



VPCs in Kubernetes

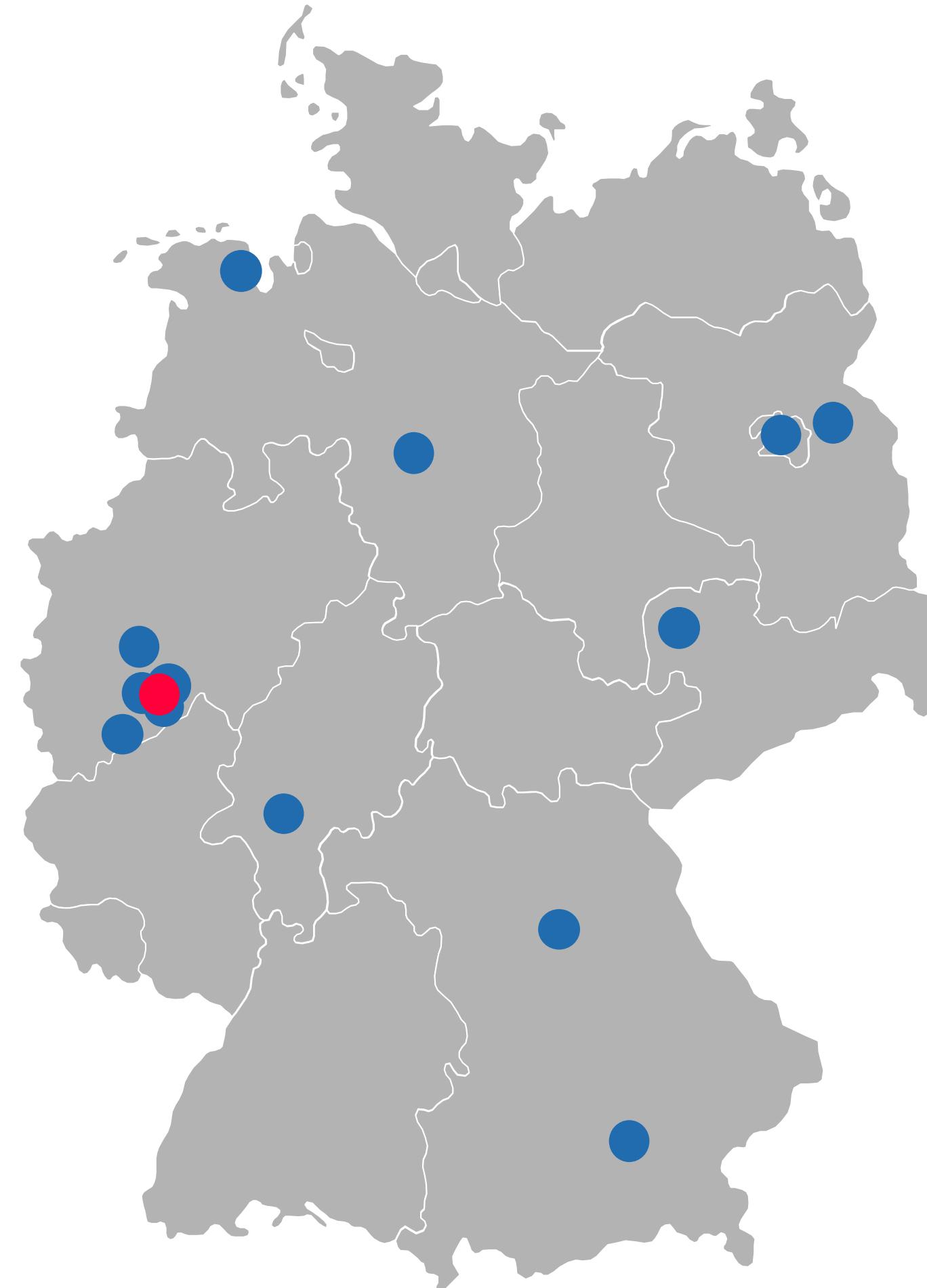
Niklas Voss



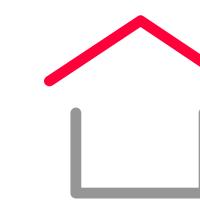
BWI is the **Bundeswehr's primary partner** of digitalization in peace, crisis, and war.



BUNDESWEHR



Regional presence with
14 major sites,
>130 locations throughout
Germany

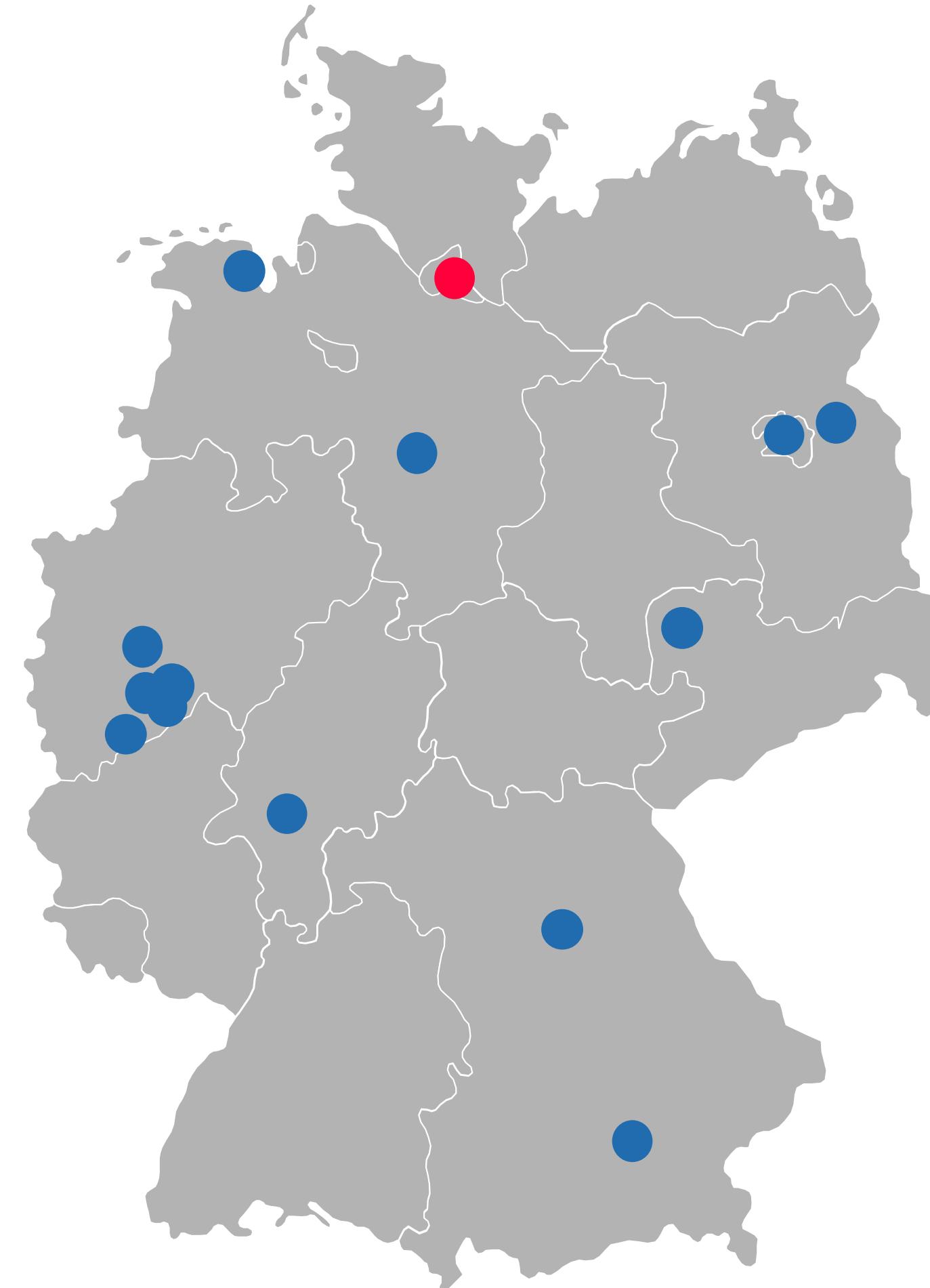


Headquarters
in Meckenheim

100% in-house company
of the federal government



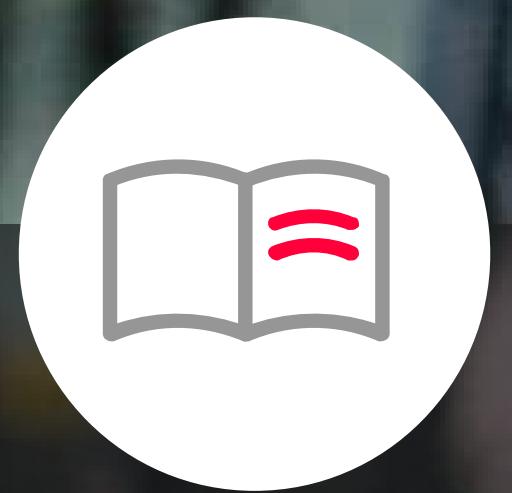
Niklas Voss
Tech Lead Private Cloud
Bundeswehr



o **Previous employers** include
Google, Finleap, ...

 **Living
in Hamburg**

Agenda



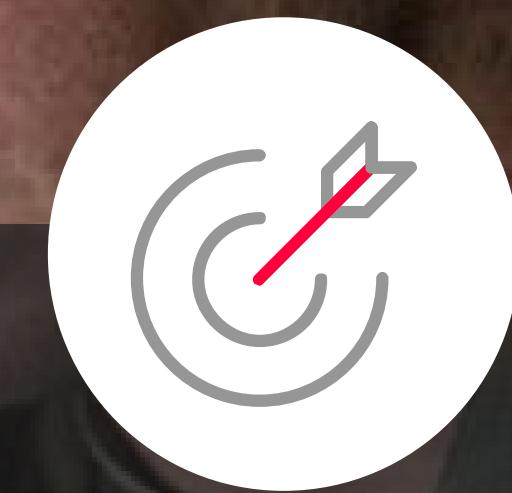
What is a VPC?



**How does the
Container Network
Interface work?**



**Let's create a CNI from
scratch!**



**What does the future
hold?**



What is a VPC?

