# SIGs Aren't Silos

**A Case Study Into Solving Inter-Domain Problems In Kubernetes Development**

KubeCon | CloudNativeCon

North America 2022

BUILDING FOR THE ROAD AHEAD

## DETROIT 2022

**October 24-28, 2021**

KubeCon | CloudNativeCon

BUILDING FOR THE ROAD AHEAD

DETROIT 2022

**Swetha Repakula**
Software Engineer
*Google*
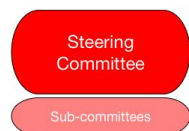
**Antonio Ojea**
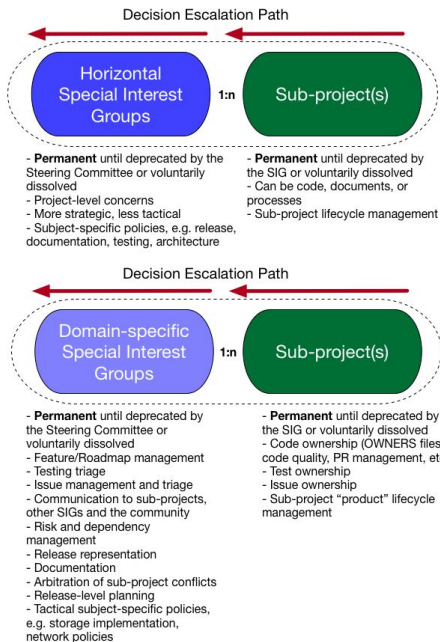Software Engineer
*Google*

# Kubernetes community: Governance



Kubernetes Community Governance Model

Decision Escalation Path
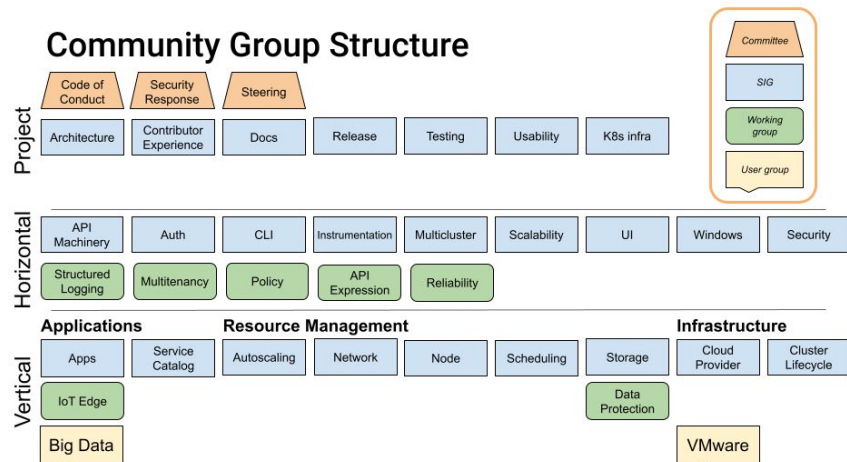
**Steering Committee**

Sub-committees

- Project-level governance
- Binding/final arbitration
- Management of sub-structures
- Security process management
- Manage project-level policies such as the Code of Conduct

**Horizontal Special Interest Groups**   1:n   **Sub-project(s)**

- **Permanent** until deprecated by the Steering Committee or voluntarily dissolved
- Project-level concerns
- More strategic, less tactical
- Subject-specific policies, e.g. release, documentation, testing, architecture

- **Permanent** until deprecated by the SIG or voluntarily dissolved
- Can be code, documents, or processes
- Sub-project lifecycle management

**Working Group**

- Organized to solve a **specific problem** then **dissolve**
- Cross-SIG collaboration around a specific effort
- Does not own code
- Can spawn sub-projects in participating SIGs
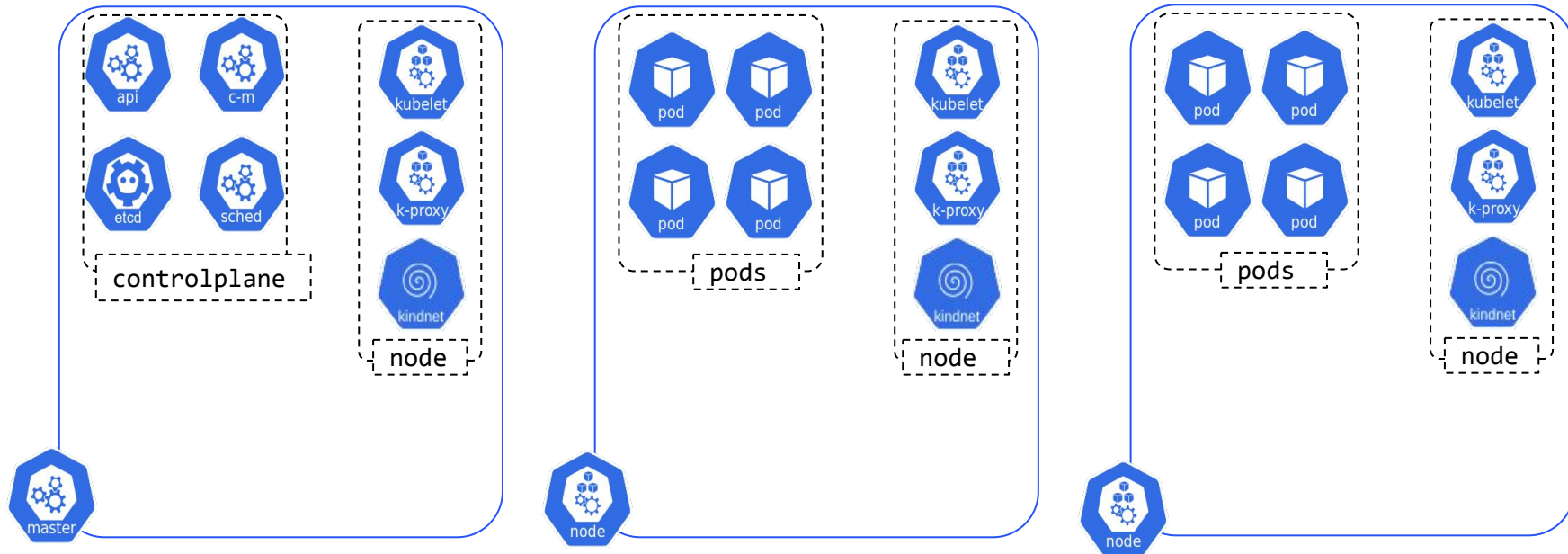- Define and address gaps

