



KubeCon



CloudNativeCon

North America 2019





KubeCon



CloudNativeCon

North America 2019

Living with the pathology of the cloud: How AWS runs lots of clusters

Micah Hausler

Sr System Development Engineer, AWS

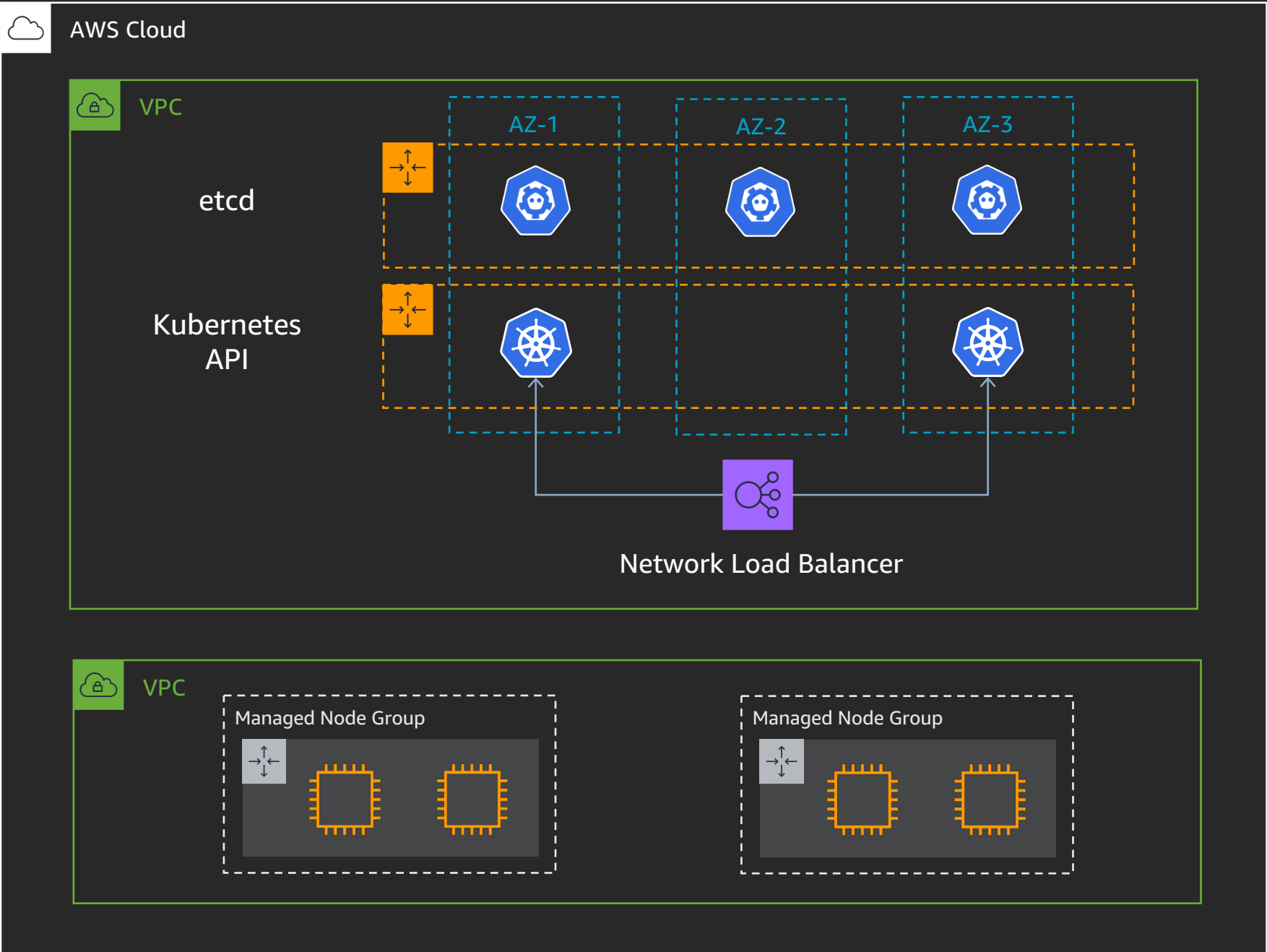


Amazon EKS

Service Priorities

Security

Operational Reliability



EKS Managed Control Plane

Customer Owned Data Plane

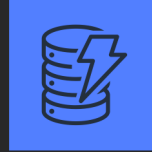
AWS Building Blocks



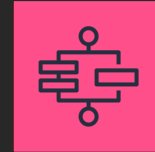
Amazon API Gateway



AWS Lambda



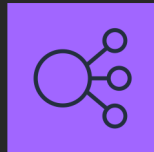
Amazon DynamoDB



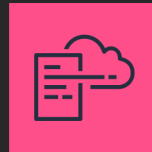