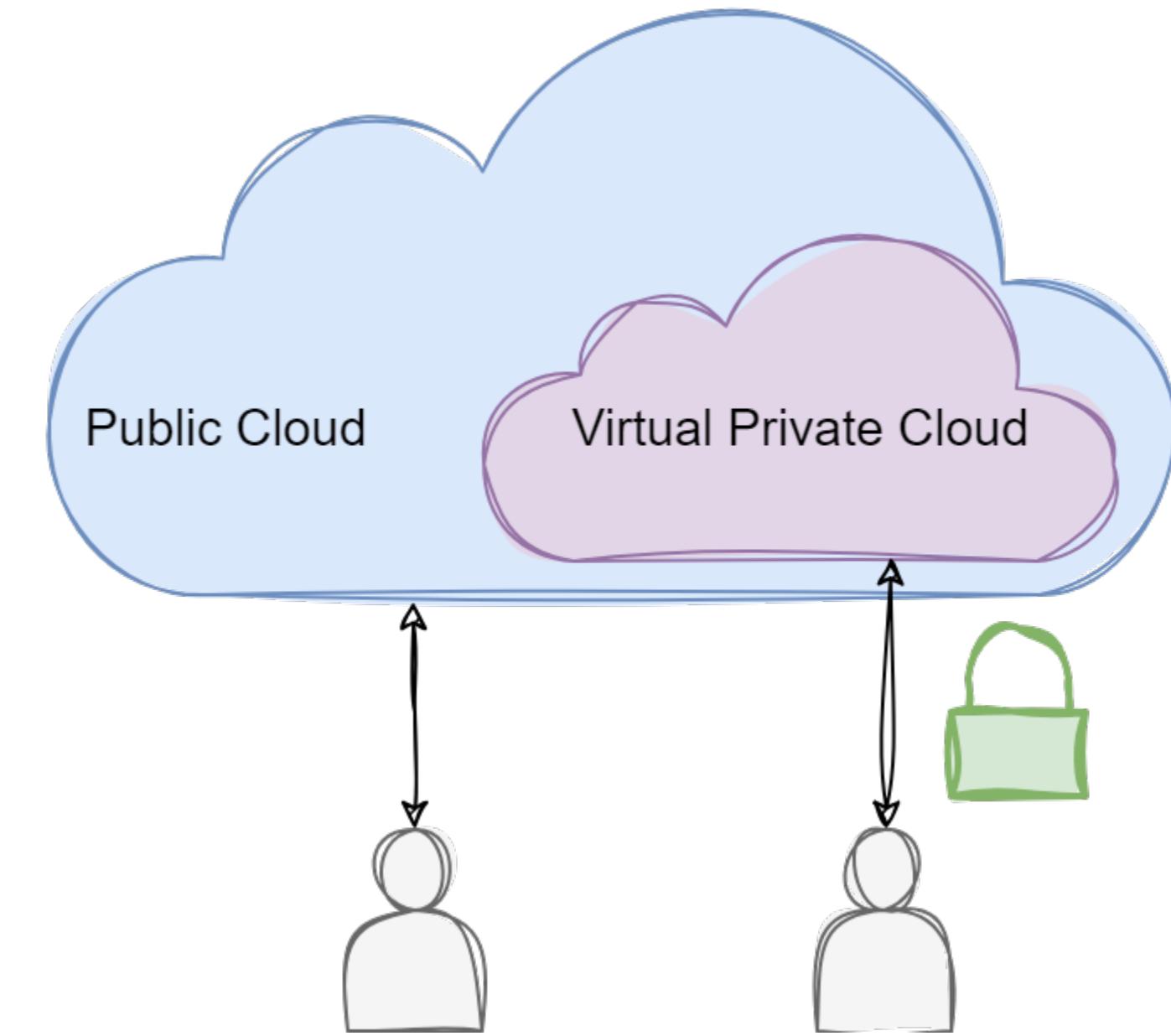


What is a VPC

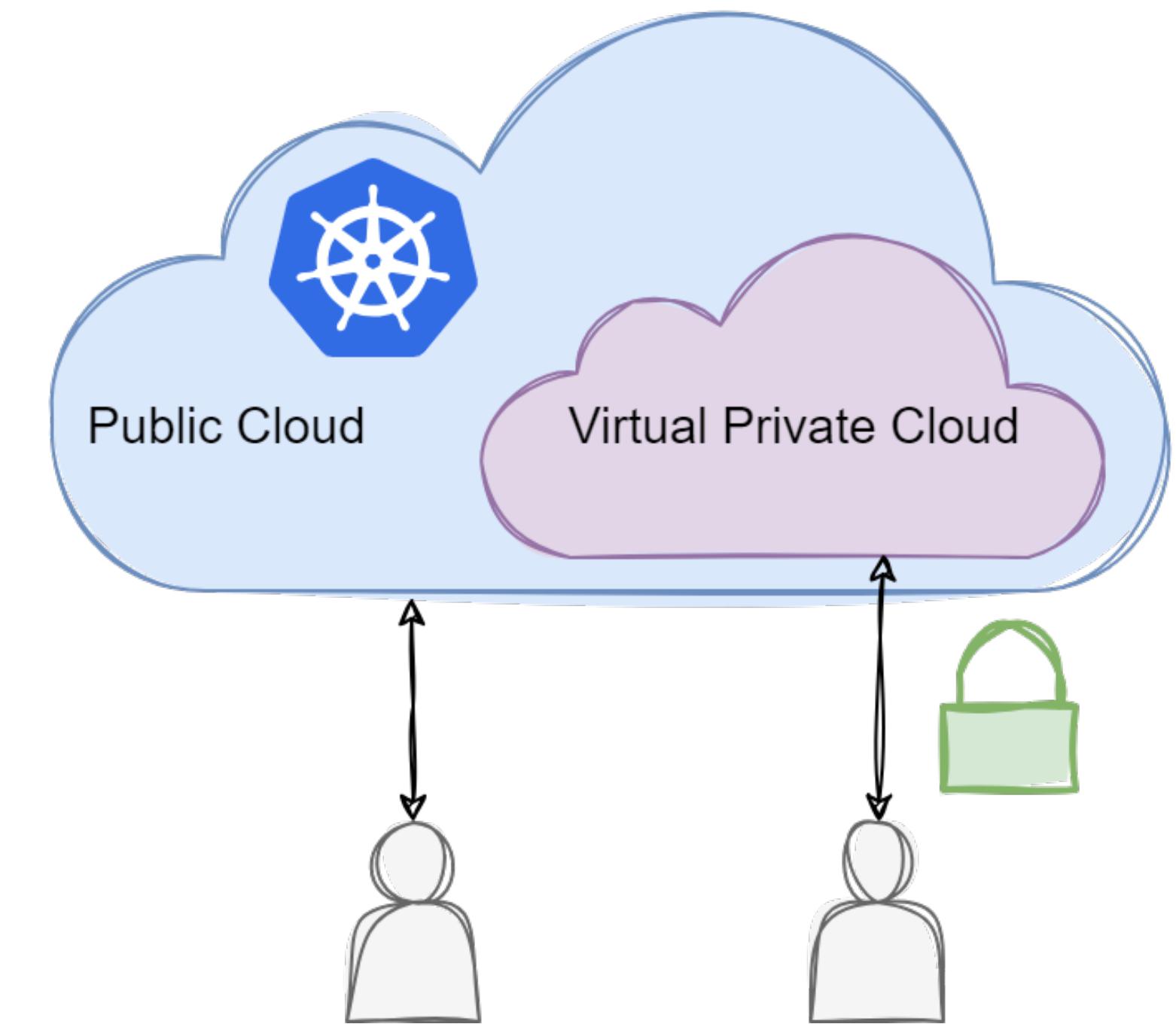
According to [Cloudflare](#):

“A virtual private cloud (VPC) is a **secure, isolated private cloud hosted within a public cloud.**”

- Strong tenancy isolation
- For Cloud Providers like AWS and Google Cloud **it is a network construct!**

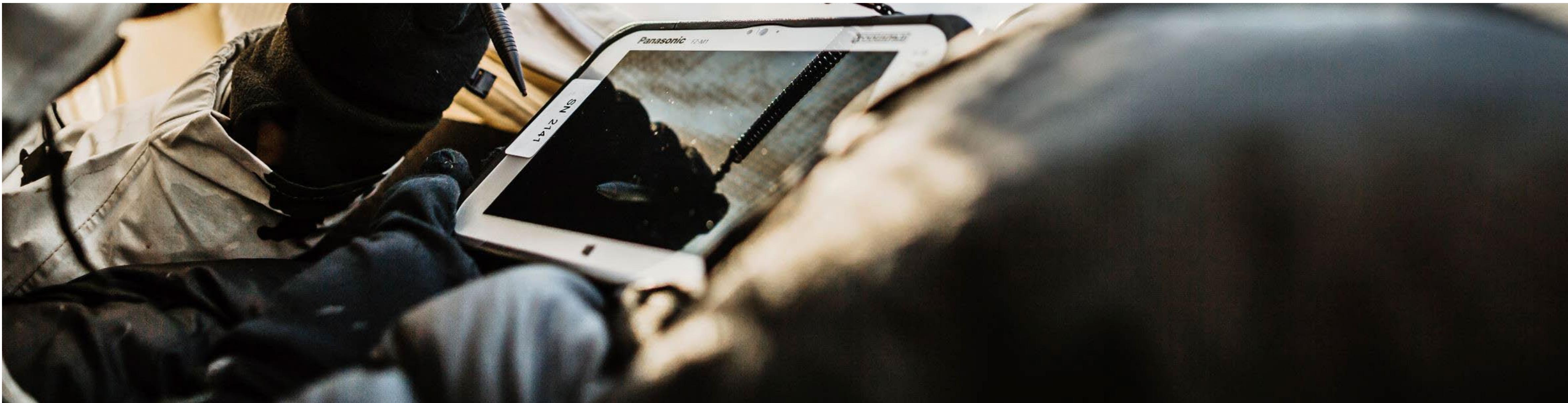


We want to build a Cloud on top of Kubernetes



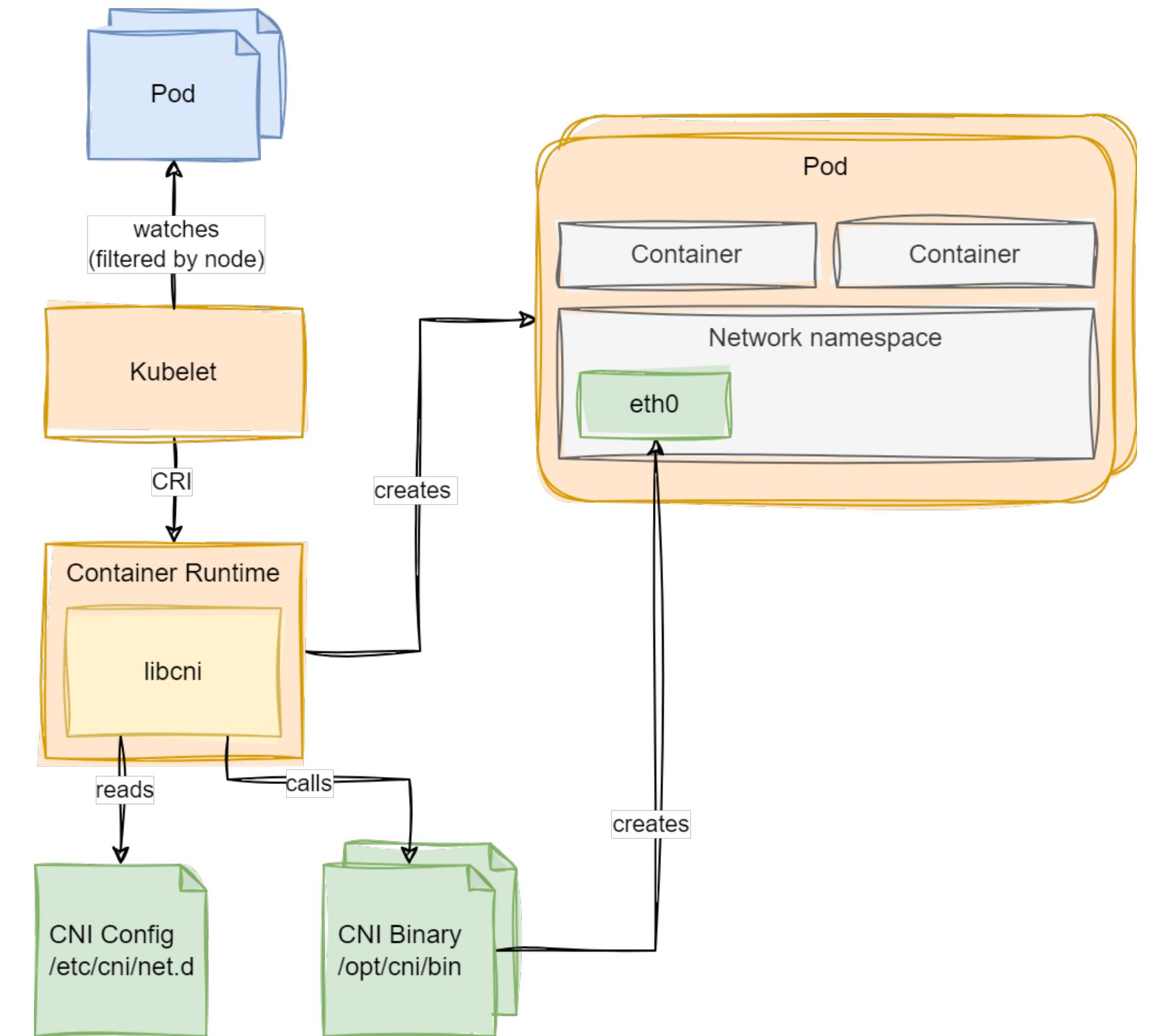


How does the Container Network Interface work?



From Kubelet to CNI

1. The **kubelet** watches pods for its node
2. Pods are **created via CRI**
3. **Container Runtime** sets up the containers and namespaces and **interacts with libcni**
4. **libcni** reads the **CNI config** and calls **CNI binary**
5. Container is started



CNI Config

CNI config is **JSON!**

WHAT!? NO YAML!?



CNI is **independent** of Kubernetes and
agnostic of container orchestrators.