Decision Escalation Path

**Domain-specific Special Interest Groups**   1:n   **Sub-project(s)**

- **Permanent** until deprecated by the Steering Committee or voluntarily dissolved
- Feature/Roadmap management
- Testing triage
- Issue management and triage
- Communication to sub-projects, other SIGs and the community
- Risk and dependency management
- Release representation
- Documentation
- Arbitration of sub-project conflicts
- Release-level planning
- Tactical subject-specific policies, e.g. storage implementation, network policies

- **Permanent** until deprecated by the SIG or voluntarily dissolved
- Code ownership (OWNERS files, code quality, PR management, etc.)
- Test ownership
- Issue ownership
- Sub-project "product" lifecycle management

## Community Group Structure

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| **Project** | Code of Conduct | Security Response | Steering | | | | | |
| | Architecture | Contributor Experience | Docs | Release | Testing | Usability | K8s infra |

Committee
SIG
Working group
User group

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| **Horizontal** | API Machinery | Auth | CLI | Instrumentation | Multicluster | Scalability | UI | Windows | Security |
| | Structured Logging | Multitenancy | Policy | API Expression | Reliability | | | |

**Applications**      **Resource Management**                    **Infrastructure**

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| **Vertical** | Apps | Service Catalog | Autoscaling | Network | Node | Scheduling | Storage | Cloud Provider | Cluster Lifecycle |
| | IoT Edge | | | | | | Data Protection |
| | Big Data | | | | | | VMware |

# Kubernetes components

# SIG Network and Node bugs

Pods with failed status IP address reused on new pods, but traffic still going to old pods across namespaces. #109414 `New issue`

`Closed` declangallagher opened this issue on Apr 11 · 11 comments · Fixed by #110255

"Terminated" pod on shutdown node listed in service edpoints. #109718 `New issue`

