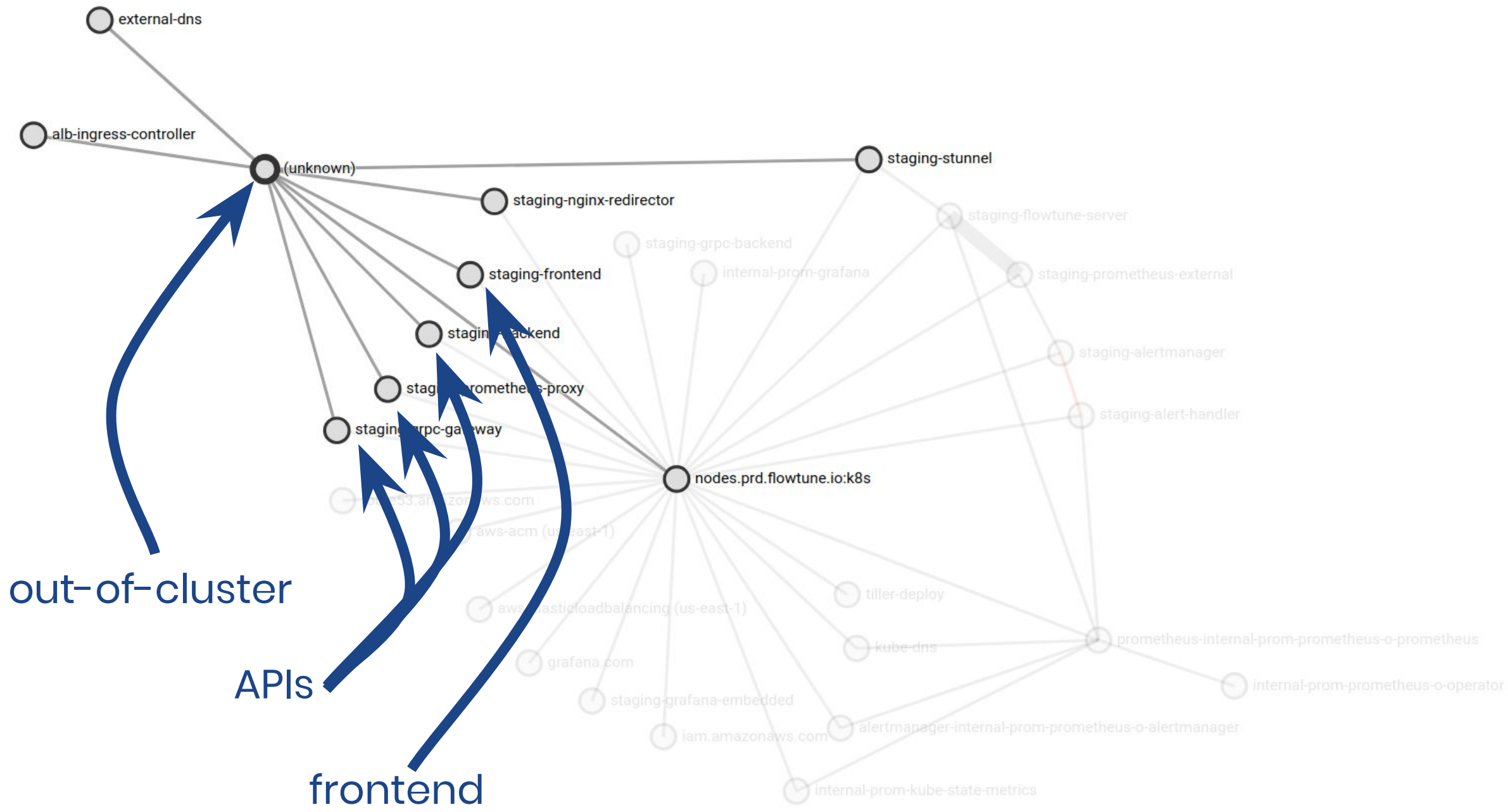


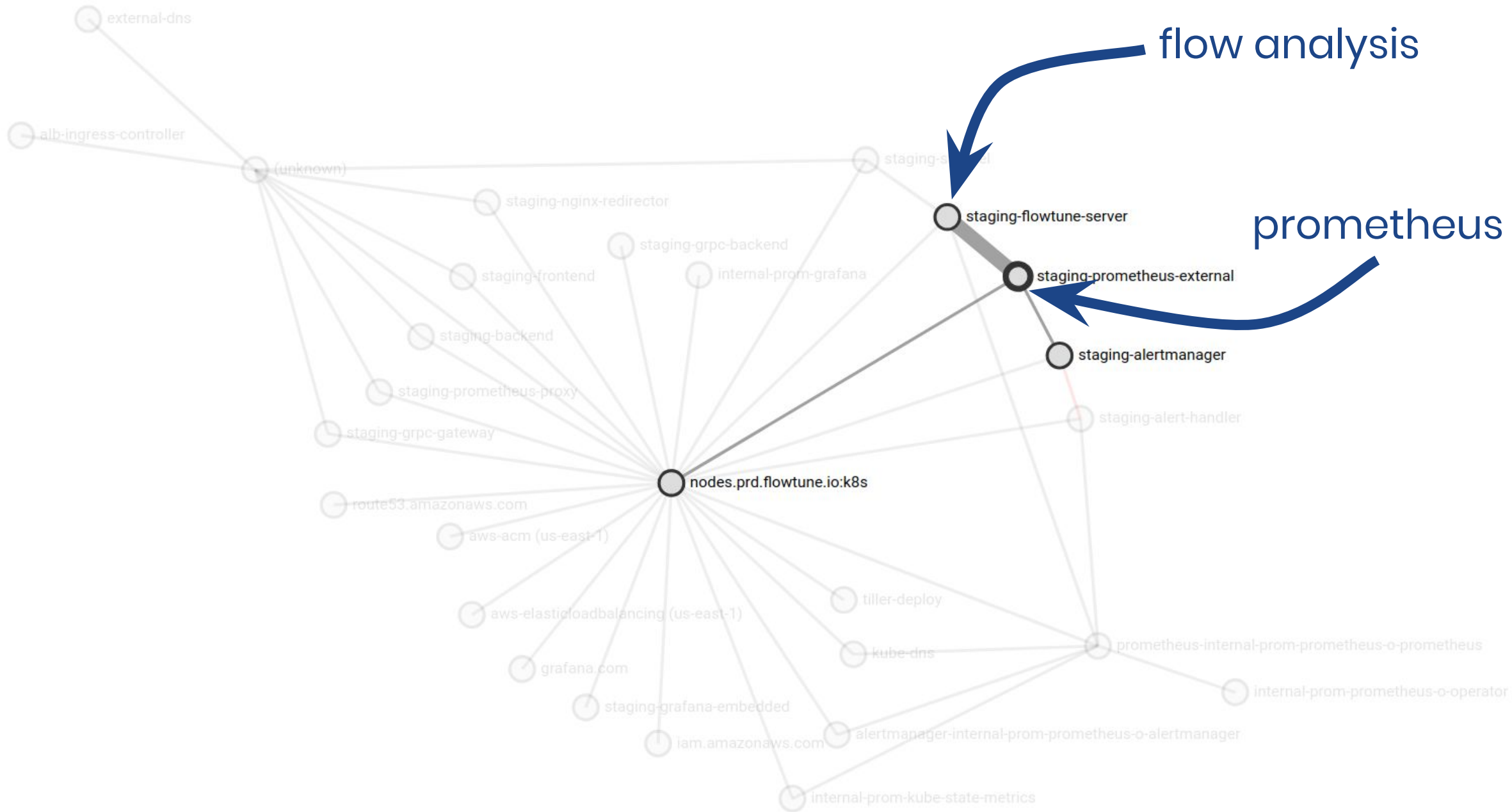
Unofficial BPF mascot by [Deirdré Straughan](#)

