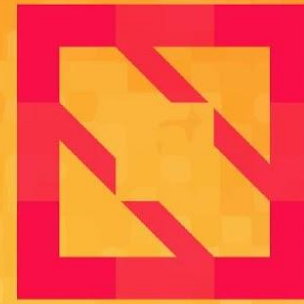




KubeCon



CloudNativeCon

North America 2019





KubeCon



CloudNativeCon

North America 2019

Making the Most Out of Kubernetes Audit Logs

Robert Boll
Laurent Bernaille

@roboll_
@lbernail



Datadog



KubeCon



CloudNativeCon

North America 2019

Monitoring service
Over 350 integrations
Over 1,200 employees
Over 8,000 customers
Runs on millions of hosts
Trillions of data points per day

10000s hosts in our infra
10s of Kubernetes clusters
Clusters from 50 to 3000 nodes
Multi-cloud
Very fast growth

Understanding what happens can be hard

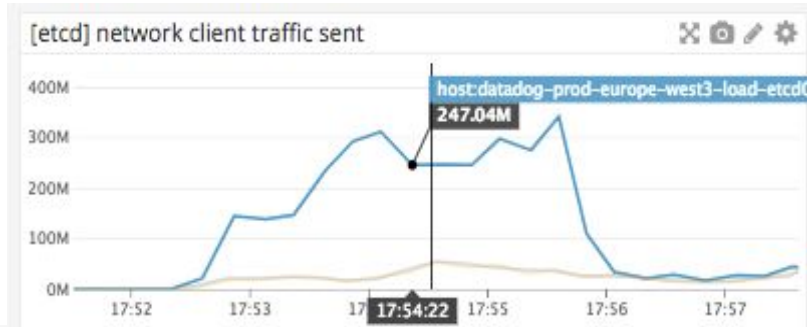


KubeCon

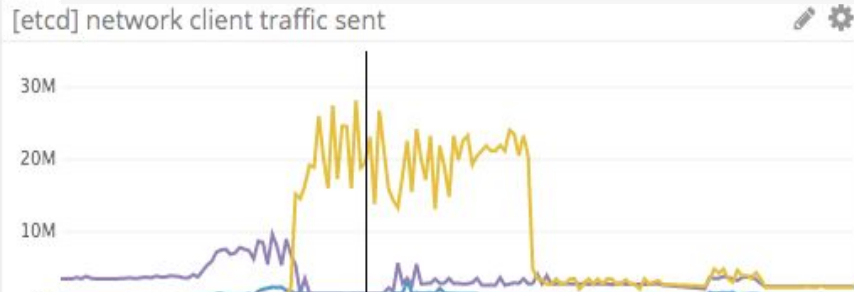
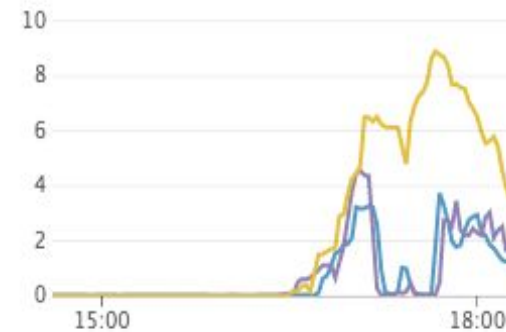


CloudNativeCon

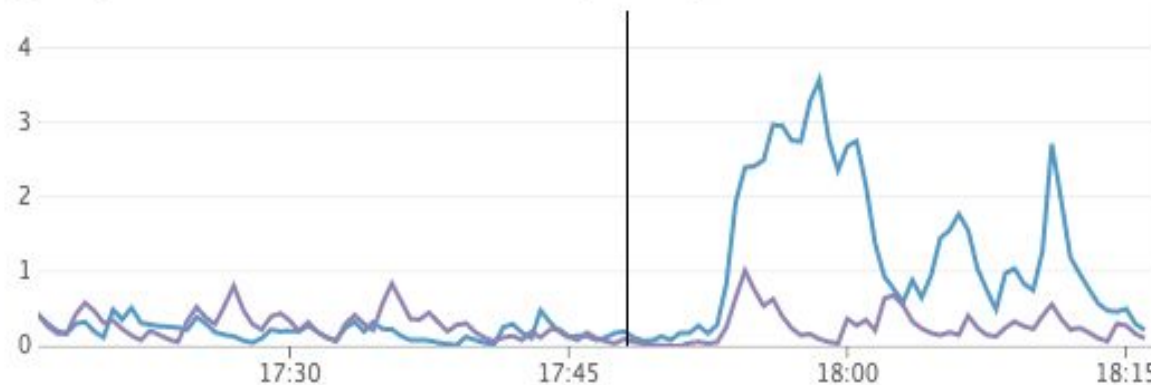
North America 2019



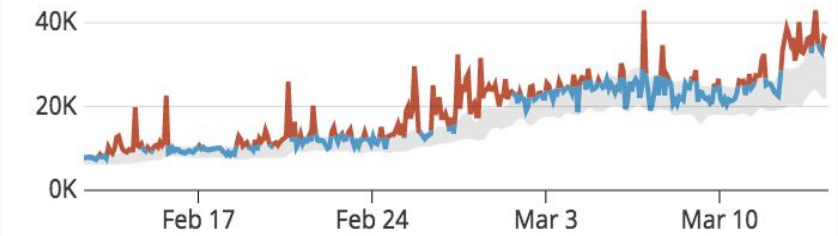
90-percentile etcd request latency



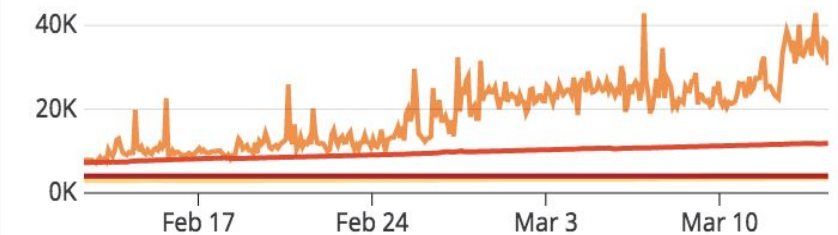
Avg of system.load.1 over \$kubernetes_cluster,role:controllers...



[apiserver] etcd objects seen



[apiserver] etcd top objects seen



Outline



KubeCon



CloudNativeCon

North America 2019

1. Background: The Kubernetes API
2. Audit Logs
3. Configuring Audit Logs
4. 10000 foot view for a large cluster
5. Understanding Kubernetes Internals
6. Troubleshooting examples

Outline



KubeCon



CloudNativeCon

North America 2019

1. **Background: The Kubernetes API**
2. Audit Logs
3. Configuring Audit Logs
4. 10000 foot view for a large cluster
5. Understanding Kubernetes Internals
6. Troubleshooting examples