`Closed` tcolgate opened this issue on Apr 29 · 35 comments

Kubernetes sending traffic to draining nodes #110195 `New issue`

`Closed` Maria-Milusheva opened this issue on May 24 · 10 comments

Kubelet Readiness Probes do not run during pod termination #110309 `New issue`

`Closed` rphillips opened this issue on May 31 · 5 comments · Fixed by #110191

# Symptoms

- Traffic was being routed to non-existent pods

- Traffic was being routed to the wrong pods

- IPs in EndpointSlices were not matching the Node Pod CIDR
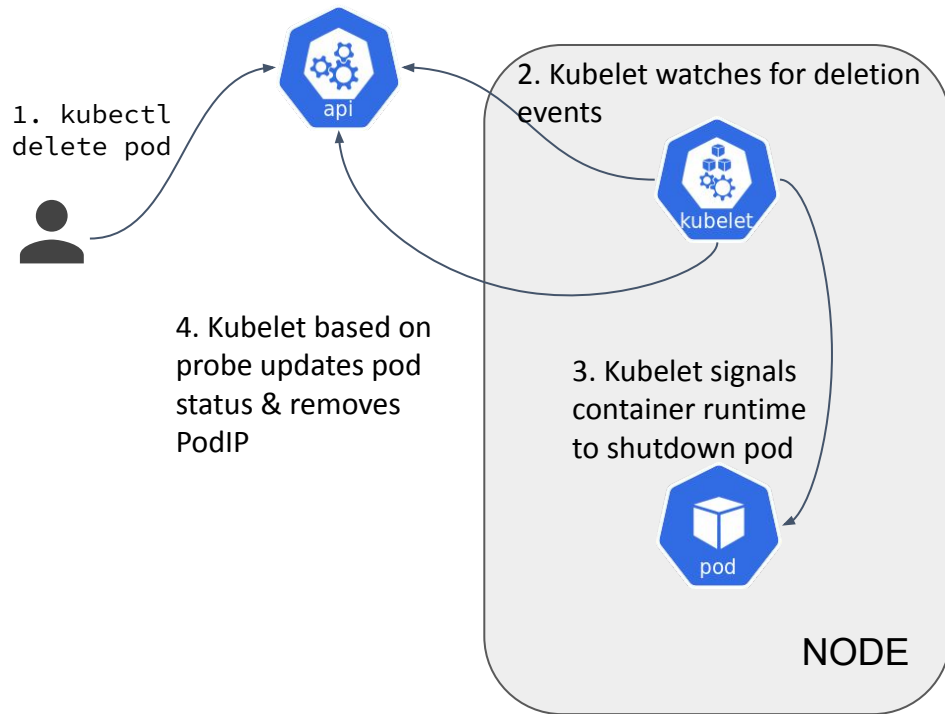
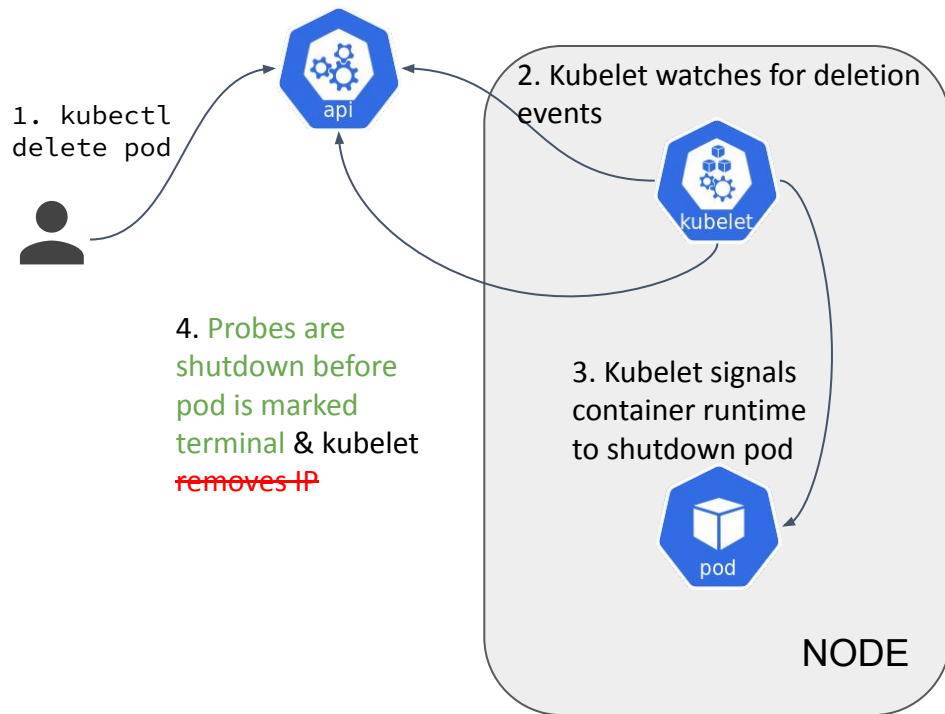# Pod Lifecycle: Pod Phases

# Pod Lifecycle: Pod Readiness

# Pod Shutdown - Before Refactor

1. `kubectl delete pod`

2. Kubelet watches for deletion events

api

kubelet

pod

4. Kubelet based on probe updates pod status & removes PodIP

3. Kubelet signals container runtime to shutdown pod

NODE

When nodes are shutdown, or if pods are evicted, pods will not be deleted until PodGC kicks in (based on threshold).

```
$ kubectl get pods -o wide
NAME       READY   STATUS       RESTARTS   AGE     IP            NODE
busybox    0/1     Completed    0          15m     10.244.0.2    121-control-plane
busybox3   0/1     StartError   0          3m56s   10.244.0.6    121-control-plane
evictme    0/1     Evicted      0          63s     <none>        121-control-plane
```
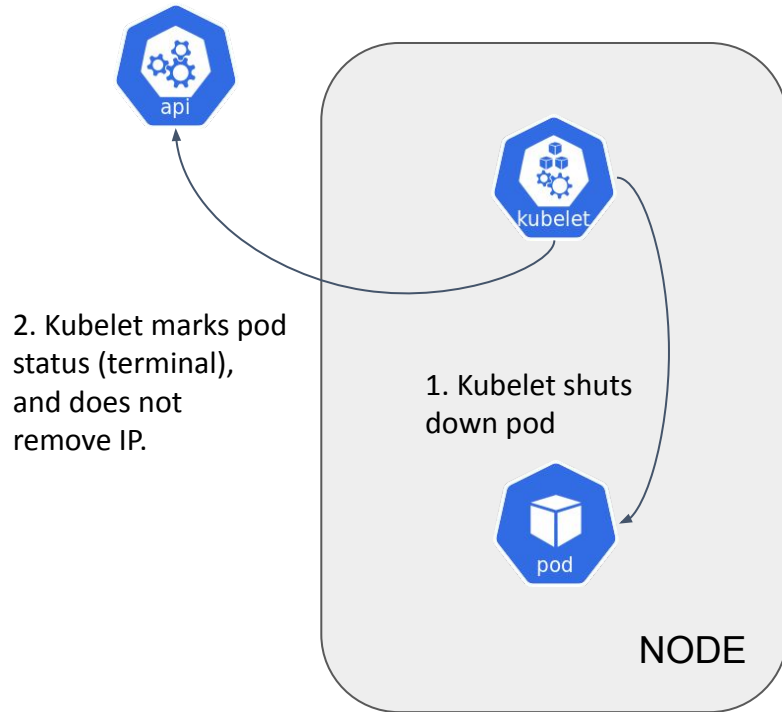
# Pod Shutdown - After Refactor

