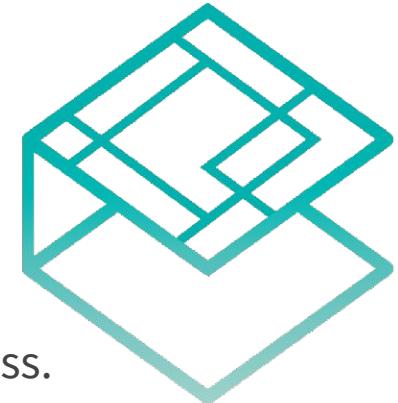




Network Services on Kubernetes on-premise

Hans Duedal
Network Services Engineer @ NPF
CTO @ Visma e-economic

Networking in Kubernetes

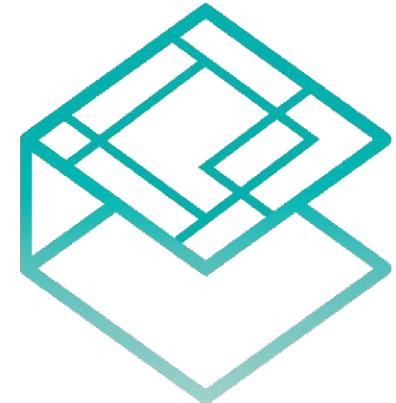


Kubernetes assumes that pods can communicate with other pods, regardless of which host they land on. Every pod gets its own IP address.

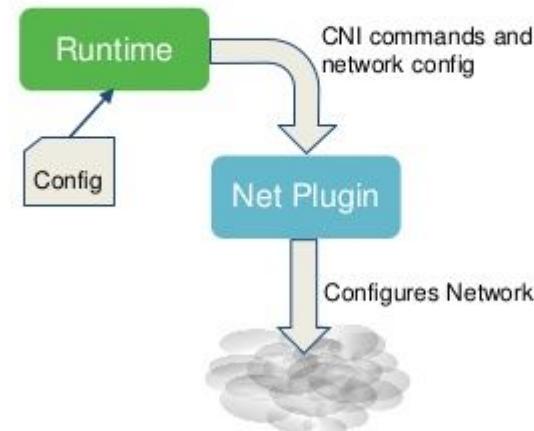
Fundamental requirements

1. all containers can communicate with all other containers without NAT
2. all nodes can communicate with all containers (and vice-versa) without NAT
3. the IP that a container sees itself as is the same IP that others see it as

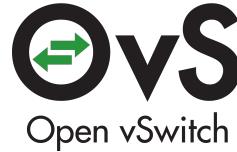
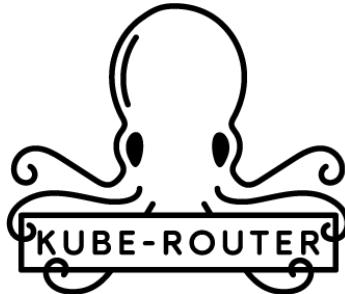
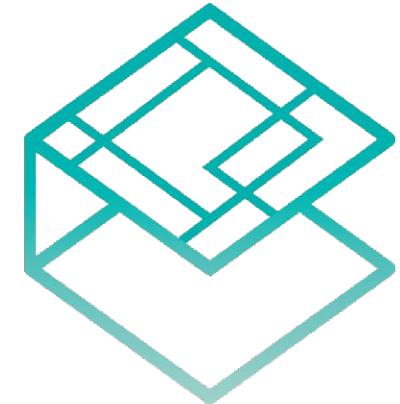
Container Network Interface



- Originated at CoreOS as part of rkt
- Cloud Native project
- Many container runtimes, ie. kubernetes, rkt, mesos
- If your network solution adheres to the CNI Spec, the runtimes can use it
- Sometimes the choice of CNI plugins are (almost) made for you, ie. on public clouds
- For an on-premise project, choosing the right CNI solution is critical to success