AWS Step Functions



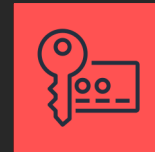
Amazon EC2



Elastic Load Balancing



AWS CloudFormation



AWS Key Management Service

Cell based architecture

Independent silos of operation

Can be logical or physical

Horizontal sharding

Scale out, not scale up

Numerous benefits

Reduced blast radius

Higher scalability

Safer deployments

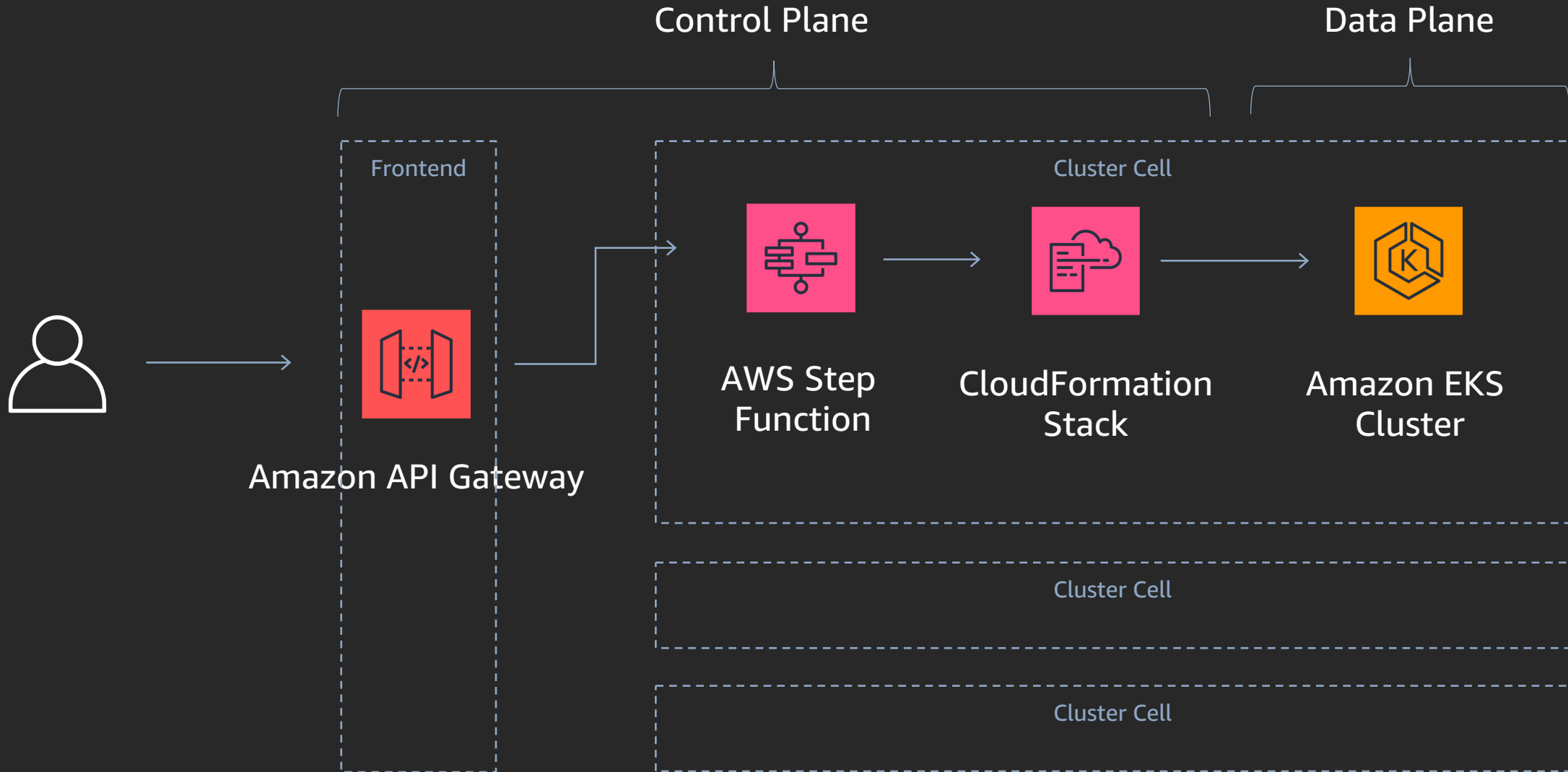
Cell based architecture

Tradeoffs

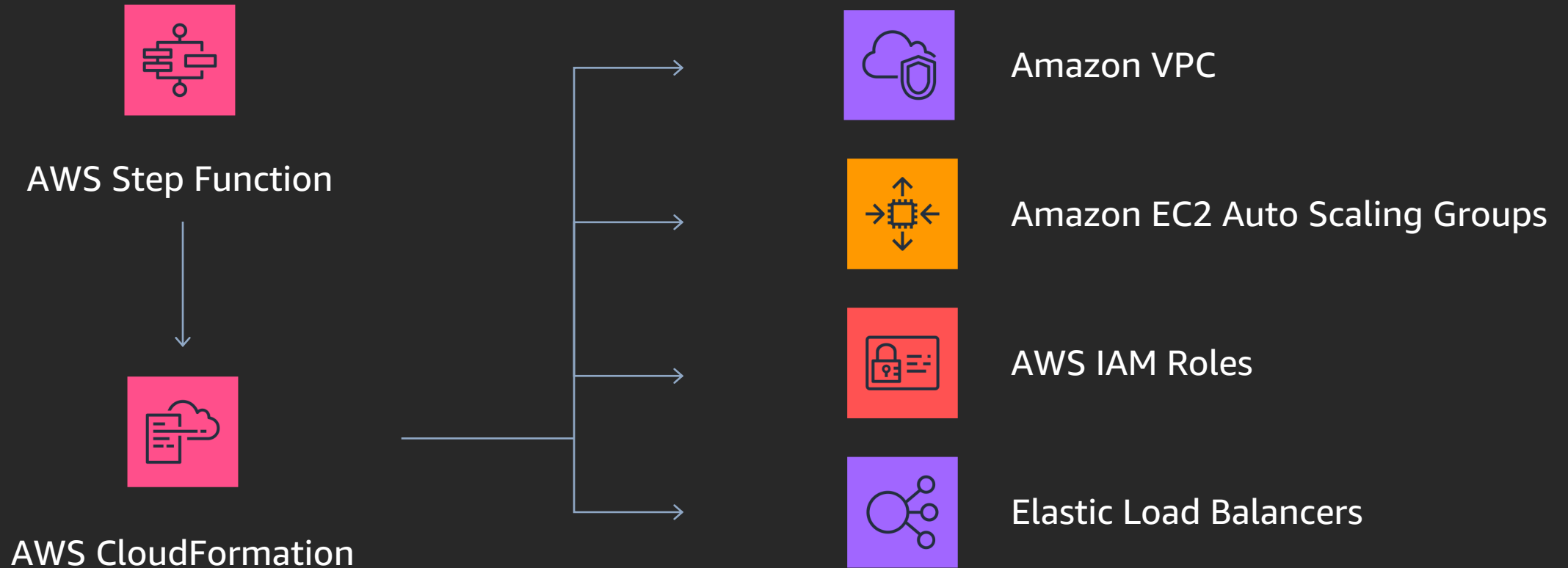
- Increased complexity

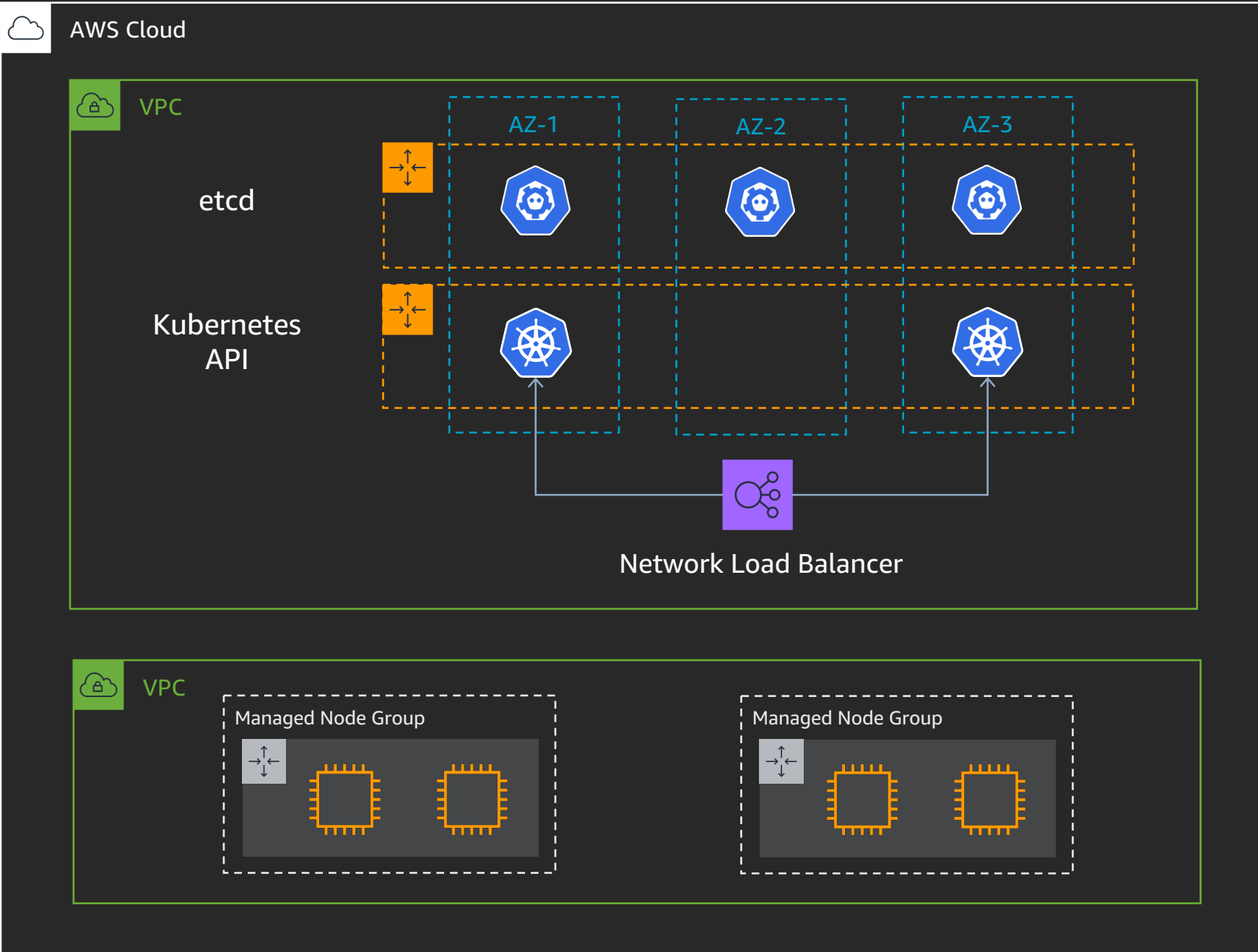
- Necessitates up front investment in tooling

How is a cluster created?



How is a cluster created?





