# Kubernetes Networking
## Seattle Kubernetes Meetup
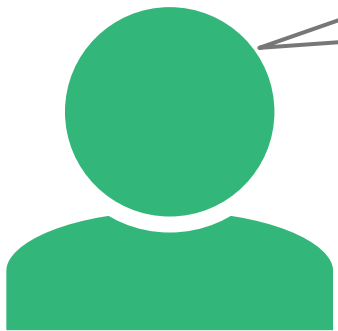
CJ Cullen <cjcullen@google.com>
Software Engineer
@cj_cullen
github.com/cjcullen

Google Cloud Platform
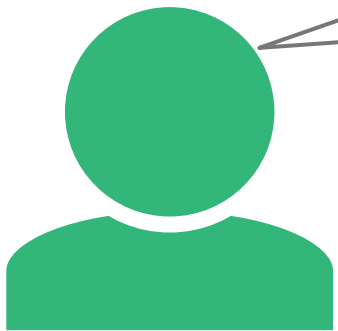
# Docker Networking

# Docker networking

# Docker networking

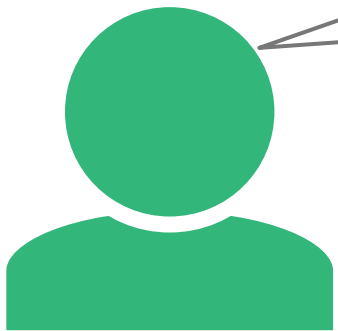# Docker networking



**docker0**                                              **172.16.1.0/24**

# Docker networking



docker0                                    172.16.1.0/24

# Docker networking
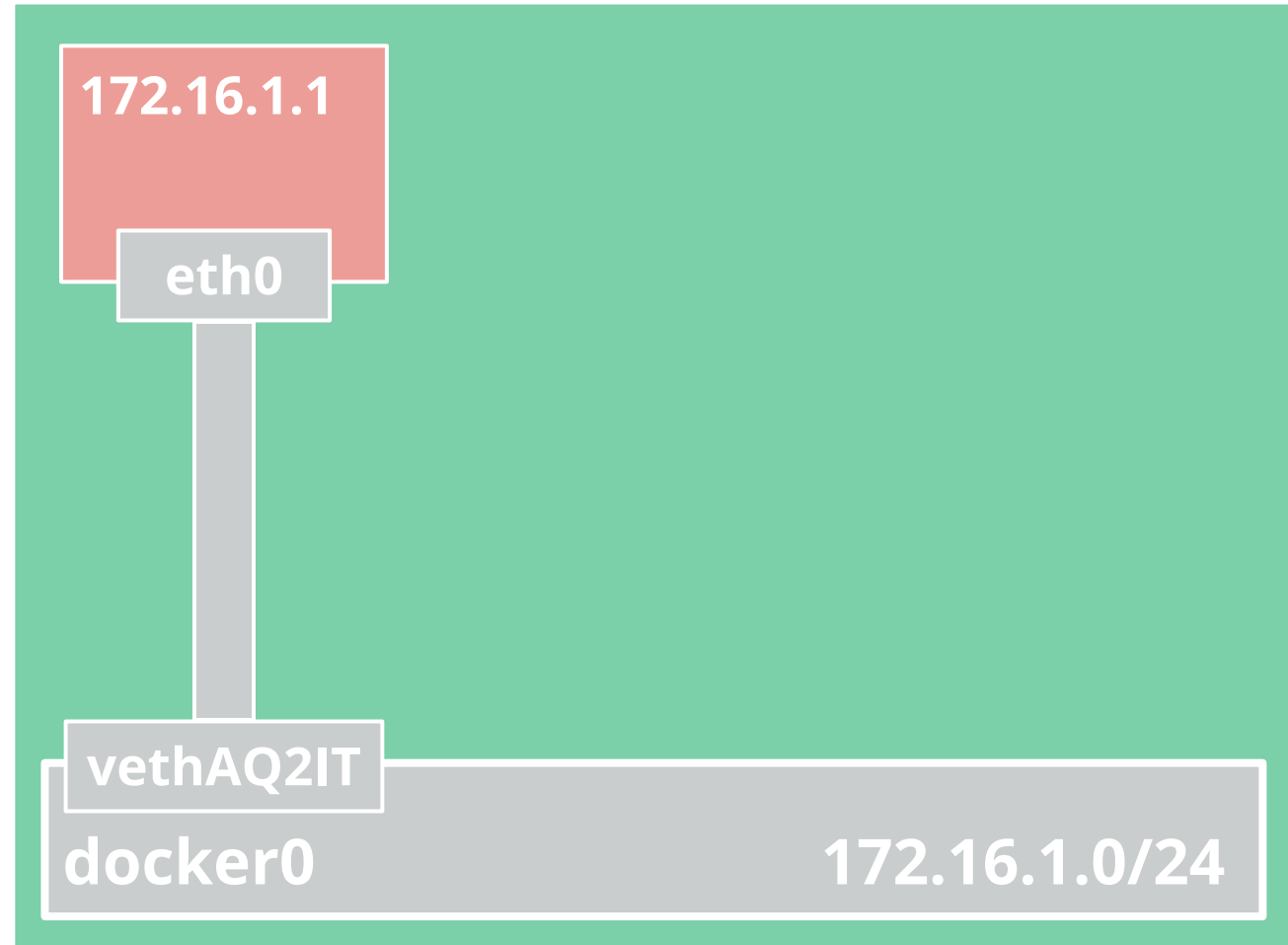


172.16.1.1

eth0

vethAQ2IT

docker0      172.16.1.0/24

# Docker networking



172.16.1.1

eth0

172.16.1.2

eth0

vethAQ2IT

vethS1LUI

docker0

172.16.1.0/24

Google Cloud Platform

# Docker networking



172.16.1.1

172.16.1.2

172.16.1.1

172.16.1.1

# Docker networking

# Host ports

Host ports

# Kubernetes Networking

# Kubernetes networking

IPs are **routable**
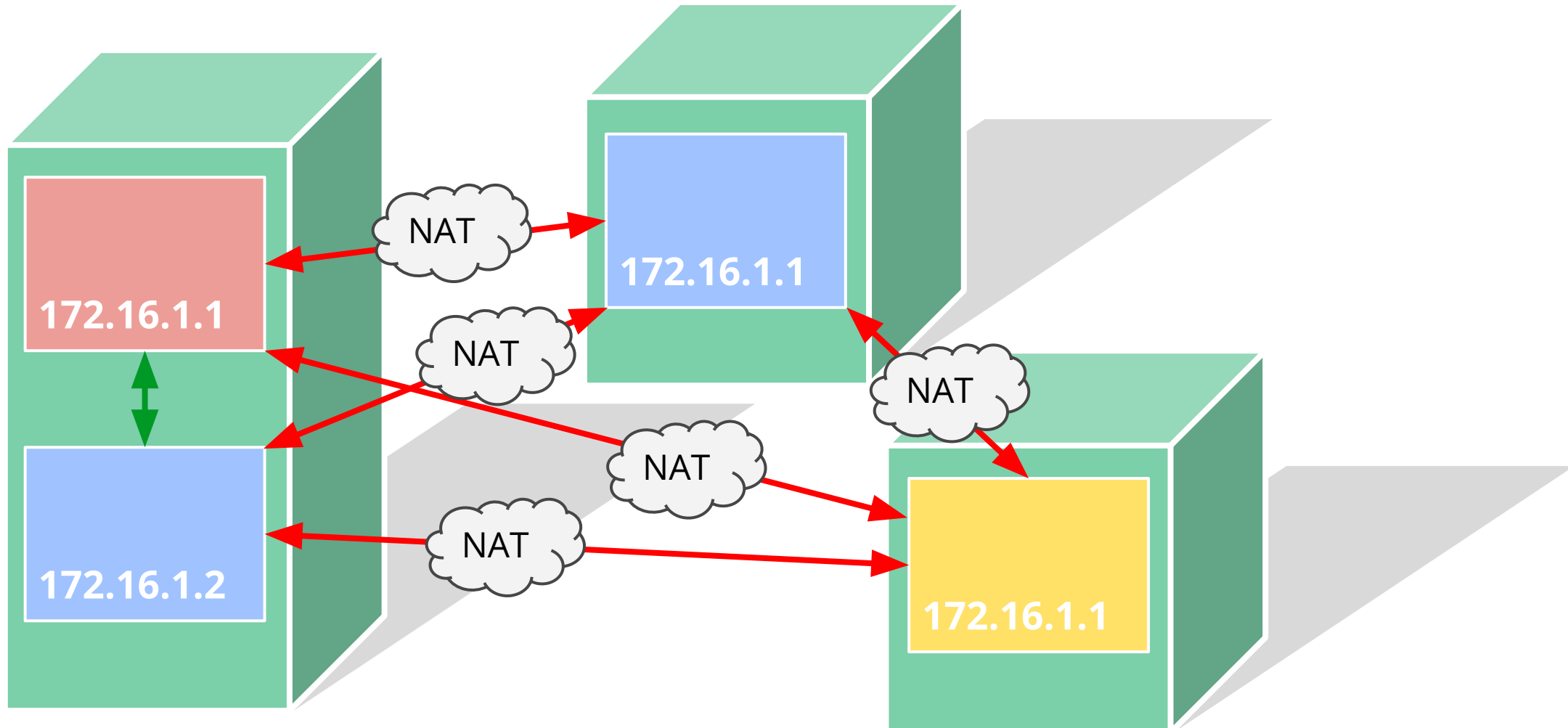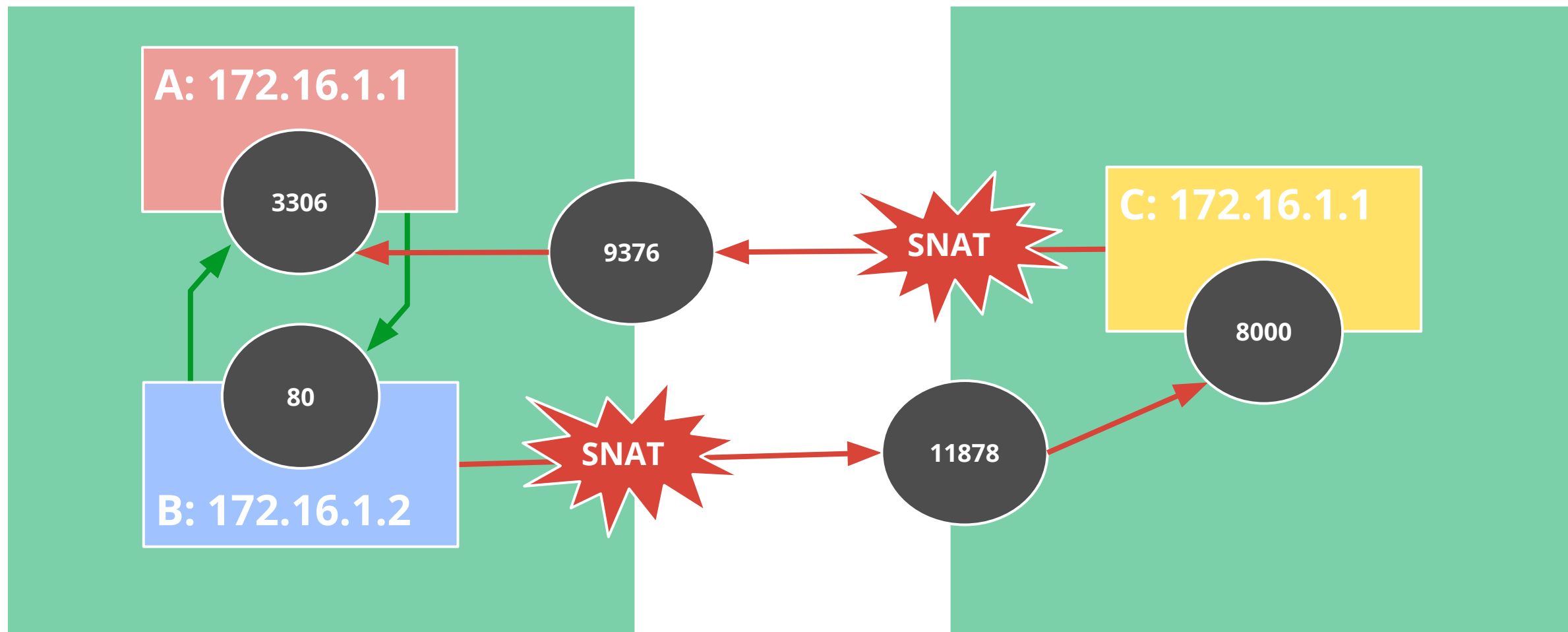- vs docker default private IP