The spec is **intentionally simple** to implement.

```

1 {
2   "cniVersion": "1.1.0",
3   "name": "mynet",
4   "plugins": [
5     {
6       "type": "bridge",
7       // plugin specific parameters
8       "bridge": "cni0",
9       "keyA": ["some more", "plugin specific", "configuration"],
10      "ipam": {
11        "type": "host-local",
12        // ipam specific
13        "subnet": "10.1.0.0/16",
14        "gateway": "10.1.0.1",
15        "routes": [
16          {"dst": "0.0.0.0/0"}
17        ]
18      },
19      "dns": {
20        "nameservers": [ "10.1.0.1" ]
21      }
22    }
23  ]
24 }
```

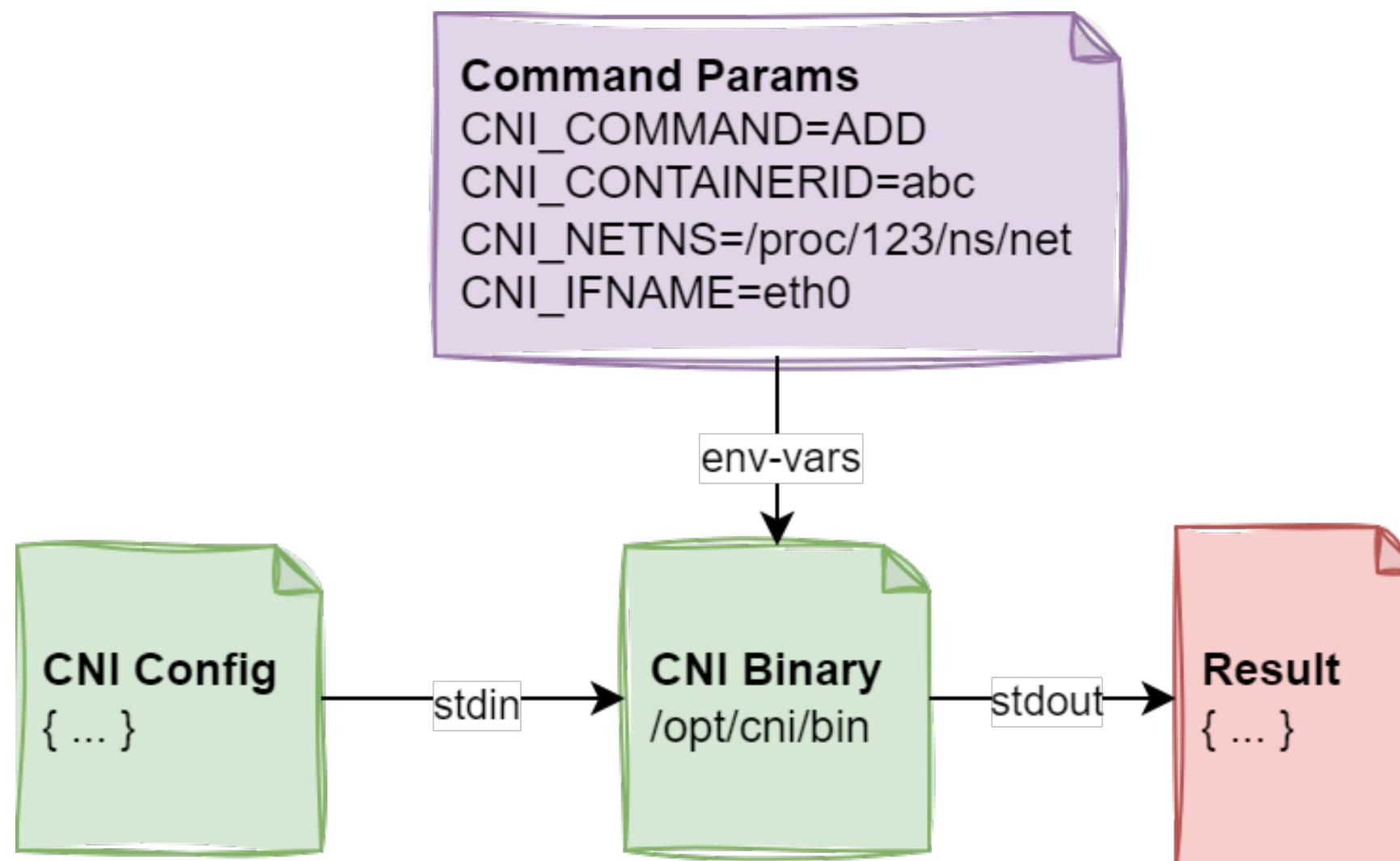
CNI Config

- CNI config is in **JSON** format
- **Versioned** CNI specification
 - So far: 0.1.0 → 0.2.0 → 0.3.1 → 0.4.0
→ 1.0.0 → 1.1.0
- Plugins **can be chained**
- Provides **sane defaults** for handling **IPAM**, **DNS** and more
- Container Runtimes **will only take first config** (first sorted by ASCII)

```

1 {
2   "cniVersion": "1.1.0",
3   "name": "mynet",
4   "plugins": [
5     {
6       "type": "bridge",
7       // plugin specific parameters
8       "bridge": "cni0",
9       "keyA": ["some more", "plugin specific", "configuration"],
10      "ipam": {
11        "type": "host-local",
12        // ipam specific
13        "subnet": "10.1.0.0/16",
14        "gateway": "10.1.0.1",
15        "routes": [
16          {"dst": "0.0.0.0/0"}
17        ]
18      },
19      "dns": {
20        "nameservers": [ "10.1.0.1" ]
21      }
22    }
23  ]
24 }
```

The CNI Binary

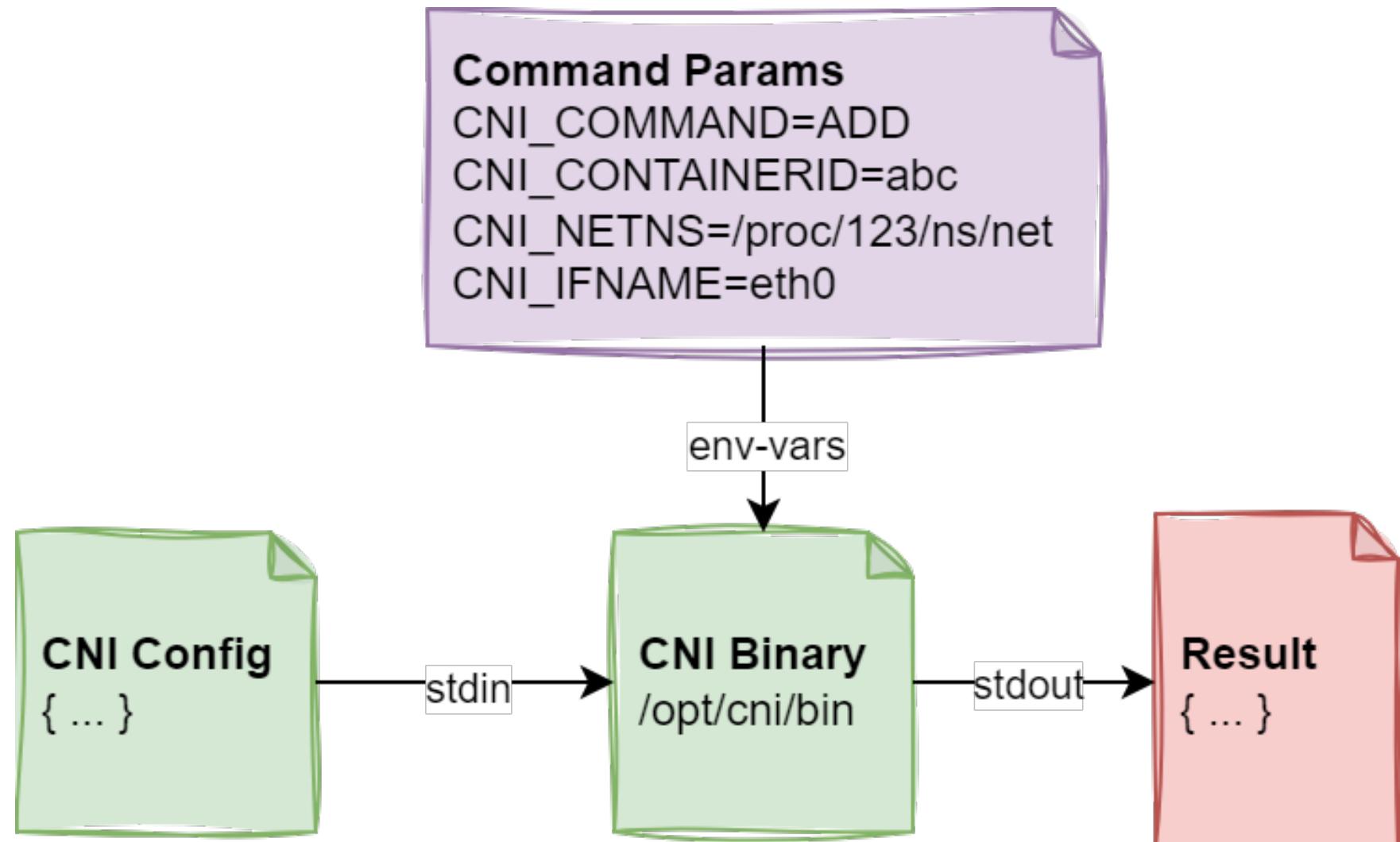


- Config provided on `stdin`
- Parameters / command arguments passed via environment variables
- Result on `stdout`

Available Commands:

- **ADD**: Add container to network, or apply modifications
- **DEL**: Remove container from network, or un-apply modifications
- **CHECK**: Check container's networking is as expected
- **STATUS**: Check plugin status
- **VERSION**: probe plugin version support
- **GC**: Clean up any stale resources