Kubernetes Network Solutions

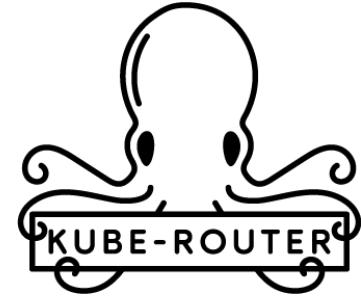


Project Calico



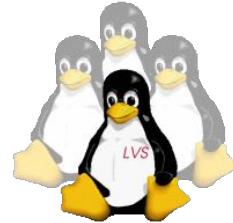
- IP-in-IP with option for BGP Based
- Independent of Kubernetes (ships with its own etcd cluster)
- Backbone of Github.com's metal cloud
<https://thenewstack.io/github-goes-kubernetes-tells/>
- Policy based access control
- Fully featured for Datacenter Level Kubernetes Networking
- Doesn't make much sense to use without BGP

Kube-router



Three components

1. IPVS/LVS (Part of Linux kernel since 1999)
2. GoBGP
3. Kuberouter itself



- BGP Based
<https://cloudnativelabs.github.io/post/2017-05-22-kube-pod-networking/>
- Heavy lifting by IPVS/LVS (in-kernel load balancer) and the routing tables
- Kubernetes specific
- Small and lightweight, with excellent community on slack
- Network Policies (but not as advanced as Calico)

- Full userspace software based solution, no hardware or kernel involved
- IP-in-IP based
- Often the goto/default choice for on-premise deployments
- Only network, no policy
- Networking: “Developer Edition”
- CoreOS project, but plays extremely well with Kubernetes
- The reference plugin for CNI

How Rancher does it



- Magically network nodes in AWS with GCE or on-premise
- It's all an illusion, it sets up a full mesh of IPSec (VPN) tunnels between each node
- High overhead of IPSec stack and encryption
- Easy to use, and most of the time “just works”

That's all fine, but you said something
about e-sports?

The logo for the National PC Federation (NPF) is located in the bottom right corner. It features a stylized blue 'G' or swirl icon followed by the acronym 'NPF' in a bold, white, sans-serif font.

NPF in Numbers

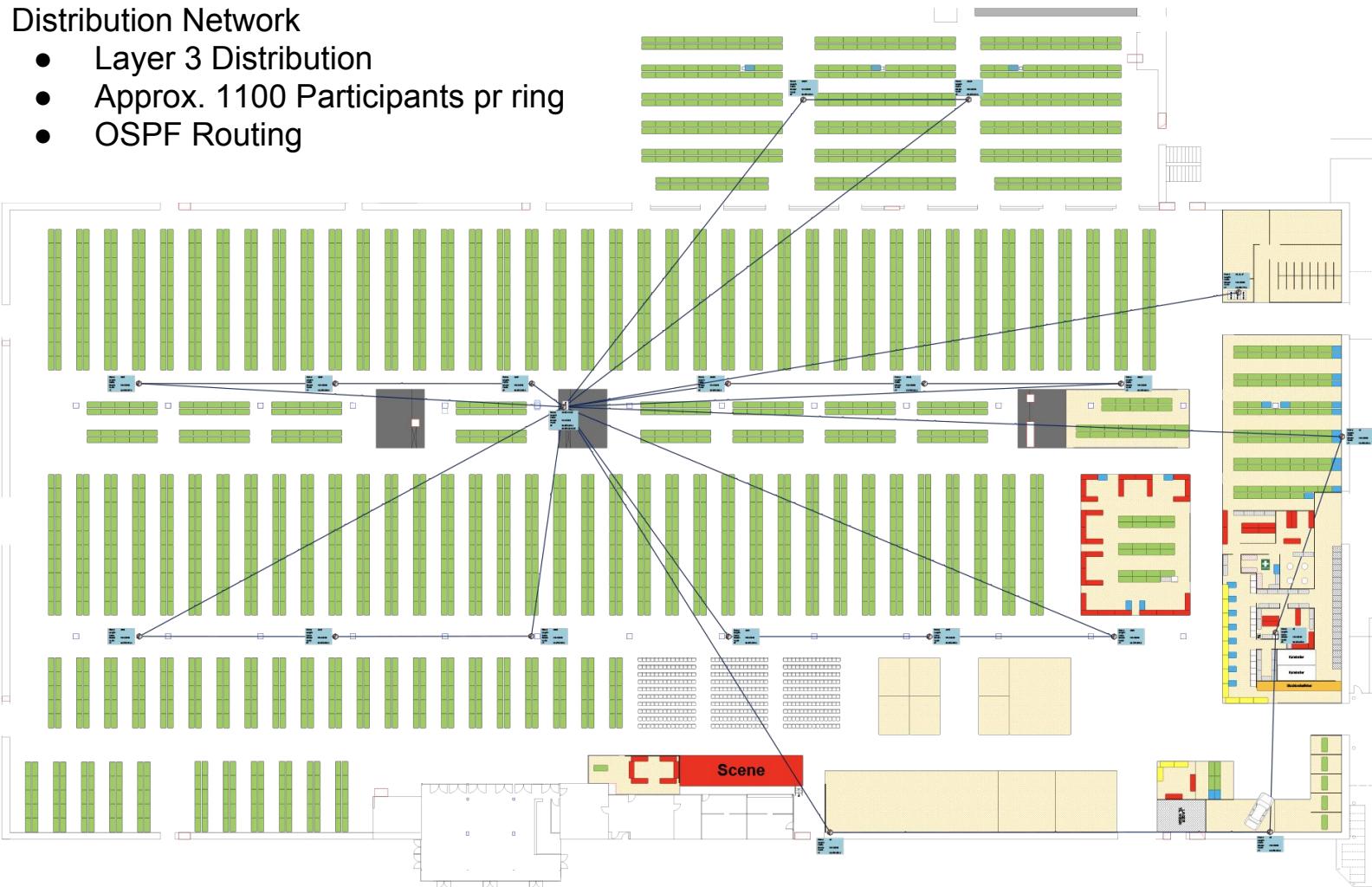
- 3600 Amps of mains, 3 phase power
- Almost 14 kilometers of network cables (2.4 km fiber, 11.5 km copper)
- 400+ volunteers
- 300+ network switches
- Built up in five days, 47 hour event, and then torn down in two days
- Almost 6000 gamers with their computers
- Internet consumption the size of Aalborg, so internet wise we become the fourth largest “city” in Denmark for a weekend