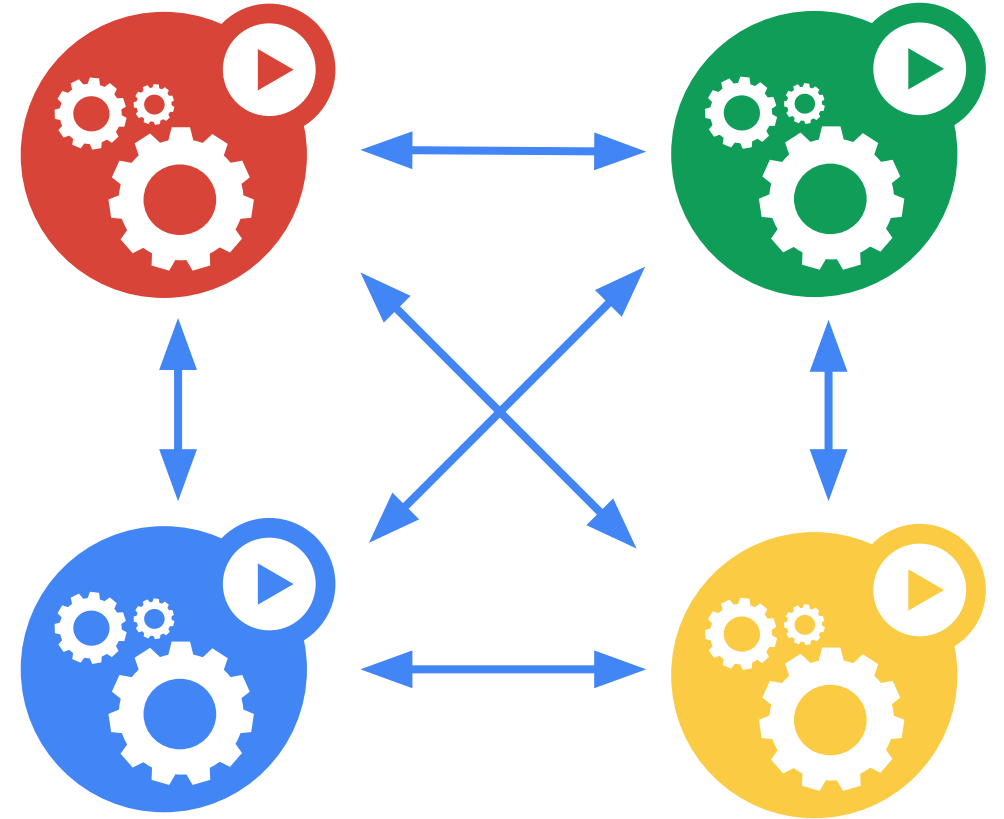Pods can reach each other without NAT
- even across nodes

**No brokering** of port numbers
- too complex, why bother?

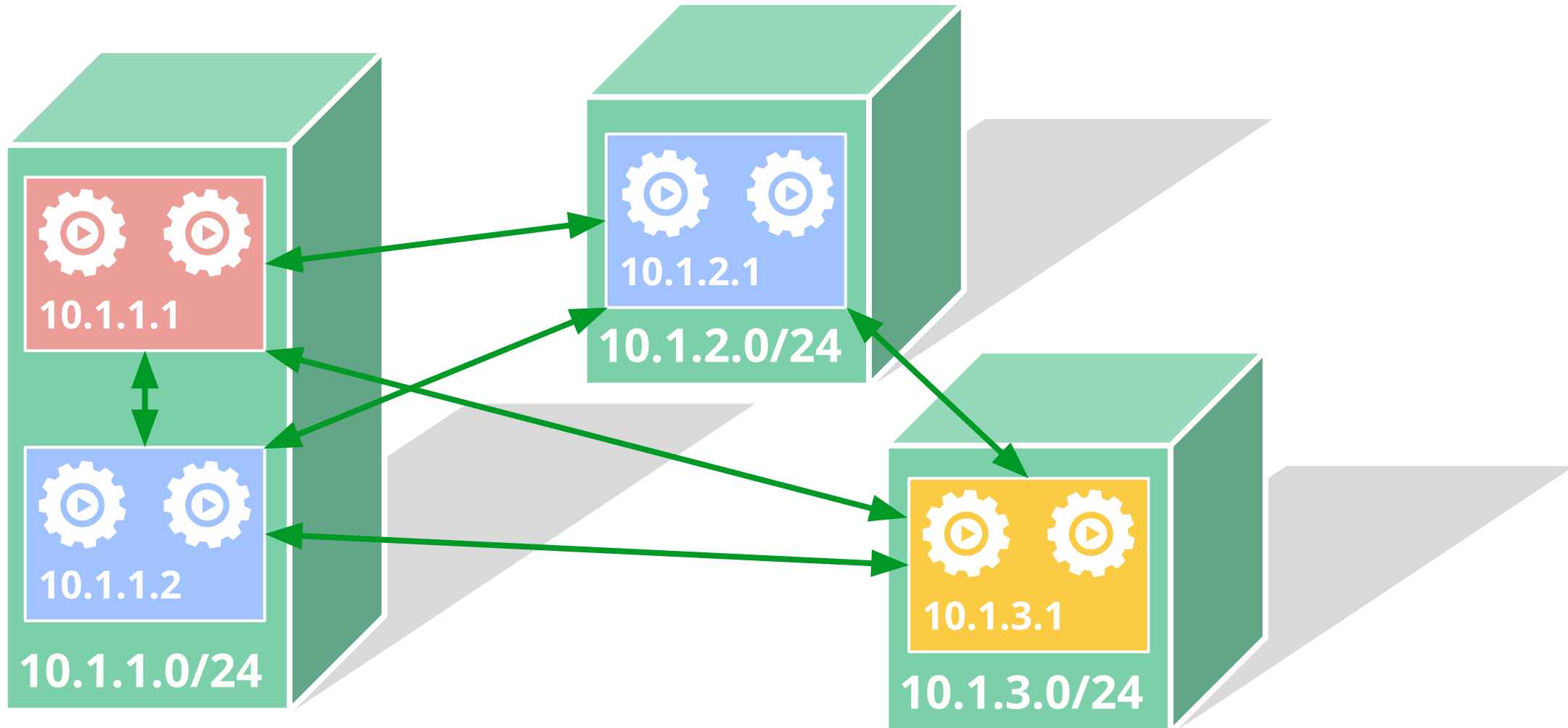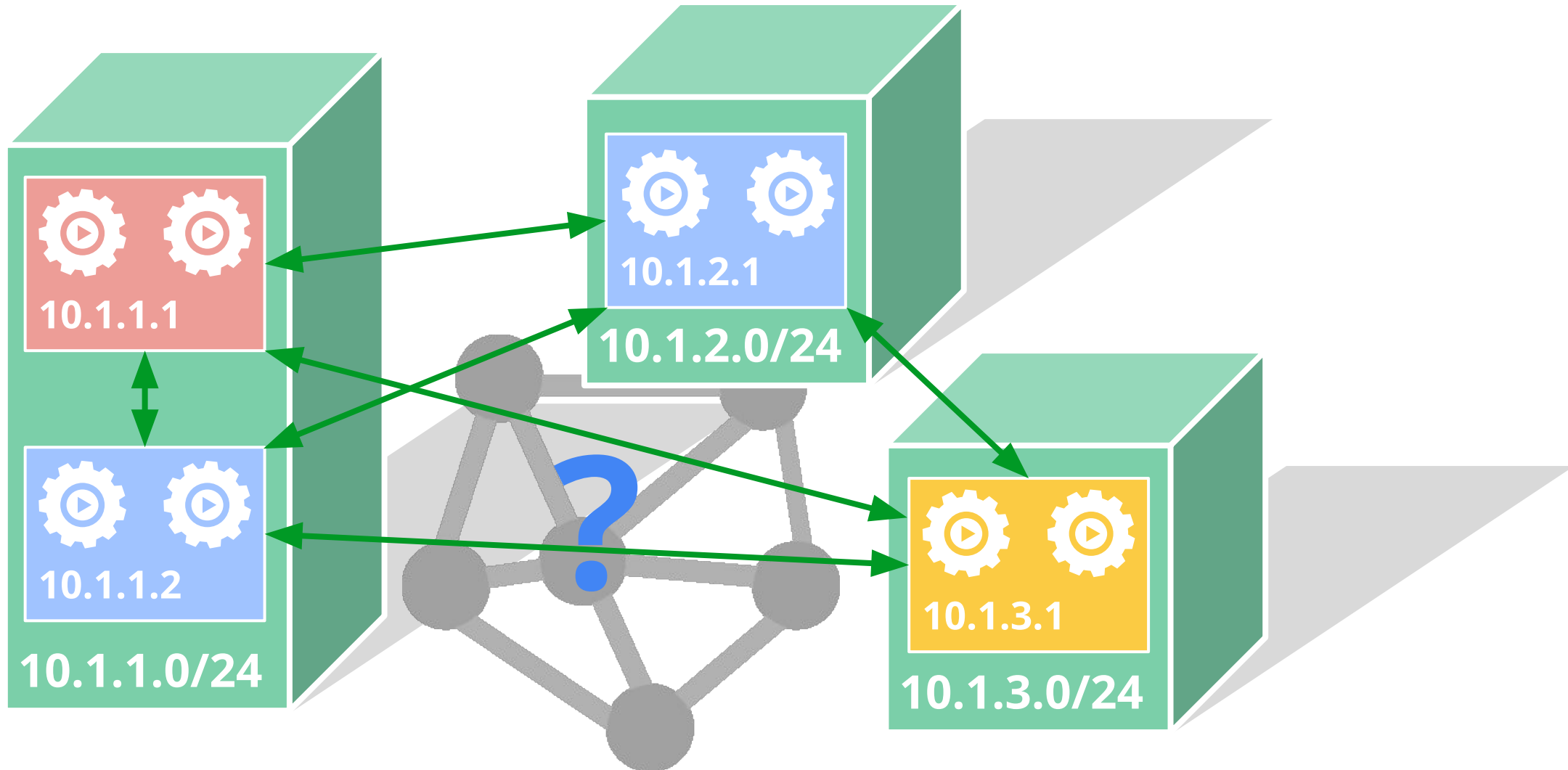**This is a fundamental requirement**
- can be L3 routed
- can be underlayed (cloud)
- can be overlayed (SDN)

Kubernetes networking

# Kubernetes networking

# Kubernetes networking

## On GCE/GKE
- GCE Advanced Routes (program the fabric)
- "Everything to 10.1.1.0/24, send to this VM"

## Plenty of other ways
- AWS: Route Tables
- Weave
- Calico
- Flannel
- OVS
- OpenContrail
- Cisco Contiv
- Others...

# Kubernetes networking

## On **GCE/GKE**

- GCE Advanced Routes (program the fabric)
- "Everything to 10.1.1.0/24, send to this VM"

## Plenty of other ways

- **AWS**: Route Tables
- Weave
- Calico
- Flannel
- OVS
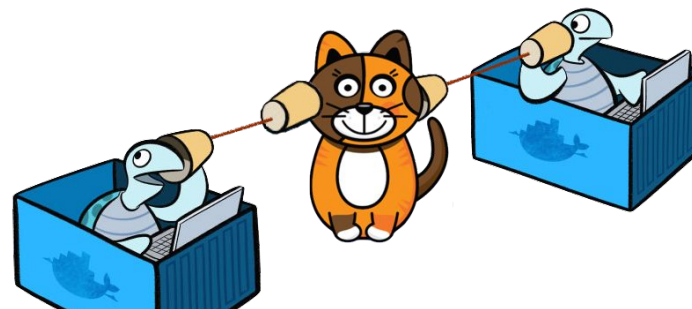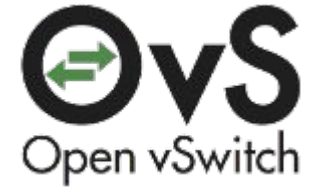- OpenContrail
- Cisco Contiv
- Others...