The CNI Binary

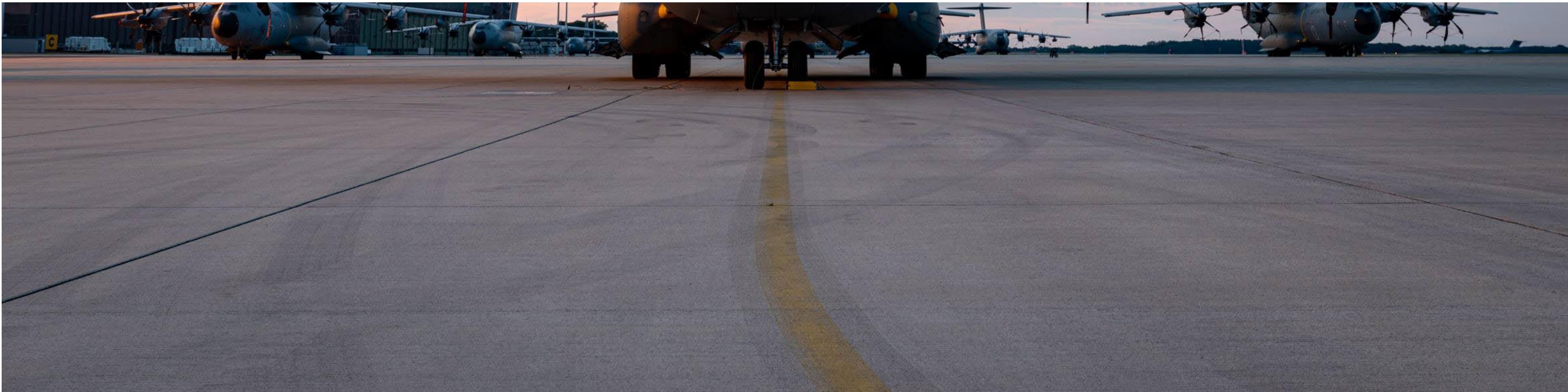


```

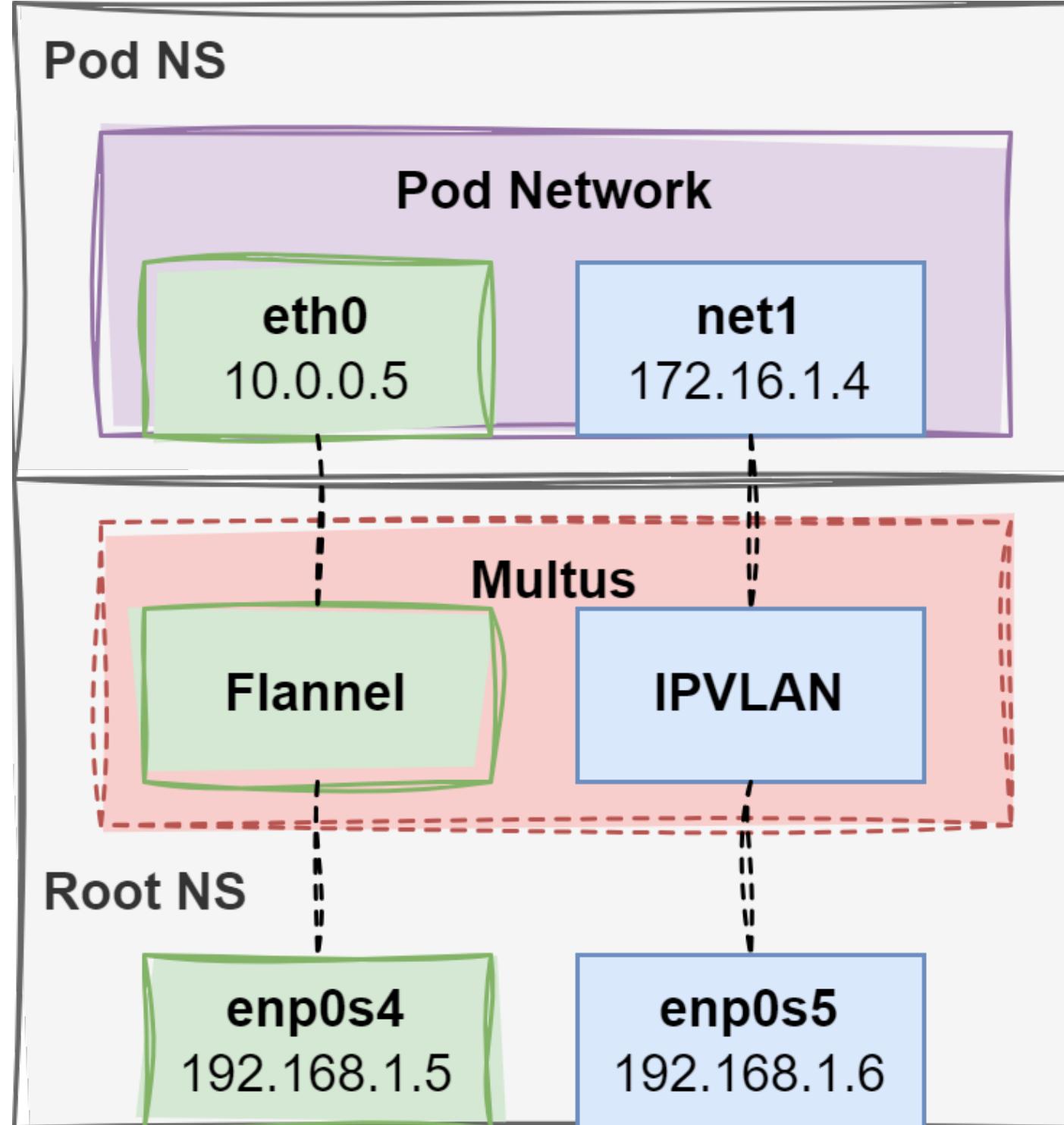
1 {
2   "ips": [
3     {
4       "address": "10.1.0.5/16",
5       "gateway": "10.1.0.1"
6     }
7   ],
8   "routes": [
9     {
10       "dst": "0.0.0.0/0"
11     }
12   ],
13   "dns": {
14     "nameservers": [ "10.1.0.1" ]
15   }
16 }
  
```



Let's create a CNI from Scratch!

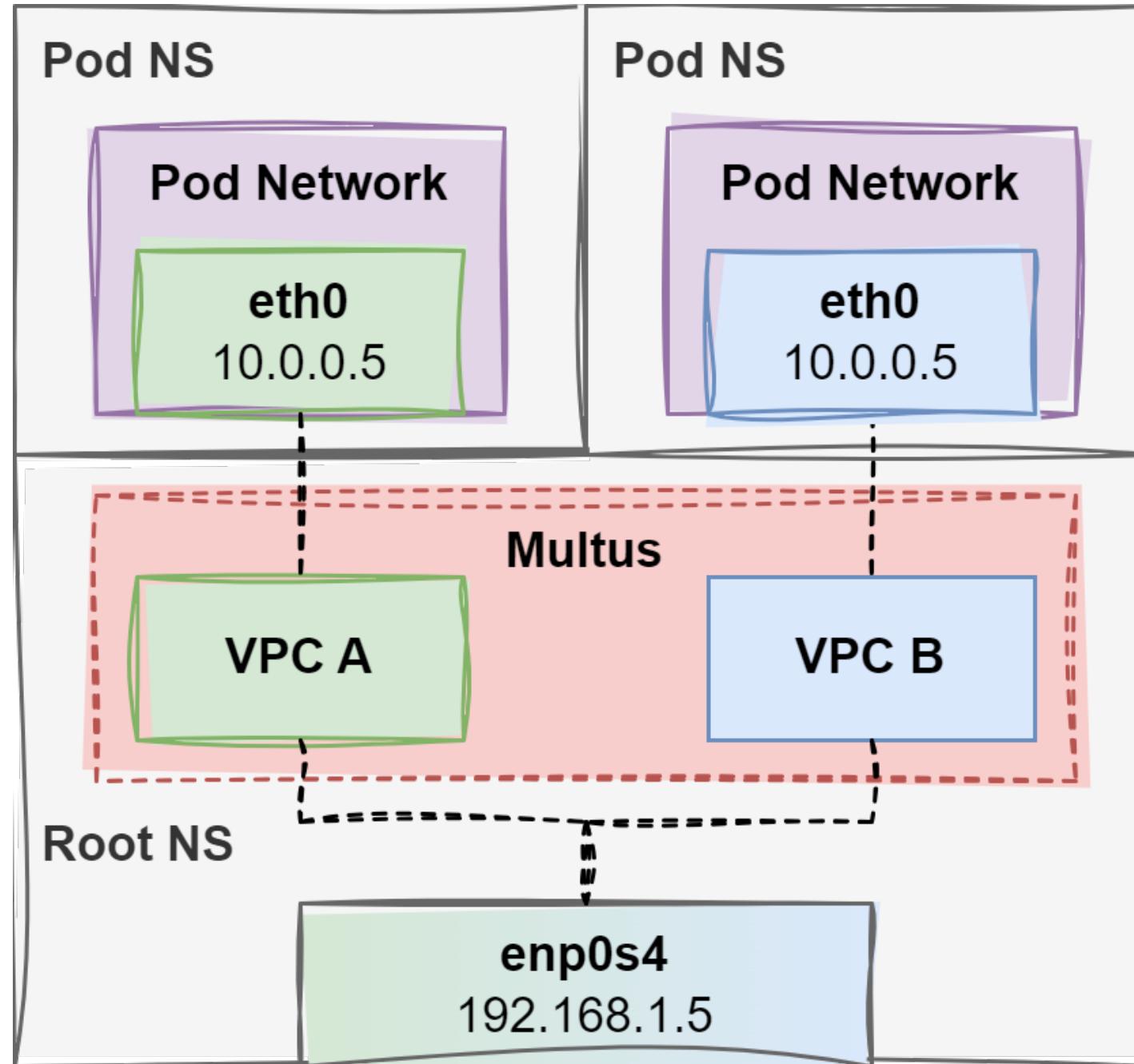


Multiple CNIs



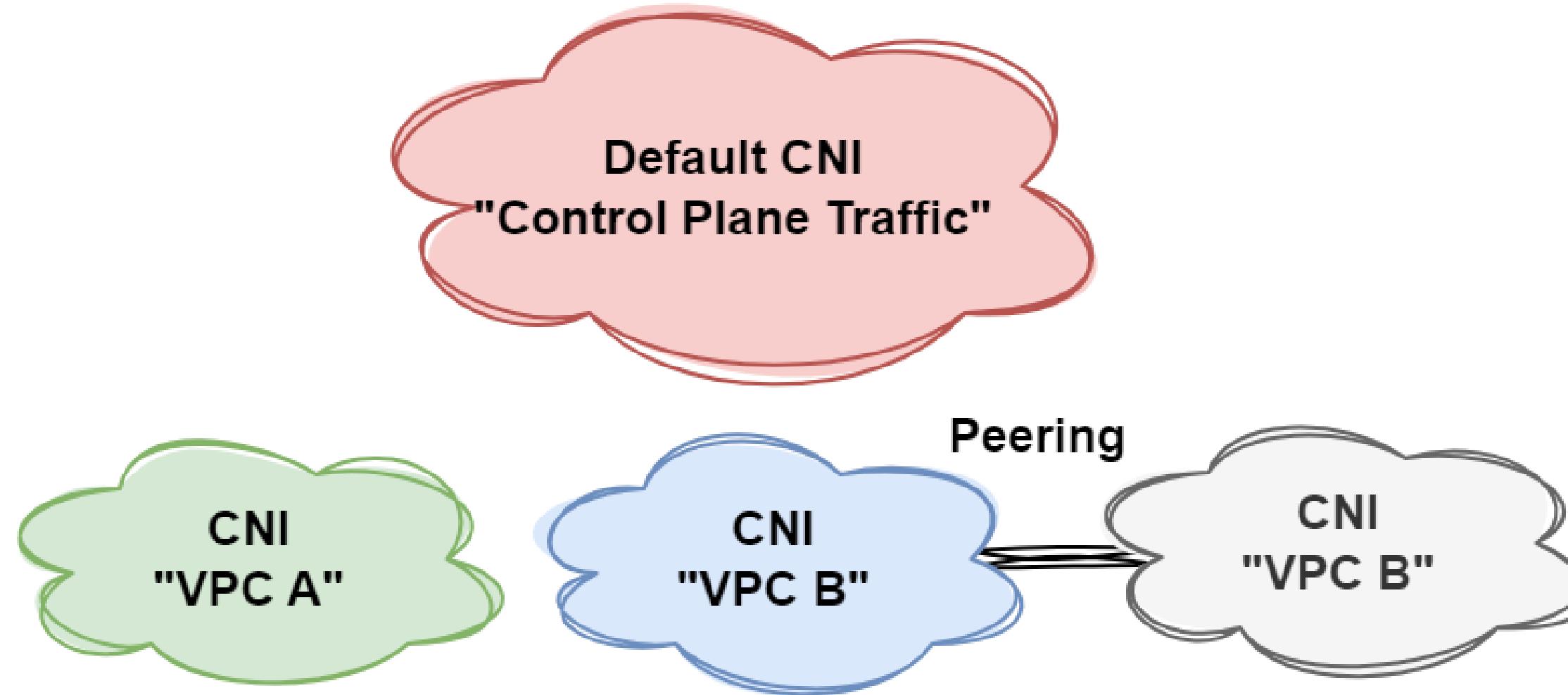
- Multus CNI is “meta-plugin”, a CNI plugin that can **call multiple other CNI plugins**
- Allows **attaching multiple network interfaces** to pods
- **Override default network**