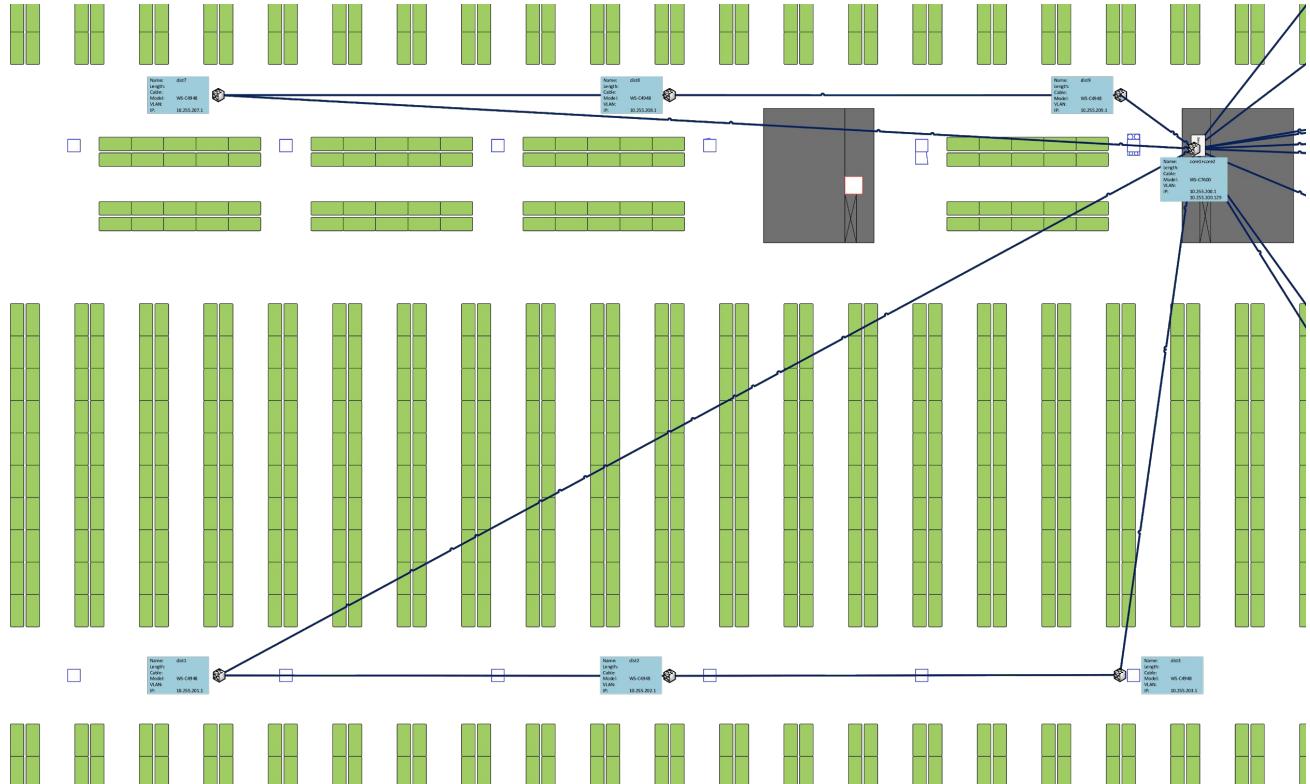
Network at NPF

Distribution Network

- Layer 3 Distribution