KubeCon



CloudNativeCon

North America 2019

Background: The Kubernetes API



Calls to the apiservers

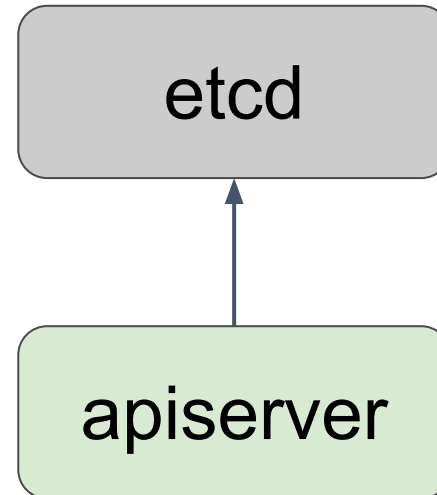


KubeCon



CloudNativeCon

North America 2019



Control plane

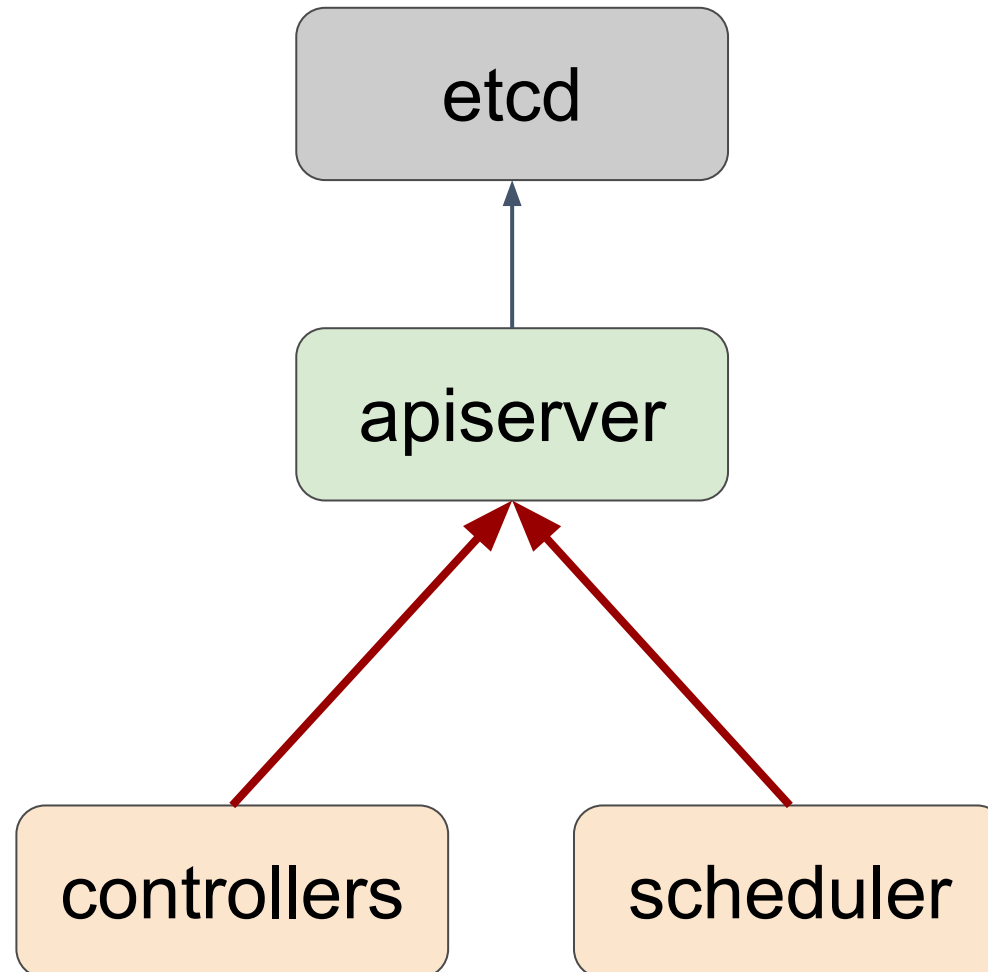


KubeCon



CloudNativeCon

North America 2019



Kubelet

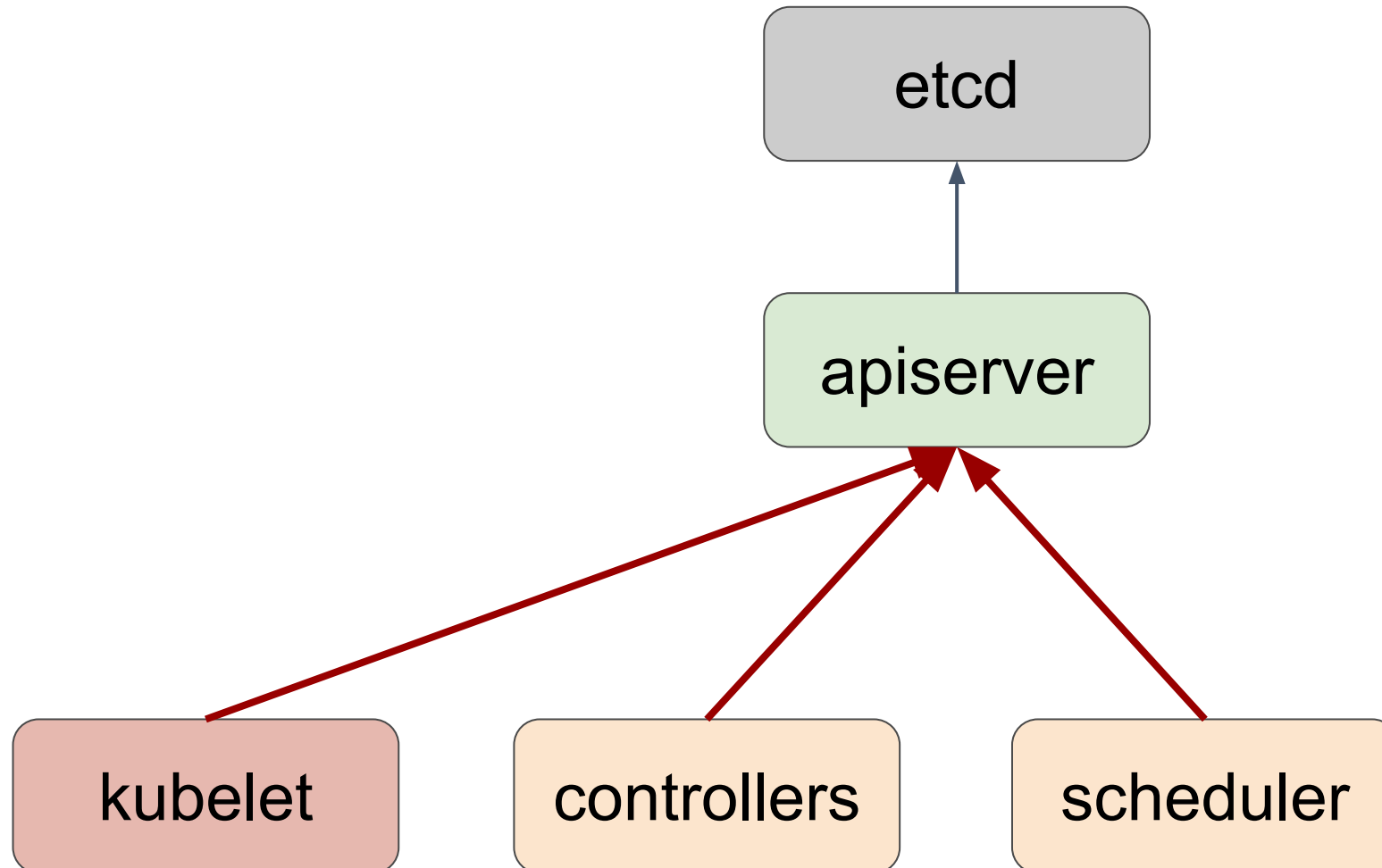


KubeCon



CloudNativeCon

North America 2019



DaemonSet: kube-proxy

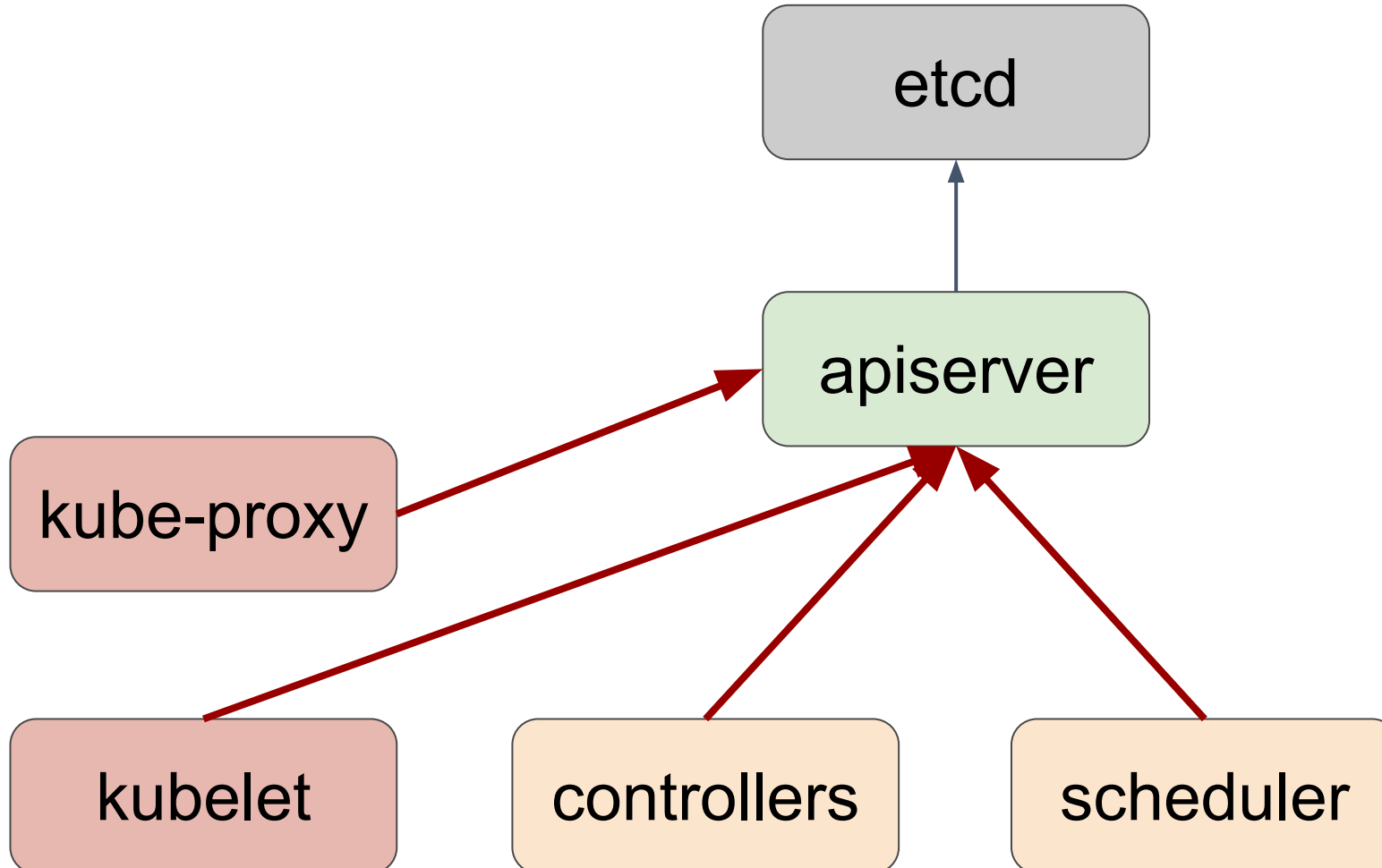


KubeCon



CloudNativeCon

North America 2019



Other DaemonSets (cni, etc.)