Multiple CNIs



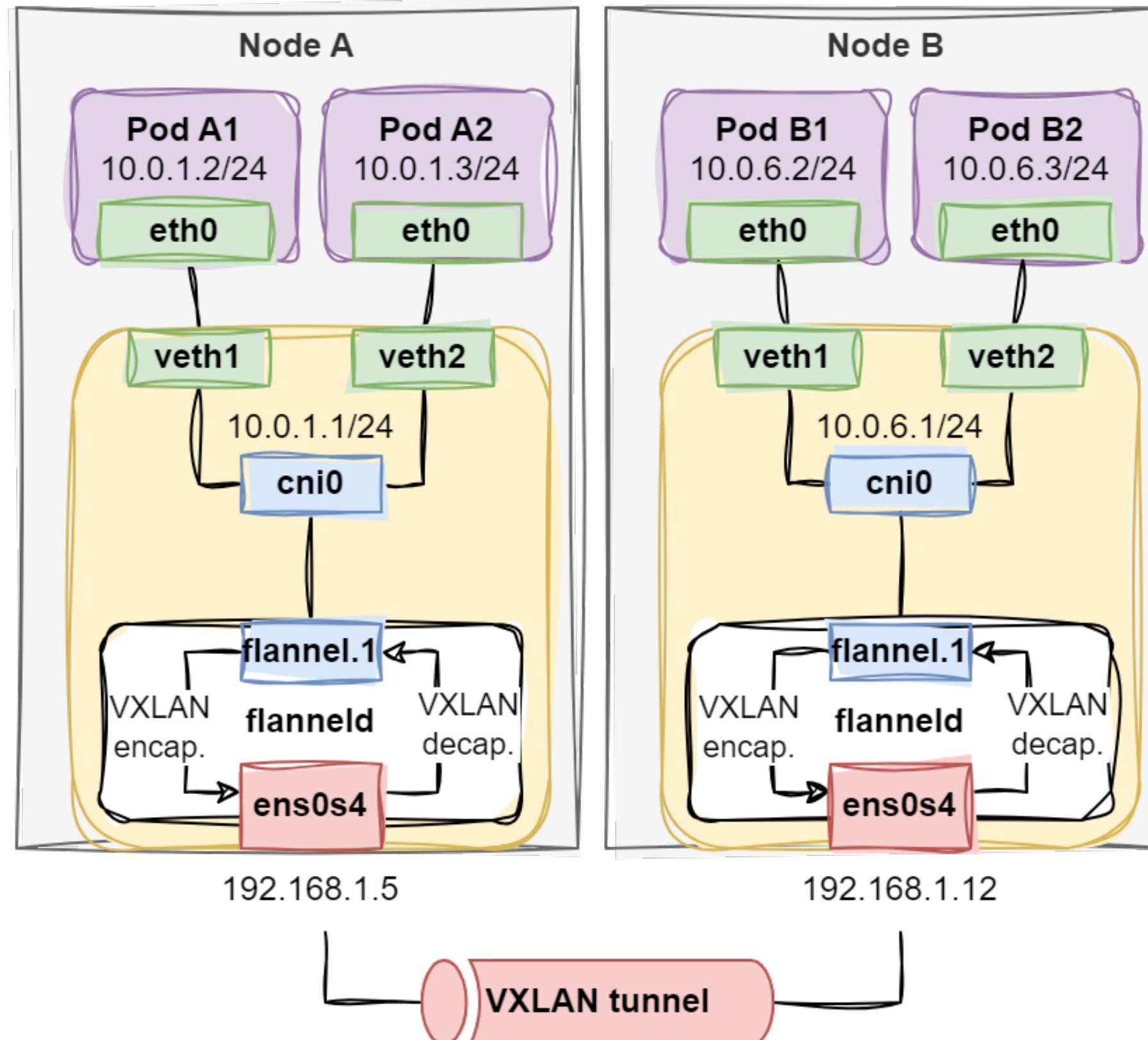
- We need to **override** the **default network**
- Every VPC is a **separate CNI config**
 - Same binary, but **different config**
 - Each CNI is an **independent overlay network**

Multiple CNIs



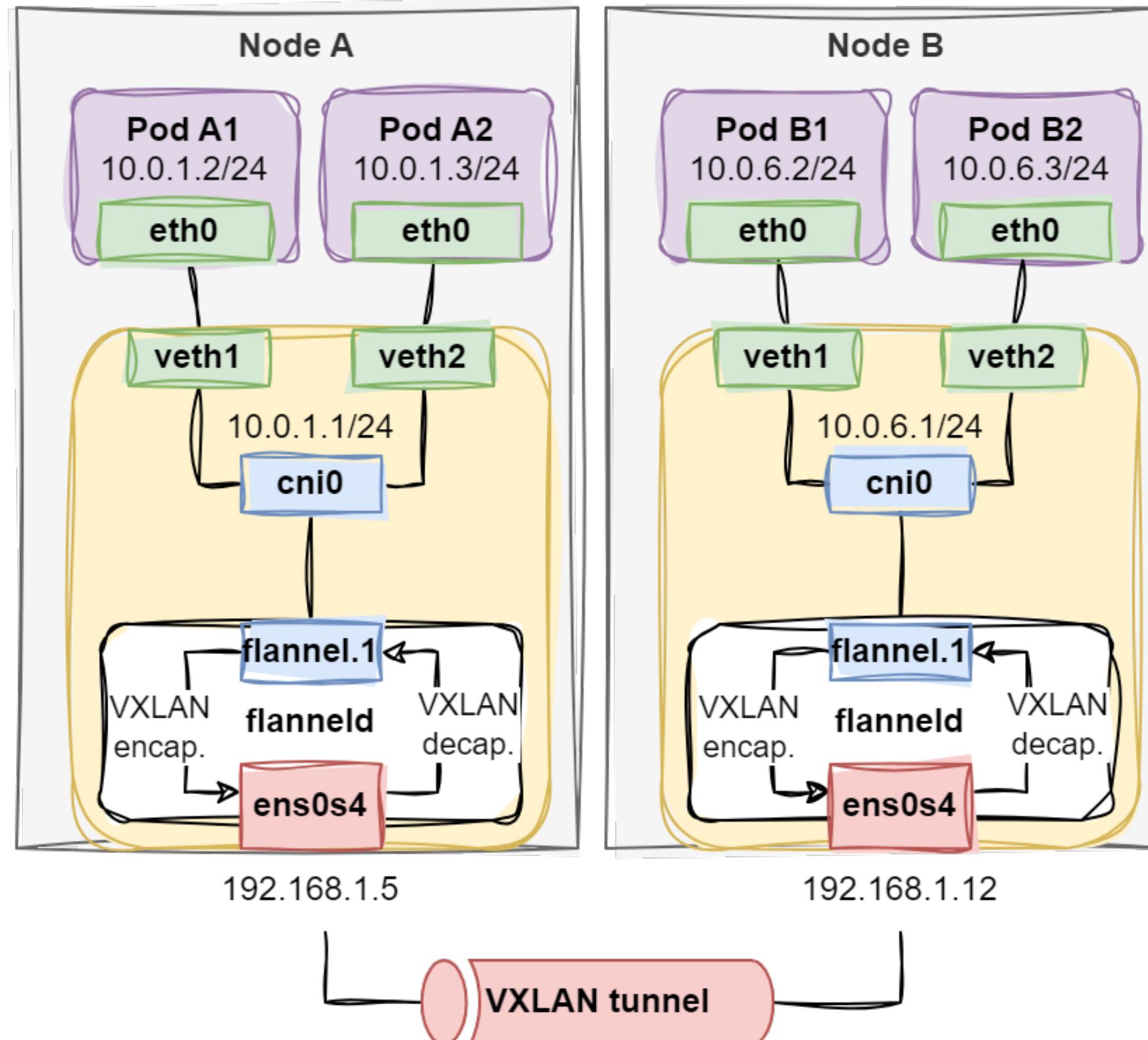
- We need to **override** the **default network**
- Every VPC is a **separate CNI config**
 - **Same binary, but different config**
 - Each CNI is an **independent overlay network**

VXLAN



- Most CNI implementations today leverage **VXLAN** (Cilium, Calico, Flannel, ...)
- Every **node gets a subnet**
- Traffic within Pod IP range is **encapsulated**

VXLAN



```

1 # Create VXLAN and bridge once per node
2 ip link add vxlan200 type vxlan \
   id 200 \
   dstport 4789 \
   local ${NODE_IP} \
   nolearning
3 ip link add br200 type bridge
4 ip link set vxlan200 master br200
5 ip link set br200 type bridge stp_state 0
6 ip link set up dev br200
7 ip link set up dev vxlan200
8
9 # Create veth pair for each pod
10 ip link add dev ${DEV} type veth peer name ${POD_DEV}
11 ip link set dev ${POD_DEV} netns ${POD_NS}
12 ip addr add ${ADDR} dev ${DEV}
13 ip link set ${DEV} master br200
14 ip link set dev ${DEV} up

```

Implementation in Go

[plugins/plugins/sample/main.go](https://github.com/main/containernetworking/plugins/blob/main/plugins/sample/main.go)
[at main · containernetworking/plugins · GitHub](https://github.com/main/containernetworking/plugins)

```

1 package main
2
3 import (
4     "encoding/json"
5
6     "github.com/containernetworking/cni/pkg/skel"
7     "github.com/containernetworking/cni/pkg/types"
8     "github.com/containernetworking/cni/pkg/version"
9     bv "github.com/containernetworking/plugins/pkg/utils/buildversion"
10 )
11
12 type PluginConf struct {
13     types.NetConf
14
15     MyAwesomeFlag     bool   `json:"myAwesomeFlag"`
16     AnotherAwesomeArg string `json:"anotherAwesomeArg"`
17 }
18
19
20 func parseConfig(stdin []byte) (*PluginConf, error) {
21     conf := PluginConf{}
22     if err := json.Unmarshal(stdin, &conf); err != nil {
23         return nil, fmt.Errorf("failed to parse network configuration: %v", err)
24     }
25     // ...
26 }
27
28 func cmdAdd(args *skel.CmdArgs) error {
29     conf, err := parseConfig(args.StdinData)
30     // ...
31
32     return types.PrintResult(result, conf.CNIVersion)
33 }
34
35 func cmdDel(args *skel.CmdArgs) error {}
36 func cmdCheck(_ *skel.CmdArgs) error {}
37
38 func main() {
39     skel.PluginMainFuncs(skel.CNIFuncs{
40         Add:    cmdAdd,
41         Check:  cmdCheck,
42         Del:    cmdDel,
43     }, version.All, bv.BuildString("mycni"))
44 }
```

Implementation in Go

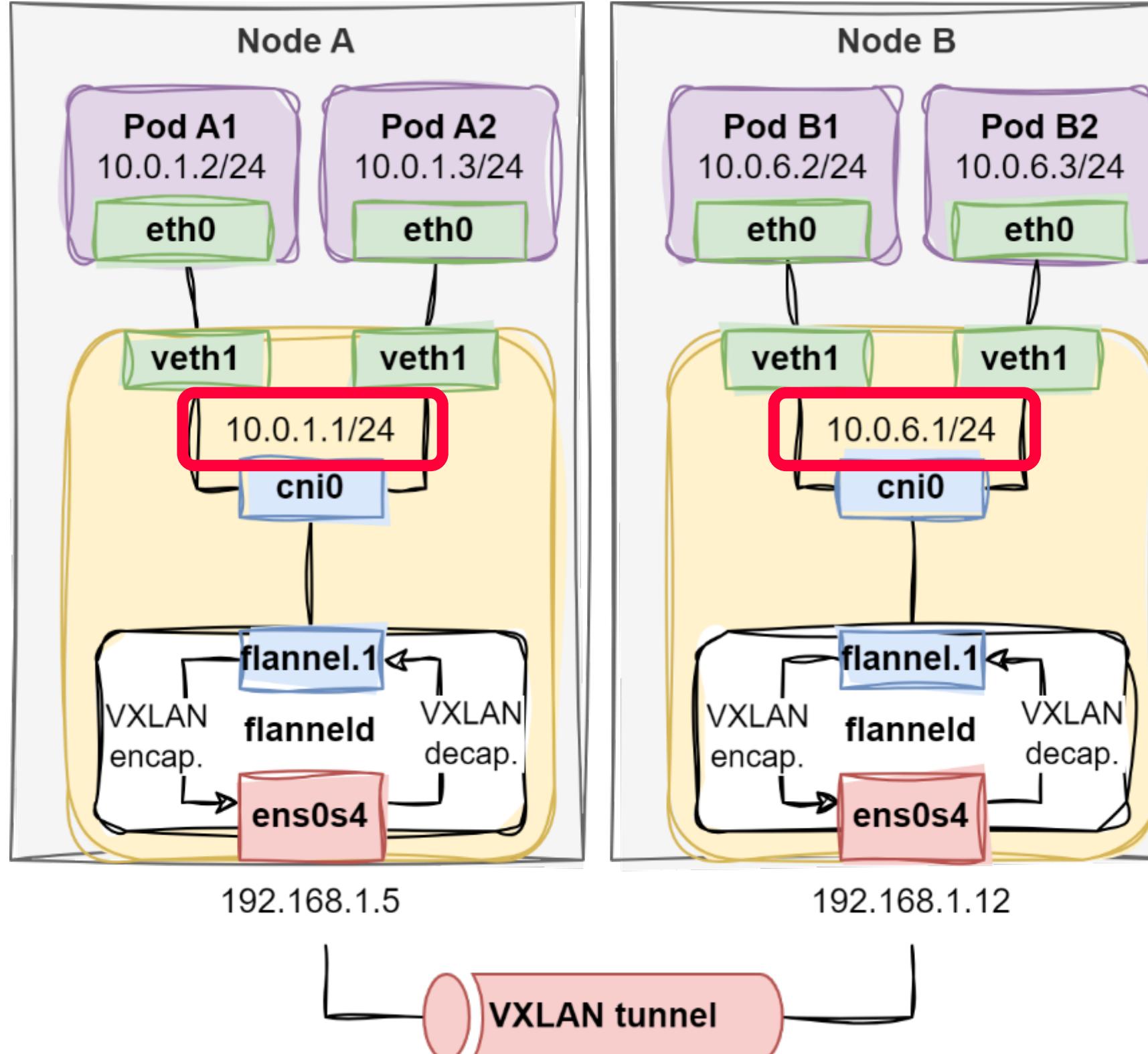
[GitHub - vishvananda/netlink:](https://github.com/vishvananda/netlink)
Simple netlink library for go.