- Approx. 1100 Participants pr ring
- OSPF Routing

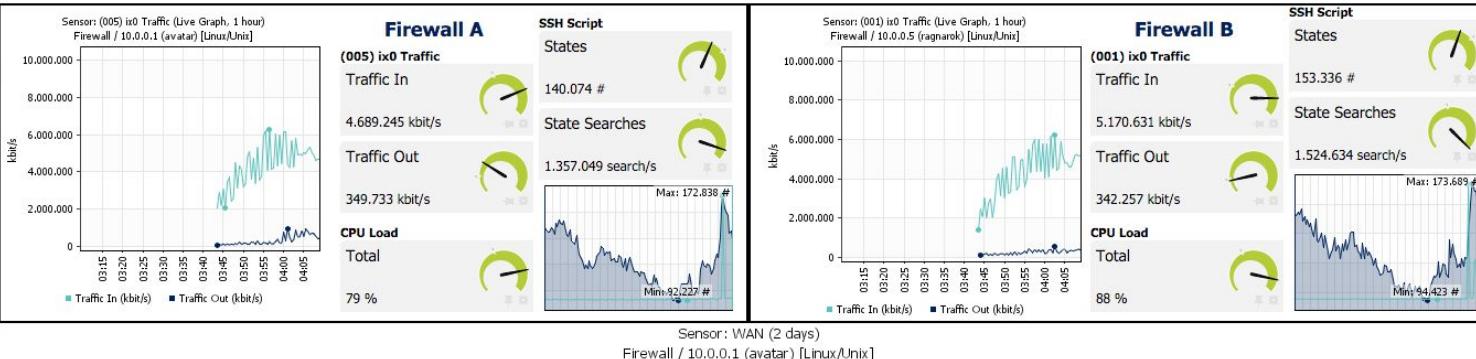
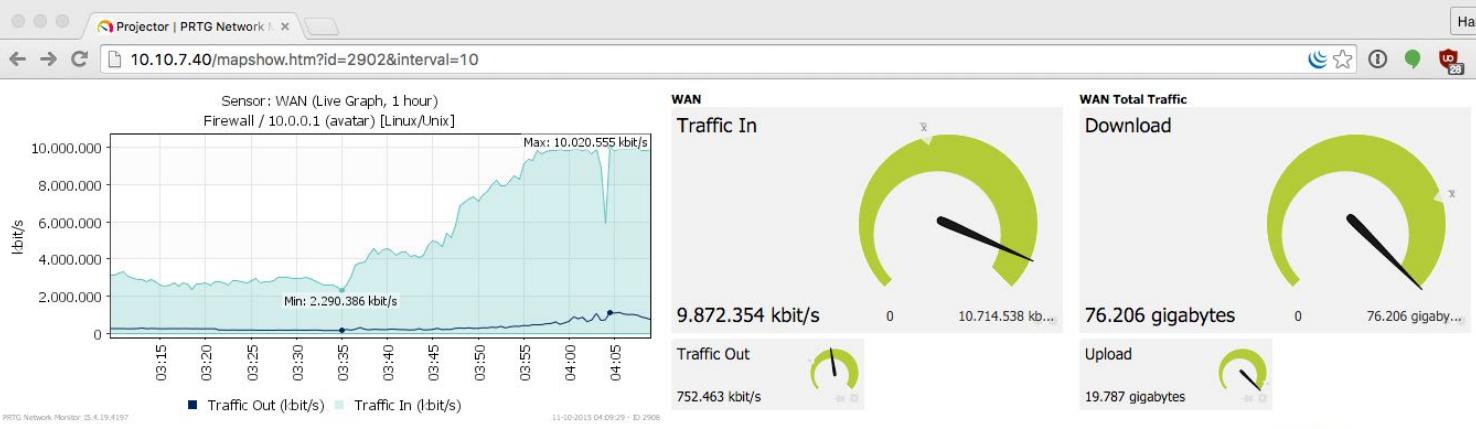