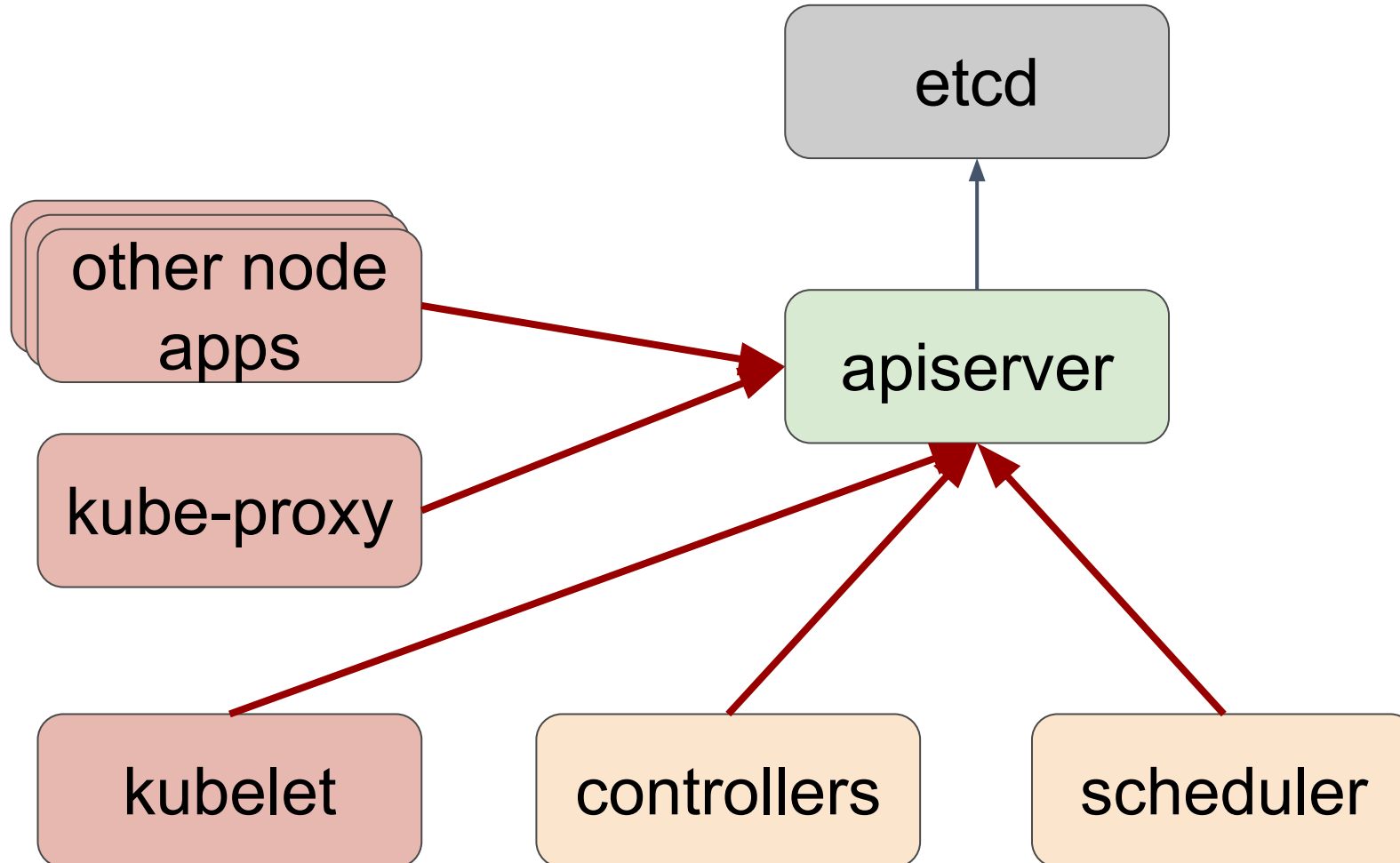


KubeCon



CloudNativeCon

North America 2019



Cluster services: DNS

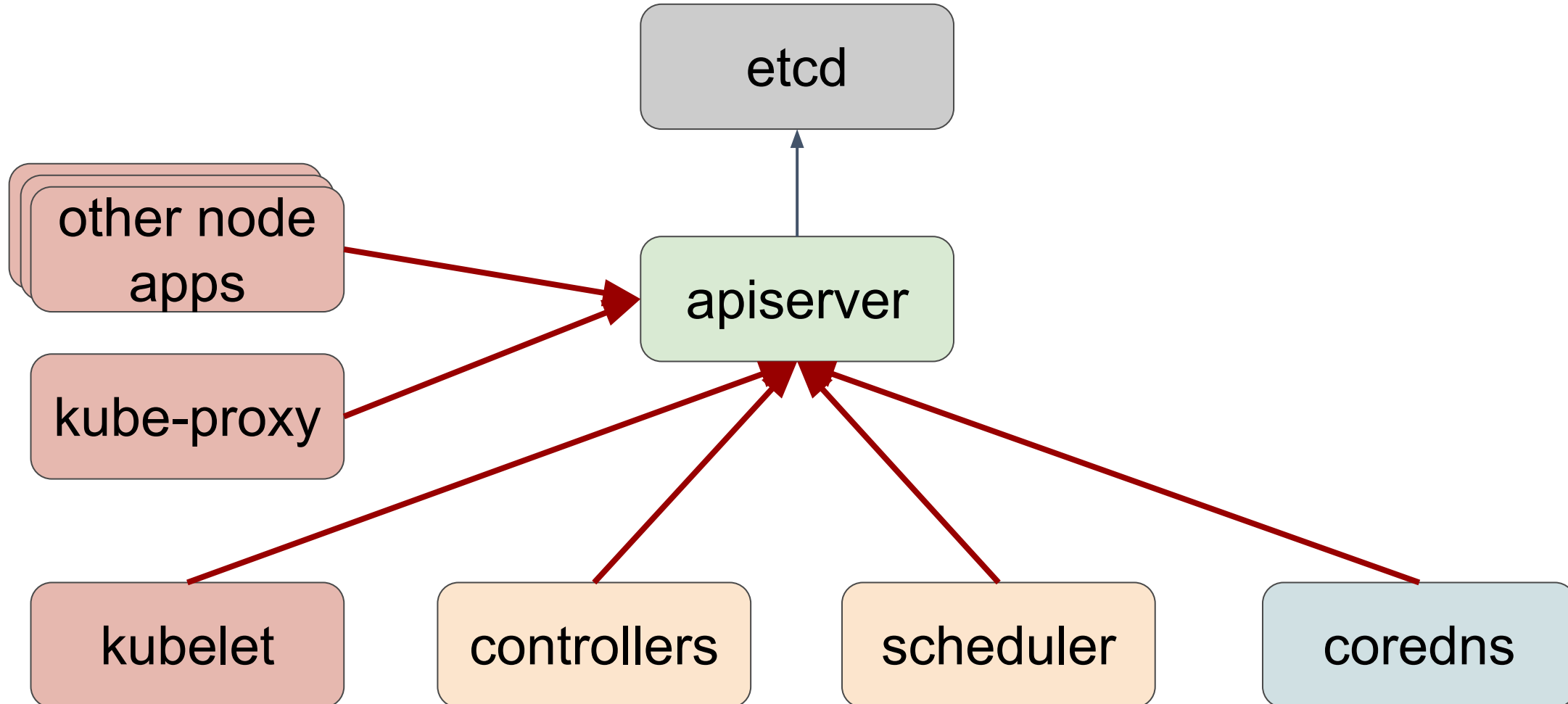


KubeCon



CloudNativeCon

North America 2019



Other cluster services

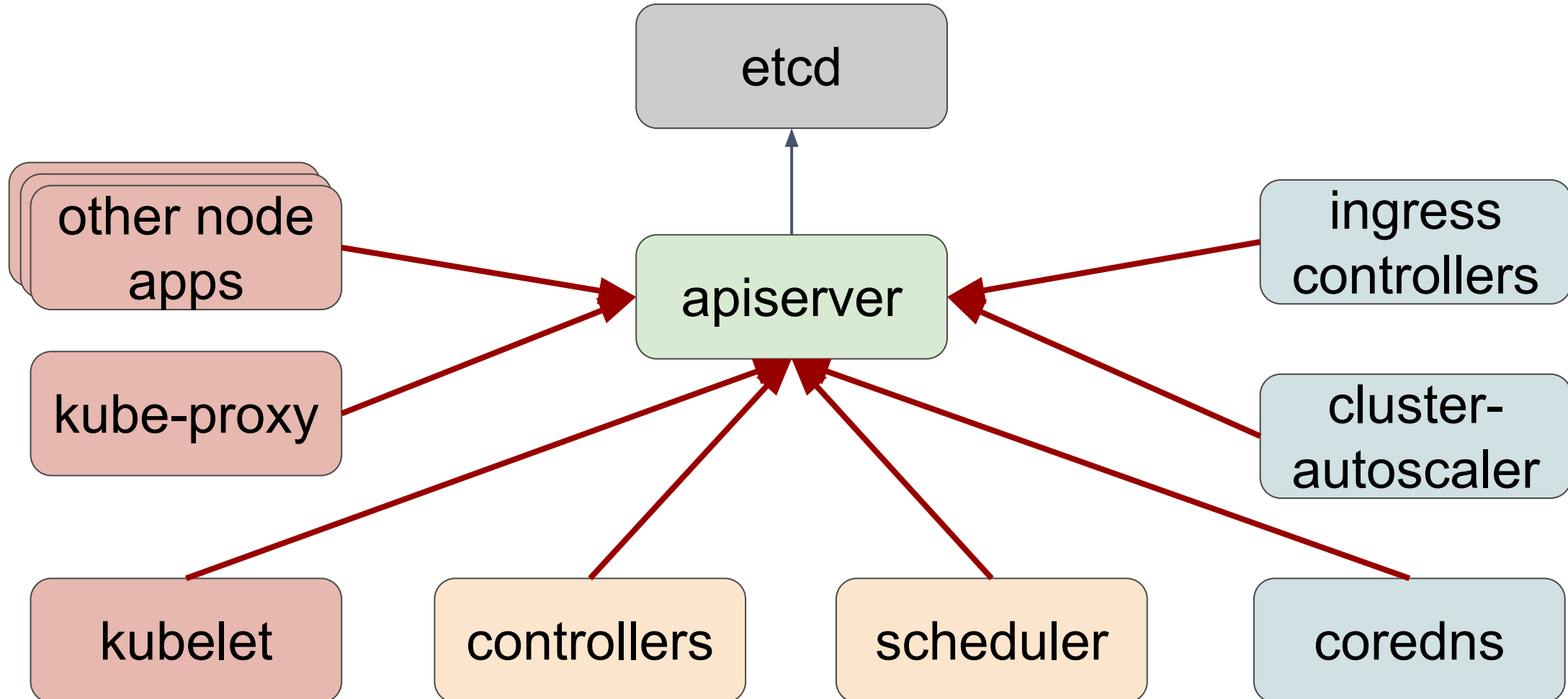


KubeCon



CloudNativeCon

North America 2019



Probably several other applications

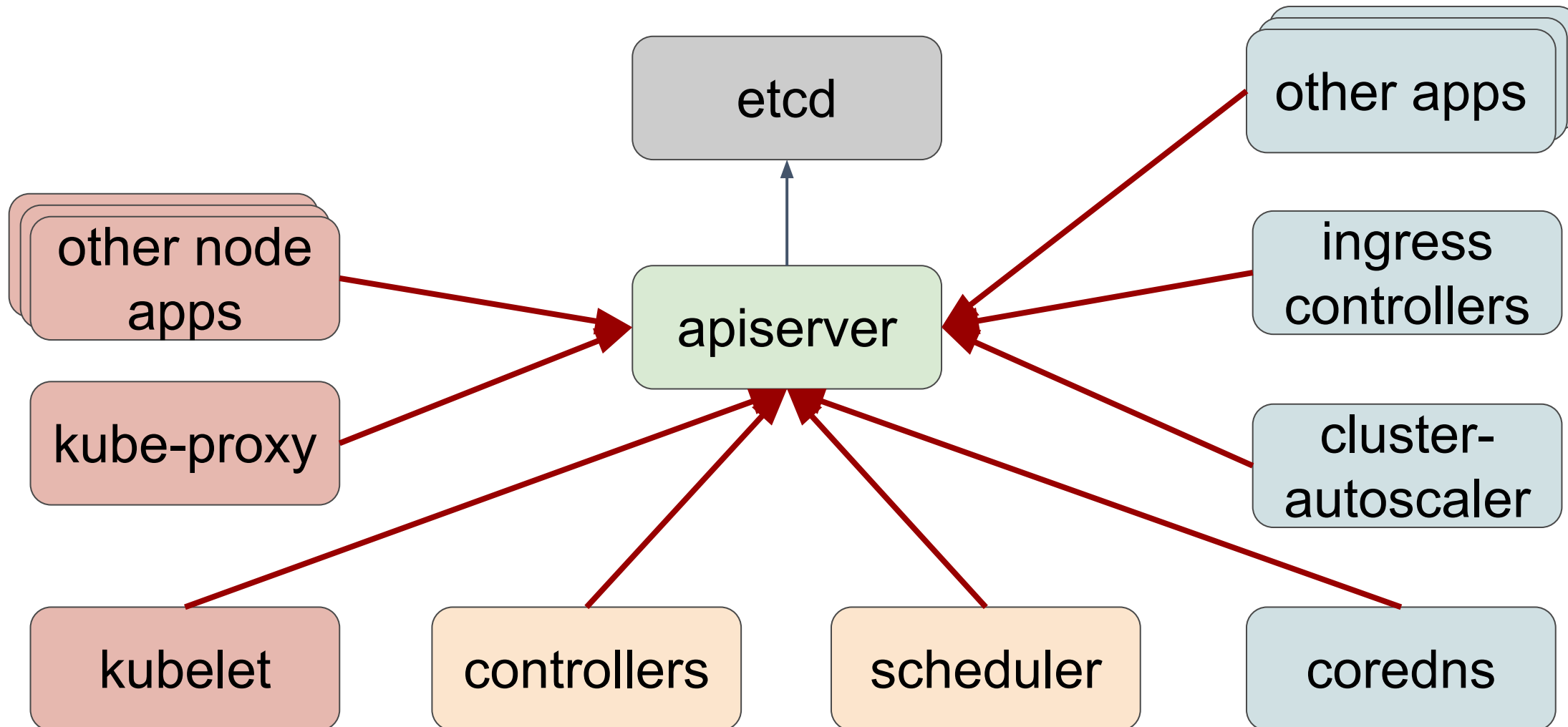


KubeCon



CloudNativeCon

North America 2019



And users, of course

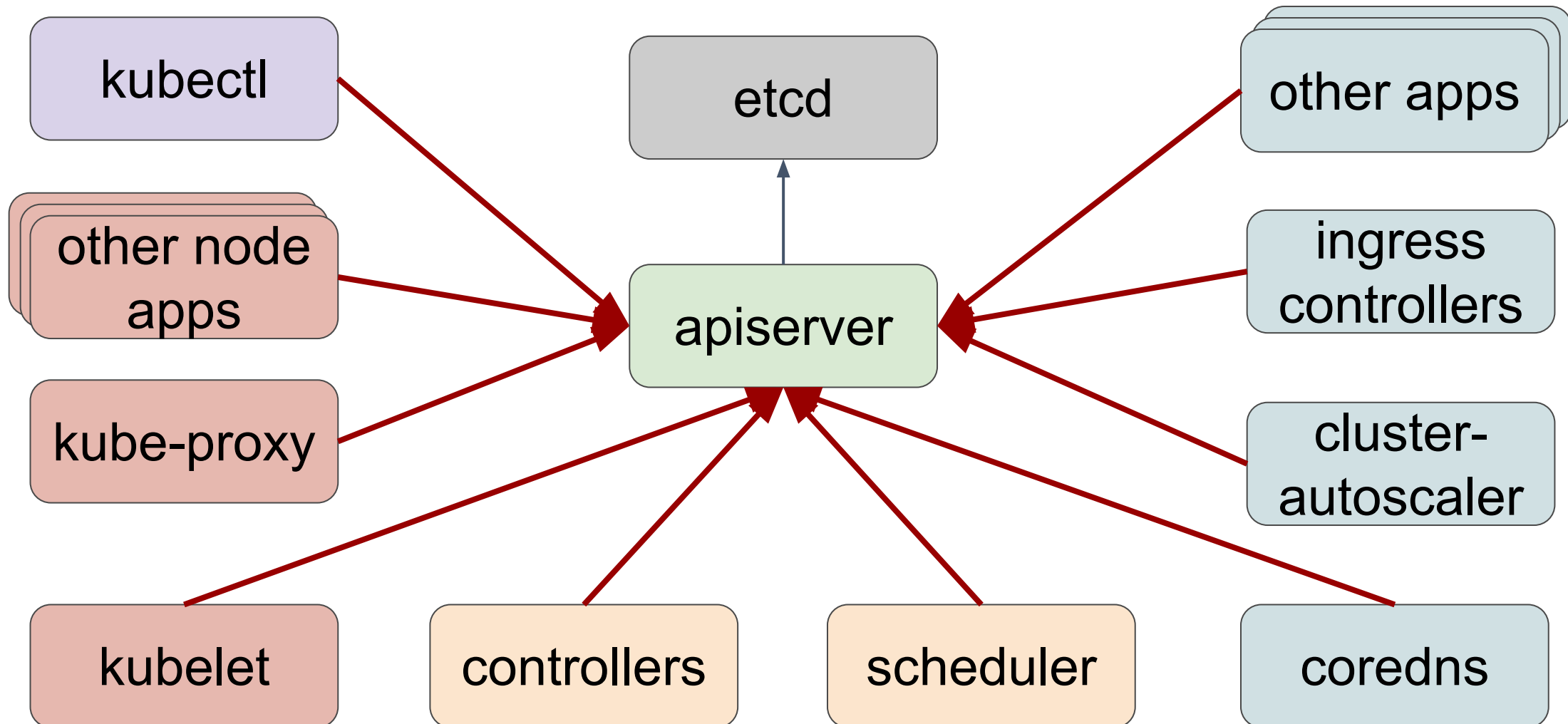


KubeCon



CloudNativeCon

North America 2019



And, surprise, the apiserver itself

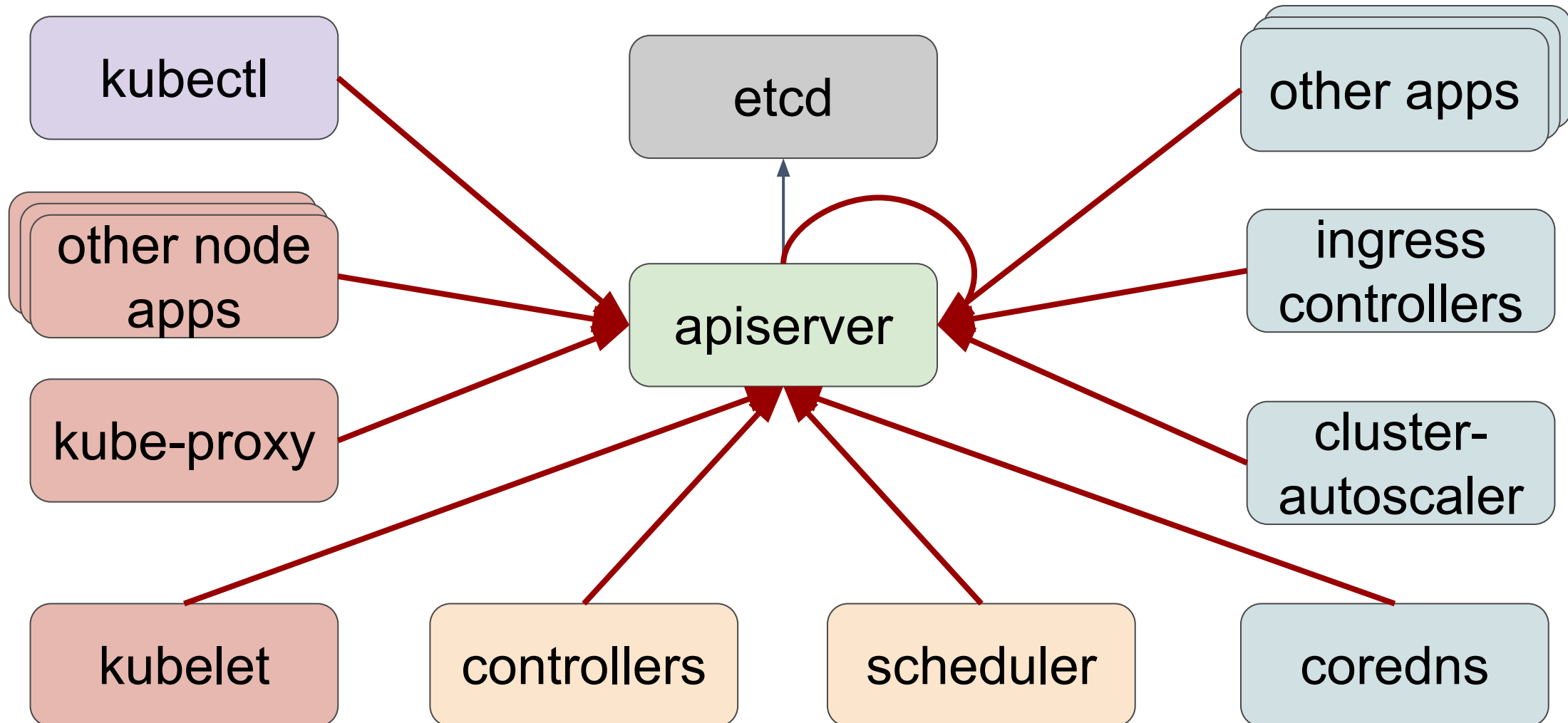


KubeCon



CloudNativeCon

North America 2019



What happens when you kubectl?



KubeCon



CloudNativeCon

North America 2019

```
$ kubectl get pod echodeploy-77cf5c6f6-brj76 -v=8  
[...]  
GET https://kubernetes.fury.us1.staging.dog/api/v1/namespaces/datadog/pods/echodeploy-77cf5c6f6-brj76
```

Let's look at details



KubeCon



CloudNativeCon

North America 2019

```
$ kubectl get pod echodeploy-77cf5c6f6-brj76 -v=8  
[...]  
GET https://kubernetes.fury.us1.staging.dog/api/v1/namespaces/datadog/pods/echodeploy-77cf5c6f6-brj76
```

apiserver

api version

namespace

resource type

resource name

A few more GET examples



KubeCon



CloudNativeCon

North America 2019

```
$ kubectl get pod echodeploy-77cf5c6f6-brj76 -v=8  
[...]  
GET https://kubernetes.fury.us1.staging.dog/api/v1/namespaces/datadog/pods/echodeploy-77cf5c6f6-brj76
```

```
$ kubectl get pods  
[...]  
GET https://kubernetes.fury.us1.staging.dog/api/v1/namespaces/datadog/pods?limit=500
```

List
(paginated)

A few more GET examples



KubeCon



CloudNativeCon

North America 2019

```
$ kubectl get pod echodeploy-77cf5c6f6-brj76 -v=8
```

```
[...]
```

```
GET https://kubernetes.fury.us1.staging.dog/api/v1/namespaces/datadog/pods/echodeploy-77cf5c6f6-brj76
```

```
$ kubectl get pods
```

```
[...]
```

```
GET https://kubernetes.fury.us1.staging.dog/api/v1/namespaces/datadog/pods?limit=500
```

```
$ kubectl get pods --watch=true -v=8
```

```
[...]
```

```
GET https://kubernetes.fury.us1.staging.dog/api/v1/namespaces/datadog/pods?limit=500
```

```
[...]
```

```
GET https://kubernetes.fury.us1.staging.dog/api/v1/namespaces/datadog/pods?resourceVersion=282725545&watch=true
```