[GitHub - trevex/homelab](https://github.com/trevex/homelab)

```

1 func ensureBridge(n *NetConf, brName string) (*netlink.Bridge, error) {
2     err := netlink.LinkAdd(&netlink.Bridge{
3         LinkAttrs: netlink.LinkAttrs{
4             Name:    brName,
5             MTU:     n.MTU,
6             TxQLen: -1,
7         },
8     })
9
10    // We ignore error if it the bridge already existed
11    if err != nil && err != syscall.EEXIST {
12        return nil, fmt.Errorf("could not add %q: %v", brName, err)
13    }
14
15    // Fetch the bridge
16    br, err := bridgeByName(brName)
17    if err != nil {
18        return nil, err
19    }
20
21    // We want to own the routes for this interface
22    _, err = sysctl.Sysctl(fmt.Sprintf("net/ipv6/conf/%s/accept_ra", brName), "0")
23    if err != nil {
24        return nil, err
25    }
26
27    // We disable Spanning Tree Protocol
28    if err := os.WriteFile(fmt.Sprintf("/sys/class/net/%s/bridge/stp_state", brName), []byte("0"), 0644); err != nil {
29        return nil, err
30    }
31
32    // And make sure it is up
33    if err := netlink.LinkSetUp(br); err != nil {
34        return nil, err
35    }
36
37    return br, nil
38 }
```

But what about routing!?

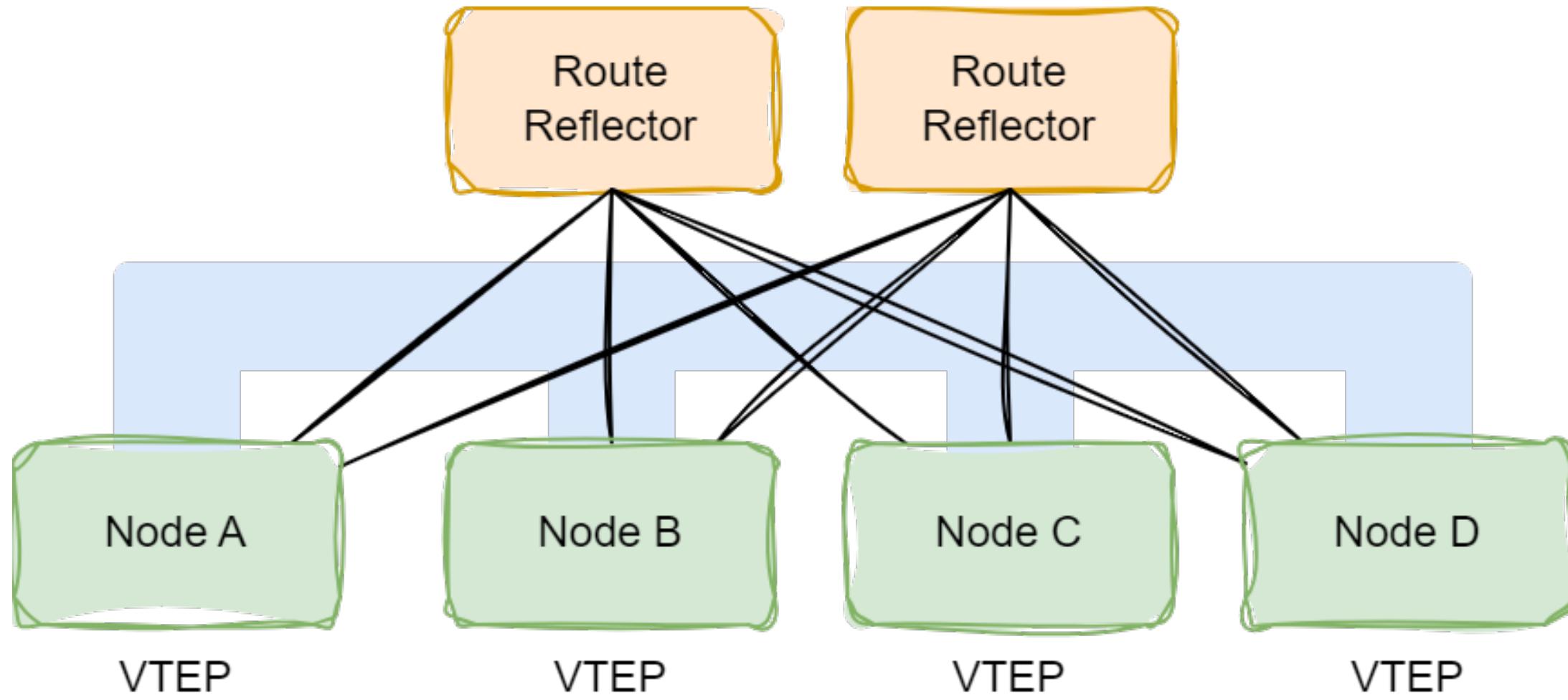


- We want to **avoid allocating a range per node**
 - VPC come **in all sizes**
 - **Not every** node will be **active for** a particular VPC
- How is overlay route management usually handled?

→ **BGP EVPN**

(Border Gateway Protocol Ethernet Virtual Private Network)

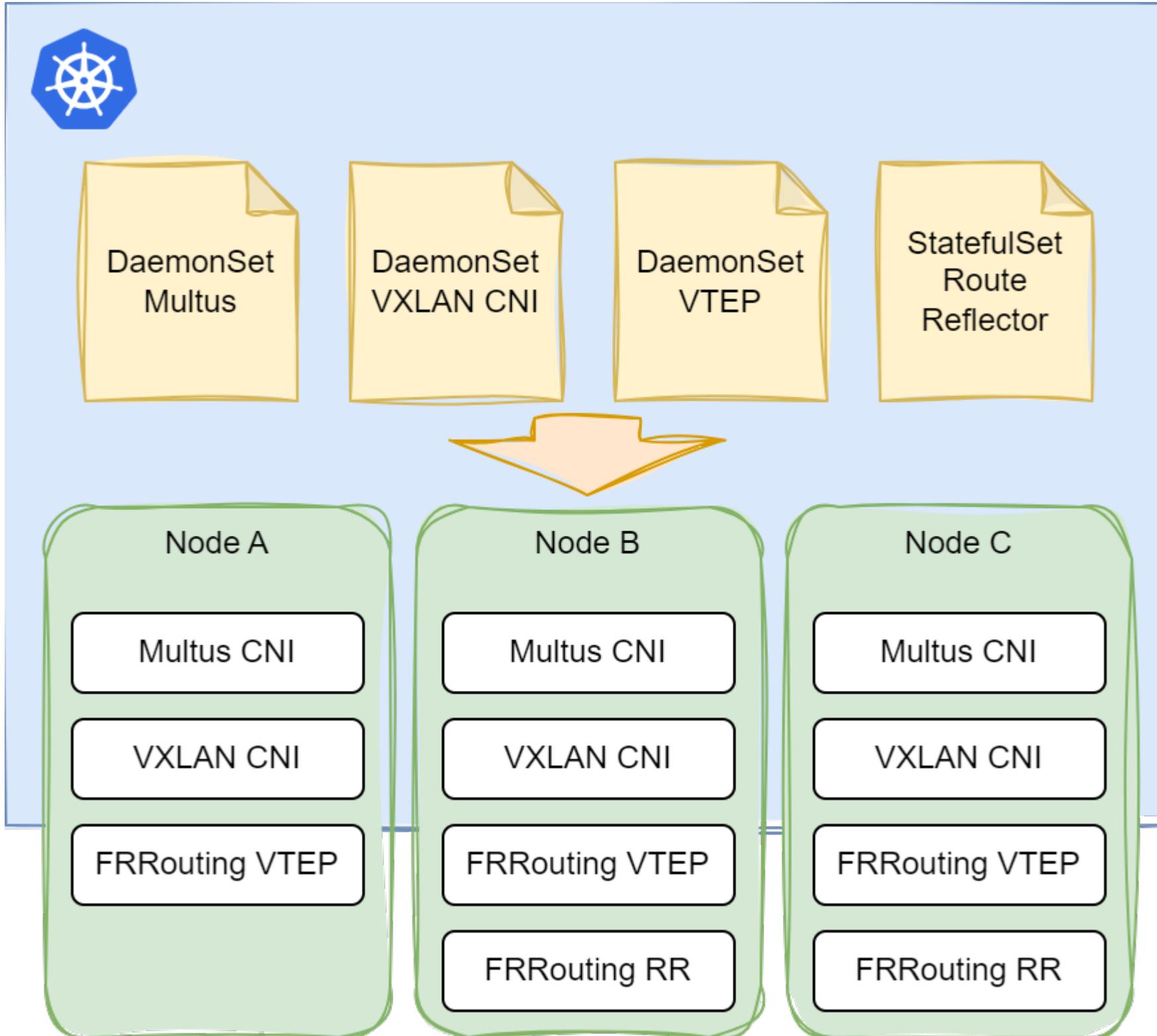
BGP EVPN



- All VXLAN Tunnel Endpoint (VTEP) routes are advertised to route reflectors (RR)
- Each VTEP receives necessary routes for local VNIs (VXLAN network identifier)
- Symmetric IRB (Integrated Bridging and Routing): MP-BGP can simultaneously advertise Layer 2 MAC addresses and Layer 3 routing information

→ Automatic discovery of routes depending on local VNIs

Deployment



1. **Multus DaemonSet**: allow secondary CNI implementation
2. **Our CNI DaemonSet**: install CNI plugin
3. **VTEP DaemonSet**: Announce endpoints
4. **RR StatefulSet**: Distribute routing information

Testing

```
1 apiVersion: "k8s.cni.cncf.io/v1"
2 kind: NetworkAttachmentDefinition
3 metadata:
4   name: pod-a
5 spec:
6   config: '{'
7     "type": "vxlan",
8     "vni": 300,
9     "addr": "172.16.32.1/24" # no ipam
10 }'
11 ---
12 apiVersion: v1
13 kind: Pod
14 metadata:
15   name: pod-a
16   annotations:
17     k8s.v1.cni.cncf.io/networks: pod-a
18 spec:
19   containers:
20     - name: test
21       image: centos/tools
22       imagePullPolicy: IfNotPresent
23       command:
24         - /sbin/init
25   nodeSelector:
26     kubernetes.io/hostname: compute-1
```

```
1 apiVersion: "k8s.cni.cncf.io/v1"
2 kind: NetworkAttachmentDefinition
3 metadata:
4   name: pod-b
5 spec:
6   config: '{'
7     "type": "vxlan",
8     "vni": 300,
9     "addr": "172.16.32.2/24" # no ipam
10 }'
11 ---
12 apiVersion: v1
13 kind: Pod
14 metadata:
15   name: pod-b
16   annotations:
17     k8s.v1.cni.cncf.io/networks: pod-b
18 spec:
19   containers:
20     - name: test
21       image: centos/tools
22       imagePullPolicy: IfNotPresent
23       command:
24         - /sbin/init
25   nodeSelector:
26     kubernetes.io/hostname: compute-2
```

Testing

```

1 apiVersion: "k8s.cni.cncf.io/v1"
2 kind: NetworkAttachmentDefinition
3 metadata:
4   name: pod-a
5 spec:
6   config: '{'
7     "type": "vxlan",
8     "vni": 300,
9     "addr": "172.16.32.1/24" # no ipam
10 }'
11 ---
12 apiVersion: v1
13 kind: Pod
14 metadata:
15   name: pod-a
16   annotations:
17     k8s.v1.cni.cncf.io/networks: pod-a
18 spec:
19   containers:
20     - name: test
21       image: centos/tools
22       imagePullPolicy: IfNotPresent
23       command:
24         - /sbin/init
25   nodeSelector:
26     kubernetes.io/hostname: compute-1

```

```

1 apiVersion: "k8s.cni.cncf.io/v1"
2 kind: NetworkAttachmentDefinition
3 metadata:
4   name: pod-b
5 spec:
6   config: '{'
7     "type": "vxlan",
8     "vni": 300,
9     "addr": "172.16.32.2/24" # no ipam
10 }'
11 ---
12 apiVersion: v1
13 kind: Pod
14 metadata:
15   name: pod-b
16   annotations:
17     k8s.v1.cni.cncf.io/networks: pod-b
18 spec:
19   containers:
20     - name: test
21       image: centos/tools
22       imagePullPolicy: IfNotPresent
23       command:
24         - /sbin/init
25   nodeSelector:
26     kubernetes.io/hostname: compute-2

```

```

1 k exec -it pod-a -- sh
2 ping ${IP_B}%net1
3 nc -l -p 3000
4 # on the other
5 k exec -it pod-b -- sh
6 nc ${IP_A}%net1 3000

```





What does the future hold?