Rings



Network “Services”?

- DHCP
- DNS
- Prometheus Monitoring
- Grafana
- SMNP Polling (Blackbox exporter)
- SNMP Traps (ELK Cluster)





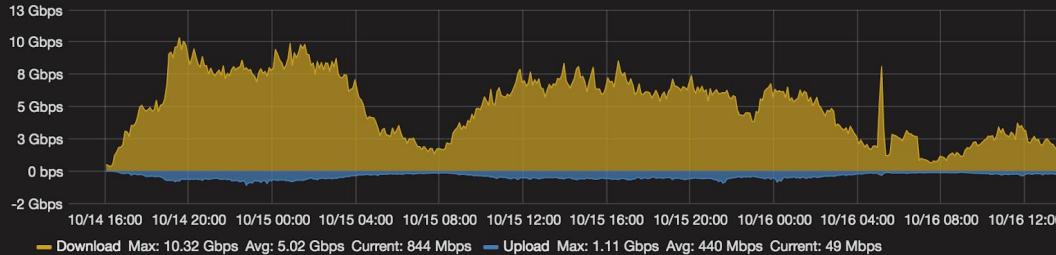
NOC



Zoom Out Last 3 hours Refresh every 10s

Firewall WAN Throughput

Last 48 hours



Balance

Download

94.78 TiB

Upload

8.38 TiB

DHCP Leases

5713

Connected devices

2575

Download hastighed

844 Mbps

Avatar LoadAvg

0.39

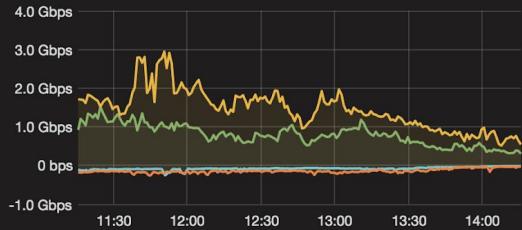
Ragnarok LoadAvg

0.65

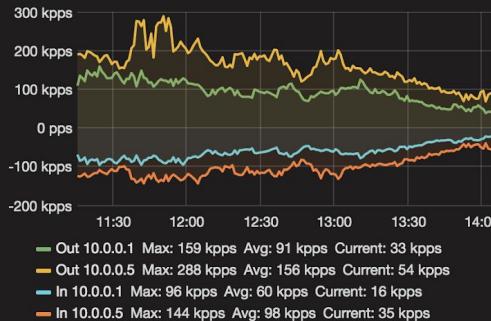
Ping offline

5

Firewall Balance



Packets



Varmeste switches

Last 12 hours



Kube-router in NPF

Key learnings

- Think in terms of ToR (top-of-rack) design
- Often it's not the new fancy kubernetes network that fails
- Kube-router toolbox is essential
 - You just kubectl exec -it into the kube-router box

Why kube-router

- We tried calico and got overwhelmed
- Very low overhead, most of it is done in-kernel
- Minimal code base, easily understood
- We run full kube-router, so no kubenet or kubeProxy on the nodes

DHCP Server Troubles

Modern DHCP server Kea failed, so we had to go with isc-dhcpd
ISC DHCPd needs to know the IP of all availability peers on startup

- ISC DHCPd can take a hostname rather than the IP
- Headless Services expose the pod ip to kube-dns
- Kubernetes DNS resolution to the rescue!
- Stateful Sets could be a possible non-hacky solution, but not sure ISC DHCPd will play nice with stateful sets.

So how did it go?



thank you!

e-economic.com/joinus

npf.dk/om-npf-crew/