1. `kubectl delete pod`

2. Kubelet watches for deletion events

3. Kubelet signals container runtime to shutdown pod

4. Probes are shutdown before pod is marked terminal & kubelet ~~removes IP~~

NODE

```
$ kubectl get pods -o wide
NAME       READY   STATUS      RESTARTS   AGE     IP           NODE
busybox    0/1     Completed   0          22m     10.244.0.5   kind-control-plane
busybox3   0/1     StartError  0          5m36s   10.244.0.6   kind-control-plane
evictme    0/1     Evicted     0          63s     10.244.0.7   kind-control-plane
```
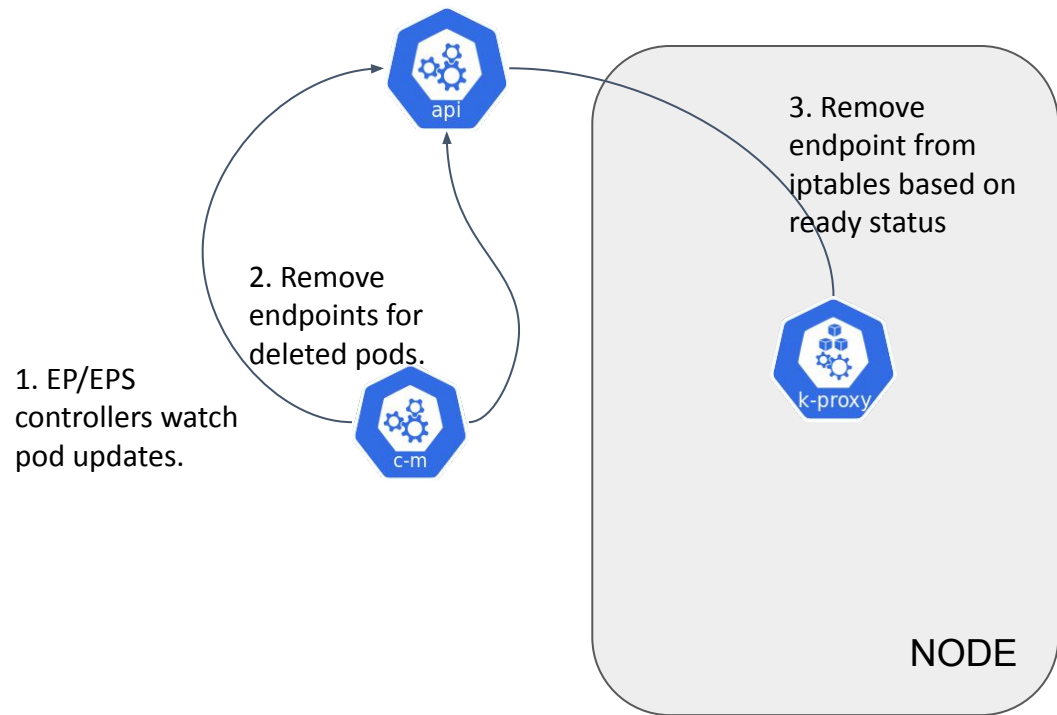
# Pod Shutdown - Pod Eviction

2. Kubelet marks pod status (terminal), and does not remove IP.

1. Kubelet shuts down pod

**NODE**

api

kubelet

pod

When nodes are shutdown, or if pods are evicted, pods will not be deleted until PodGC kicks in (based on threshold).

# Service/Endpoints Removal



3. Remove endpoint from iptables based on ready status

2. Remove endpoints for deleted pods.

1. EP/EPS controllers watch pod updates.

NODE

# The bug …

```go
        for _, pod := range pods {
                if len(pod.Status.PodIP) == 0 {
                        klog.V(5).Infof("Failed to find an IP for pod %s/%s", pod.Namespace, pod.Name)
                        continue
                }
                if !tolerateUnreadyEndpoints && pod.DeletionTimestamp != nil {
                        klog.V(5).Infof("Pod is being deleted %s/%s", pod.Namespace, pod.Name)
                        continue
                }