# Kubernetes networking

## On GCE/GKE
- GCE Advanced Routes (program the fabric)
- "Everything to 10.1.1.0/24, send to this VM"

## Plenty of other ways
- AWS: Route Tables
- Weave
- Calico
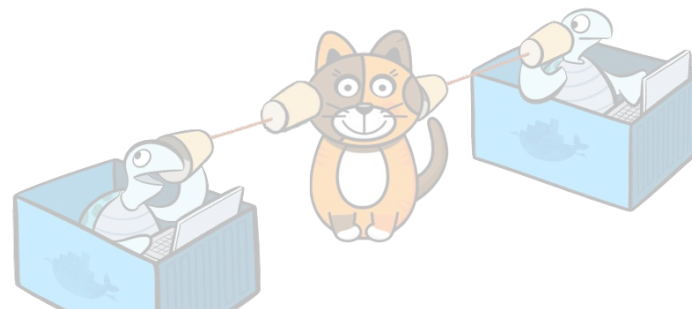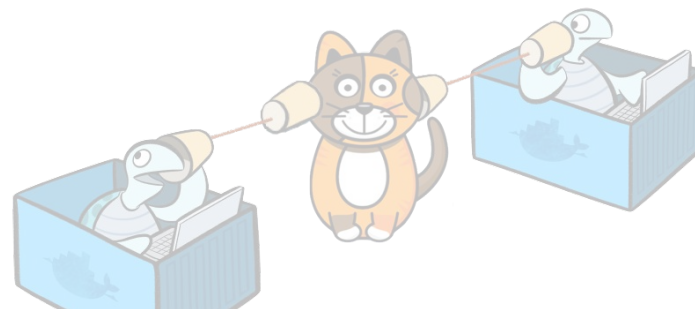- **Flannel**
- OVS
- OpenContrail
- Cisco Contiv
- Others…

# Pods

# Pods

**Small group** of containers & volumes

**Tightly** coupled

The atom of scheduling & placement

Shared namespace
- share IP address & localhost
- share IPC, etc.

Managed lifecycle
- bound to a node, restart in place
- can die, cannot be reborn with same ID

**Example: data puller & web server**

## Pods

**Small group** of containers & volumes

**Tightly** coupled
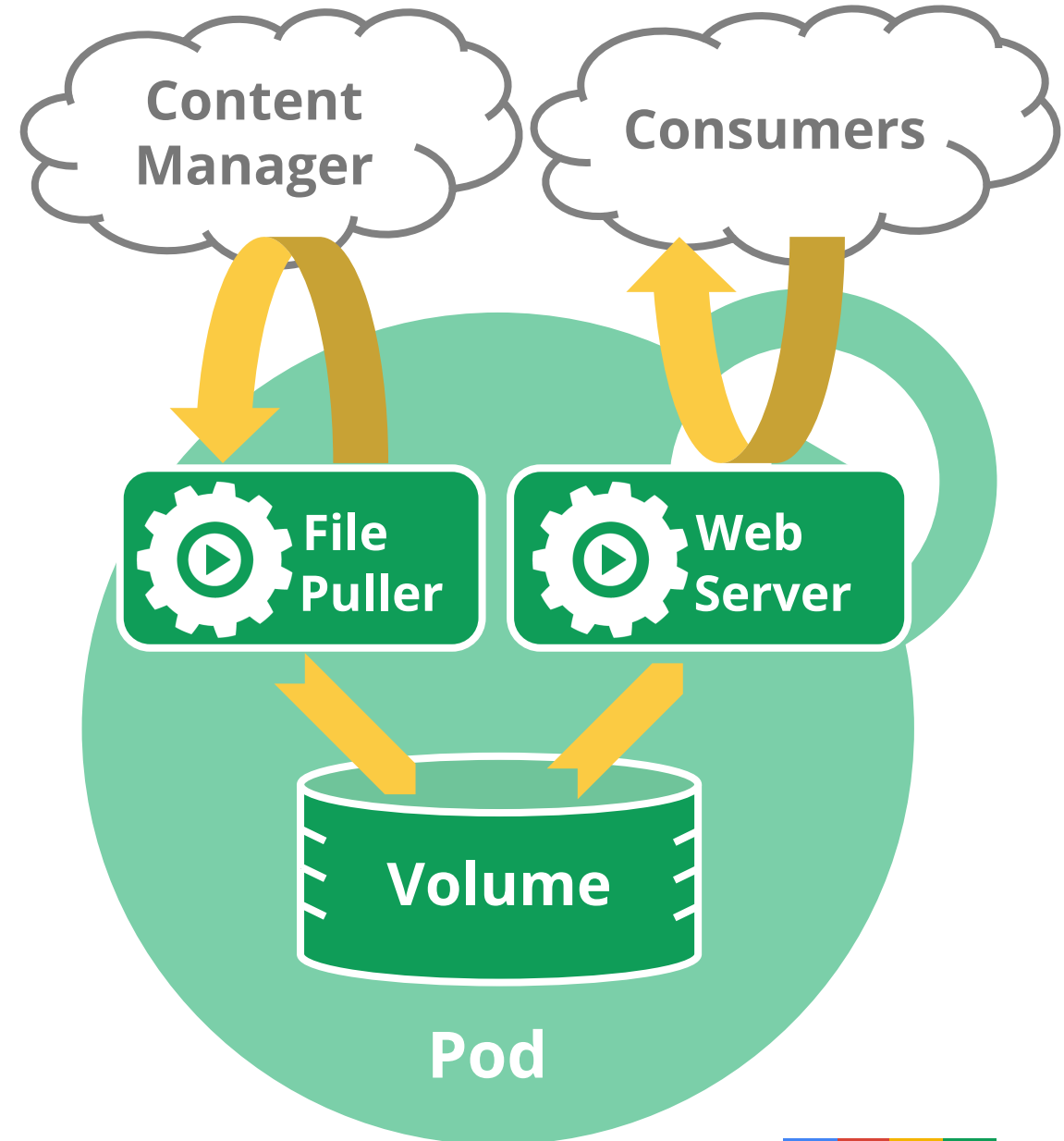
The atom of scheduling & placement

Shared namespace
- share IP address & localhost
- share IPC, etc.

Managed lifecycle
- bound to a node, restart in place
- can die, cannot be reborn with same ID

**Example: data puller & web server**

**10.1.1.2**

# Pods

**Small group** of containers & volumes

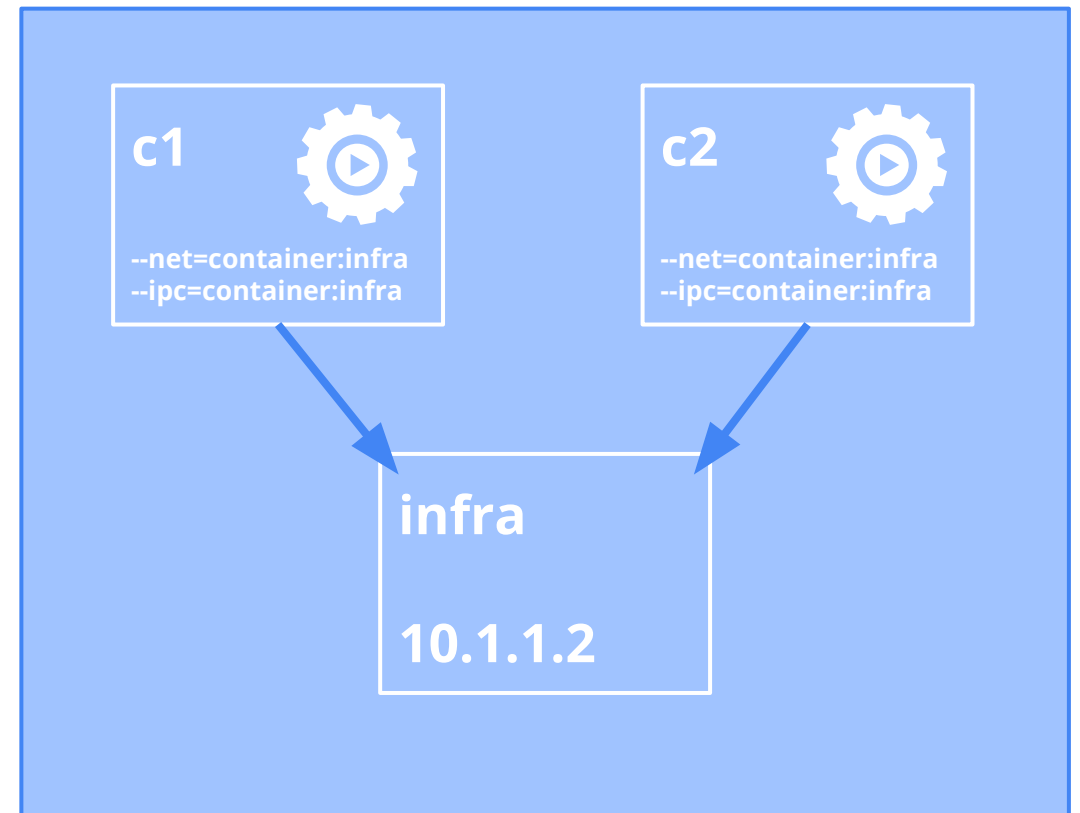**Tightly** coupled

The atom of scheduling & placement

**Shared namespace**
- **share IP address & localhost**
- **share IPC, etc.**

Managed lifecycle
- bound to a node, restart in place
- can die, cannot be reborn with same ID

**Example: data puller & web server**



c1
--net=container:infra
--ipc=container:infra

c2
--net=container:infra
--ipc=container:infra

**infra**

**10.1.1.2**

# Services

# Services

A group of pods that **work together**
- grouped by a selector

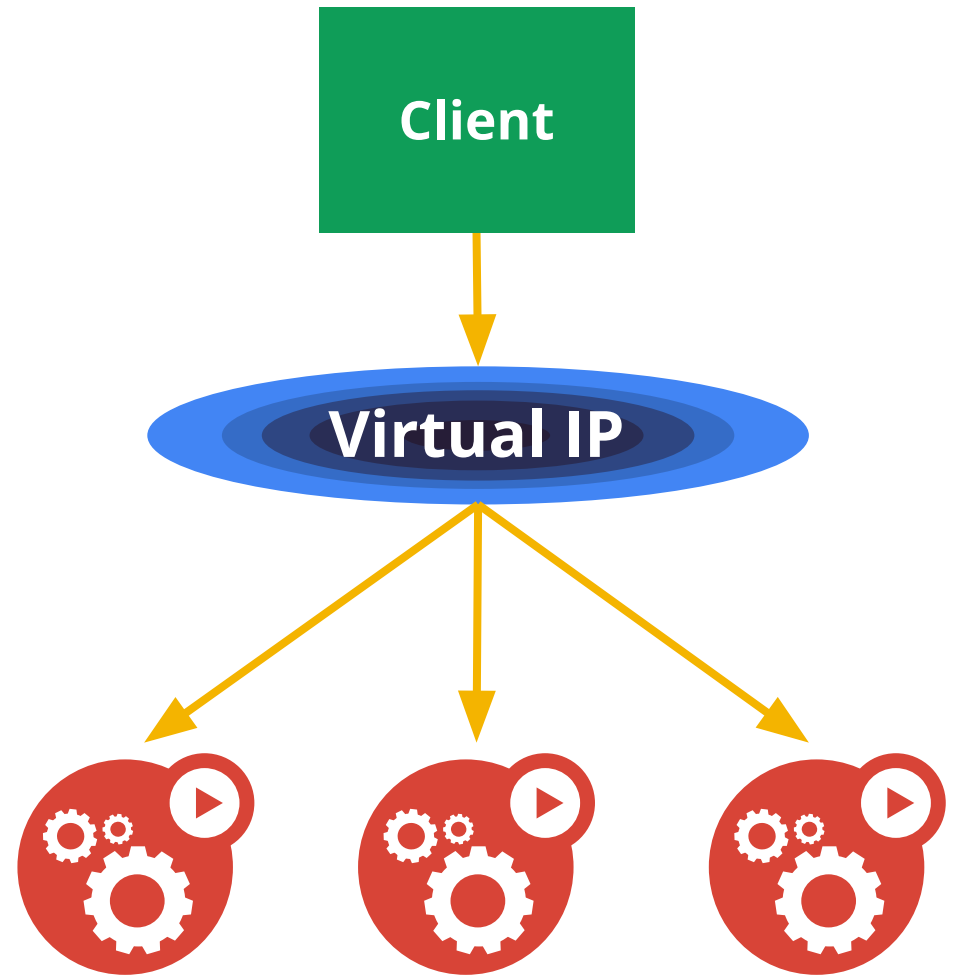Defines access policy
- "load balanced" or "headless"