EKS Managed Control Plane

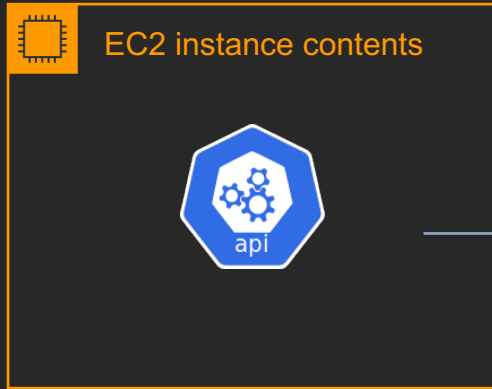
Customer Owned Data Plane

Failure Stories

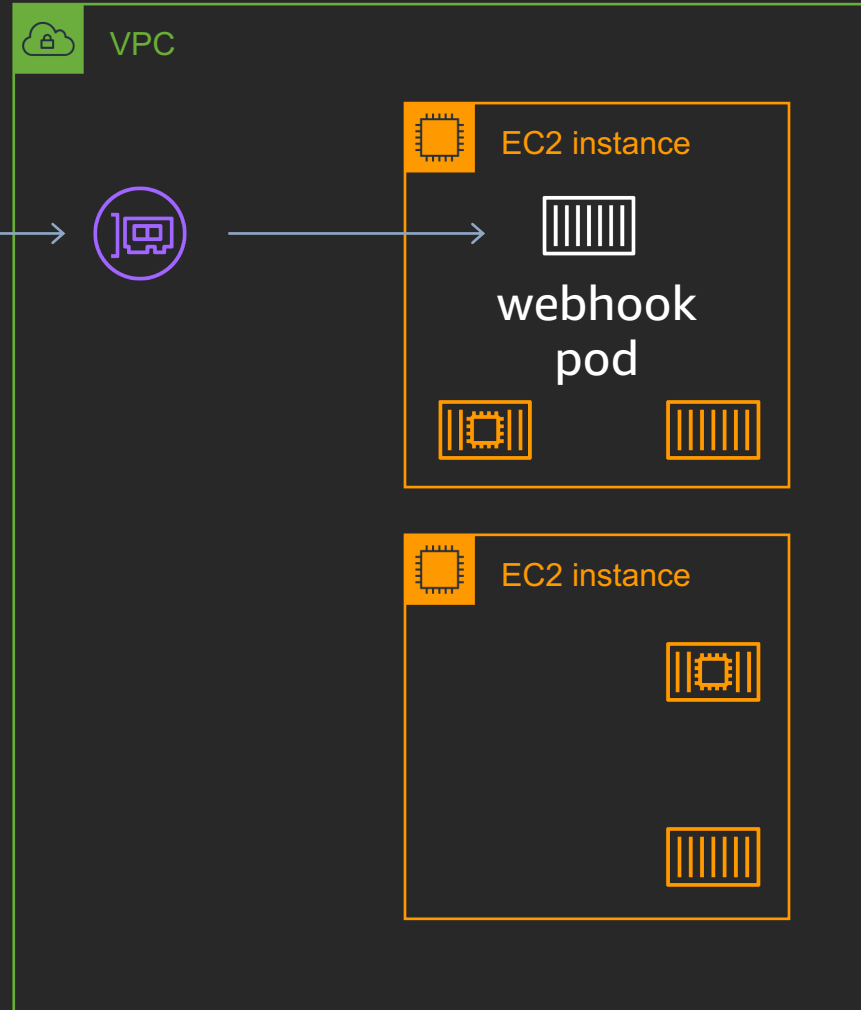
kube-apiserver fails to connect to new webhook pod

- Customer is running an mutating admission webhook in a pod
- EC2 node the pod is on is terminated (no FIN)
- New webhook pod comes up almost immediately
- kube-apiserver doesn't reconnect to new pod for 15 minutes