```

pkg/controller/endpoint/endpoints_controller.go (v1.22.0)
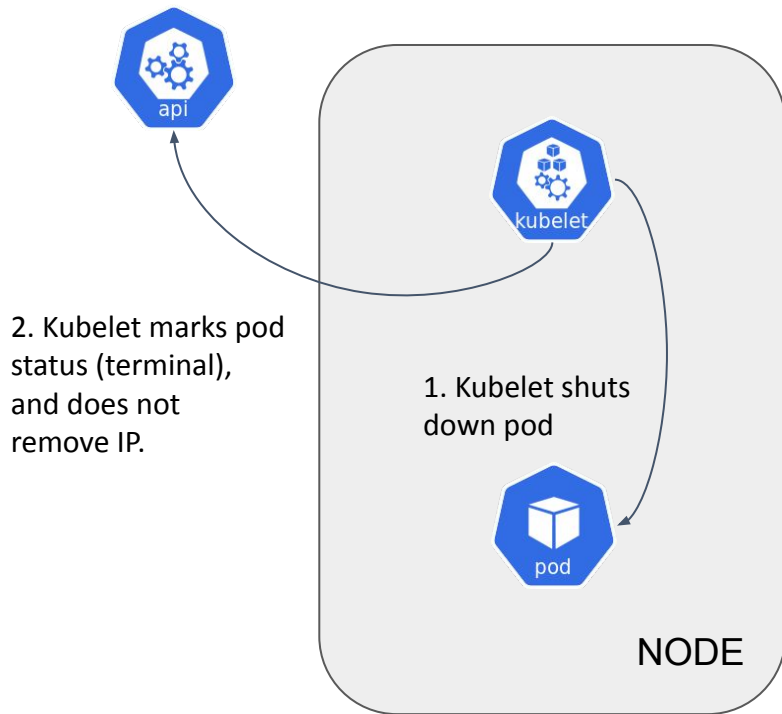
```
448          for _, pod := range pods {
449              if len(pod.Status.PodIP) == 0 {
450                  klog.V(5).Infof("Failed to find an IP for pod %s/%s", pod.Namespace, pod.Name)
451                  continue
452              }
453              if !tolerateUnreadyEndpoints && pod.DeletionTimestamp != nil {
454                  klog.V(5).Infof("Pod is being deleted %s/%s", pod.Namespace, pod.Name)
455                  continue
456              }
457
```

pkg/controller/endpoint/endpoints_controller.go (v1.22.0)

# Pod Shutdown - Pod Eviction

2. Kubelet marks pod status (terminal), and does not remove IP.

1. Kubelet shuts down pod

NODE

When nodes are shutdown, or if pods are evicted, pods will not be deleted until PodGC kicks in (based on threshold).

```
448        for _, pod := range pods {
449            if len(pod.Status.PodIP) == 0 {
450                klog.V(5).Infof("Failed to find an IP for pod %s/%s", pod.Namespace, pod.Name)
451                continue
452            }
453            if !tolerateUnreadyEndpoints && pod.DeletionTimestamp != nil {
454                klog.V(5).Infof("Pod is being deleted %s/%s", pod.Namespace, pod.Name)
455                continue
456            }
457
```

pkg/controller/endpoint/endpoints_controller.go (v1.22.0)

# The bug …

```go
622  func addEndpointSubset(subsets []v1.EndpointSubset, pod *v1.Pod, epa v1.EndpointAddress,
623           epp *v1.EndpointPort, tolerateUnreadyEndpoints bool) ([]v1.EndpointSubset, int, int) {
624           var readyEps int
625           var notReadyEps int
626           ports := []v1.EndpointPort{}
627           if epp != nil {
628                   ports = append(ports, *epp)
629           }
630           if tolerateUnreadyEndpoints || podutil.IsPodReady(pod) {
631                   subsets = append(subsets, v1.EndpointSubset{
632                           Addresses: []v1.EndpointAddress{epa},
633                           Ports:     ports,
634                   })
635                   readyEps++
636           } else if shouldPodBeInEndpoints(pod) {
637                   klog.V(5).Infof("Pod is out of service: %s/%s", pod.Namespace, pod.Name)
638                   subsets = append(subsets, v1.EndpointSubset{
639                           NotReadyAddresses: []v1.EndpointAddress{epa},
640                           Ports:             ports,
641                   })
642                   notReadyEps++
643           }
644           return subsets, readyEps, notReadyEps
645  }
```