Gets a stable **virtual IP** and port
- sometimes called the service *portal*
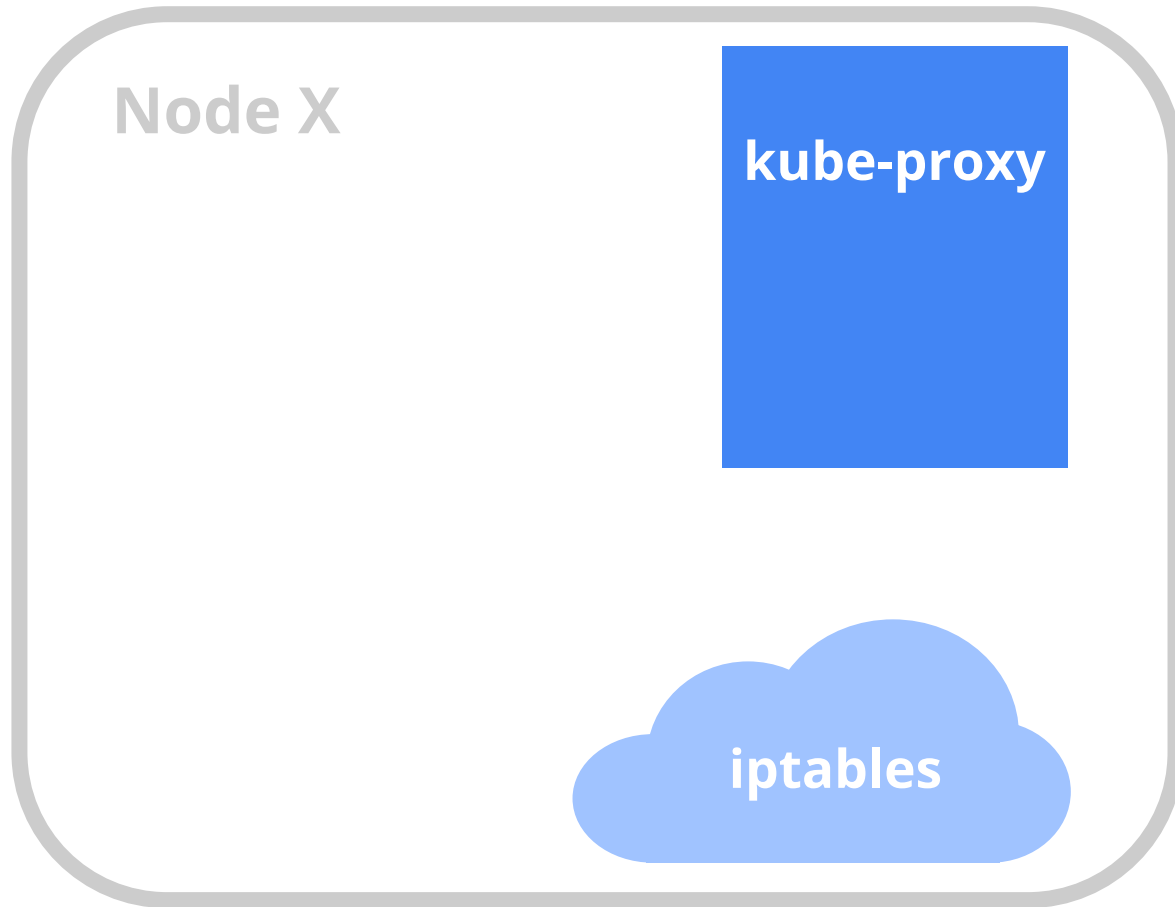- also a DNS name

VIP is managed by *kube-proxy*
- watches all services
- updates iptables when backends change

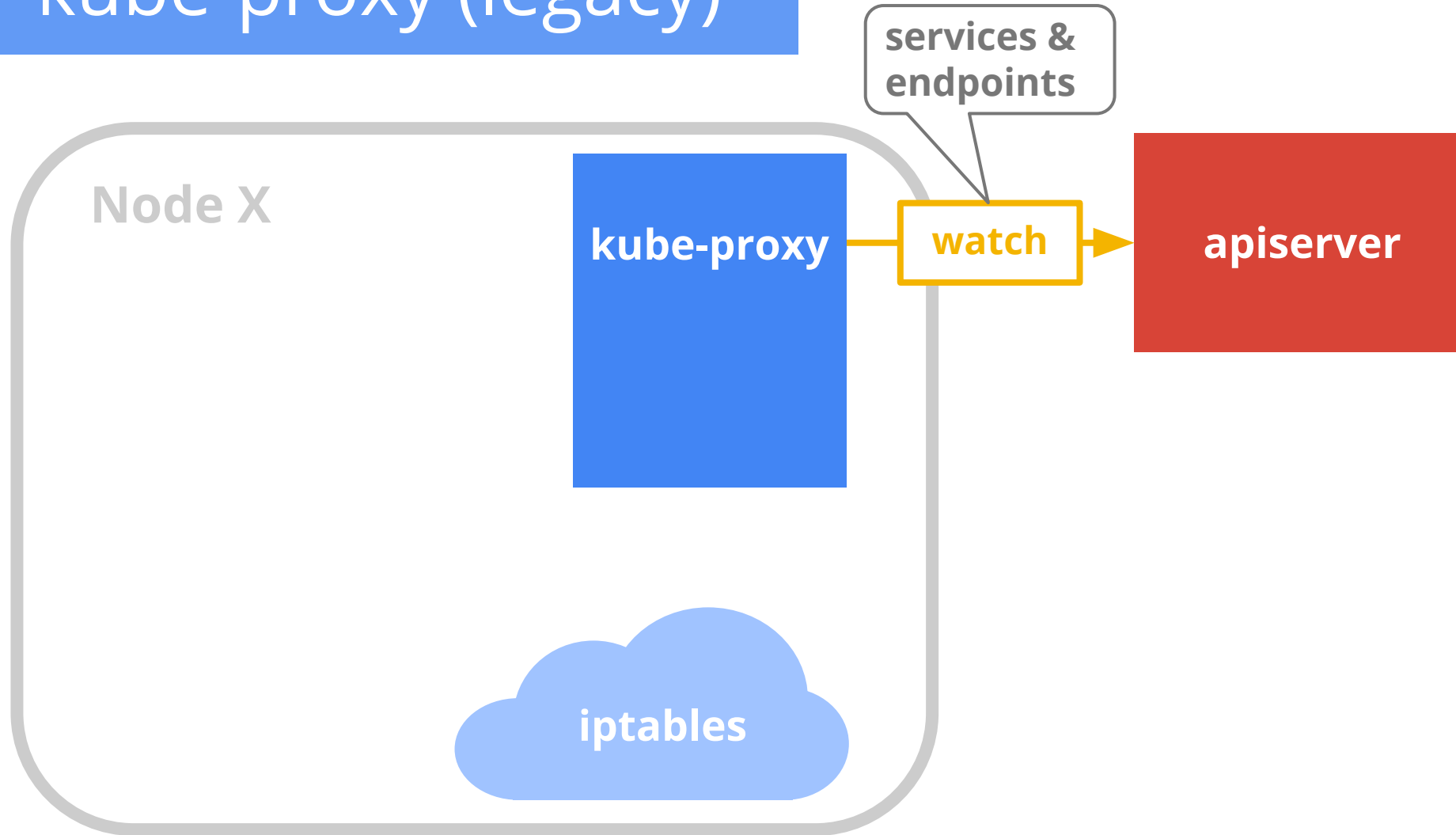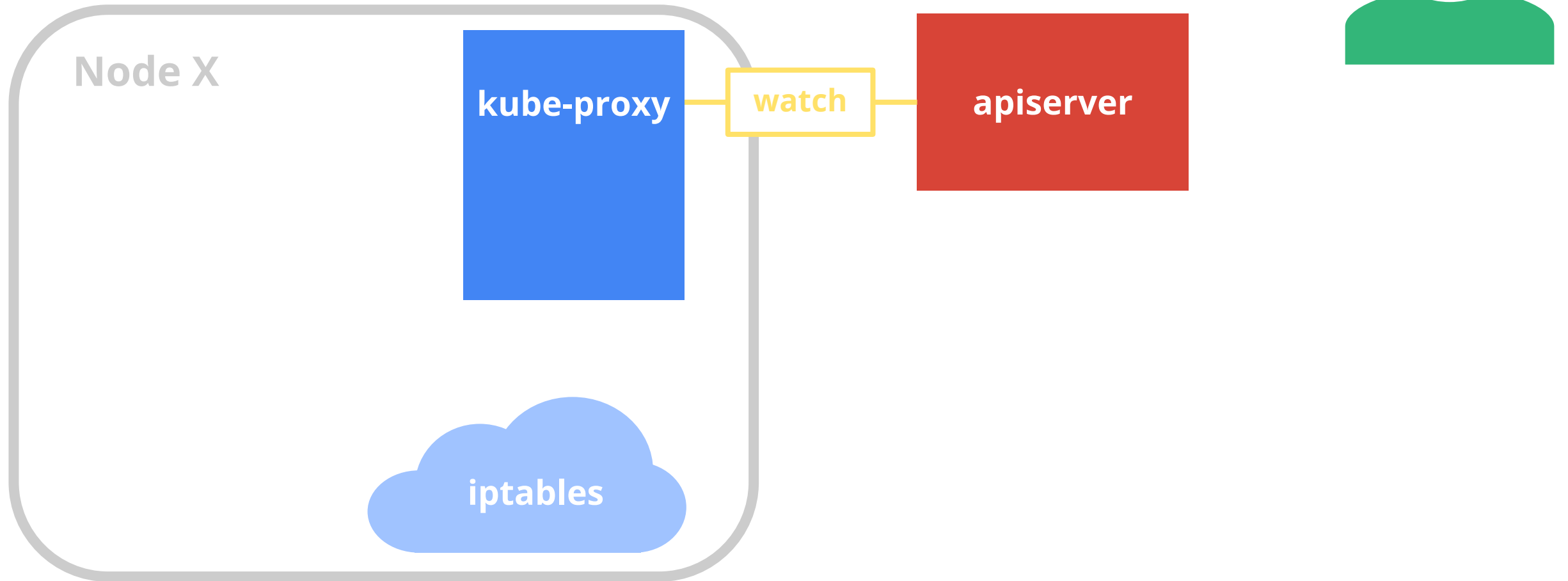Hides complexity - ideal for non-native apps
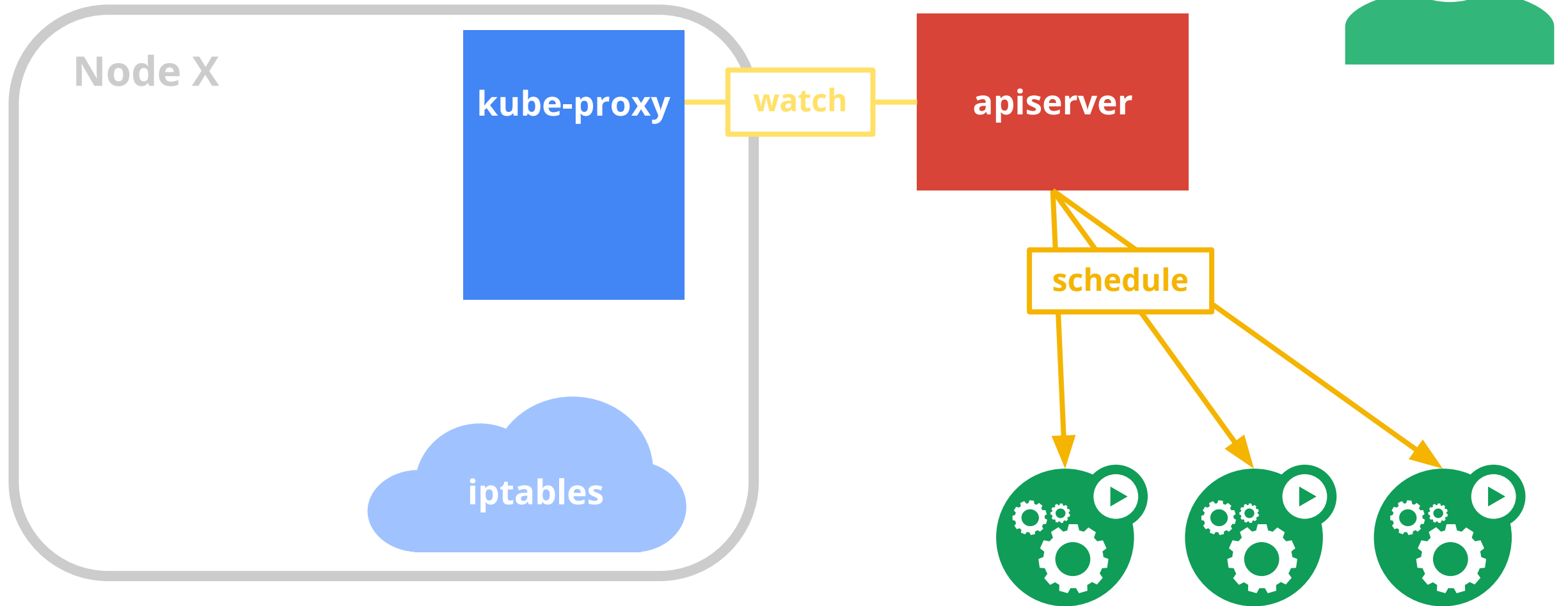
Client