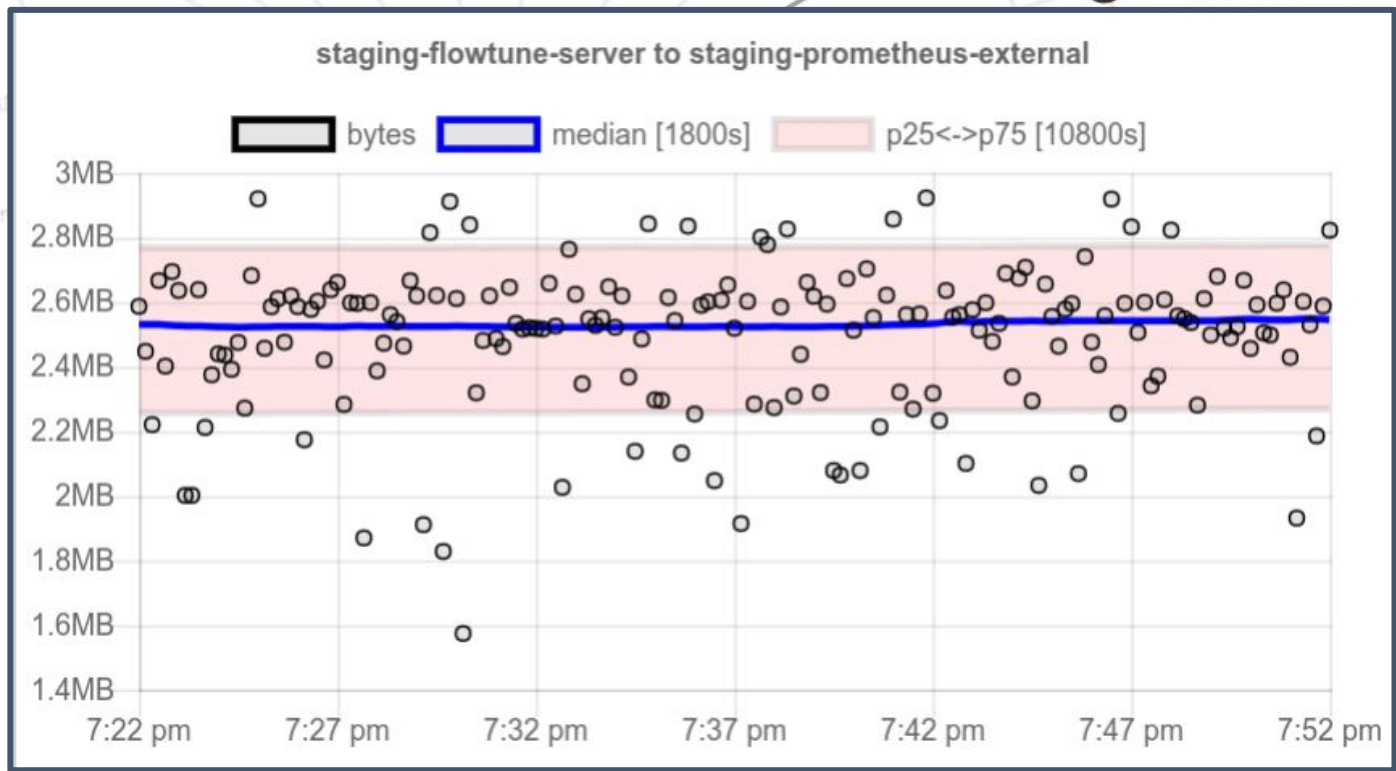
→ **100% coverage + no app changes + low overhead ftw!**











■ Productizing: Challenges

- **CPU overhead:** profile + iterate → 0.1% CPU
- **Network overhead:** encode efficiently and compress
- **Security:** TLS, OAuth everywhere
- **Real-time:** stream, don't batch → 2 second latency
- Pre-aggregate to manage **cardinality**
- Workload **baselining** for **automatic alerting**

Jonathan Perry <jperry@flowmill.com>
www.flowmill.com