pkg/controller/endpoint/endpoints_controller.go (v1.22.0)

```
622  func addEndpointSubset(subsets []v1.EndpointSubset, pod *v1.Pod, epa v1.EndpointAddress,
623          epp *v1.EndpointPort, tolerateUnreadyEndpoints bool) ([]v1.EndpointSubset, int, int) {
624          var readyEps int
625          var notReadyEps int
626          ports := []v1.EndpointPort{}
627          if epp != nil {
628                  ports = append(ports, *epp)
629          }
630          if tolerateUnreadyEndpoints || podutil.IsPodReady(pod) {
631                  subsets = append(subsets, v1.EndpointSubset{
632                          Addresses: []v1.EndpointAddress{epa},
633                          Ports:     ports,
634                  })
635                  readyEps++
636          } else if shouldPodBeInEndpoints(pod) {
637                  klog.V(5).Infof("Pod is out of service: %s/%s", pod.Namespace, pod.Name)
638                  subsets = append(subsets, v1.EndpointSubset{
639                          NotReadyAddresses: []v1.EndpointAddress{epa},
640                          Ports:             ports,
641                  })
642                  notReadyEps++
643          }
644          return subsets, readyEps, notReadyEps
645  }
```
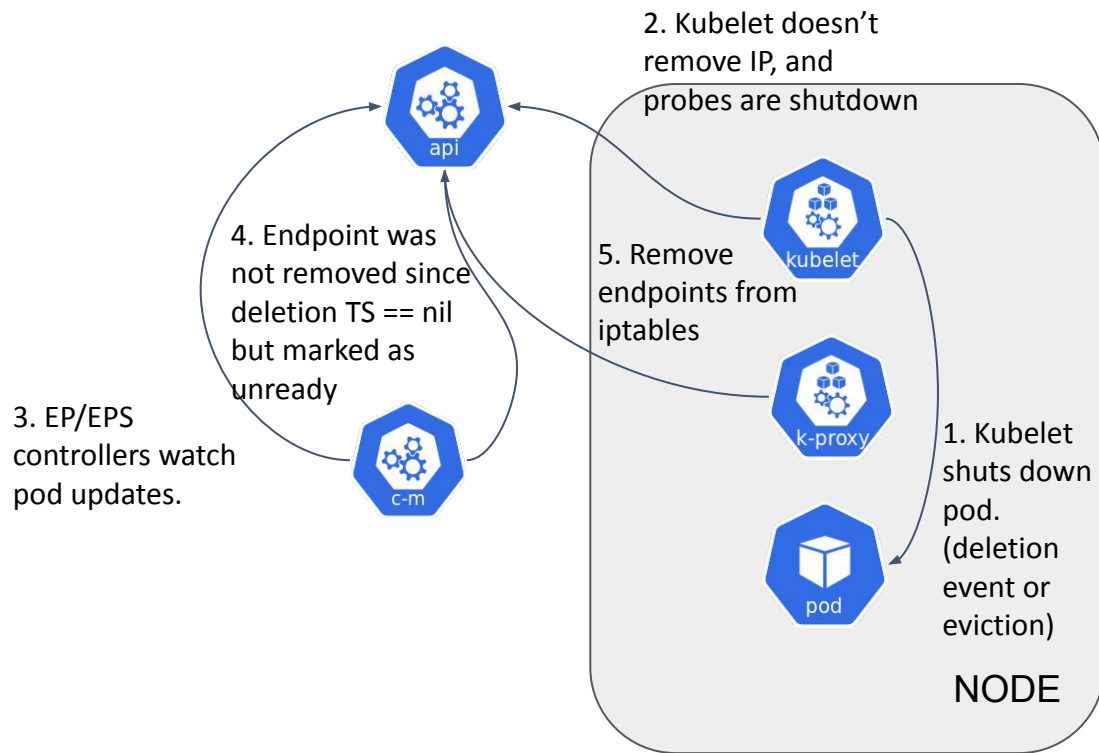
pkg/controller/endpoint/endpoints_controller.go (v1.22.0)

```
622  func addEndpointSubset(subsets []v1.EndpointSubset, pod *v1.Pod, epa v1.EndpointAddress,
623          epp *v1.EndpointPort, tolerateUnreadyEndpoints bool) ([]v1.EndpointSubset, int, int) {
624          var readyEps int
625          var notReadyEps int
626          ports := []v1.EndpointPort{}
627          if epp != nil {
628                  ports = append(ports, *epp)
629          }
630          if tolerateUnreadyEndpoints || podutil.IsPodReady(pod) {
631                  subsets = append(subsets, v1.EndpointSubset{
632                          Addresses: []v1.EndpointAddress{epa},
633                          Ports:     ports,
634                  })
635                  readyEps++
636          } else if shouldPodBeInEndpoints(pod) {
637                  klog.V(5).Infof("Pod is out of service: %s/%s", pod.Namespace, pod.Name)
638                  subsets = append(subsets, v1.EndpointSubset{
639                          NotReadyAddres
640                          Ports:
641                  })
642                  notReadyEps++
643          }
644          return subsets, readyEps, notR
645  }
```
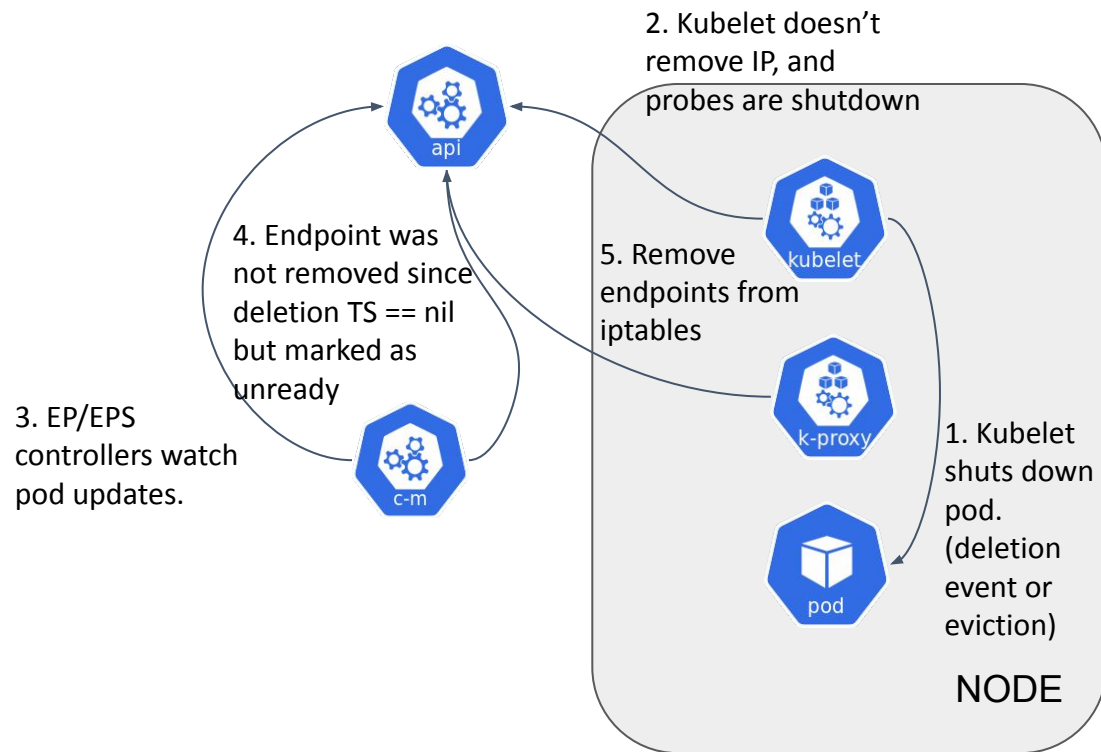
```
647  func shouldPodBeInEndpoints(pod *v1.Pod) bool {
648          switch pod.Spec.RestartPolicy {
649          case v1.RestartPolicyNever:
650                  return pod.Status.Phase != v1.PodFailed && pod.Status.Phase != v1.PodSucceeded
651          case v1.RestartPolicyOnFailure:
652                  return pod.Status.Phase != v1.PodSucceeded
653          default:
654                  return true
655          }
656  }
```