Virtual IP

# kube-proxy

# kube-proxy (legacy)

Node X

**kube-proxy**

**iptables**

**apiserver**

Google Cloud Platform

# kube-proxy (legacy)



**Node X**

**kube-proxy**

**watch**

services & endpoints

**apiserver**

**iptables**

# kube-proxy (legacy)

kubectl run ...

**Node X**

**kube-proxy** — watch — **apiserver**

**iptables**

Google Cloud Platform

# kube-proxy (legacy)



Node X

**kube-proxy** — **watch** — **apiserver**

**iptables**

**schedule**

# kube-proxy (legacy)

**Node X**

**kube-proxy**

listen

watch

**apiserver**

**iptables**

# kube-proxy (legacy)

**Node X**

**kube-proxy**

**listen**

**watch**

**apiserver**

**iptables**

# kube-proxy (legacy)

**Node X**

**kube-proxy** —— watch —— **apiserver**

**configure**

**iptables**

kube-proxy (legacy)

# kube-proxy (legacy)



Node X

Client

VIP

kube-proxy

watch

apiserver

iptables

Google Cloud Platform

# kube-proxy (legacy)



Node X

Client

VIP

kube-proxy

iptables

watch

apiserver

# kube-proxy (legacy)

**Node X**

**kube-proxy** — **watch** — **apiserver**

**Client**

**VIP**

**iptables**

Google Cloud Platform

kube-proxy (legacy)

# kube-proxy (legacy)

Userspace proxy isn't ideal

Burns CPU copying bytes
- "Proxy" is just parallel copy loops.
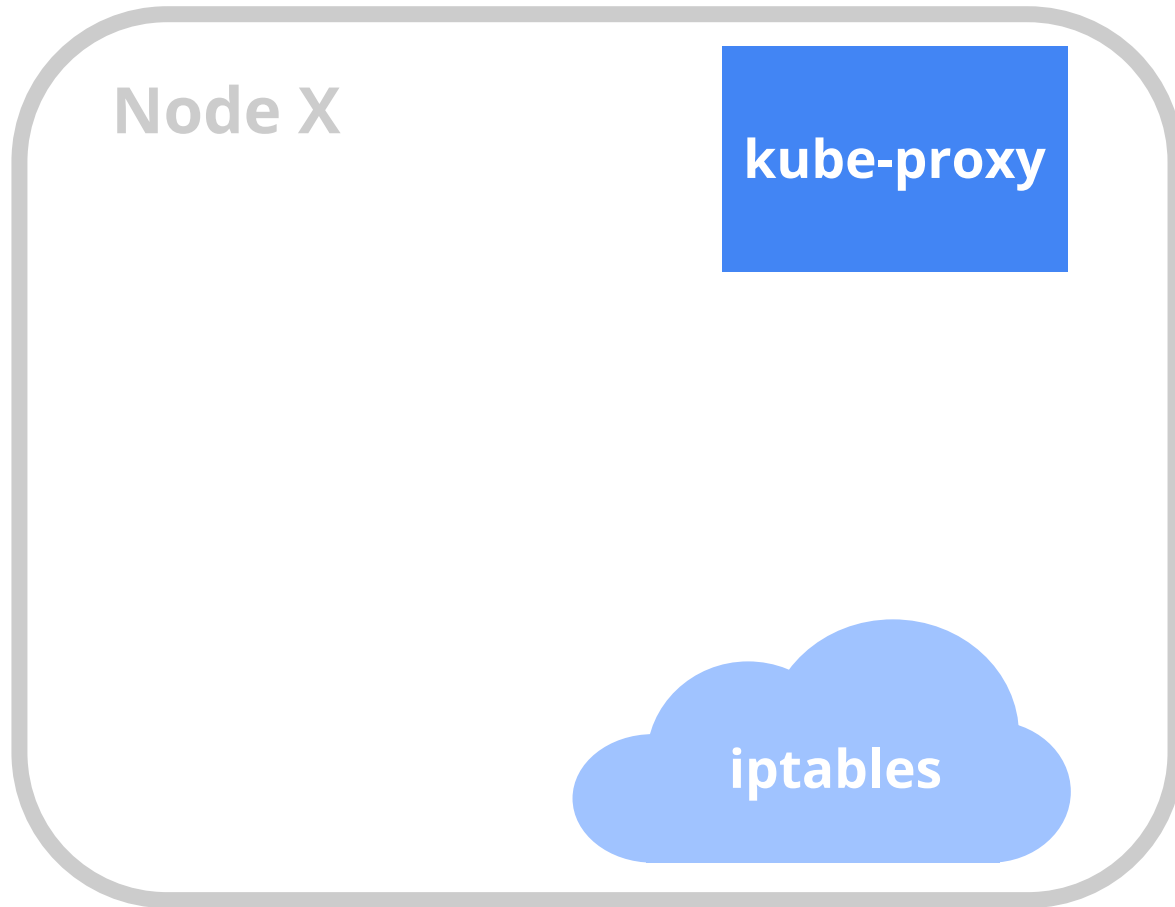
Loses source IP
- Everything looks like it's from the node IP.

Userspace TCP listening = higher latency

# iptables kube-proxy

# iptables kube-proxy

Node X

**kube-proxy**

**apiserver**

**iptables**

Google Cloud Platform

# iptables kube-proxy

services &
endpoints

| Node X | |
|---|---|
| **kube-proxy** | **watch** → **apiserver** |

**iptables**

Google Cloud Platform

# iptables kube-proxy

**Node X**

**kube-proxy** — watch — **apiserver**

**iptables**

**schedule**

# iptables kube-proxy

**Node X**

**kube-proxy** — watch — **apiserver**

configure

**iptables**

# iptables kube-proxy

Node X

kube-proxy — watch — apiserver

VIP
iptables

# iptables kube-proxy

# iptables kube-proxy

**Node X**

**kube-proxy** — watch — **apiserver**

**VIP**

**20**

**iptables**

# iptables kube-proxy



Node X

Client

kube-proxy — watch — apiserver

VIP

20

iptables

# iptables kube-proxy



**Node X**

**kube-proxy** — watch — **apiserver**

**Client**

**VIP**

20

**iptables**

Google Cloud Platform

# iptables kube-proxy

Node X

Client → VIP → **20** (die)

kube-proxy — watch — apiserver

iptables

# iptables kube-proxy

# iptables kube-proxy

## Mean Latency

`contrib/for-tests/netperf-tester --number=1000`



**iptables** kube-proxy

**legacy** kube-proxy

Mean Latency Microseconds

# Services

Services are just an abstraction
- Only requirement: route (and maybe load balance) a virtual IP to a set of backends.

Kube-proxy is an implementation
- Kube-proxy watches apiserver.
- iptables is re-configured on changes.

There could be other ways
- Userspace, iptables, IP Virtual Servers?

## Run SkyDNS as a pod in the cluster
- kube2sky bridges Kubernetes API -> SkyDNS
- Tell kubelets about it (static service IP)

## Strictly optional, but practically required
- LOTS of things depend on it
- Probably will become more integrated

## Or plug in your own!

```
kubernetes

kubernetes.default

kubernetes.default.svc.cluster.local

foo.my-namespace.svc.cluster.local
```

# DNS

Run SkyDNS as a pod in the cluster
* kube2sky bridges Kubernetes API -> SkyDNS
* Tell kubelets about it (static service IP)

Strictly optional, but practically required
* LOTS of things depend on it
* Probably will become more integrated

Or plug in your own!

kube-dns-qxin

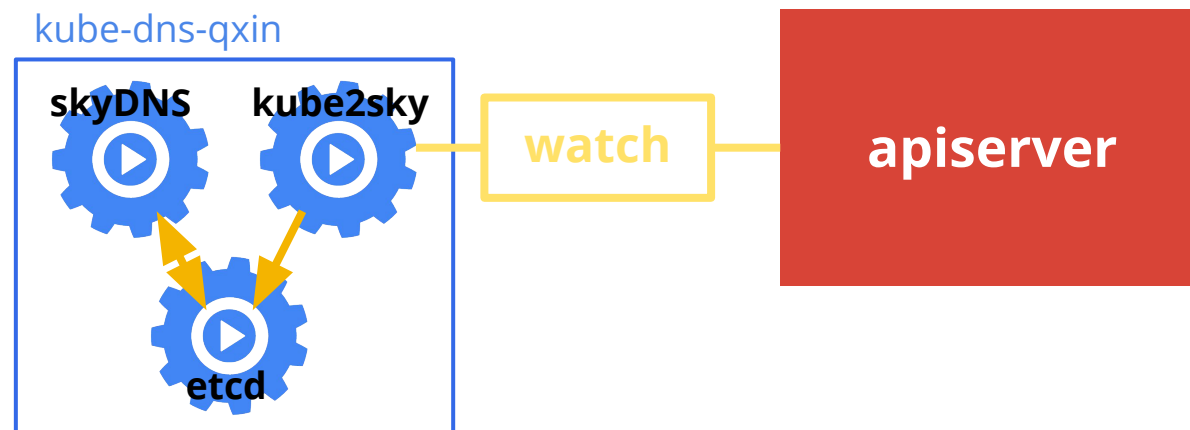skyDNS    kube2sky

etcd

watch

apiserver

# DNS

Run SkyDNS as a pod in the cluster
- kube2sky bridges Kubernetes API -> SkyDNS
- Tell kubelets about it (static service IP)

Strictly optional, but practically required
- LOTS of things depend on it
- Probably will become more integrated

Or plug in your own!