List &
Watch

Describe resource



KubeCon



CloudNativeCon

North America 2019

```
kubectl describe pod echodeploy-77cf5c6f6-5wmw9 -v=8
[...]
GET https://kubernetes.fury.us1.staging.dog/api/v1/namespaces/datadog/pods/echodeploy-77cf5c6f6-5wmw9
[...]
GET https://kubernetes.fury.us1.staging.dog/api/v1/namespaces/datadog/events?
      fieldSelector=involvedObject.name=echodeploy-77cf5c6f6-5wmw9,
      involvedObject.namespace=datadog,
      involvedObject.uid=770b3a5e-0631-11ea-bc60-12d7306f3c0c
[...]
ResponseBody
{
  "kind": "EventList",
  "items": [
    {
      "involvedObject": { "kind": "Pod", "namespace": "datadog", "name": "echodeploy-77cf5c6f6-5wmw9"},
      "reason": "Scheduled",
      "message": "Successfully assigned echodeploy-77cf5c6f6-5wmw9 to ip-10-128-205-156.ec2.internal",
      "source": {
        "component": "default-scheduler"
      },
    },
  ],
}
```

Describe resource



KubeCon



CloudNativeCon

North America 2019

```
kubectl describe pod echodeploy-77cf5c6f6-5wmw9 -v=8
```

```
[...]
```

```
GET https://kubernetes.fury.us1.staging.dog/api/v1/namespaces/datadog/pods/echodeploy-77cf5c6f6-5wmw9
```

Get resource

```
[...]
```

```
GET https://kubernetes.fury.us1.staging.dog/api/v1/namespaces/datadog/events?
```

```
  fieldSelector=involvedObject.name=echodeploy-77cf5c6f6-5wmw9,
```

```
    involvedObject.namespace=datadog,
```

```
    involvedObject.uid=770b3a5e-0631-11ea-bc60-12d7306f3c0c
```

```
[...]
```

```
ResponseBody
```

```
{
```

```
  "kind": "EventList",
```

```
  "items": [
```

```
    {
```

```
      "involvedObject": { "kind": "Pod", "namespace": "datadog", "name": "echodeploy-77cf5c6f6-5wmw9"},
```

```
      "reason": "Scheduled",
```

```
      "message": "Successfully assigned echodeploy-77cf5c6f6-5wmw9 to ip-10-128-205-156.ec2.internal",
```

```
      "source": {
```

```
        "component": "default-scheduler"
```

```
    },
```

```
  ],
```

```
}
```

```
}
```

Describe resource



KubeCon



CloudNativeCon

North America 2019

```
kubectl describe pod echodeploy-77cf5c6f6-5wmw9 -v=8
```

```
[...]
```

```
GET https://kubernetes.fury.us1.staging.dog/api/v1/namespaces/datadog/pods/echodeploy-77cf5c6f6-5wmw9
```

```
[...]
```

```
GET https://kubernetes.fury.us1.staging.dog/api/v1/namespaces/datadog/events?
    fieldSelector=involvedObject.name=echodeploy-77cf5c6f6-5wmw9,
    involvedObject.namespace=datadog,
    involvedObject.uid=770b3a5e-0631-11ea-bc60-12d7306f3c0c
```

```
[...]
```

```
ResponseBody
```

```
{
  "kind": "EventList",
  "items": [
    {
      "involvedObject": { "kind": "Pod", "namespace": "datadog", "name": "echodeploy-77cf5c6f6-5wmw9"},
      "reason": "Scheduled",
      "message": "Successfully assigned echodeploy-77cf5c6f6-5wmw9 to ip-10-128-205-156.ec2.internal",
      "source": {
        "component": "default-scheduler"
      },
    },
  ],
}
```

Get resource

Get **events** associated with resource

A few other examples



KubeCon



CloudNativeCon

North America 2019

```
$ kubectl delete pod echodeploy-77cf5c6f6-brj76 -v=8
```

```
[...]
```

```
DELETE https://kubernetes.fury.us1.staging.dog/api/v1/namespaces/datadog/pods/echodeploy-77cf5c6f6-brj76
```

```
GET https://kubernetes.fury.us1.staging.dog/api/v1/namespaces/datadog/pods?fieldSelector=metadata.name%3Dechodeploy-77cf5c6f6-brj76
```

```
GET https://kubernetes.fury.us1.staging.dog/api/v1/namespaces/datadog/pods?fieldSelector=metadata.name%3Dechodeploy-77cf5c6f6-brj76&resourceVersion=282733788&watch=true
```

Delete
+
List &
Watch

A few other examples



KubeCon



CloudNativeCon

North America 2019

```
$ kubectl delete pod echodeploy-77cf5c6f6-brj76 -v=8
[...]
DELETE https://kubernetes.fury.us1.staging.dog/api/v1/namespaces/datadog/pods/echodeploy-77cf5c6f6-brj76
GET https://kubernetes.fury.us1.staging.dog/api/v1/namespaces/datadog/pods?fieldSelector=metadata.name%3Dechodeploy-77cf5c6f6-brj76
GET https://kubernetes.fury.us1.staging.dog/api/v1/namespaces/datadog/pods?fieldSelector=metadata.name%3Dechodeploy-77cf5c6f6-brj76&resourceVersion=282733788&watch=true
```

```
$ kubectl create deployment test --image=busybox -v=8
```

Request Body:

```
{"apiVersion":"apps/v1","kind":"Deployment","metadata":{"creationTimestamp":null,"labels":{"app":"test"},"name":"test"},"spec":{"replicas":1,"selector":{"matchLabels":{"app":"test"},"strategy":{"template":{"metadata":{"creationTimestamp":null,"labels":{"app":"test"},"spec":{"containers":[{"image":"busybox","name":"busybox","resources":{}}]}}}}},"status":{}}
```

Minimal
deployment
spec

```
POST https://kubernetes.fury.us1.staging.dog/apis/apps/v1/namespaces/datadog/deployments
```

POST call

A few other examples



KubeCon



CloudNativeCon

North America 2019

```
$ kubectl delete pod echodeploy-77cf5c6f6-brj76 -v=8
[...]
DELETE https://kubernetes.fury.us1.staging.dog/api/v1/namespaces/datadog/pods/echodeploy-77cf5c6f6-brj76
GET https://kubernetes.fury.us1.staging.dog/api/v1/namespaces/datadog/pods?fieldSelector=metadata.name%3Dechodeploy-77cf5c6f6-brj76
GET https://kubernetes.fury.us1.staging.dog/api/v1/namespaces/datadog/pods?fieldSelector=metadata.name%3Dechodeploy-77cf5c6f6-brj76&resourceVersion=282733788&watch=true
```

```
$ kubectl create deployment test --image=busybox -v=8
```

Request Body:

```
{"apiVersion":"apps/v1","kind":"Deployment","metadata":{"creationTimestamp":null,"labels":{"app":"test"},"name":"test"},"spec":{"replicas":1,"selector":{"matchLabels":{"app":"test"},"strategy":{"template":{"metadata":{"creationTimestamp":null,"labels":{"app":"test"},"spec":{"containers":[{"image":"busybox","name":"busybox","resources":{}}]}}}}},"status":{}}
```

```
POST https://kubernetes.fury.us1.staging.dog/apis/apps/v1/namespaces/datadog/deployments
```

```
$ kubectl scale deploy test --replicas=2 -v=8
```

```
GET https://kubernetes.fury.us1.staging.dog/apis/extensions/v1beta1/namespaces/datadog/deployments/test
```

GET current

```
Request Body: {"spec":{"replicas":2}}
```

```
PATCH https://kubernetes.fury.us1.staging.dog/apis/extensions/v1beta1/namespaces/datadog/deployments/test/scale
```

PATCH body
+call

Takeaways



KubeCon



CloudNativeCon

North America 2019

- A lot of components are making calls
 - Control plane: controllers, scheduler
 - Node daemons: kubelet, kube-proxy
 - Other controllers: autoscaler, ingress
- “Simple” user ops translate to **many** API calls

How can we understand what is going on?

Outline



KubeCon



CloudNativeCon

North America 2019

1. Background: The Kubernetes API
- 2. Audit Logs**
3. Configuring Audit Logs
4. 10000 foot view for a large cluster
5. Understanding Kubernetes Internals
6. Troubleshooting examples



KubeCon



CloudNativeCon

North America 2019

Audit Logs



What are Audit Logs?



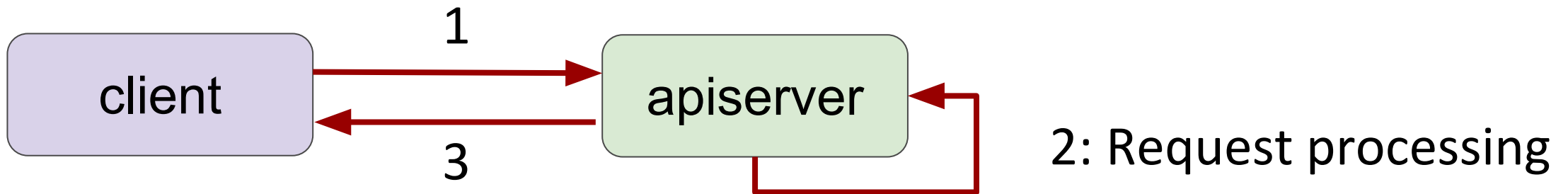
KubeCon



CloudNativeCon

North America 2019

- **Rich Structured json logs** output by the apiserver
- **Configurable Verbosity** for each resource
- Logging can happen at **different processing stages**



1: Apiserver receives request, Stage: RequestReceived

2: Apiserver processes request

3: Apiserver answers, Stage: ResponseComplete/ResponseStarted

Content of Audit Logs



KubeCon



CloudNativeCon

North America 2019

- What happened?
- Who initiated it?
- Why was it authorized?
- When did it happen?
- From where?
- Depending on verbosity, Request/Response

GET example from earlier



KubeCon



CloudNativeCon

North America 2019

```
$ kubectl get pod echodeploy-77cf5c6f6-5wmw9 -v=8
```

```
GET https://kubernetes.fury.us1.staging.dog/api/v1/namespaces/datadog/pods/echodeploy-77cf5c6f6-5wamw9
```

```
{
  apiVersion: audit.k8s.io/v1beta1,
  auditID: 0fe92a67-248e-461e-b5ec-af6e90f81c49,
  http: {
    status_code: 200
  },
  kind: Event,
  level: Metadata,
  metadata: {
    creationTimestamp: 2019-11-19T13:49:30Z
  },
  objectRef: {
    apiVersion: v1,
    name: echodeploy-77cf5c6f6-5wmw9,
    namespace: datadog,
    resource: pods
  },
  requestReceivedTimestamp: 2019-11-19T13:49:30.882458Z,
  requestURI: /api/v1/namespaces/datadog/pods/echodeploy-77cf5c6f6-5wmw9,
  responseStatus: {
    code: 200
  },
  sourceIPs: [
    10.128.215.38
  ],
  stage: ResponseComplete,
  stageTimestamp: 2019-11-19T13:49:30.885163Z,
  timestamp: 2019-11-19T13:49:30Z,
  user: {
    groups: [
      datadoghq.com,
      system:authenticated
    ],
    username: laurent.bernaille@datadoghq.com
  },
  verb: get
}
```

Structured JSON log
A lot of information

What happened?



KubeCon



CloudNativeCon

North America 2019

```
$ kubectl get pod echodeploy-77cf5c6f6-5wmw9 -v=8
```

```
GET https://kubernetes.fury.us1.staging.dog/api/v1/namespaces/datadog/pods/echodeploy-77cf5c6f6-5wmw9
```

requestURI

/api/v1/namespaces/datadog/pods/echodeploy-77cf5c6f6-5wmw9

verb

get

```
objectRef {
```

```
  apiVersion  v1
```

```
  name        echodeploy-77cf5c6f6-5wmw9
```

```
  namespace   datadog
```

```
  resource    pods
```

```
}
```

```
responseStatus {
```

```
  code 200
```

```
}
```

Who initiated it?



KubeCon



CloudNativeCon

North America 2019

```
$ kubectl get pod echodeploy-77cf5c6f6-5wmw9 -v=8
```

```
GET https://kubernetes.fury.us1.staging.dog/api/v1/namespaces/datadog/pods/echodeploy-77cf5c6f6-5wamw9
```

```
user {  
  groups [  
    datadoghq.com,  
    system:authenticated  
  ]  
  username laurent.bernaille@datadoghq.com  
}
```

User was **laurent.bernaille@datadoghq.com** and mapped to groups

- **datadoghq.com**
- **system:authenticated**

Why was it authorized?



KubeCon



CloudNativeCon

North America 2019

```
$ kubectl get pod echodeploy-77cf5c6f6-5wmw9 -v=8
```

```
GET https://kubernetes.fury.us1.staging.dog/api/v1/namespaces/datadog/pods/echodeploy-77cf5c6f6-5wamw9
```

```
authorization {
```

```
  k8s {
```

```
    io/decision allow
```

```
    io/reason RBAC: allowed by ClusterRoleBinding "datadoghq:cluster-admin-binding" of ClusterRole "datadoghq:cluster-user" to Group "datadoghq.com"
```

```
  }
```

```
}
```

```
}
```

It was authorized because group **datadoghq.com** is bound to role **datadoghq:cluster-user** by ClusterRoleBinding **datadoghq:cluster-admin-binding** (and this role has the required permissions)

When, and from where?