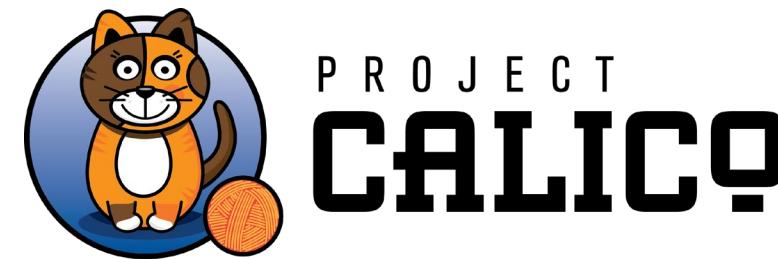


Current developments



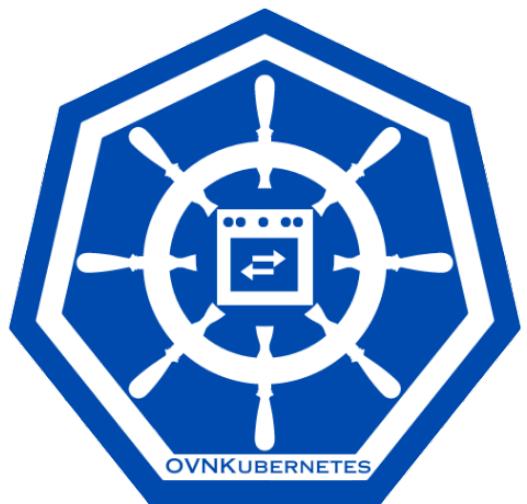
Multi-Pool (Beta)

The Multi-Pool IPAM mode supports allocating PodCIDRs from multiple different IPAM pools, depending on properties of the workload defined by the user, e.g. annotations.



IP pools

IP pools are ranges of IP addresses from which Calico assigns pod IPs. [...] You can control which pool Calico uses for each pod using node selectors, or annotations on the pod or the pod's namespace.



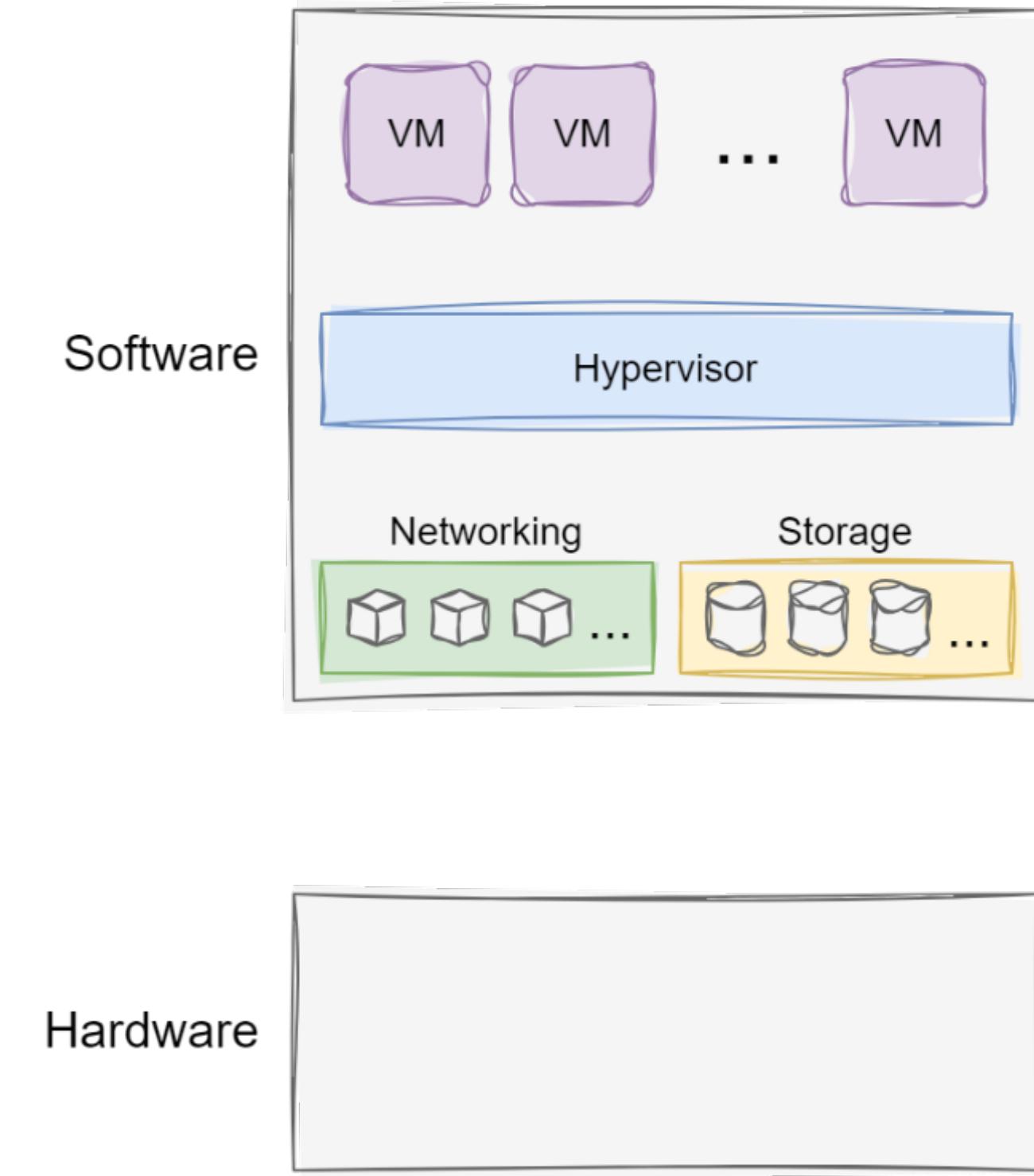
OVN-Kubernetes (Open Virtual Networking - Kubernetes) is an open-source project that provides a robust networking solution for Kubernetes clusters with OVN (Open Virtual Networking) and Open vSwitch (Open Virtual Switch) at its core. It is a Kubernetes networking conformant plugin written according to the CNI (Container Network Interface) specifications.

What is (Hyper-)Converged Infrastructure?

According to [Wikipedia](#):

"Hyper-converged infrastructure (HCI) is a **software-defined IT infrastructure** that virtualizes all of the elements of conventional "hardware-defined" systems. HCI includes, at a minimum, **virtualized computing** (a **hypervisor**), **software-defined storage**, and **virtualized networking (software-defined networking)**. HCI typically runs on **commercial off-the-shelf (COTS) servers**."

→ **Everything** is software defined

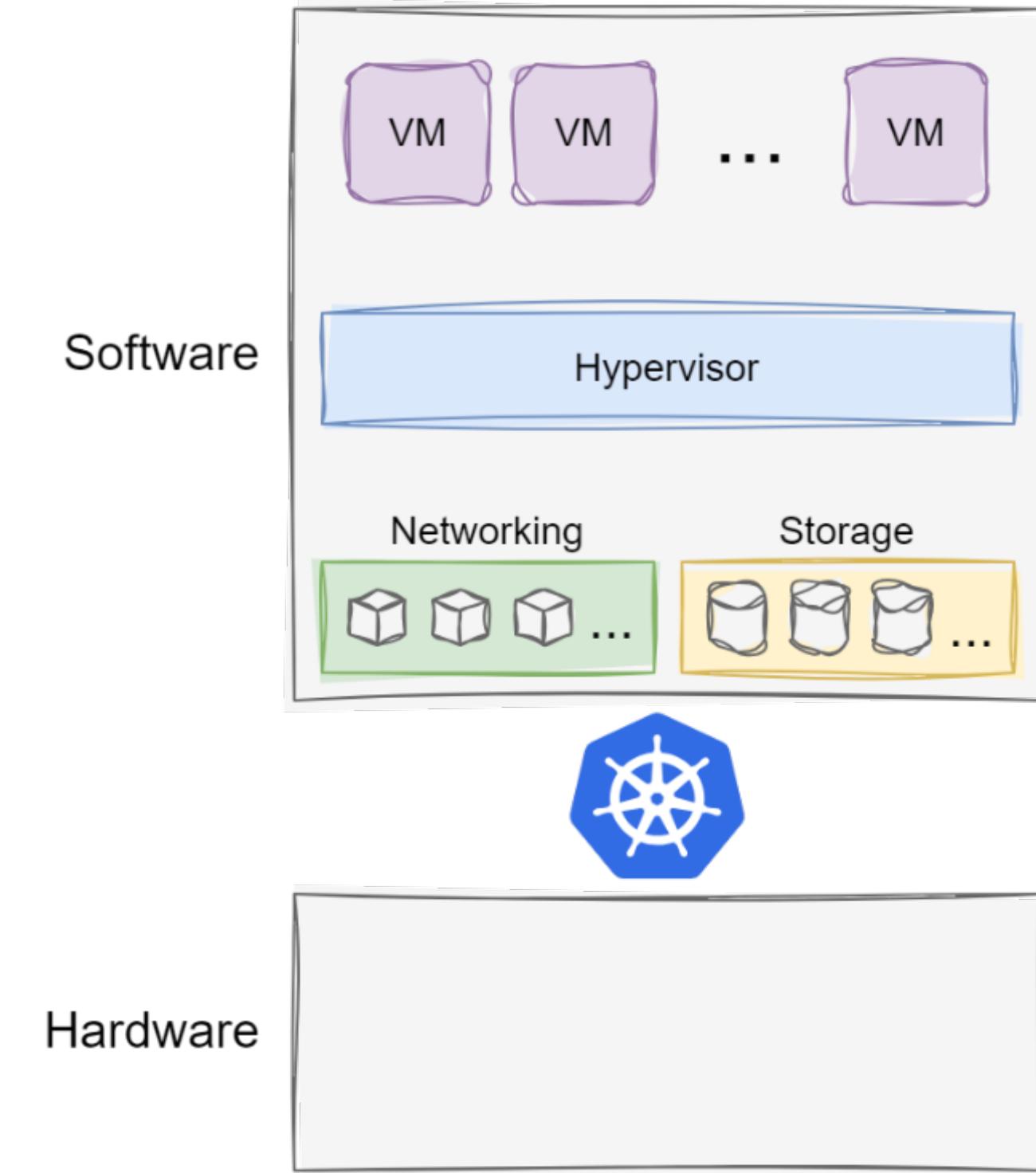


What is (Hyper-)Converged Infrastructure?