/etc/resolv.conf

nameserver 10.0.0.10
...

kube-dns-qxin

skyDNS    kube2sky

watch    apiserver

etcd

# DNS

Run SkyDNS as a pod in the cluster
- kube2sky bridges Kubernetes API -> SkyDNS
- Tell kubelets about it (static service IP)

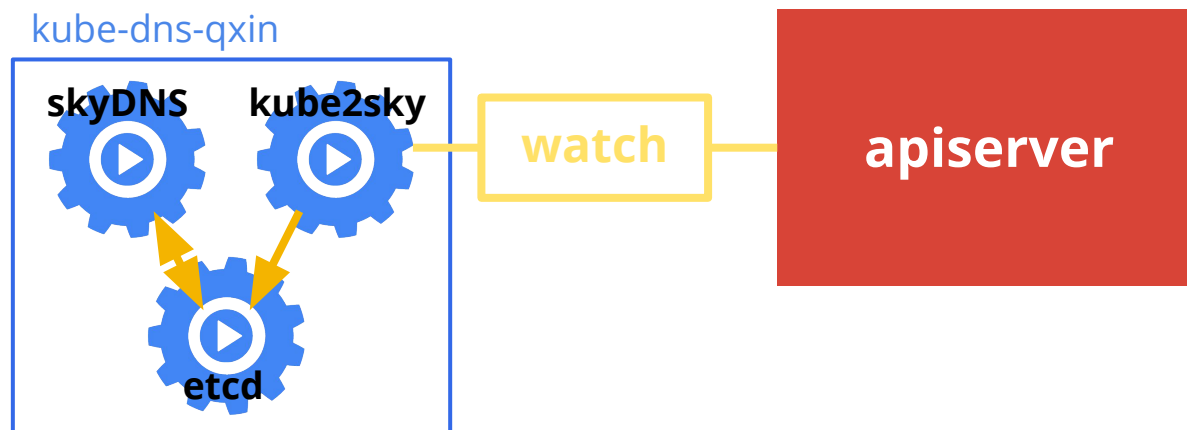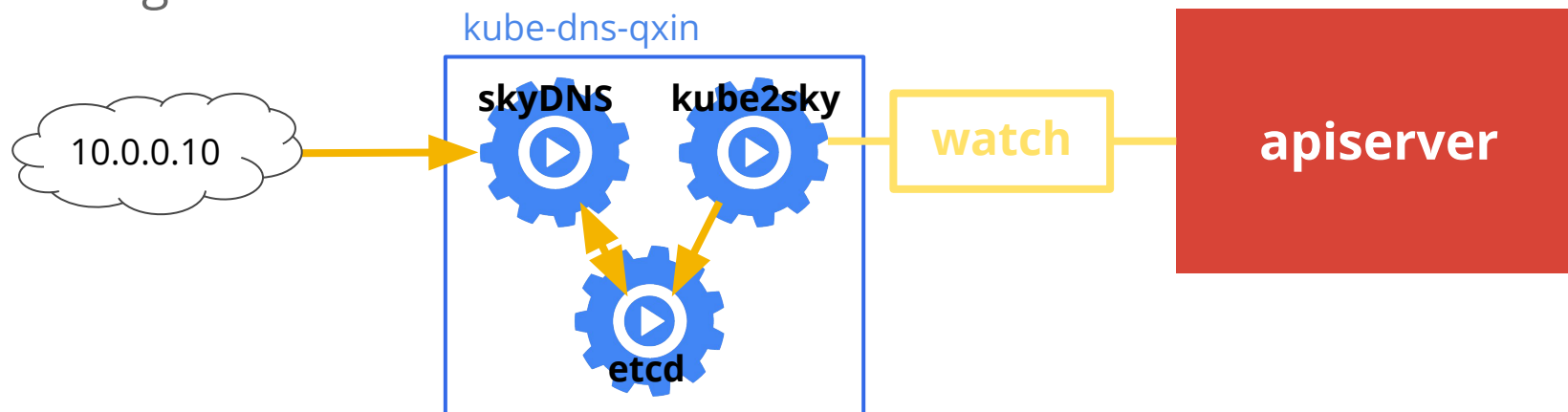Strictly optional, but practically required
- LOTS of things depend on it
- Probably will become more integrated

Or plug in your own!

/etc/resolv.conf

nameserver 10.0.0.10
...

kube-dns-qxin

skyDNS    kube2sky

10.0.0.10

etcd

watch

apiserver

# Putting it Together

What happens when I...

```
$ curl foo.my-namespace
```

# Putting it Together

What happens when I...

`$ curl foo.my-namespace`

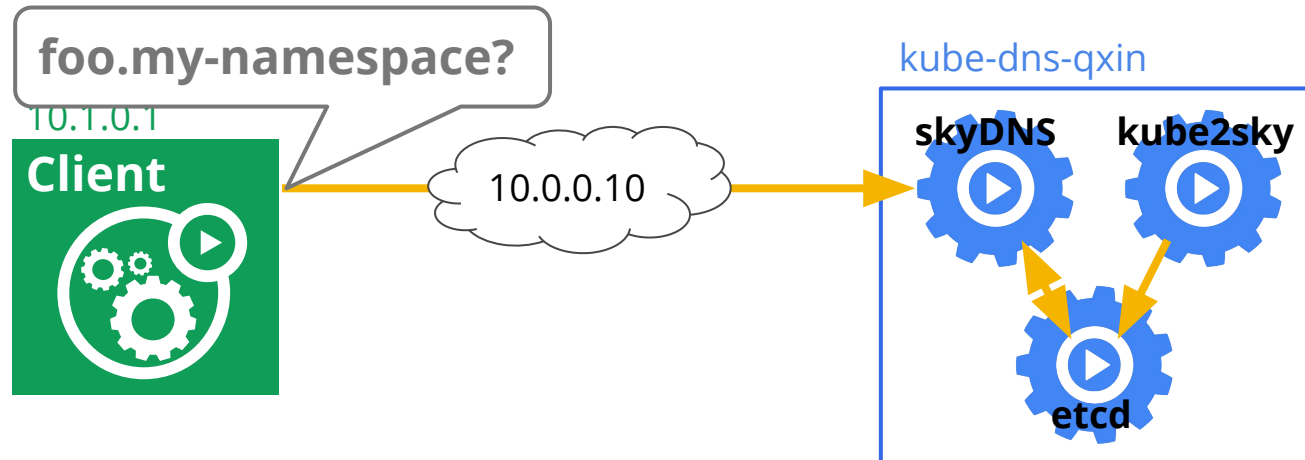/etc/resolv.conf

nameserver 10.0.0.10
...