KubeCon



CloudNativeCon

North America 2019

```
$ kubectl get pod echodeploy-77cf5c6f6-5wmw9 -v=8
```

```
GET https://kubernetes.fury.us1.staging.dog/api/v1/namespaces/datadog/pods/echodeploy-77cf5c6f6-5wamw9
```

```
requestReceivedTimestamp 2019-11-13T20:33:26.757736Z
```

```
stage ResponseComplete
```

```
stageTimestamp 2019-11-13T20:33:26.771303Z
```

```
sourceIPs [
```

```
  10.1 . . . . . 74
```

```
]
```

Request received at **20:33:26.757**

Response completed at **20:33:26:771**

Duration: **14ms**

From IP: **10.X.Y.74**

Another GET call



KubeCon



CloudNativeCon

North America 2019

```
$ kubectl get pods -v=8
```

```
GET https://kubernetes.fury.us1.staging.dog/api/v1/namespaces/datadog/pods?limit=500
```

```
requestURI      /api/v1/namespaces/datadog/pods?limit=500
```

```
verb            list
```

```
objectRef {
```

```
  apiVersion    v1
```

```
  namespace     datadog
```

```
  resource      pods
```

```
}
```

GET is mapped to different verbs (get/list)

Watches



KubeCon



CloudNativeCon

North America 2019

```
$ kubectl get pods -v=8 -w
```

```
GET https://kubernetes.fury.us1.staging.dog/api/v1/namespaces/datadog/pods?limit=500
```

```
GET https://kubernetes.fury.us1.staging.dog/api/v1/namespaces/datadog/pods?resourceVersion=288656279&watch=true
```

```
objectRef {
  apiVersion v1
  namespace datadog
  resource pods
}
requestReceivedTimestamp 2019-11-19T14:17:06.162487Z
requestURI /api/v1/namespaces/datadog/pods?limit=500
+ responseStatus {...}
+ sourceIPs {...}
stage ResponseComplete
stageTimestamp 2019-11-19T14:17:06.165555Z
timestamp 2019-11-19T14:17:06Z
+ user {...}
verb list
```

Call 1 : list

stage: ResponseComplete

```
objectRef {
  apiVersion v1
  namespace datadog
  resource pods
}
requestReceivedTimestamp 2019-11-19T14:17:06.301721Z
requestURI /api/v1/namespaces/datadog/pods?resourceVersion=288656279&watch=true
+ responseStatus {...}
+ sourceIPs {...}
stage ResponseStarted
stageTimestamp 2019-11-19T14:17:06.302044Z
timestamp 2019-11-19T14:17:06Z
+ user {...}
verb watch
```

Call 2 : watch

get + watch parameters

136ms later

stage: ResponseStarted



@lbernail @roboll_

Create call



KubeCon



CloudNativeCon

North America 2019

```
$ kubectl create deployment test --image=busybox -v=8  
POST https://kubernetes.fury.us1.staging.dog/apis/apps/v1/namespaces/datadog/deployments
```

```
objectRef {  
  apiGroup  apps  
  apiVersion v1  
  name      test  
  namespace datadog  
  resource  deployments  
}  
+ requestObject {...}  
  requestReceivedTimestamp 2019-11-19T14:39:04.850905Z  
  requestURI                /apis/apps/v1/namespaces/datadog/deployments  
- responseStatus {  
  code 201  
}  
+ sourceIPs [...]  
  stage           ResponseComplete  
  stageTimestamp  2019-11-19T14:39:04.870988Z  
  timestamp       2019-11-19T14:39:04Z  
+ user {...}  
  verb create
```

Create call



KubeCon



CloudNativeCon

North America 2019

```
$ kubectl create deployment test --image=busybox -v=8  
POST https://kubernetes.fury.us1.staging.dog/apis/apps/v1/namespaces/datadog/deployments
```

```
objectRef {  
  apiGroup  apps  
  apiVersion v1  
  name      test  
  namespace datadog  
  resource  deployments  
}
```

```
+ requestObject {...}
```

```
requestReceivedTimestamp 2019-11-19T14:39:04.850905Z  
requestURI                /apis/apps/v1/namespaces/datadog/deployments
```

```
responseStatus {  
  code 201  
}
```

```
+ sourceIPs [...]
```

```
stage           ResponseComplete  
stageTimestamp  2019-11-19T14:39:04.870988Z  
timestamp       2019-11-19T14:39:04Z
```

```
+ user {...}
```

```
verb           create
```

```
requestObject {  
  apiVersion apps/v1  
  kind       Deployment  
  metadata {  
    creationTimestamp null  
    labels {...}  
    name             test  
  }  
  spec {  
    progressDeadlineSeconds 600  
    replicas                 1  
    revisionHistoryLimit    10  
    selector {...}  
    strategy {...}  
    template {  
      metadata {  
        creationTimestamp null  
        labels {  
          app test  
        }  
      }  
      spec {  
        containers [  
          {"image":"busybox","imagePullPolicy":"Always","terminationMessagePolicy":"File","terminationMessagePath":"/dev/termination-log","name":"busybox","resources":{}}  
        ]  
        dnsPolicy           ClusterFirst  
        restartPolicy       Always  
        schedulerName       default-scheduler  
        terminationGracePeriodSeconds 30  
      }  
    }  
  }  
}
```

Takeaways



KubeCon



CloudNativeCon

North America 2019

Audit logs contain information on all API calls

- What happened?
- Who initiated it?
- Why was it authorized?
- When did it happen?
- From where?
- Depending on verbosity, Request/Response

OK, how do I get them?

Outline



KubeCon



CloudNativeCon

North America 2019

1. Background: The Kubernetes API
2. Audit Logs
- 3. Configuring Audit Logs**
4. 10000 foot view for a large cluster
5. Understanding Kubernetes Internals
6. Troubleshooting examples



KubeCon



CloudNativeCon

North America 2019

Configuring Audit Logs



Apiserver configuration



KubeCon



CloudNativeCon

North America 2019

Minimum configuration

```
kube-apiserver
```

```
[...]
```

```
--audit-log-path=/var/log/kubernetes/apiserver/audit.log
```

Where to store them

```
--audit-policy-file=/etc/kubernetes/audit-policies/policy.yaml
```

What to collect

Advanced

- Rotation parameters (max size, backup options)
- Alternative backend (webhook)
- Batching mode

Audit policy: what to log?



KubeCon



CloudNativeCon

North America 2019

```
apiVersion: audit.k8s.io/v1 kind: Policy
```

```
rules:
```

```
# Log pod changes at RequestResponse level  
- level: RequestResponse  
  omitStages:  
  - "RequestReceived"  
  resources:  
  - group: "" # core API group  
    resources: ["pods"]  
  verbs: ["create", "patch", "update", "delete"]
```

```
# Log "pods/log", "pods/status" at Metadata level  
- level: Metadata  
  omitStages:  
  - "RequestReceived"  
  resources:  
  - group: ""  
    resources: ["pods/log", "pods/status"]
```

Set of rules

Audit policy: what to log?



KubeCon



CloudNativeCon

North America 2019

```
apiVersion: audit.k8s.io/v1 kind: Policy
```

```
rules:
```

```
# Log pod changes at RequestResponse level
```

```
- level: RequestResponse
```

```
omitStages:
```

```
- "RequestReceived"
```

```
resources:
```

```
- group: "" # core API group
```

```
resources: ["pods"]
```

```
verbs: ["create", "patch", "update", "delete"]
```

```
# Log "pods/log", "pods/status" at Metadata level
```

```
- level: Metadata
```

```
omitStages:
```

```
- "RequestReceived"
```

```
resources:
```

```
- group: ""
```

```
resources: ["pods/log", "pods/status"]
```

Rules match api call

- api group / version
- resource
- verbs

> Similar to RBAC

Audit policy: when to log?



KubeCon



CloudNativeCon

North America 2019

```
apiVersion: audit.k8s.io/v1 kind: Policy
```

```
rules:
```

```
# Log pod changes at RequestResponse level
```

```
- level: RequestResponse
```

```
  omitStages:
```

```
    - "RequestReceived"
```

```
  resources:
```

```
    - group: "" # core API group
```

```
      resources: ["pods"]
```

```
    verbs: ["create", "patch", "update", "delete"]
```

```
# Log "pods/log", "pods/status" at Metadata level
```

```
- level: Metadata
```

```
  omitStages:
```

```
    - "RequestReceived"
```

```
  resources:
```

```
    - group: ""
```

```
      resources: ["pods/log", "pods/status"]
```

For matching API calls

- Which verbosity?
- When? (stage)

Gotchas



KubeCon



CloudNativeCon

North America 2019

```
apiVersion: audit.k8s.io/v1 kind: Policy

rules:

# Log pod changes at RequestResponse level
- level: RequestResponse
  omitStages:
  - "RequestReceived"
  resources:
  - group: "" # core API group
    resources: ["pods"]
  verbs: ["create", "patch", "update", "delete"]

# Log "pods/log", "pods/status" at Metadata level
- level: Metadata
  omitStages:
  - "RequestReceived"
  resources:
  - group: ""
    resources: ["pods/log", "pods/status"]
```



Rules are evaluated in order
First matching rule sets level

Request/RequestResponse
> contain payload data
Careful with security implications
ex: tokenreviews calls

group: "" means core API only
Don't forget to add

- 3rd party apiservices
- your apiservices

Recommendations



KubeCon



CloudNativeCon

North America 2019

- Ignore **RequestReceived** stage
- Use at least **Metadata** level for almost everything
 - Possibly ignore healthz, metrics
- Use **Request/Response** level for critical resource/verbs
 - Very valuable for retroactive debugging
 - Careful for large/sensitive request/response bodies
- Very complete example in GKE
 - <https://github.com/kubernetes/kubernetes/blob/master/cluster/gce/gci/configure-helper.sh>
- Documentation
 - <https://kubernetes.io/docs/tasks/debug-application-cluster/audit/#audit-policy>

Takeaways



KubeCon



CloudNativeCon

North America 2019

- Getting audit logs is “simple”: 2 flags
- Getting policies right is harder
- You will get **a lot** of logs
- Requires a solution to analyze them

Let's look at an overview on a real large cluster

Outline



KubeCon



CloudNativeCon

North America 2019

1. Background: The Kubernetes API
2. Audit Logs
3. Configuring Audit Logs
4. **10000 foot view for a large cluster**
5. Understanding Kubernetes Internals
6. Troubleshooting examples



KubeCon



CloudNativeCon

North America 2019

10000 foot view for a large cluster



Total number of API calls



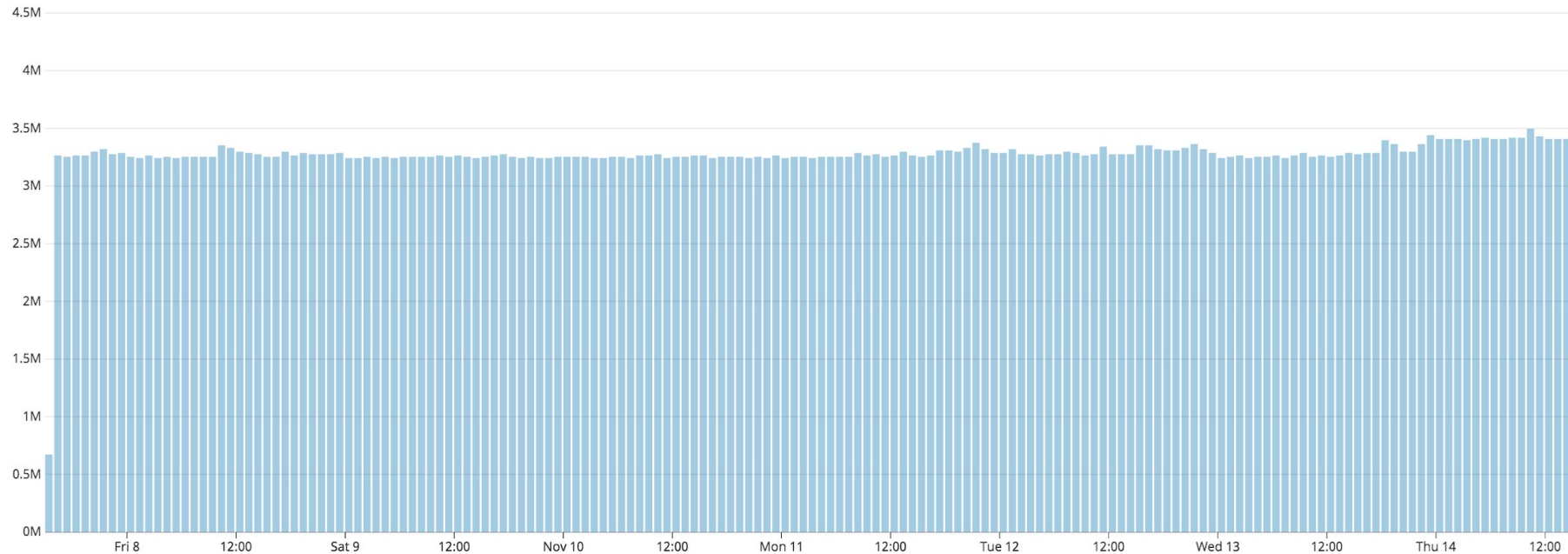
KubeCon



CloudNativeCon

North America 2019

Number of audit logs per hour



~900 calls/second on this 2500 nodes cluster

Top API users?