EKS managed kube-apiserver

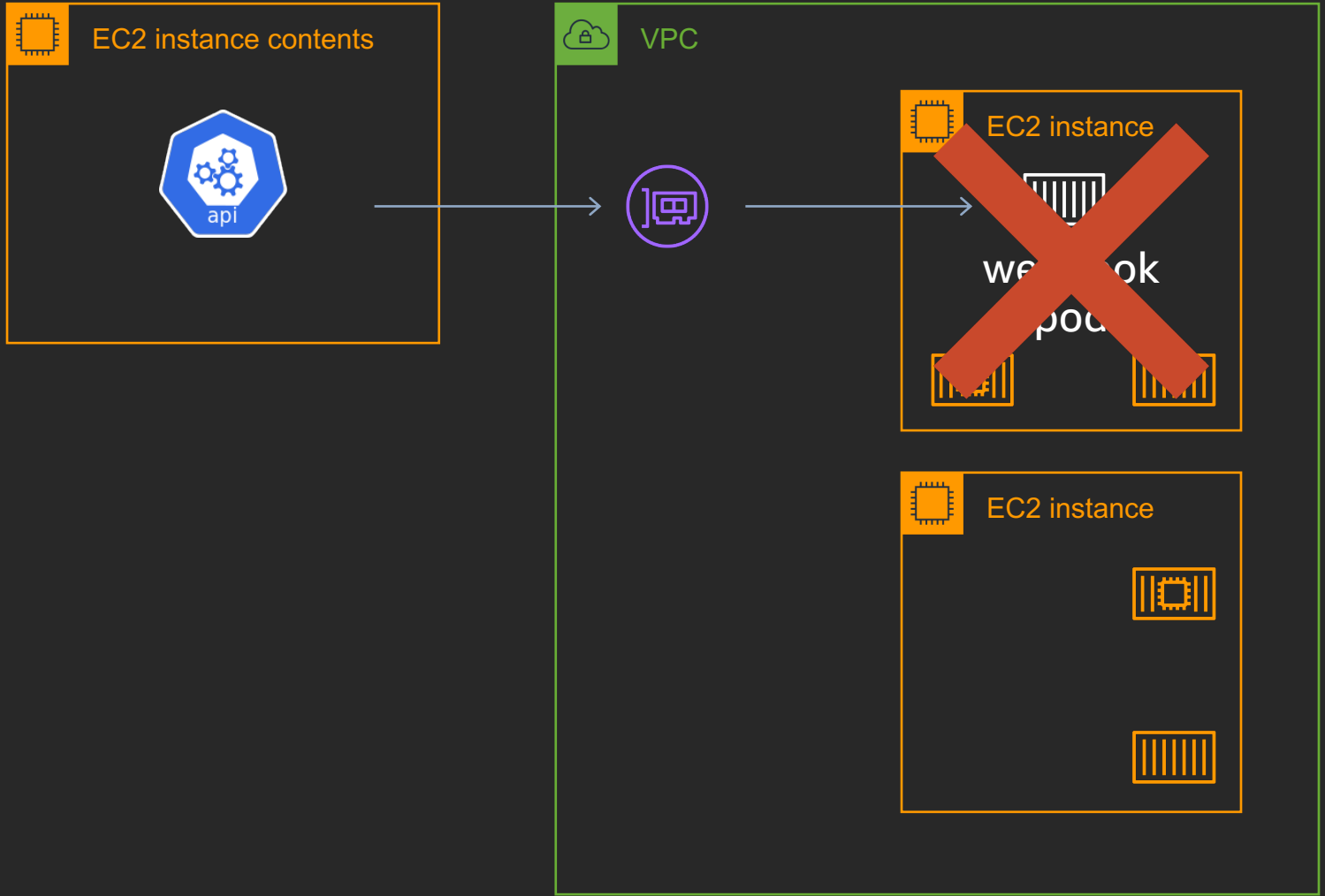


Customer VPC



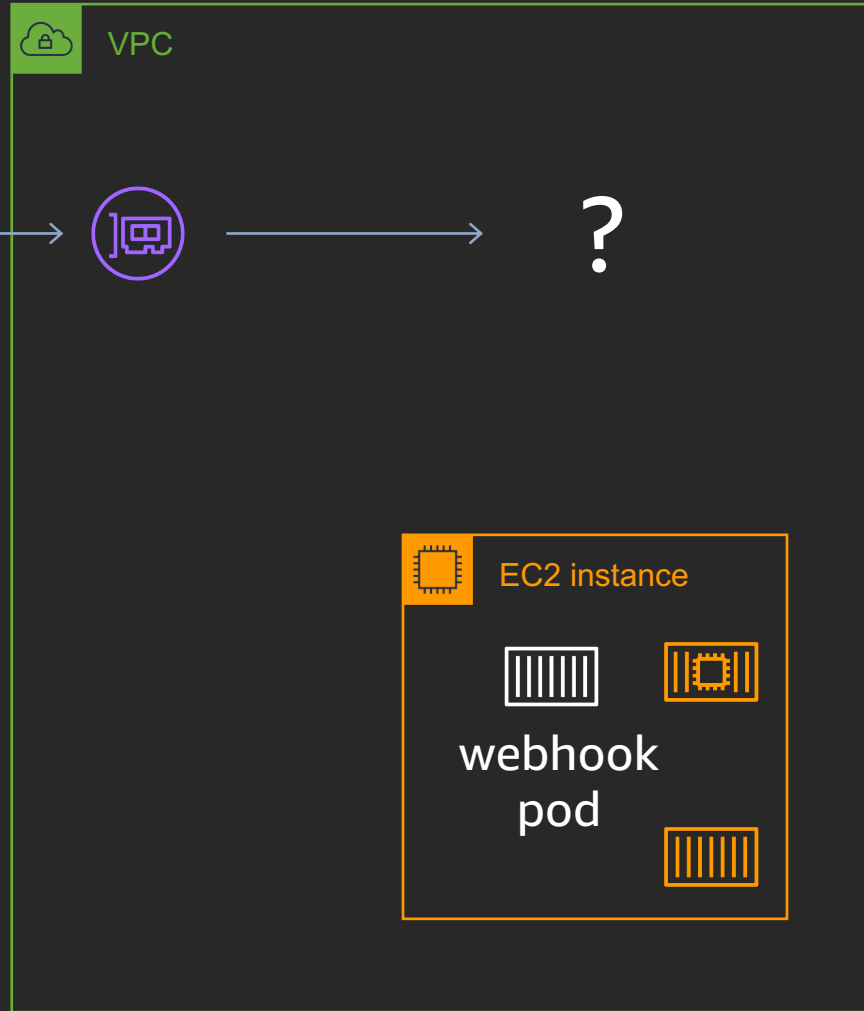
EKS managed kube-apiserver

Customer VPC



EKS managed kube-apiserver

Customer VPC




```
$ netstat -n
```

```
Active Internet connections (w/o servers)
```

Proto	Recv-Q	Send-Q	Local Address	Foreign Address	State
tcp	0	0	10.19.0.12:59824	10.20.0.13:443	ESTABLISHED

After some digging around

`tcp_retries2` - INTEGER

This value influences the timeout of an alive TCP connection, when RTO retransmissions remain unacknowledged.

Given a value of N, a hypothetical TCP connection following exponential backoff with an initial RTO of `TCP_RTO_MIN` would retransmit N times before killing the connection at the (N+1)th RTO.

The default value of 15 yields a hypothetical timeout of 924.6 seconds and is a lower bound for the effective timeout.

TCP will effectively time out at the first RTO which exceeds the hypothetical timeout.

RFC 1122 recommends at least 100 seconds for the timeout, which corresponds to a value of at least 8.

kubernetes/kubernetes

- [#80313](#) - Admission webhooks affected by dead tcp connections ([# 65012](#), [# 75791](#))
- [golang/go# 31643](#)
- Go's net/http in HTTP2 doesn't implement PING frames



http - The Go Programming Lan



golang.org/pkg/net/http/

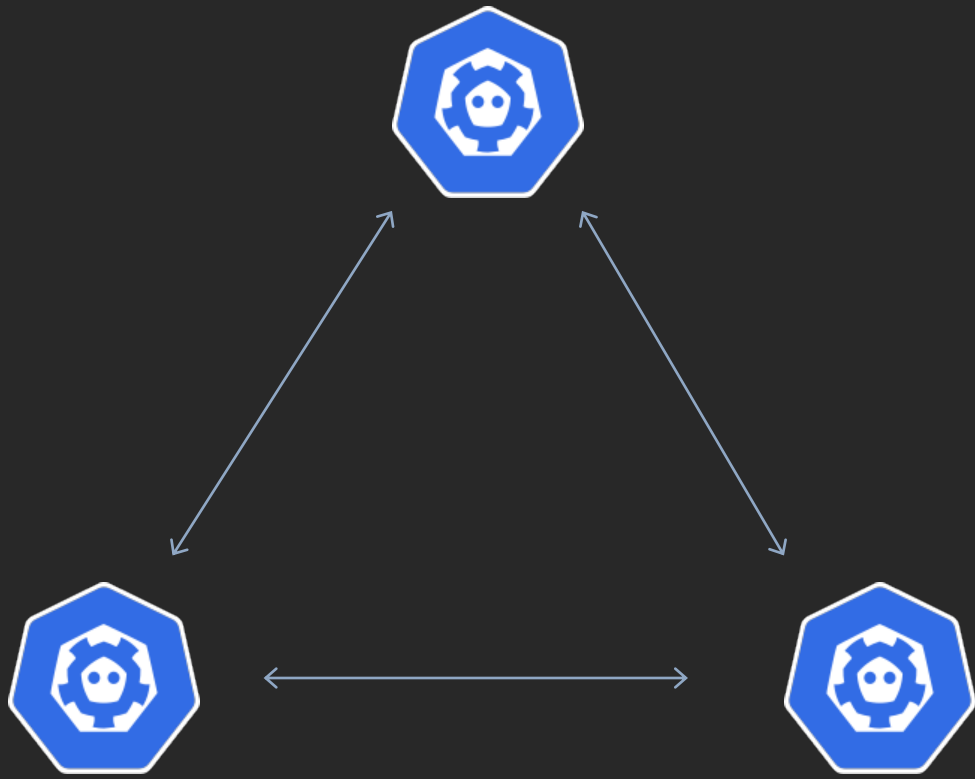


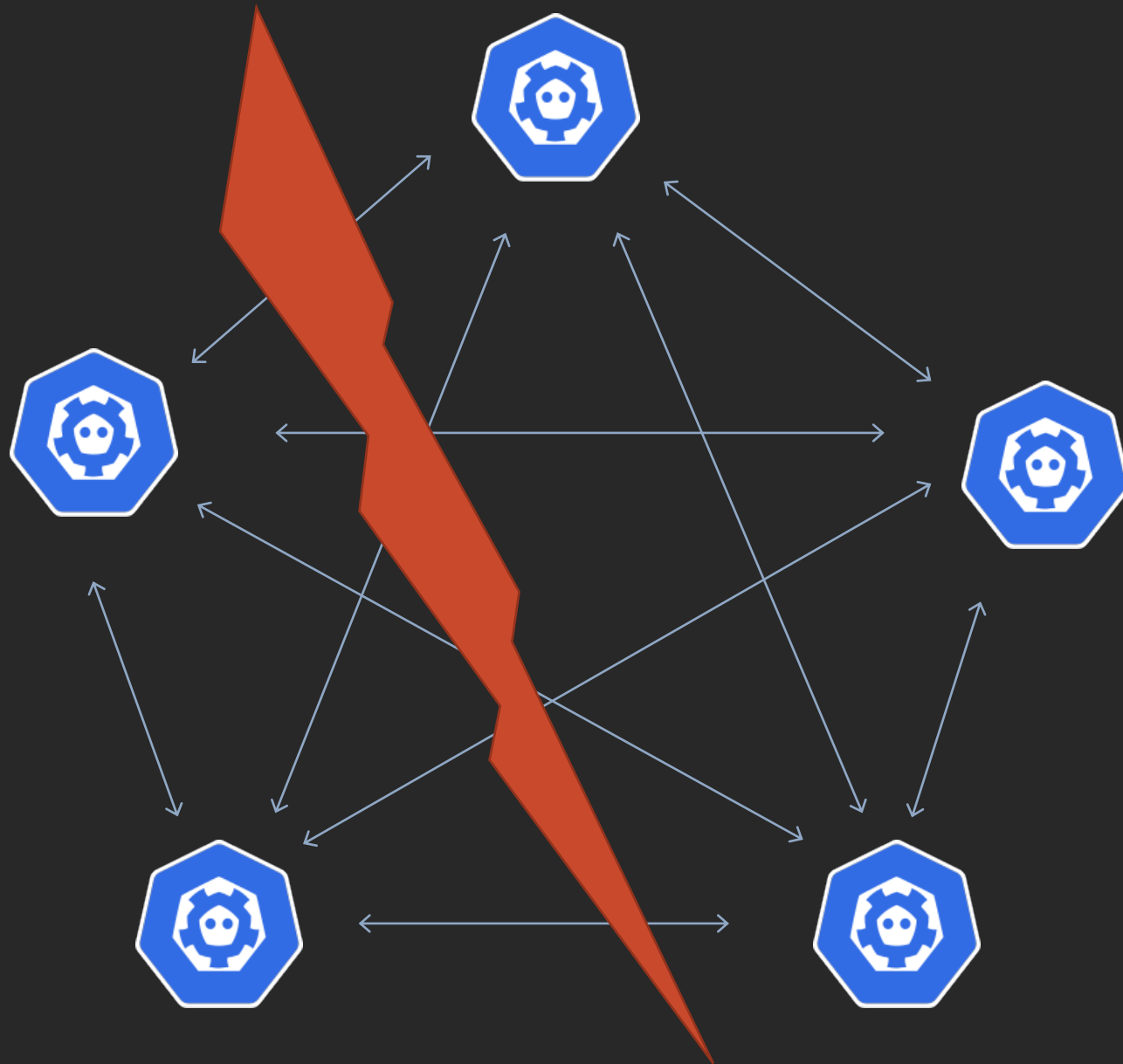
Starting with Go 1.6, the `http` package has transparent support for the HTTP/2 protocol when using HTTPS. Programs that must disable HTTP/2 can do so by setting `Transport.TLSNextProto` (for clients) or `Server.TLSNextProto` (for servers) to a non-nil, empty map. Alternatively, the following GODEBUG environment variables are currently supported:

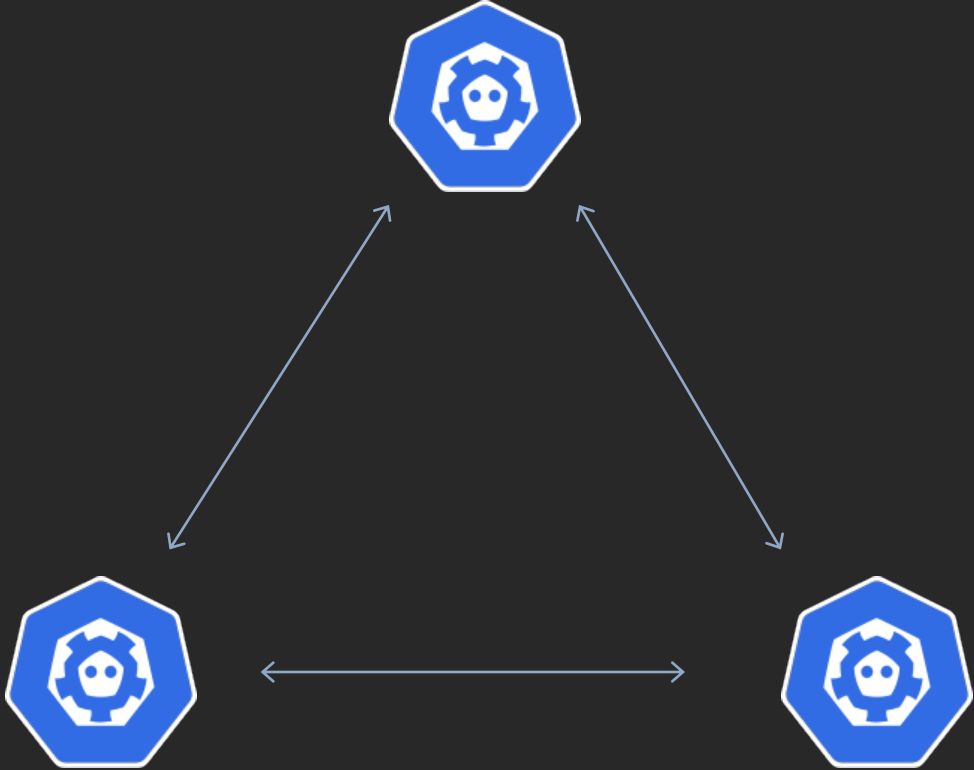
```
GODEBUG=http2client=0 # disable HTTP/2 client support
GODEBUG=http2server=0 # disable HTTP/2 server support
GODEBUG=http2debug=1  # enable verbose HTTP/2 debug logs
GODEBUG=http2debug=2  # ... even more verbose, with frame dumps
```

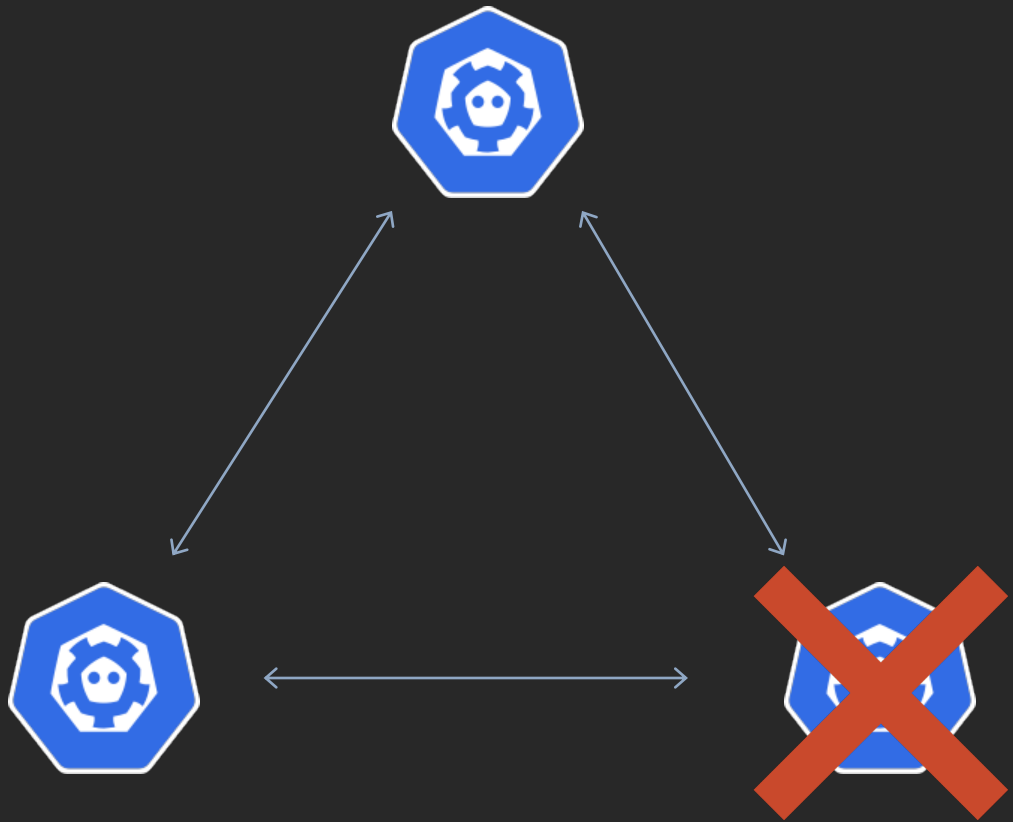
The GODEBUG variables are not covered by Go's API compatibility promise. Please report any issues before disabling HTTP/2 support: <https://golang.org/s/http2bug>