According to [Wikipedia](#):

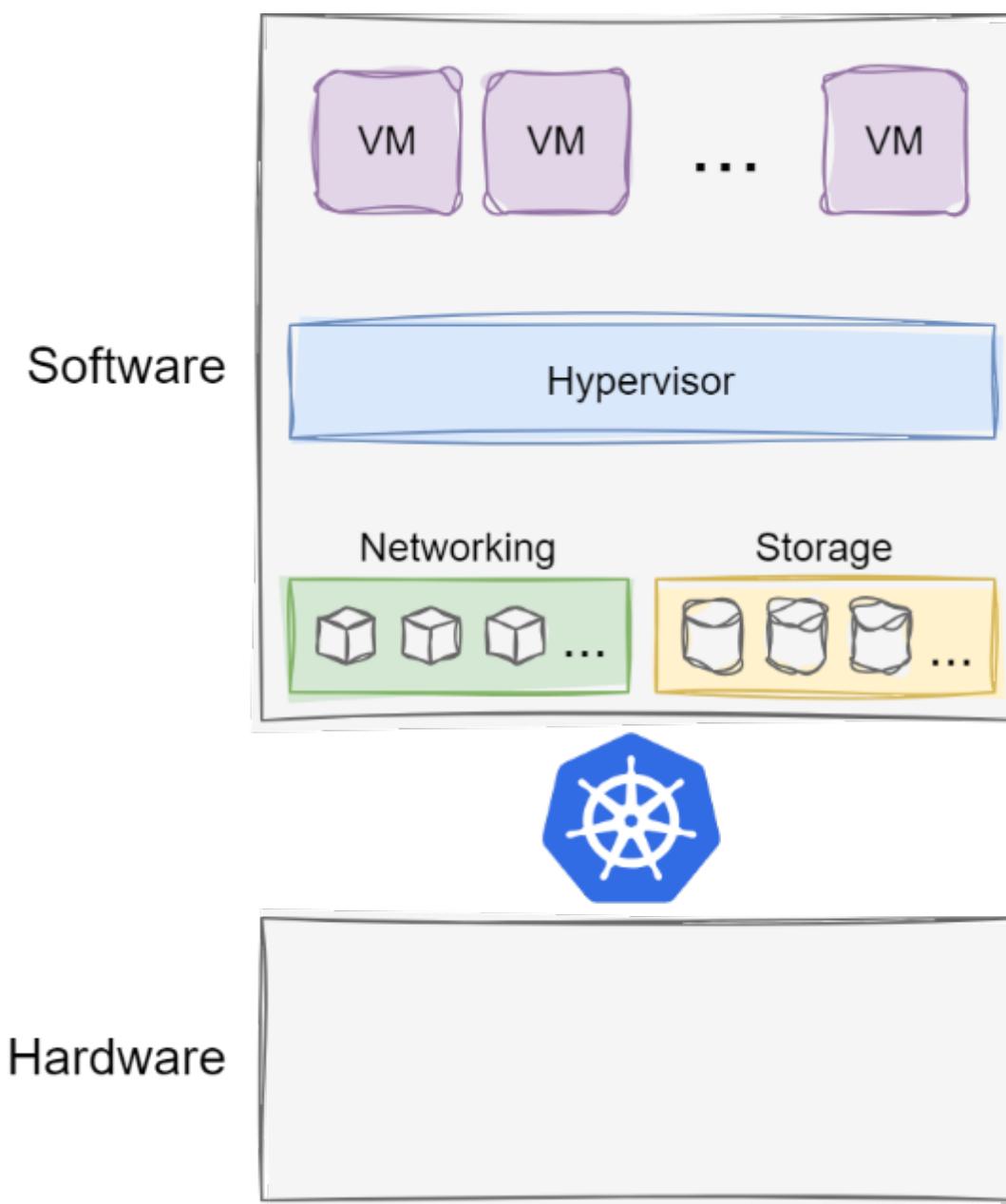
"Hyper-converged infrastructure (HCI) is a software-defined IT infrastructure that virtualizes all of the elements of conventional "hardware-defined" systems. HCI includes, at a minimum, virtualized computing (a hypervisor), software-defined storage, and virtualized networking (software-defined networking). HCI typically runs on commercial off-the-shelf (COTS) servers."

→ Everything is **software defined** in **Kubernetes**

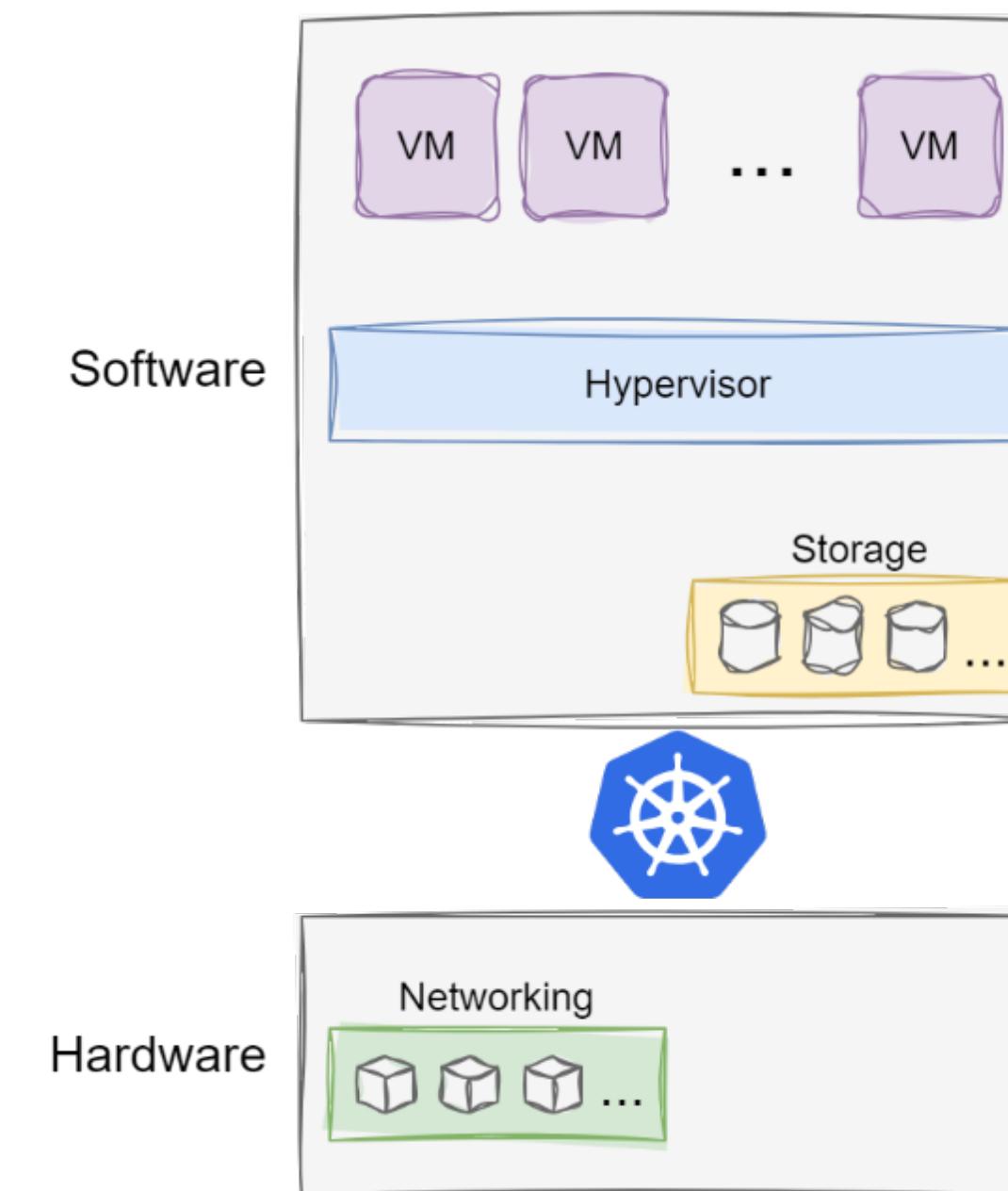


Current trends

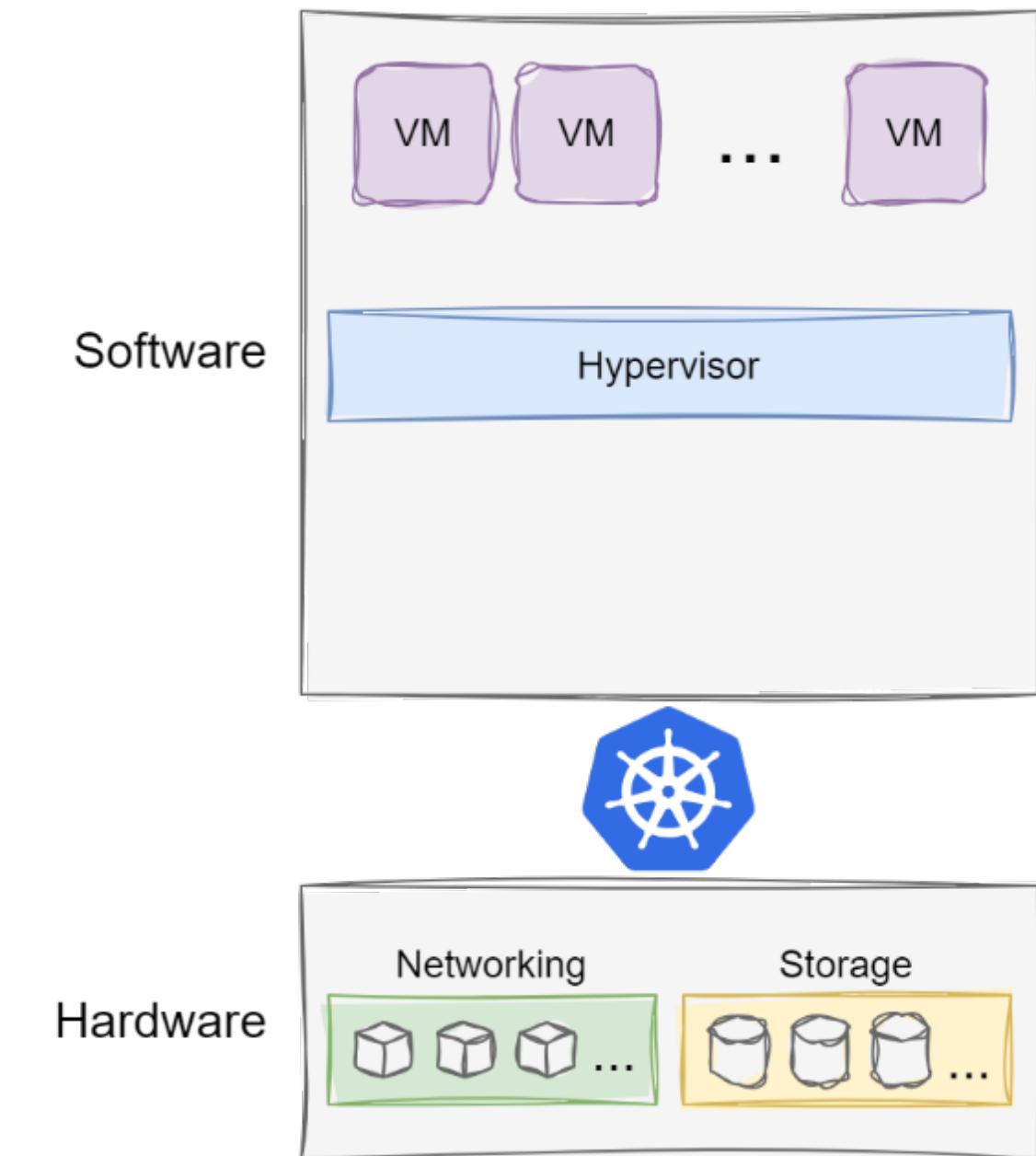
Converged 1.0



Converged 2.0



Converged 3.0



Current trends

- Offloading to DPU/IPUs is becoming popular, particularly at scale
 - AWS Nitro System
 - Microsoft acquired Fungible (DPU vendor)
 - Google partnering with Intel to acquire E2100 IPUs
- Offloading support is gaining traction in the CNCF ecosystem
 - [SR-IOV VMs in Kubevirt](#)
 - [Kube OVN DPDK Kubevirt fork](#)
 - [IronCore by SAP](#)
 - [Userspace CNI by Intel](#)

→ With more FOSS developments in this area everyone can benefit

Conclusion

- Converged infrastructure brings new requirements into the CNCF ecosystem
- Networking in Kubernetes is flexible
- The CNI spec is approachable
- Ideas how to leverage Kubernetes in new scenarios

Questions?

Code



Slides