KubeCon

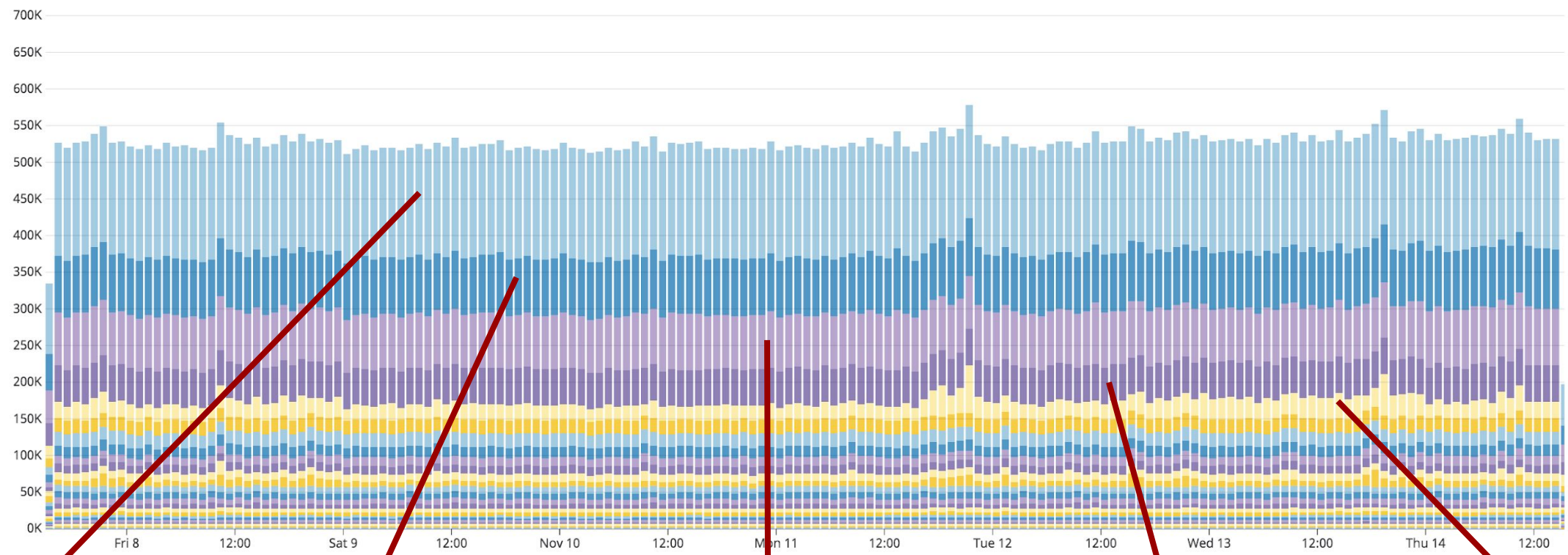


CloudNativeCon

North America 2019

Count * group by User Username X roll up every 1h limit to top 25

Hide Controls Timeseries Display: Bars Color: Classic Stacked as: Value



apiserver: 40 rps

kube-proxy: 20 rps

local-volume-provisioner:20 rps

cronjob-controller: 15 rps

spinnaker: 5-10 rps

@lbernail @roboll_

Top list, missing “small” users



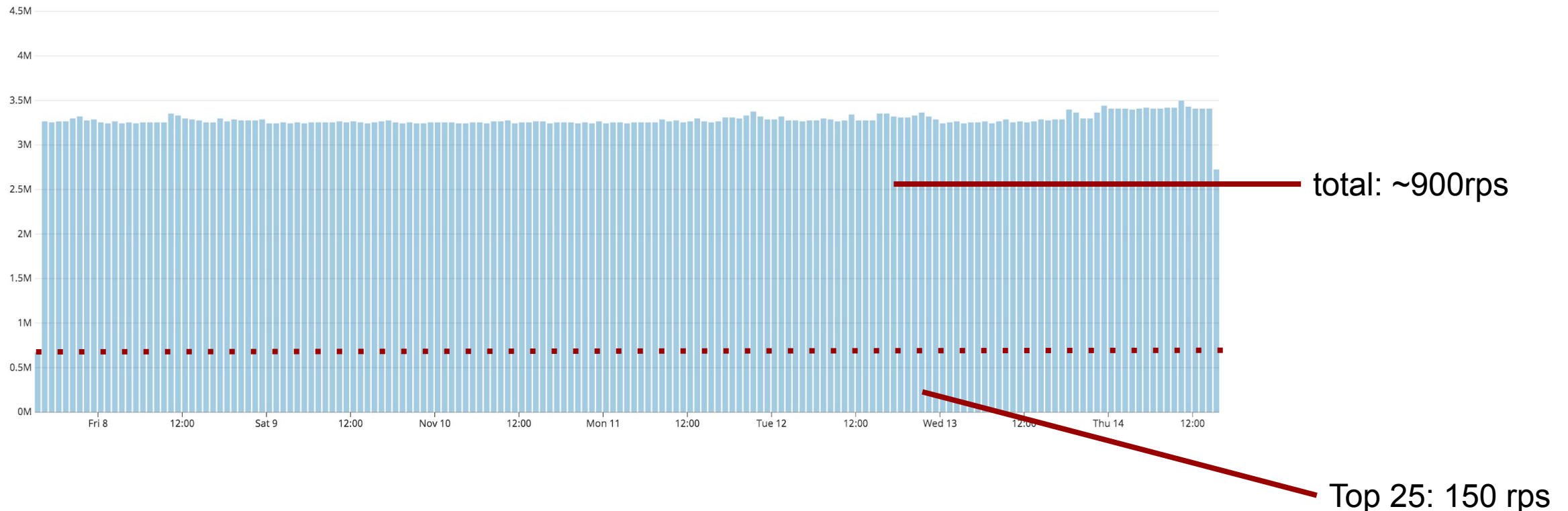
KubeCon



CloudNativeCon

North America 2019

Total number of API calls



900 calls/second on this 2500 nodes cluster
What is doing ~80% of API calls?

Grouping by users is not helping



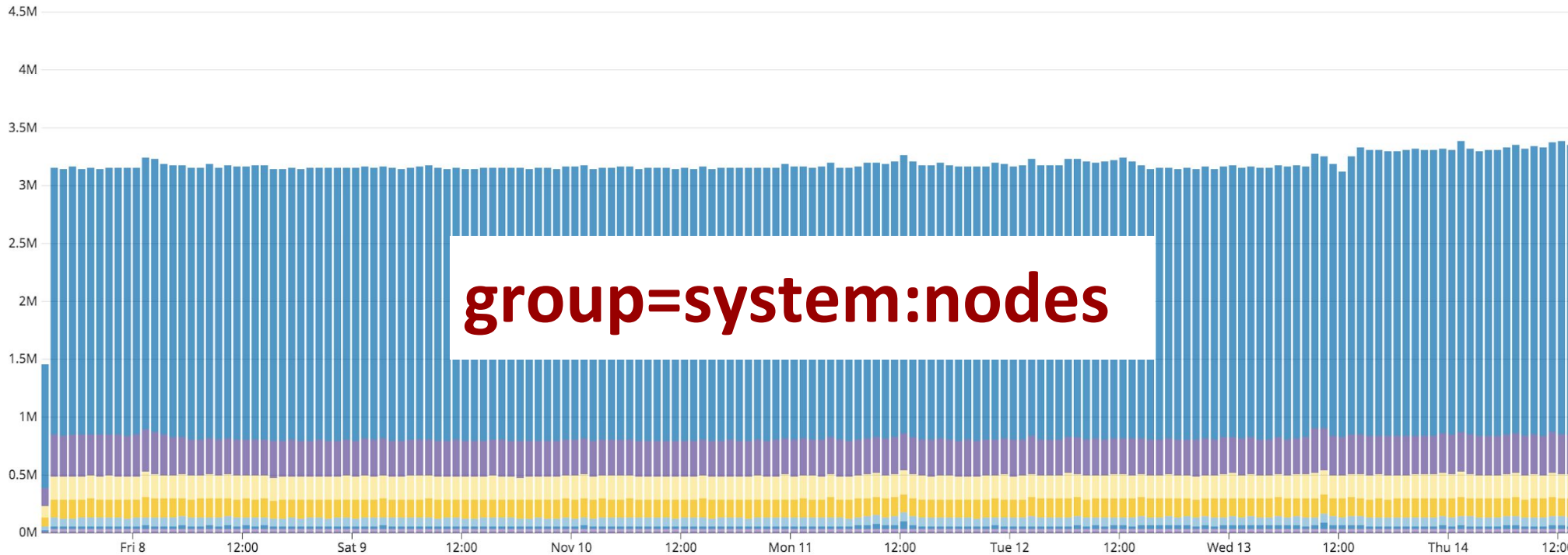
KubeCon



CloudNativeCon

North America 2019

Calls by user group



Calls from users in group “system:nodes”: 750 rps (~80% of API calls)

In this 2500-nodes cluster, this means **0.3 rps per node!**

Why is “system:nodes” so high?



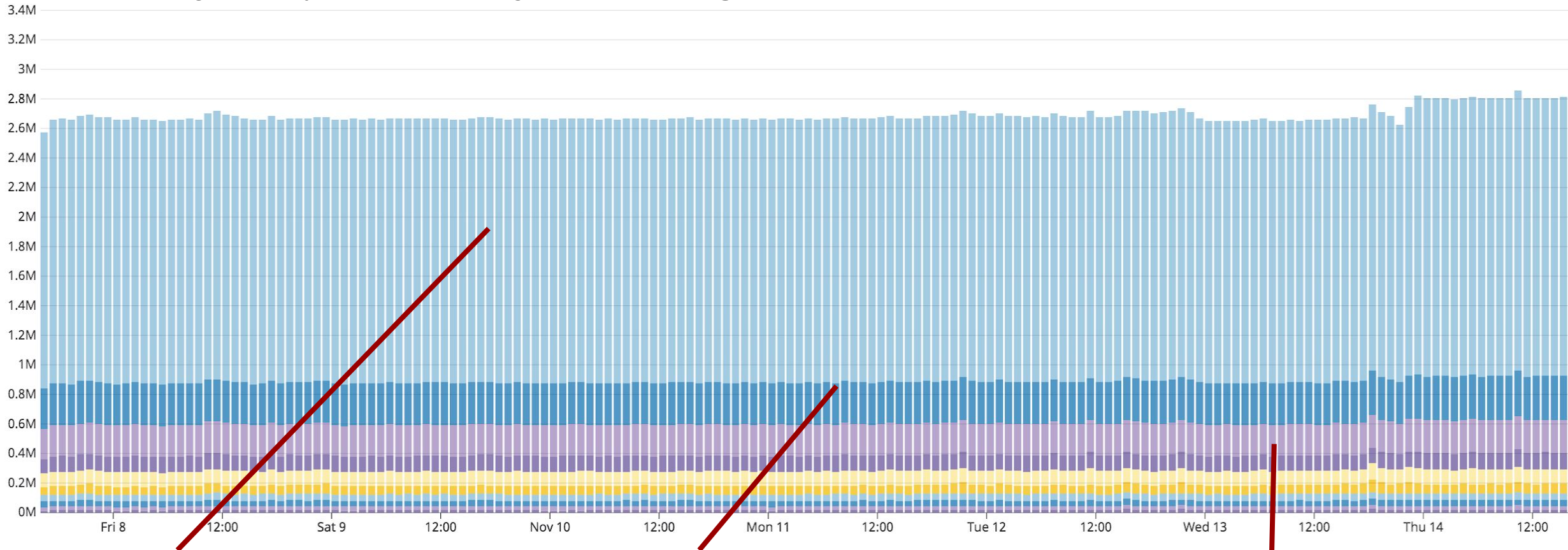
KubeCon



CloudNativeCon

North America 2019

Calls from group “system:nodes” by resource targeted



nodes: 500 rps

configmaps: 75 rps

secrets: 60 rps

Verbs on “node” for a kubelet

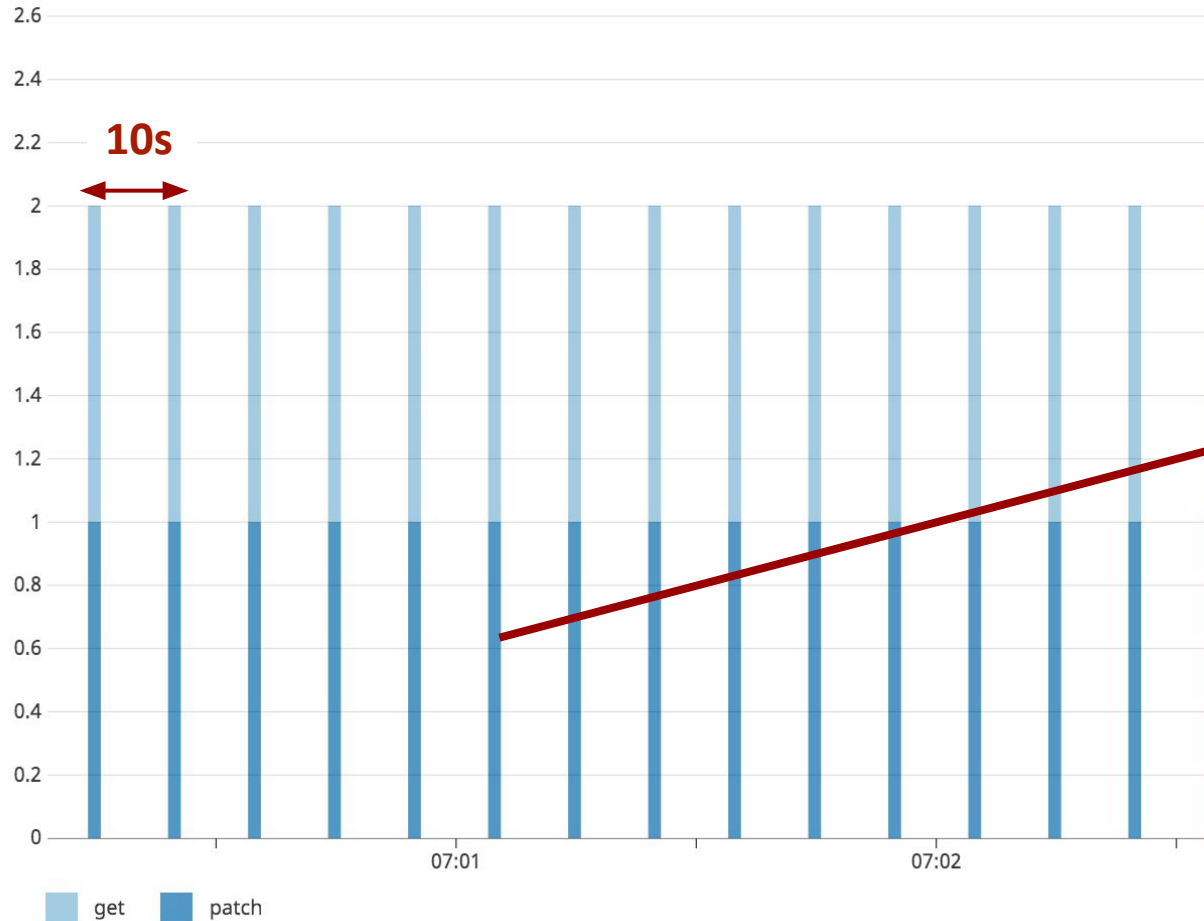


KubeCon



CloudNativeCon

North America 2019



```
user {
  groups [
    system:nodes,
    system:authenticated
  ]
  username system:node:eu1-prod-dog-app3-k8s-nodepool-vm-c602a136b569f8c0-b0j7
}
kind Event
level Metadata
metadata {
  creationTimestamp 2019-11-19T15:03:04Z
}
objectRef {
  apiVersion v1
  name eu1-prod-dog-app3-k8s-nodepool-vm-c602a136b569f8c0-b0j7
  resource nodes
  subresource status
}
verb patch
```

Each node update its status every 10s

Verbs on “configmaps”