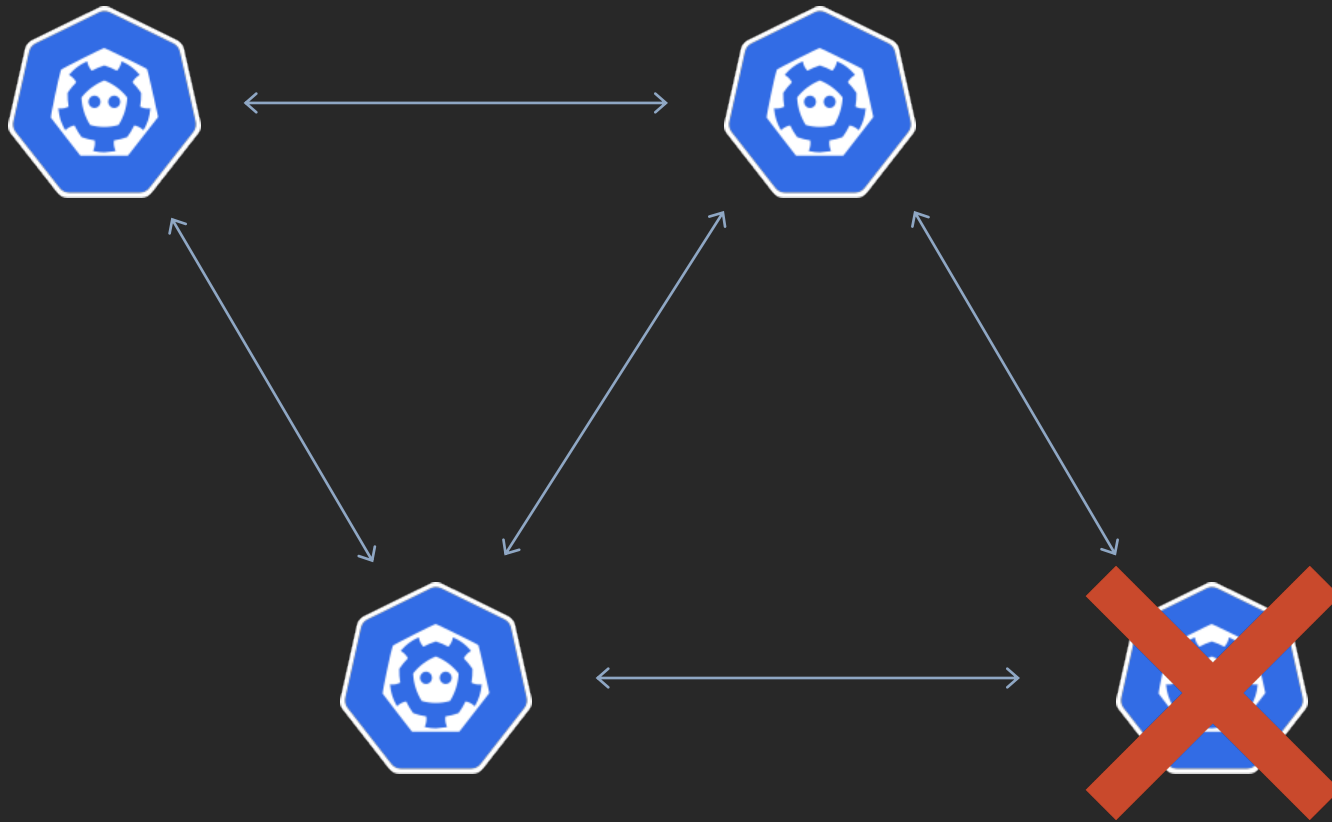
etcd

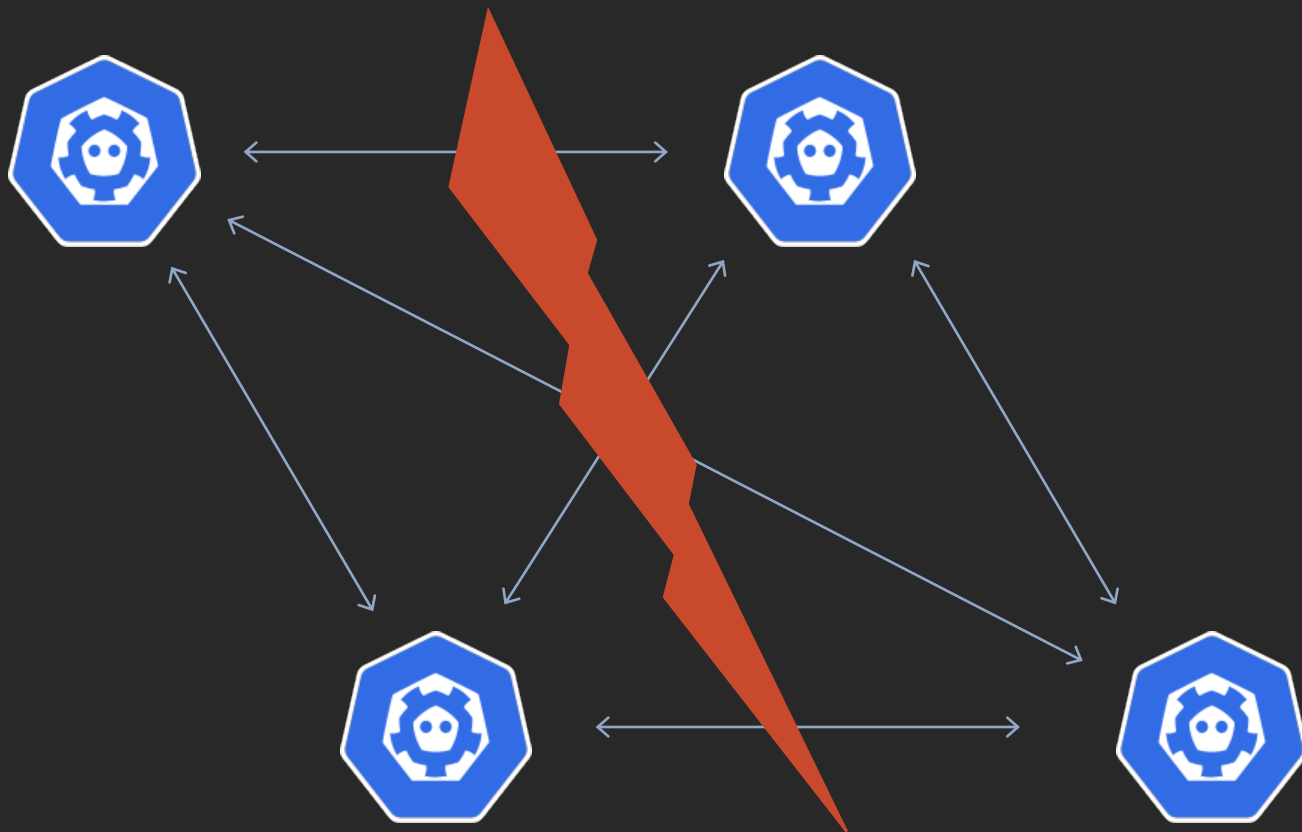


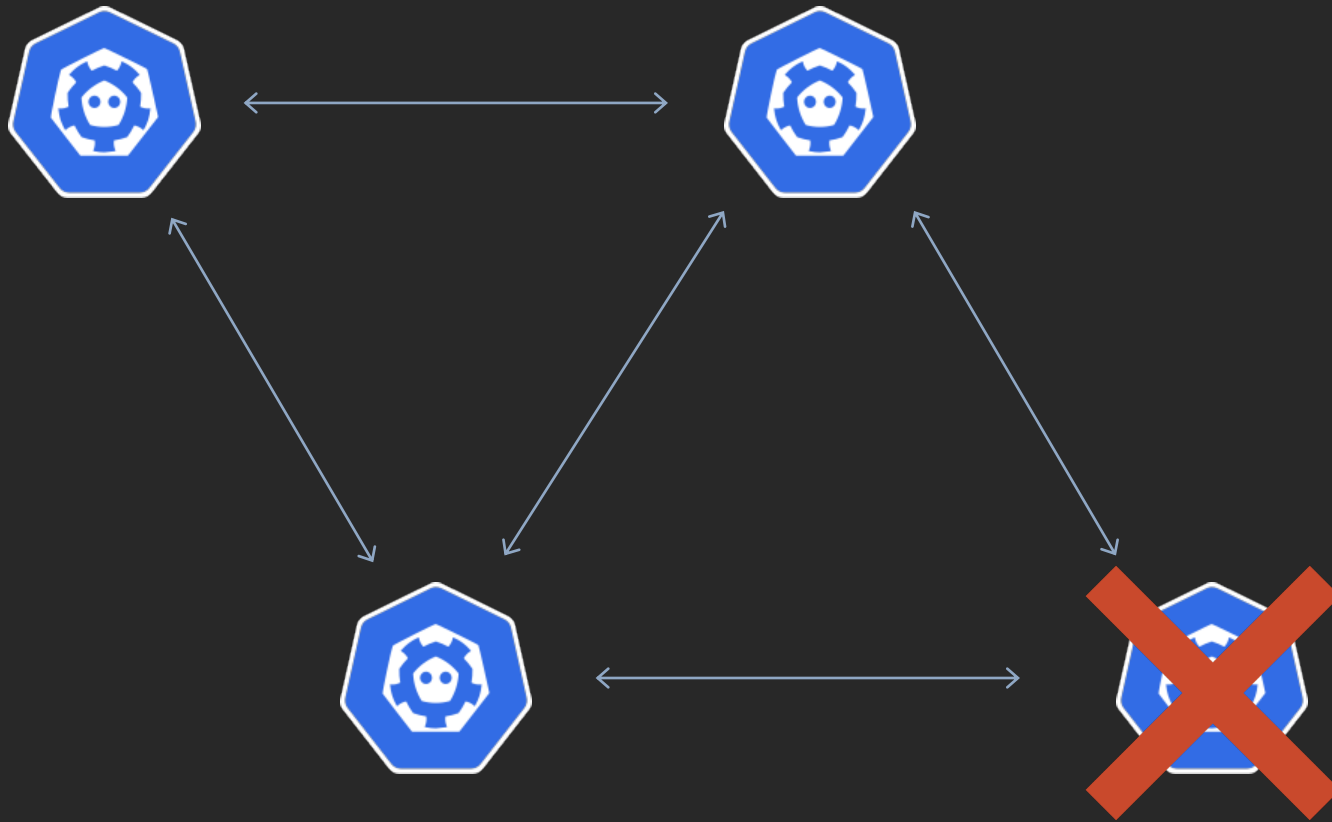












Lessons learned

- Keep backups of etcd
- Monitor quorum size and membership
- Check your dependencies are telling the truth