10.1.0.1
**Client**

# Putting it Together
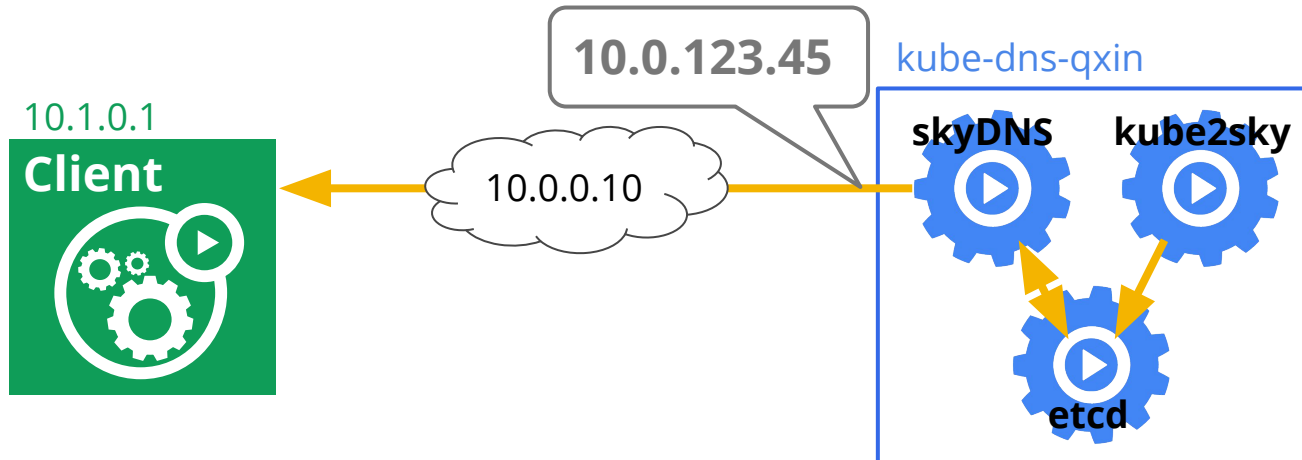
What happens when I...

```
$ curl foo.my-namespace
```

# Putting it Together

What happens when I...

```
$ curl foo.my-namespace
```

# Putting it Together

What happens when I...

```
$ curl foo.my-namespace
```
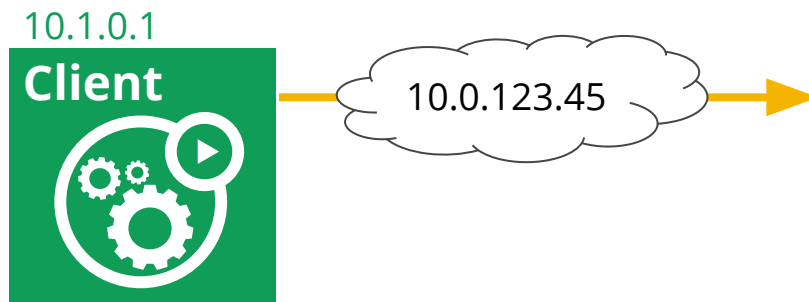
# Putting it Together
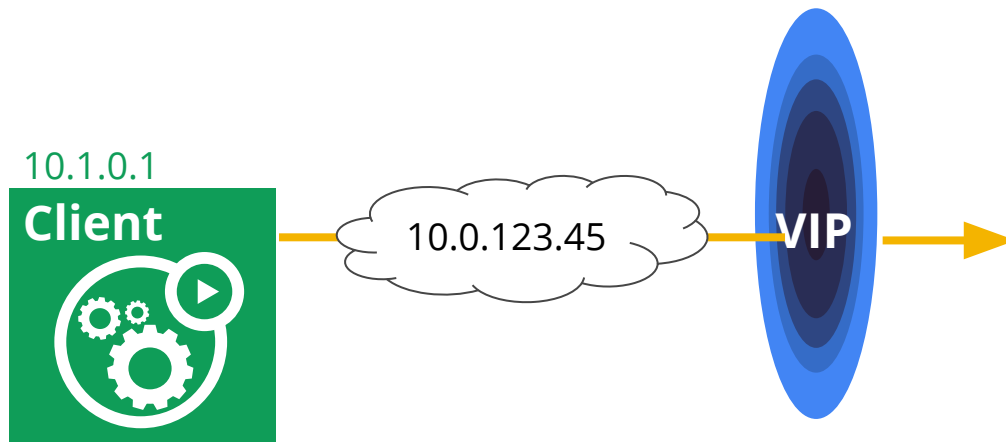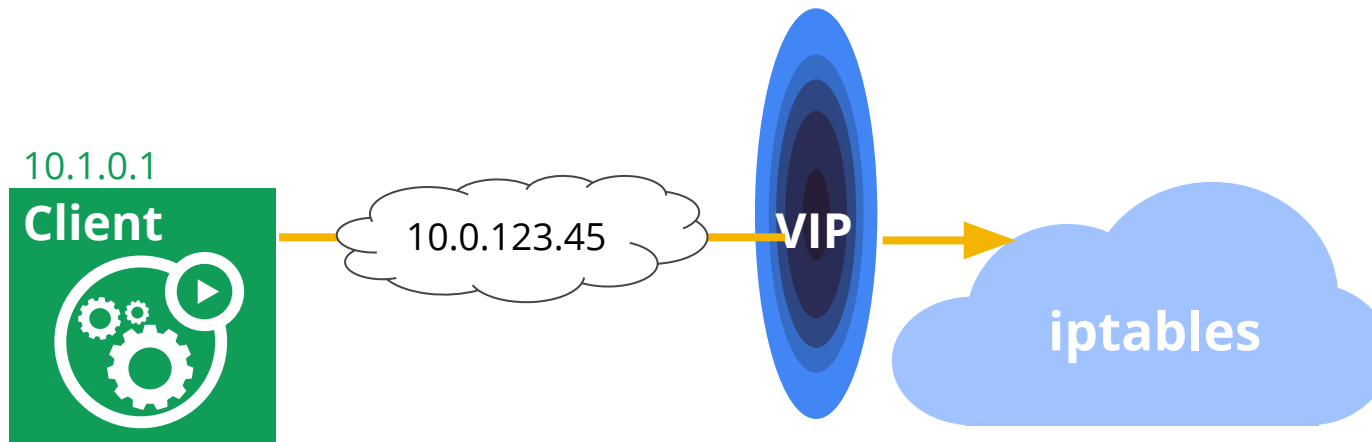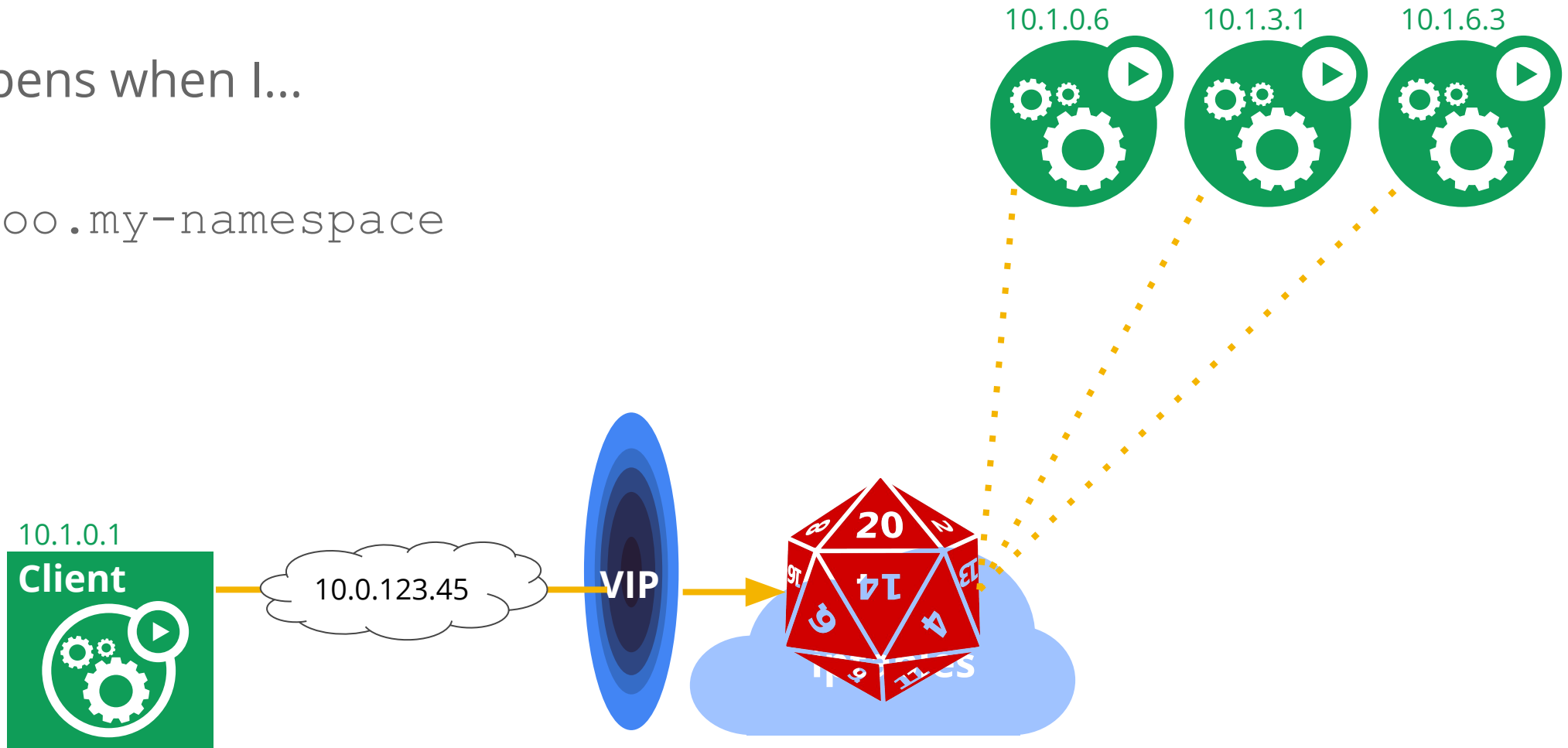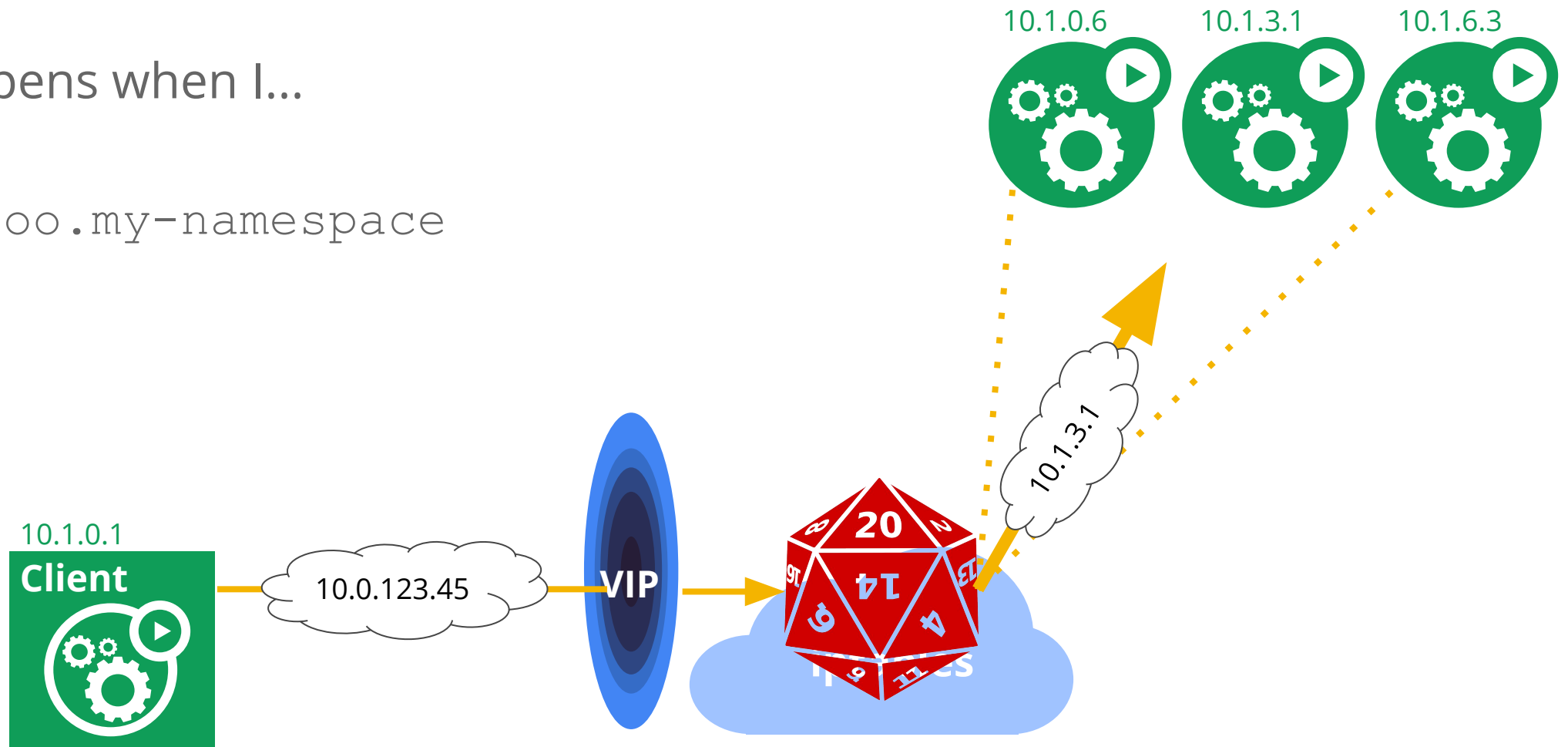
What happens when I...

```
$ curl foo.my-namespace
```

# Putting it Together

What happens when I...

```
$ curl foo.my-namespace
```

# Putting it Together

What happens when I...

```
$ curl foo.my-namespace
```

10.1.0.6    10.1.3.1    10.1.6.3

10.1.0.1
**Client**

10.0.123.45

**VIP**

**20**

# Putting it Together

What happens when I...

```
$ curl foo.my-namespace
```

10.1.0.6    10.1.3.1    10.1.6.3

10.1.3.1

10.1.0.1
**Client**

10.0.123.45

**VIP**

20

iptables

Google Cloud Platform

# Putting it Together

What happens when I...

`$ curl foo.my-namespace`

**10.1.0.6**   **10.1.3.1**   **10.1.6.3**

10.1.3.0/24 -> Node X

10.1.3.1

**10.1.0.1**

**Client**

10.0.123.45

**VIP**

**20**

# Putting it Together

What happens when I...

`$ curl foo.my-namespace`



10.1.0.6   10.1.3.1   10.1.6.3

10.1.0.1
**Client**

10.0.123.45

**VIP**

10.1.3.1

**20**

# Putting it Together

What happens when I...

$ curl foo.my-namespace

# What about external?

# External Services

Services IPs are only available **inside** the cluster

Need to receive traffic from "the outside world"

Builtin: Service "type"
- nodePort: expose on a port on every node
- loadBalancer: provision a cloud load-balancer
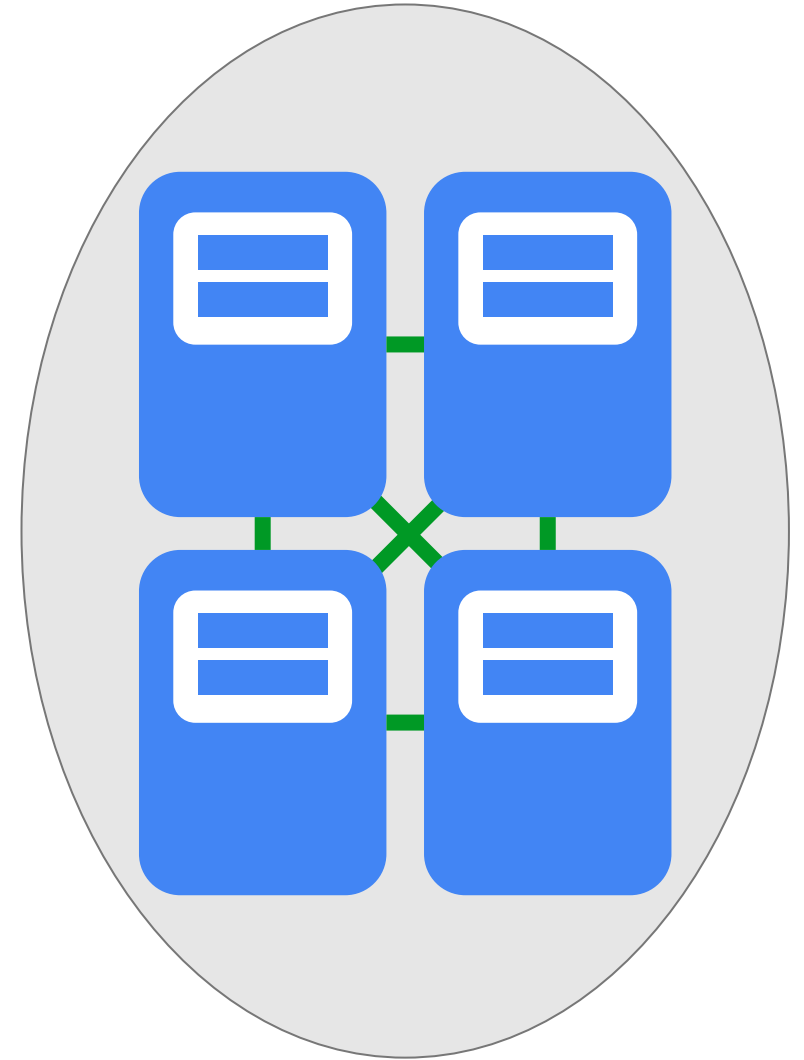
DiY load-balancer solutions
- socat (for nodePort remapping)
- haproxy
- nginx