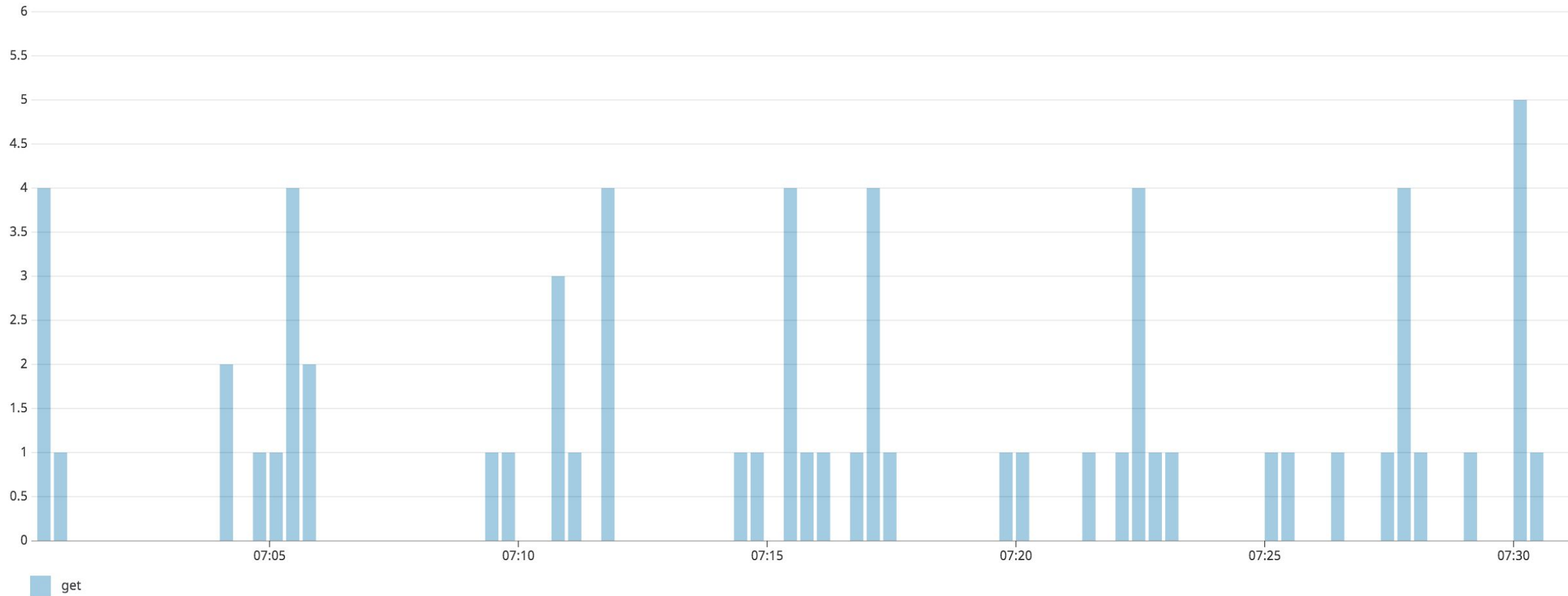


KubeCon



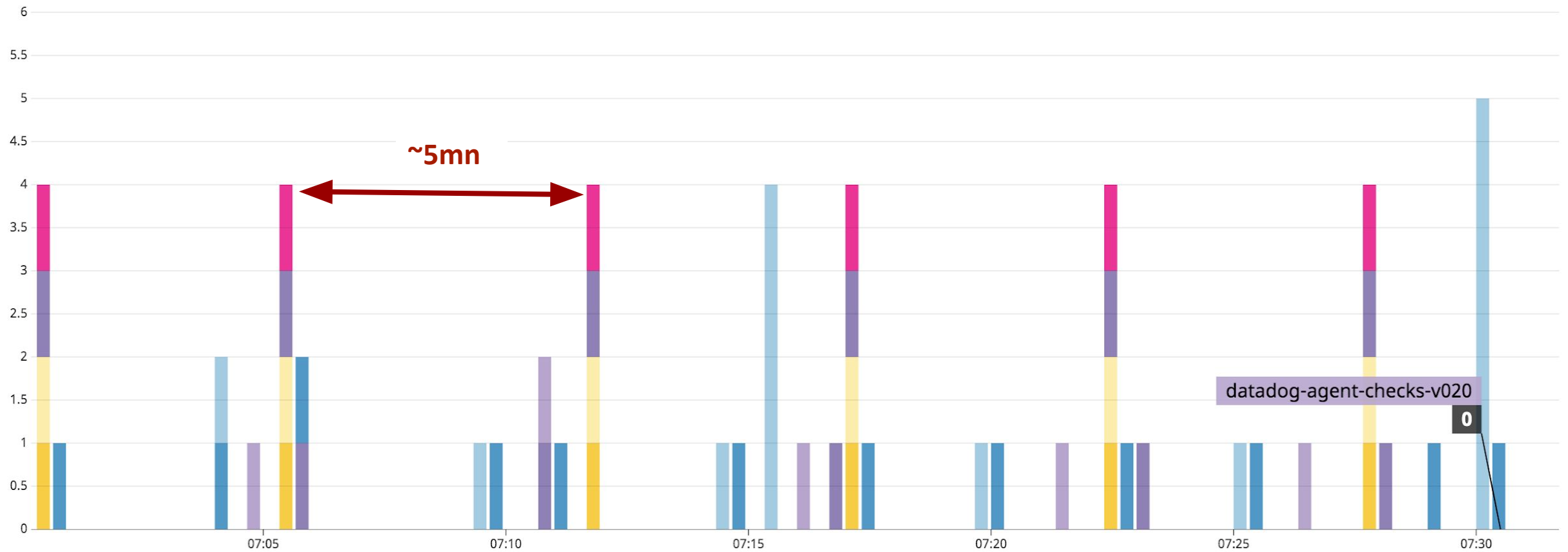
CloudNativeCon

North America 2019



Only GETs
Regularity but not clear pattern

Group by resource name



Each configmap is refreshed every ~5mn => GET call
Similar for secrets

List latency by resource



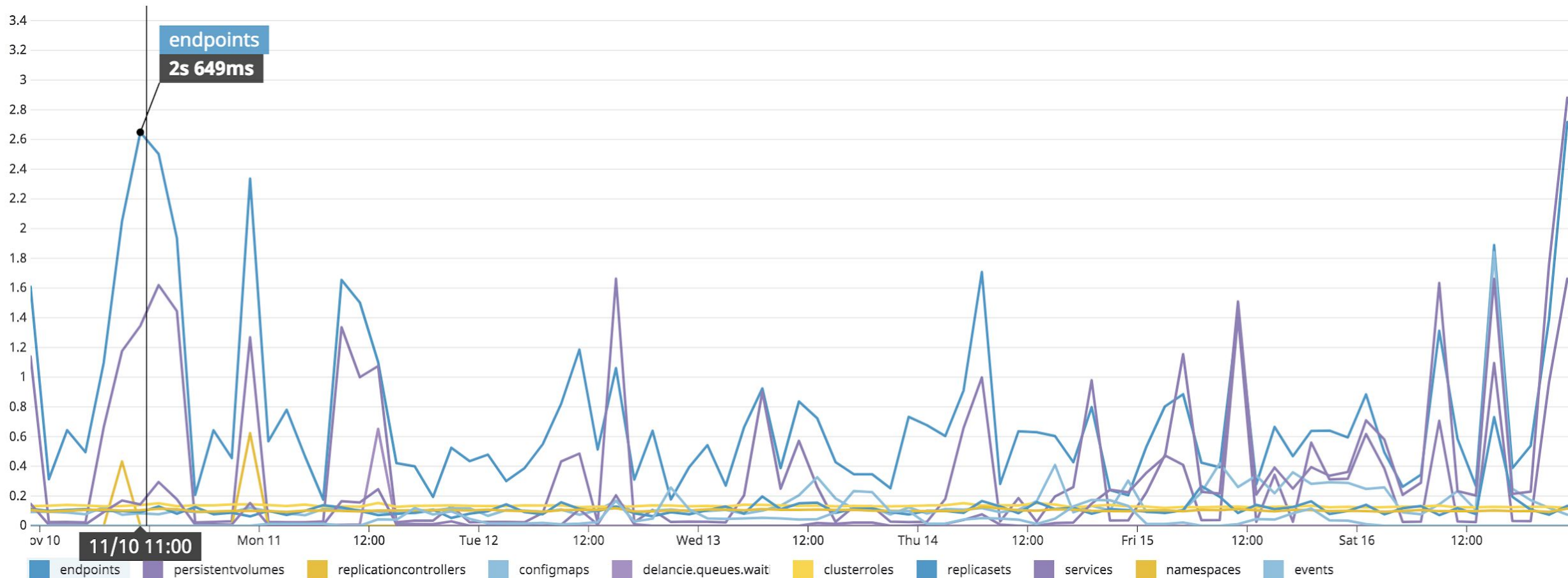
KubeCon



CloudNativeCon

North America 2019

Latency for list by resource



List latency by resource



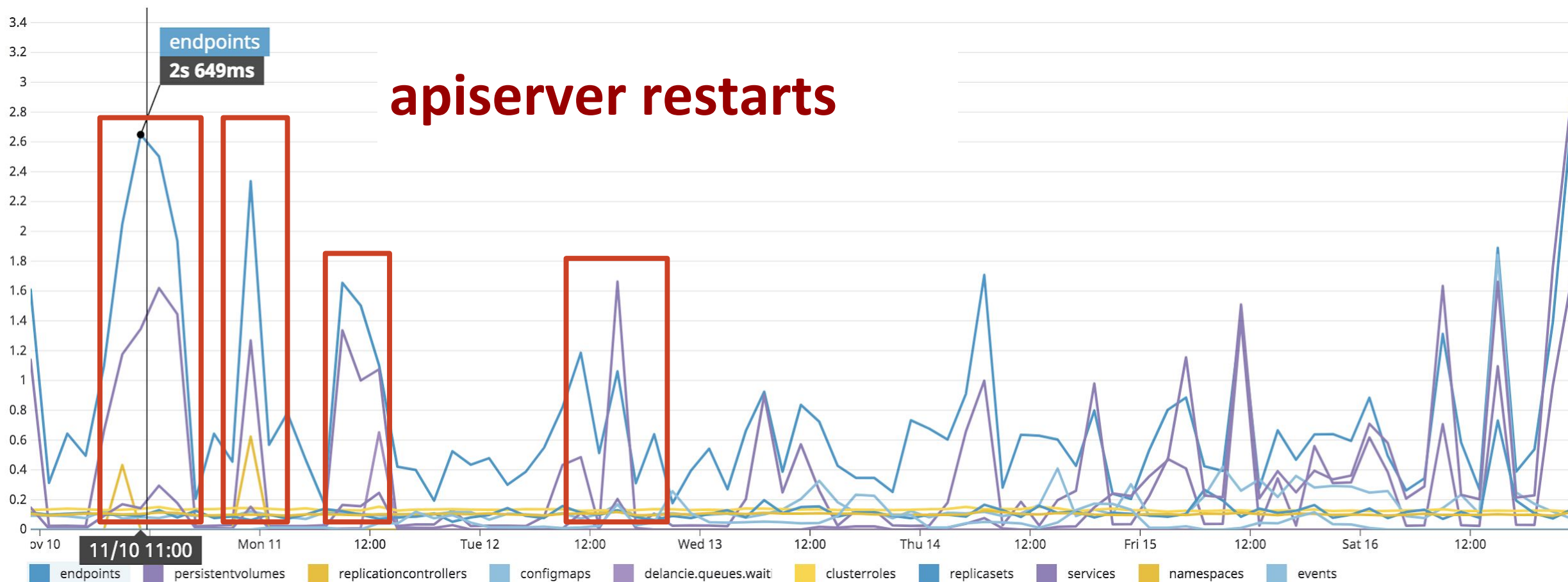
KubeCon



CloudNativeCon

North America 2019

Latency for list by resource



Compare cluster performances



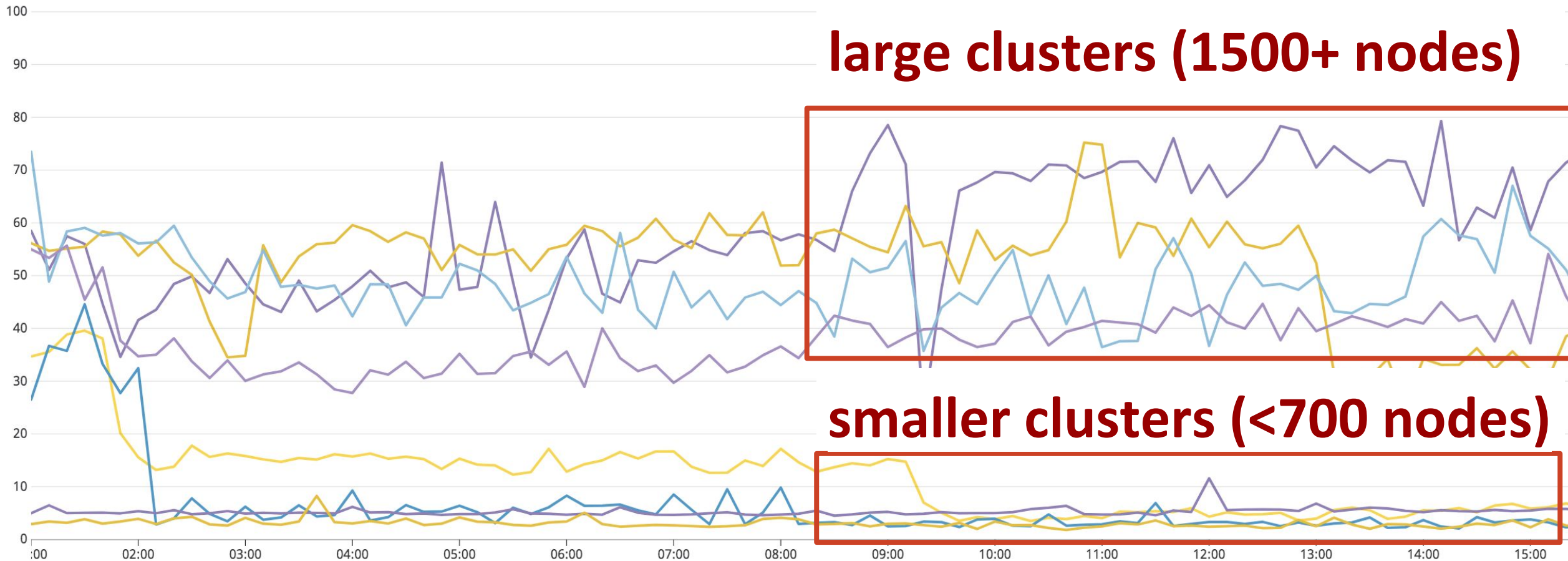
KubeCon



CloudNativeCon

North America 2019

Latency for get pod by cluster (ms)



large clusters (1500+ nodes)

smaller clusters (<700 nodes)