IAM for Service Accounts

Mutating Webhook

```
apiVersion: v1
kind: ServiceAccount
metadata:
  name: my-serviceaccount
  namespace: default
  annotations:
    eks.amazonaws.com/role-arn: "arn:aws:iam::111122223333:role/s3-reader"
```

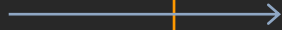
```
apiVersion: v1
kind: Pod
metadata:
  name: my-pod
spec:
  serviceAccountName: my-serviceaccount
  containers:
  - name: container-name
    image: container-image:version
```



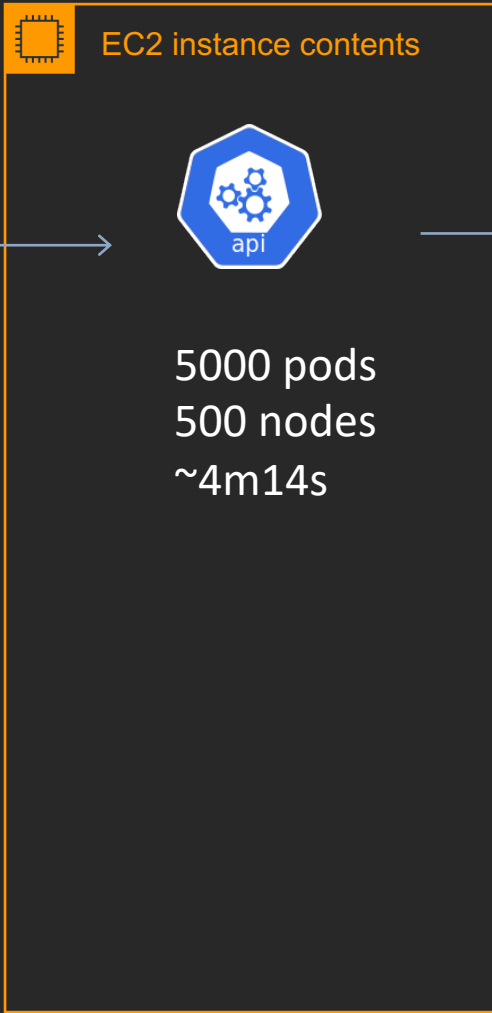
```
apiVersion: v1
kind: Pod
metadata:
  name: my-pod
spec:
  serviceAccountName: my-serviceaccount
  containers:
  - name: container-name
    image: container-image:version
  ### Everything below is added by the webhook ###
  env:
  - name: AWS_ROLE_ARN
    value: "arn:aws:iam::111122223333:role/s3-reader"
  - name: AWS_WEB_IDENTITY_TOKEN_FILE
    value: "/var/run/secrets/eks.amazonaws.com/serviceaccount/token"
  volumeMounts:
  - mountPath: "/var/run/secrets/eks.amazonaws.com/serviceaccount/"
    name: aws-token
  volumes:
  - name: aws-token
    projected:
      sources:
      - serviceAccountToken:
          audience: "sts.amazonaws.com"
          expirationSeconds: 86400
          path: token
```

EKS managed kube-apiserver

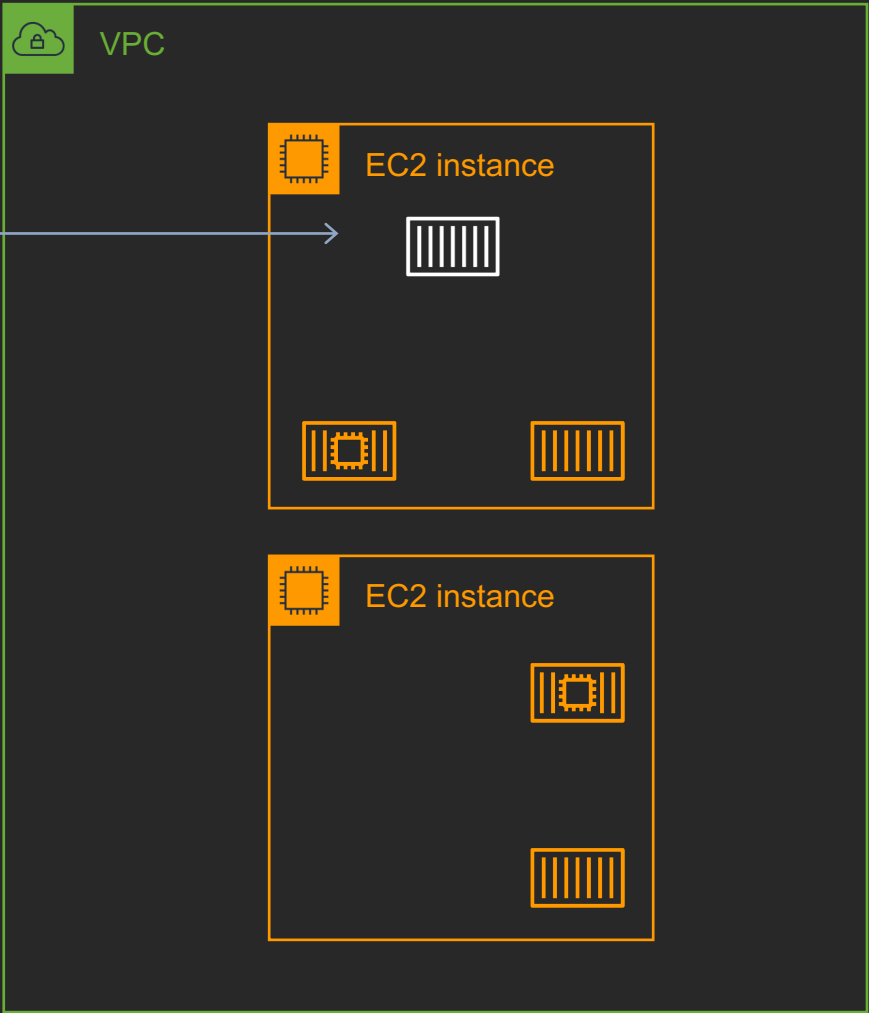
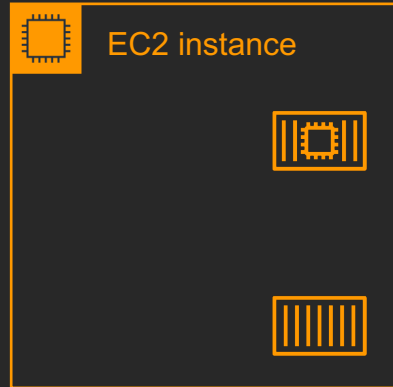
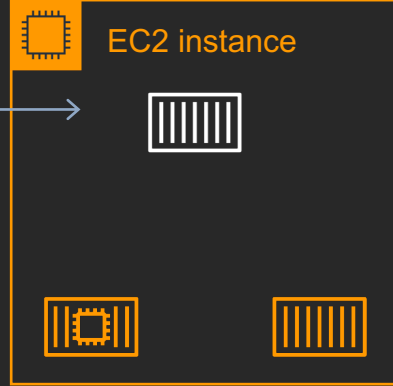
Customer VPC