# The Bleeding Edge

# Ingress (L7)

Services are assumed L3/L4

Lots of apps want HTTP/HTTPS

Ingress maps incoming traffic to backend services

- by HTTP host headers
- by HTTP URL paths

HAProxy and GCE implementations

No SSL yet

Status: **BETA** in Kubernetes v1.1

**Client**

**URL Map**

# Ingress (L7)

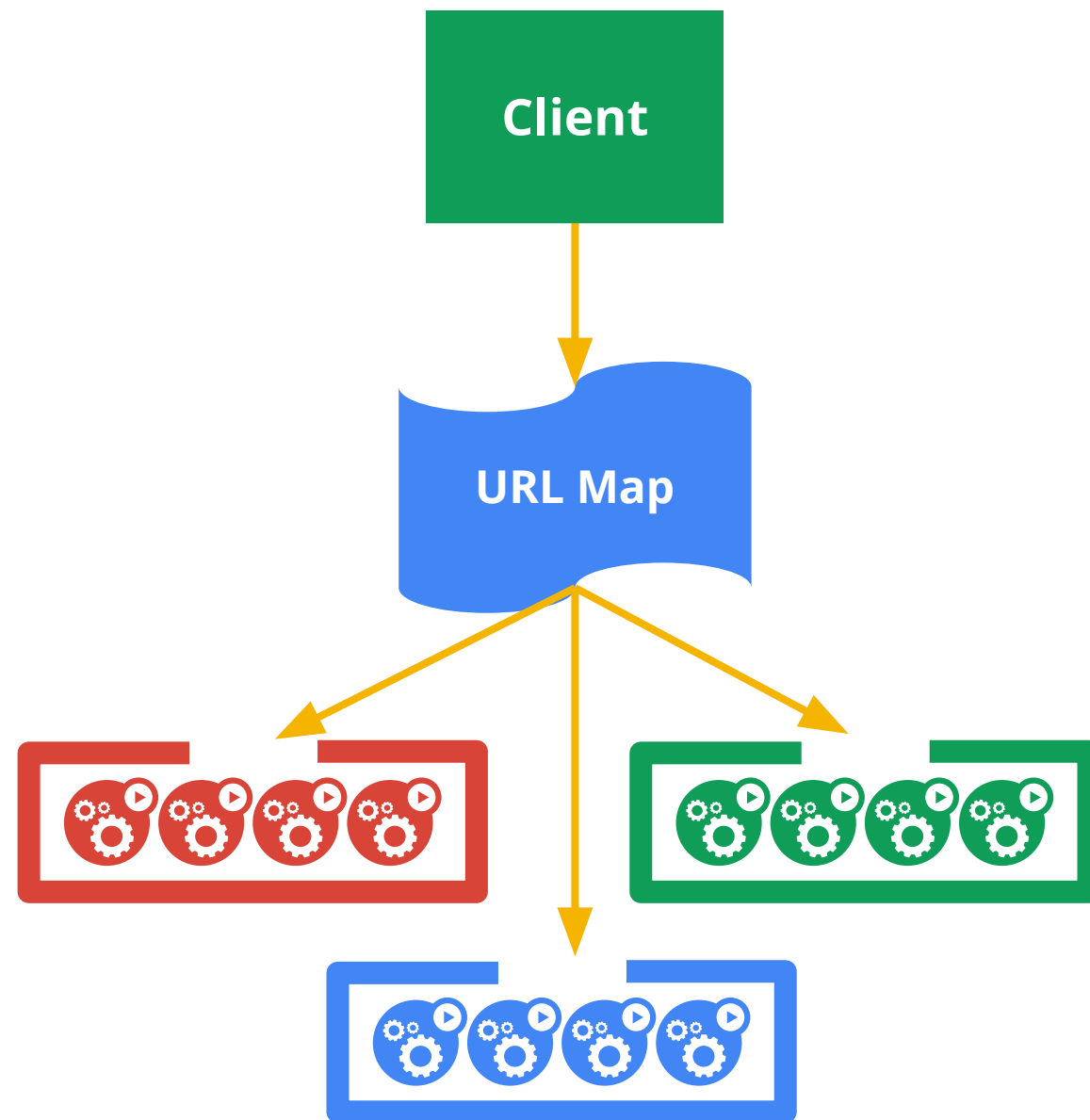Services are assumed L3/L4
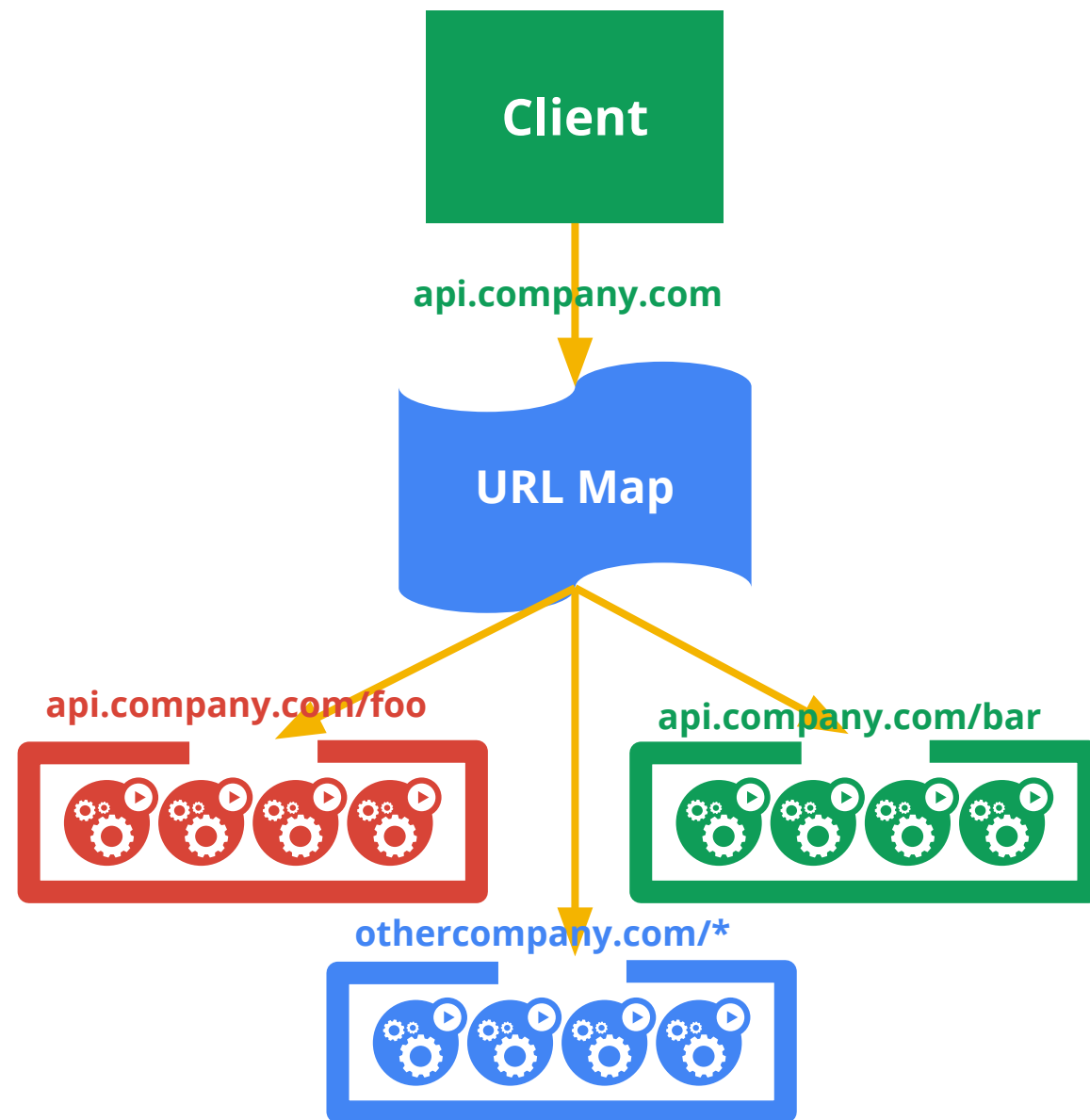
Lots of apps want HTTP/HTTPS

Ingress maps incoming traffic to backend services

- by HTTP host headers
- by HTTP URL paths

HAProxy and GCE implementations

No SSL yet

Status: **BETA** in Kubernetes v1.1

**Client**

api.company.com

**URL Map**

api.company.com/foo

api.company.com/bar

othercompany.com/*

# Network Plugins

# Network Plugins

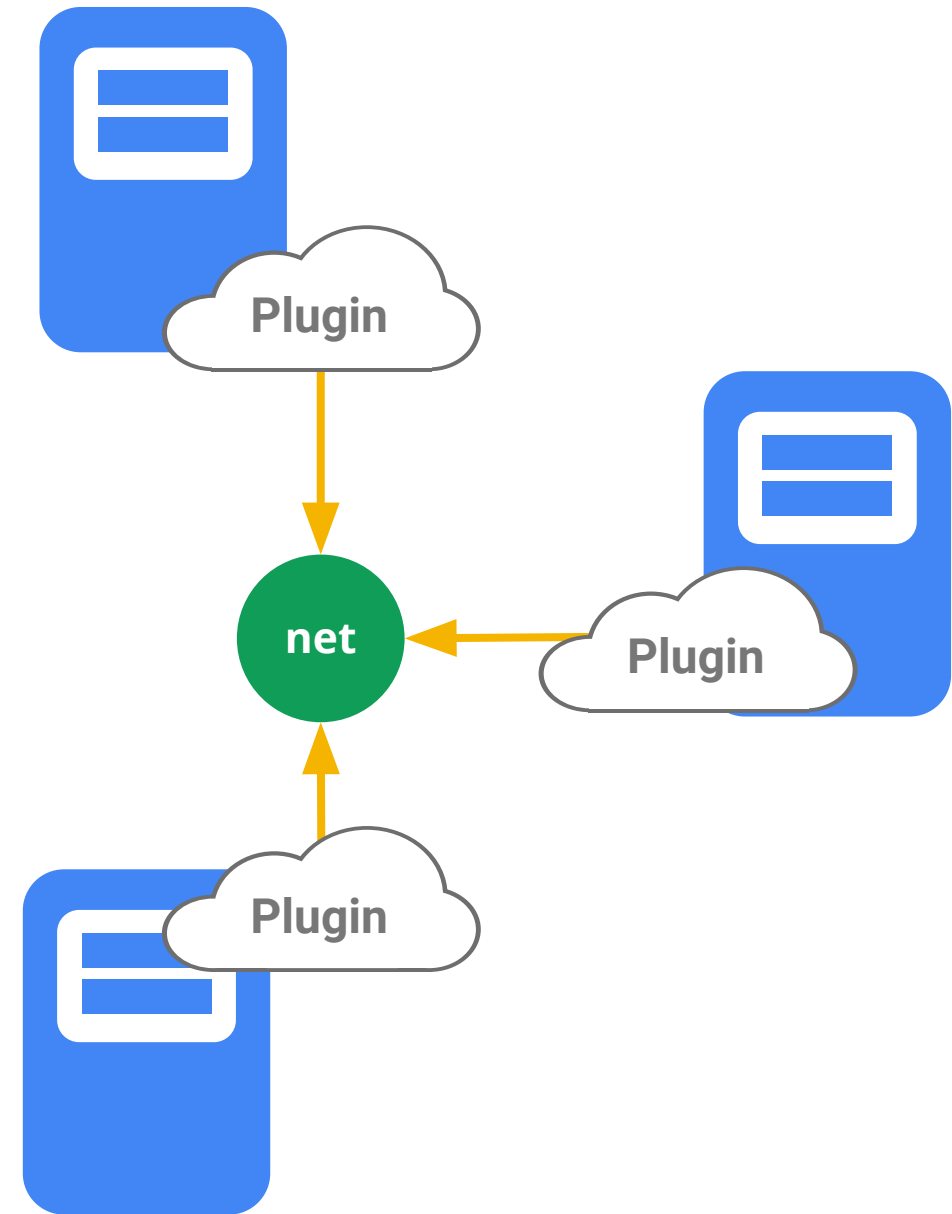Introduced in Kubernetes v1.0
- VERY experimental

Uses **CNI** (CoreOS) in v1.1
- Simple exec interface
- Not using Docker libnetwork
  - but can defer to Docker for networking

Cluster admins can customize their installs
- DHCP, MACVLAN, Flannel, custom

# Kubernetes is **Open**

- open community
- open design
- open source
- open to ideas

# Networking is **Hard**

- help guide us!

http://kubernetes.io
https://github.com/kubernetes/kubernetes

slack: kubernetes                              twitter: @kubernetesio