))

# The bug …



2. Kubelet doesn't remove IP, and probes are shutdown

4. Endpoint was not removed since deletion TS == nil but marked as unready

5. Remove endpoints from iptables

3. EP/EPS controllers watch pod updates.

1. Kubelet shuts down pod. (deletion event or eviction)

NODE

# … and its consequences

2. Kubelet doesn't remove IP, and probes are shutdown

api

4. Endpoint was not removed since deletion TS == nil but marked as unready

5. Remove endpoints from iptables

kubelet

3. EP/EPS controllers watch pod updates.

c-m

k-proxy

1. Kubelet shuts down pod. (deletion event or eviction)

pod
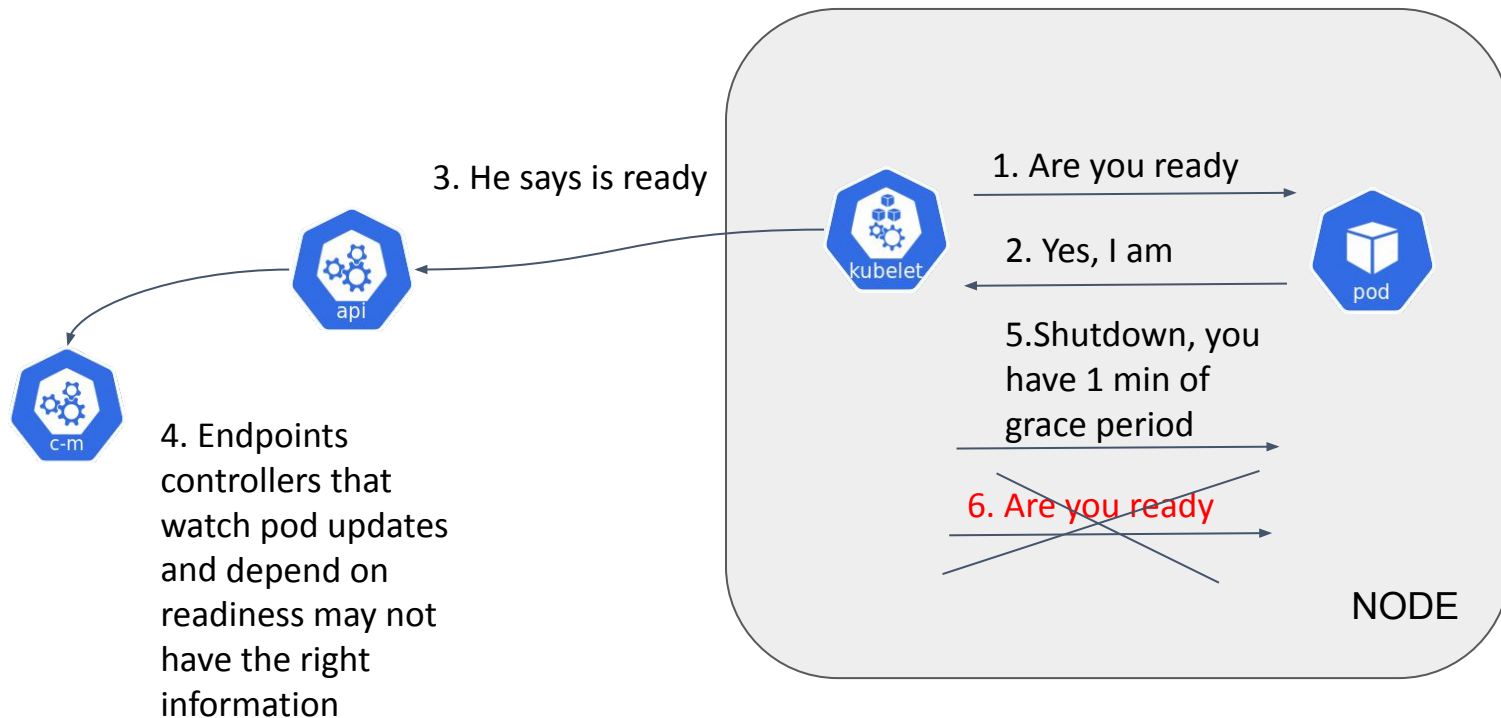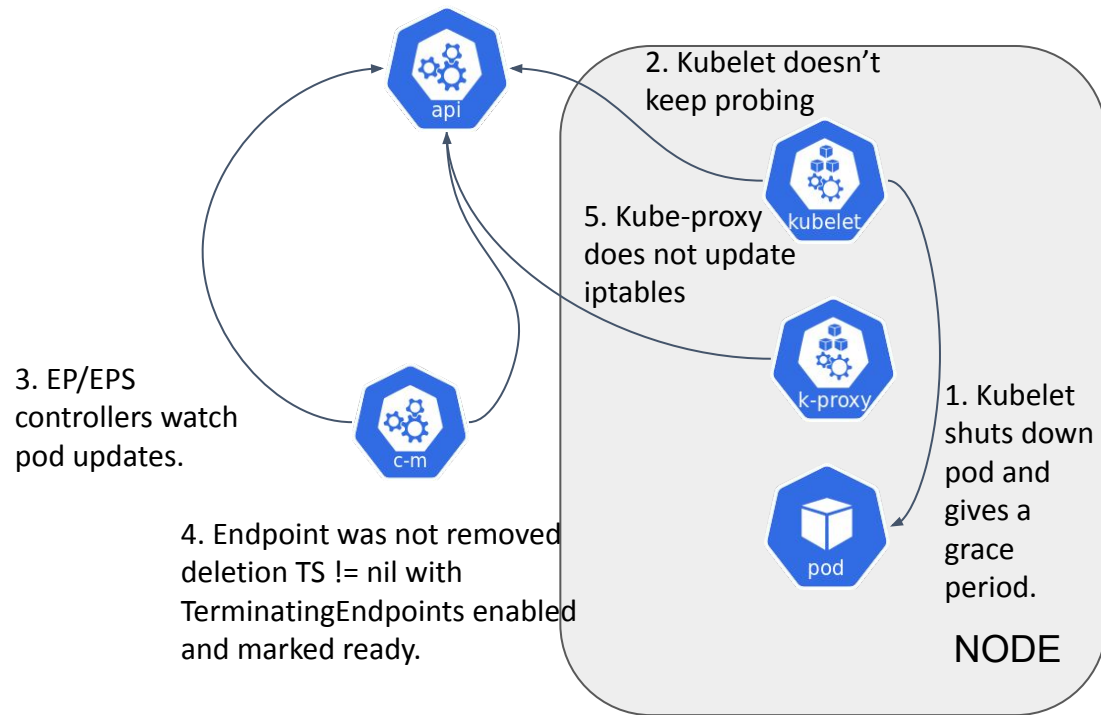
NODE

Endpoints and EndpointSlice resources have an endpoint that doesn't point to a current pod but marked as unready, so kube-proxy will removed from iptables.

# The other bug … TerminatingEndpoints

# … and its consequences

2. Kubelet doesn't keep probing

5. Kube-proxy does not update iptables

3. EP/EPS controllers watch pod updates.

4. Endpoint was not removed deletion TS != nil with TerminatingEndpoints enabled and marked ready.

1. Kubelet shuts down pod and gives a grace period.

NODE

Traffic could be routed pods that area no longer ready. If pod readiness is not checked during the termination grace period, the pod keeps receiving traffic.