5000 pods
500 nodes
~4m14s

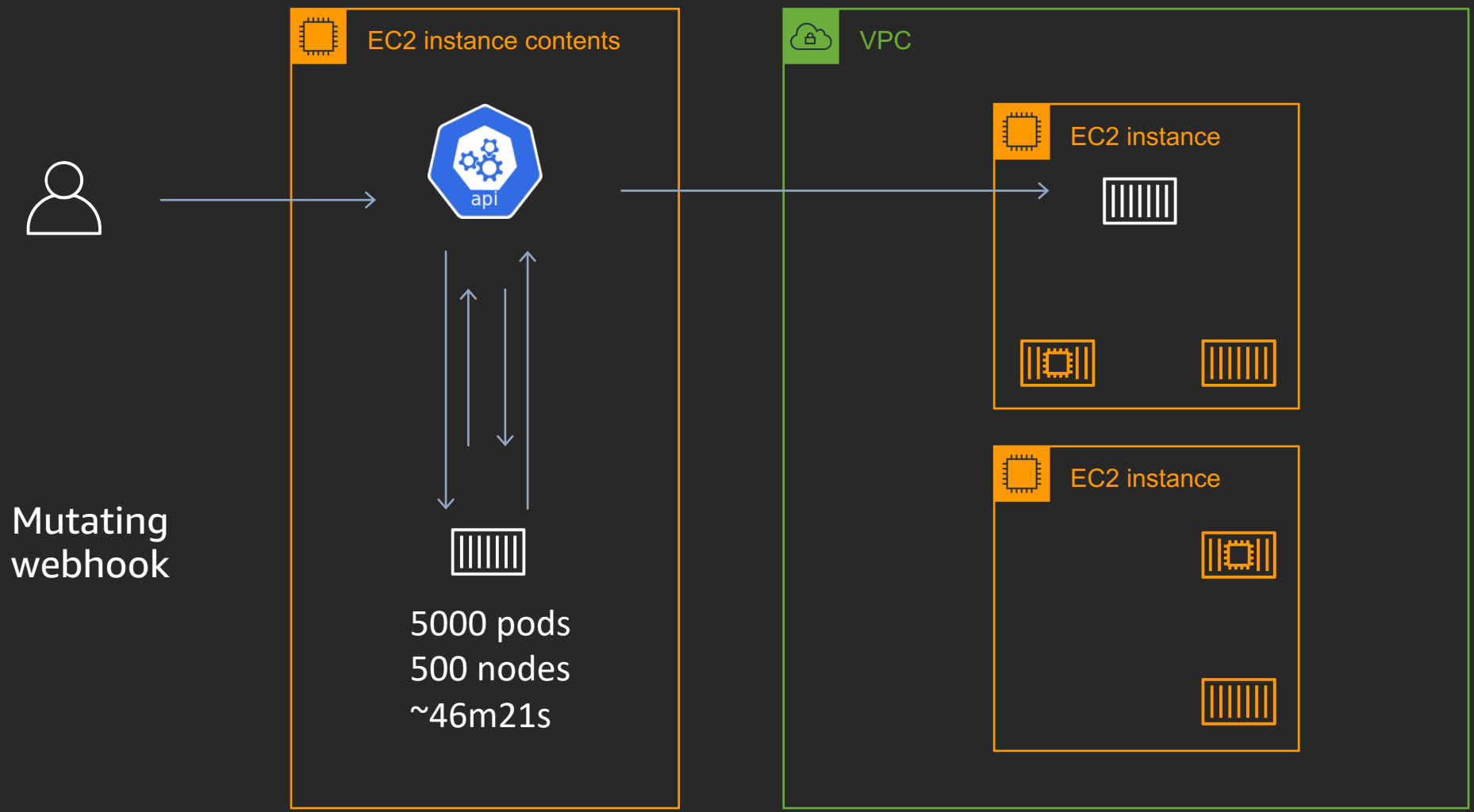


VPC



EKS managed kube-apiserver

Customer VPC



Mutating
webhook

EC2 instance contents



5000 pods
500 nodes
~46m21s

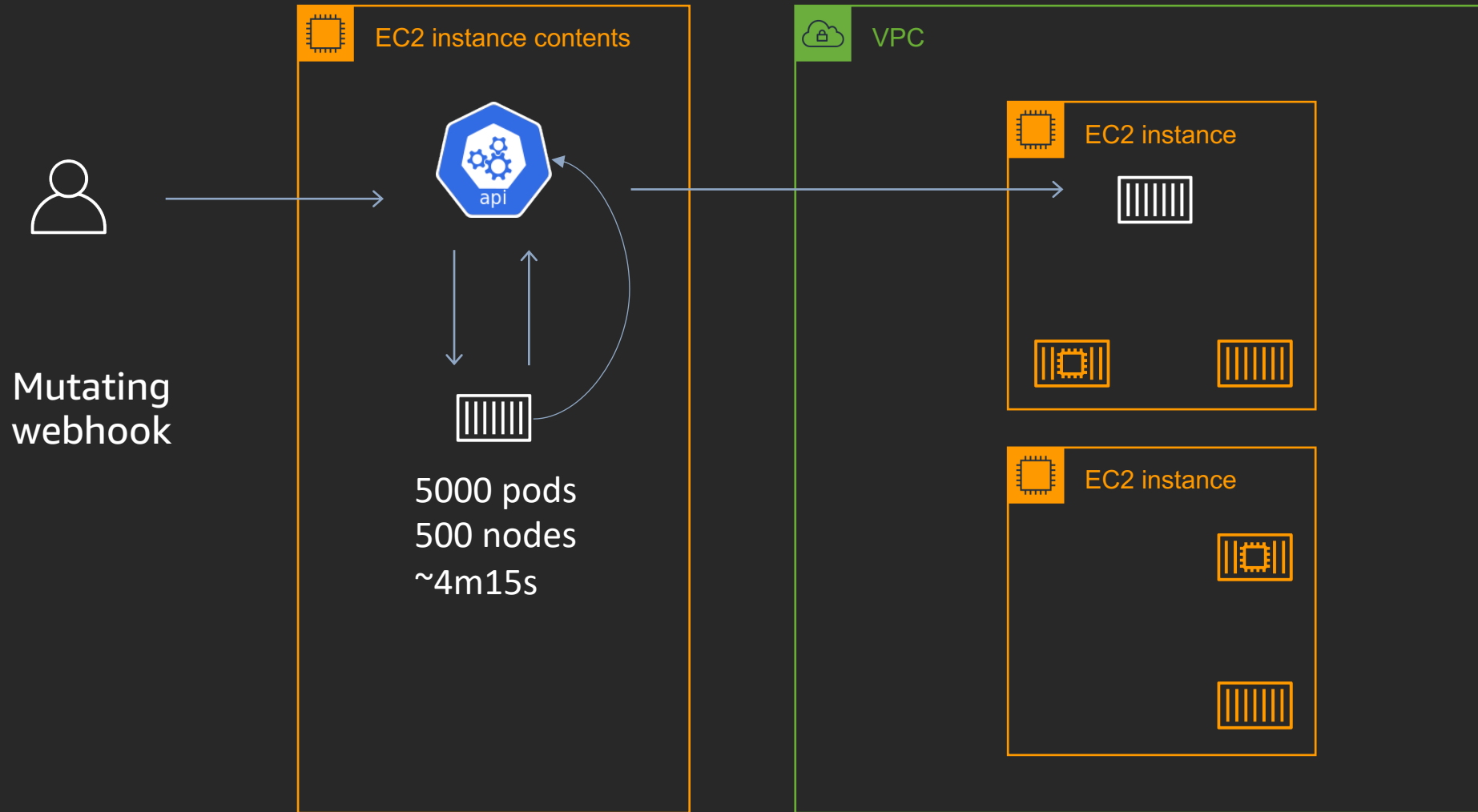
VPC

EC2 instance

EC2 instance

EKS managed kube-apiserver

Customer VPC



Lessons learned

- Keep webhooks as stateless as possible
- When not possible, add a cache
- Always measure your changes



KubeCon



CloudNativeCon

North America 2019

Questions?





KubeCon



CloudNativeCon

North America 2019

Thank you!

Micah Hausler
Sr System Development Engineer, AWS
@micahhausler