# The bugfix



**kubelet: only shutdown probes for pods that are terminated**

This fixes a bug where terminating pods would not run their readiness probes. Terminating pods are found within the possiblyRunningPods map.

< Prev    Next >

master (#110191)
v1.26.0-alpha.2 ... v1.25.0-alpha.1

rphillips committed on Jun 6                                    commit f25ca15e1c47ebd8036c1fb5c06fe8fe6714807a

4 ■■■■ pkg/kubelet/kubelet_pods.go                                                             ...

```
       @@ -1105,8 +1105,8 @@ func (kl *Kubelet) HandlePodCleanups() error {
1105 1105          }
1106 1106
1107 1107          // Stop probing pods that are not running
1108    -          klog.V(3).InfoS("Clean up probes for terminating and terminated pods")
1109    -          kl.probeManager.CleanupPods(runningPods)
     1108 +          klog.V(3).InfoS("Clean up probes for terminated pods")
     1109 +          kl.probeManager.CleanupPods(possiblyRunningPods)
```

# The lessons learned

- All behavior changes have potential to break others

- Networking and Pod/Node lifecycle are tightly coupled

- Behavior Contracts need to be enforced with testing

- "Broken telephone" bugs are hard to solve

- Bugs can become a feature

# Thank You!

Thank you!