Takeaways



KubeCon



CloudNativeCon

North America 2019

- Biggest users are the ones running on each node
 - kubelet, daemonsets (kube-proxy)
 - A lot of effort upstream to reduce their load
 - Be extra careful of daemonsets doing API calls
- Audit logs structure allow to filter and slice & dice
- Audit logs are verbose (1000 logs/s in our example)

Outline



KubeCon



CloudNativeCon

North America 2019

1. Background: The Kubernetes API
2. Audit Logs
3. Configuring Audit Logs
4. 10000 foot view for a large cluster
- 5. Understanding Kubernetes Internals**
6. Troubleshooting examples



KubeCon



CloudNativeCon

North America 2019

Understanding Kubernetes Internals



Creating a simple deployment



KubeCon



CloudNativeCon

North America 2019

```
apiVersion: apps/v1
kind: Deployment
metadata:
  name: nginx
  labels:
    app: nginx

spec:
  replicas: 2
  selector:
    matchLabels:
      app: nginx
  template:
    metadata:
      labels:
        app: nginx

spec:
  containers:
    - name: nginx
      image: nginx
```

Sequence

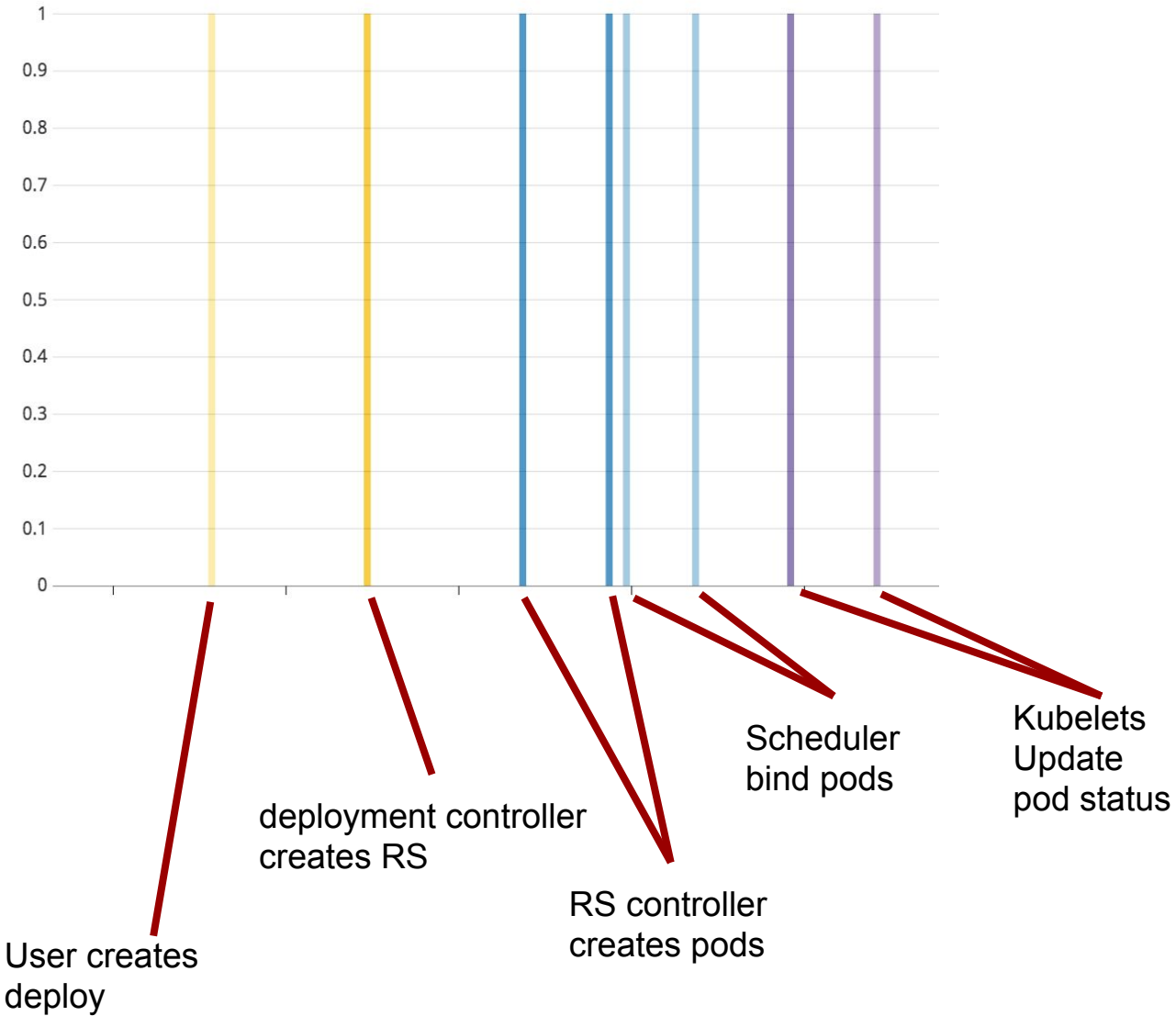


KubeCon



CloudNativeCon

North America 2019



Sequence

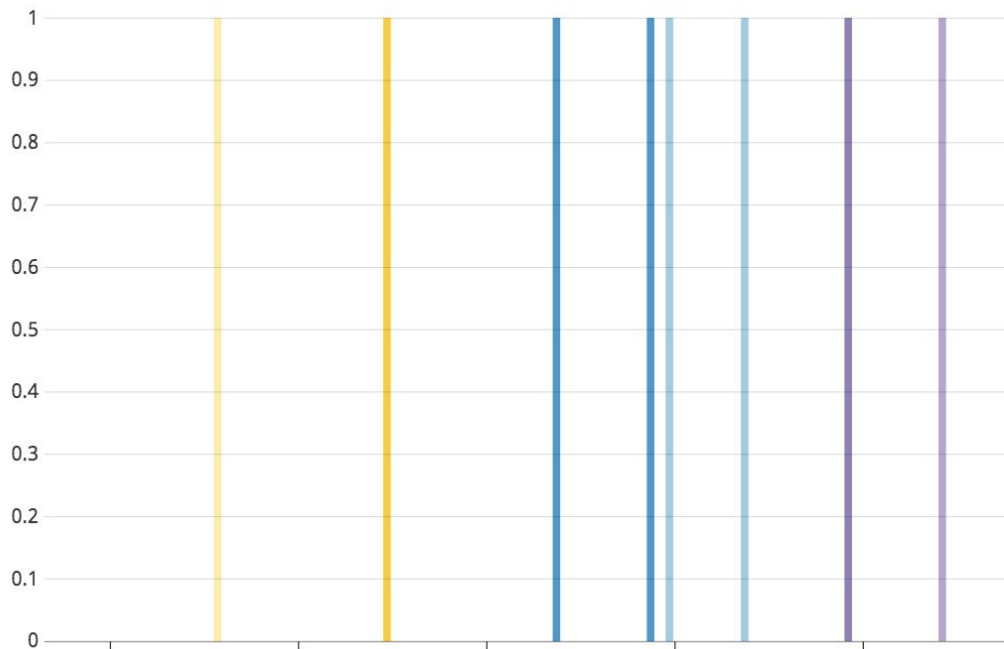


KubeCon



CloudNativeCon

North America 2019



Scheduler
bind pods

Create

```
http {
  method    create
  status_code 201
  url_details {
    path /api/v1/namespaces/demo/pods/nginx-68c5f9f877-29qsb/binding
  }
}
```

Scheduler call

Binding
For nginx pod

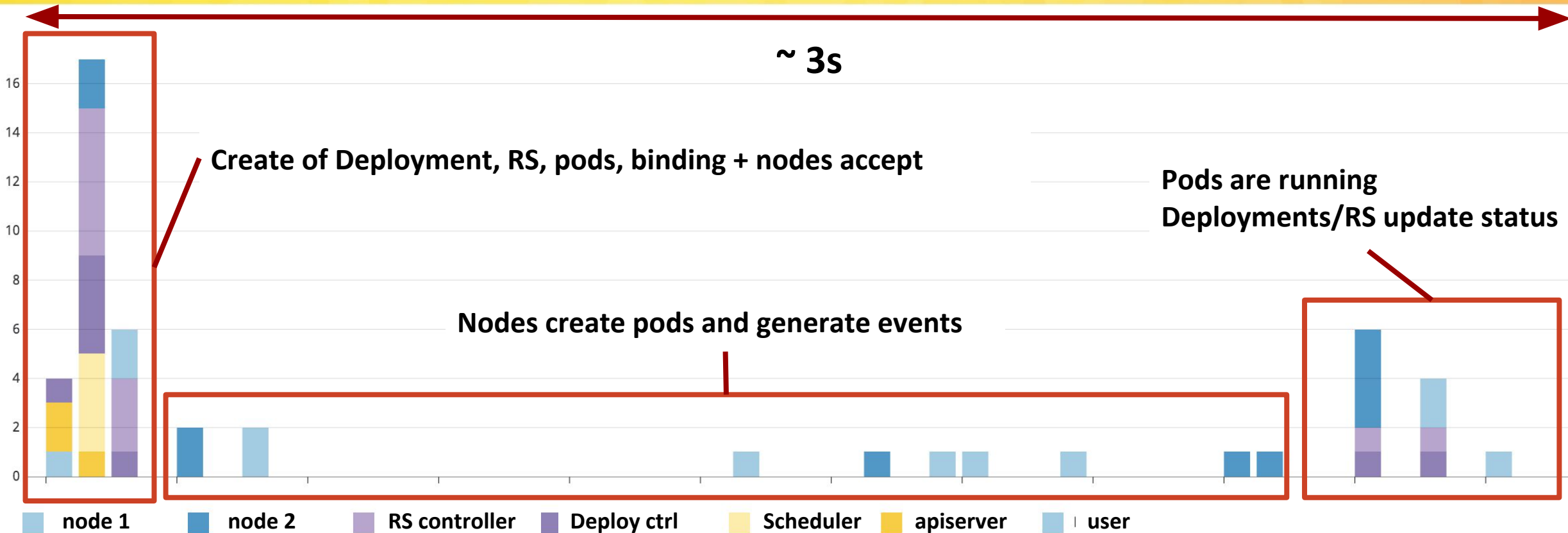
```
objectRef {
  apiVersion v1
  name      nginx-68c5f9f877-29qsb
  namespace demo
  resource  pods
  subresource binding
  uid      43a0691b-07a8-11ea-8101-12d7306f3c0c
}

requestObject {
  apiVersion v1
  kind      Binding
  metadata {
    creationTimestamp null
    name             nginx-68c5f9f877-29qsb
    namespace        demo
    uid              43a0691b-07a8-11ea-8101-12d7306f3c0c
  }
}
```

To node
ip-10-x-y-123

```
target {
  kind Node
  name ip-10- -123.ec2.internal
}
```

Actually a lot more



Additions

- Apiserver verifies Quotas
- Components also get/list
- Creation of events + Update of status fields
- Complete node workflow

Node API calls



KubeCon



CloudNativeCon

North America 2019

```
Nov 15 14:03:07.904 > {"duration":4000000,"objectRef":{"resource":"pods","subresource":"status","namespace":"demo","name":"nginx-68c5f9f877-d7qqg"},'
Nov 15 14:03:07.900 > {"duration":2000000,"objectRef":{"resource":"pods","namespace":"demo","name":"nginx-68c5f9f877-d7qqg","apiVersion":"v1"},'
Nov 15 14:03:07.393 > {"duration":5000000,"objectRef":{"resource":"events","namespace":"demo","name":"nginx-68c5f9f877-d7qqg.15d757ab26df72e4"},'
Nov 15 14:03:07.248 > {"duration":4000000,"objectRef":{"resource":"events","namespace":"demo","name":"nginx-68c5f9f877-d7qqg.15d757ab1e1587b1"},'
Nov 15 14:03:07.162 > {"duration":4000000,"objectRef":{"resource":"events","namespace":"demo","name":"nginx-68c5f9f877-d7qqg.15d757ab19061b3d"},'
Nov 15 14:03:06.897 > {"duration":3000000,"objectRef":{"resource":"events","namespace":"demo","name":"nginx-68c5f9f877-d7qqg.15d757ab09445624"},'
Nov 15 14:03:06.111 > {"duration":4000000,"objectRef":{"resource":"events","namespace":"demo","name":"nginx-68c5f9f877-d7qqg.15d757aada77ac19"},'
Nov 15 14:03:06.100 > {"duration":2000000,"objectRef":{"resource":"secrets","namespace":"demo","name":"default-token-dqmqz","apiVersion":"v1"},'
Nov 15 14:03:05.908 > {"duration":3000000,"objectRef":{"resource":"pods","subresource":"status","namespace":"demo","name":"nginx-68c5f9f877-d7qqg"},'
Nov 15 14:03:05.900 > {"duration":6000000,"objectRef":{"resource":"pods","namespace":"demo","name":"nginx-68c5f9f877-d7qqg","apiVersion":"v1"},'
```

Initial calls

Get pod

Update pod/status: ContainerCreating

Get service account token

Node API calls



KubeCon



CloudNativeCon

North America 2019

```
Nov 15 14:03:07.904 > {"duration":4000000,"objectRef":{"resource":"pods","subresource":"status","namespace":"demo","name":"nginx-68c5f9f877-d7qqg"},'
Nov 15 14:03:07.900 > {"duration":2000000,"objectRef":{"resource":"pods","namespace":"demo","name":"nginx-68c5f9f877-d7qqg","apiVersion":"v1"},'
Nov 15 14:03:07.393 > {"duration":5000000,"objectRef":{"resource":"events","namespace":"demo","name":"nginx-68c5f9f877-d7qqg.15d757ab26df72e4"},'
Nov 15 14:03:07.248 > {"duration":4000000,"objectRef":{"resource":"events","namespace":"demo","name":"nginx-68c5f9f877-d7qqg.15d757ab1e1587b1"},'
Nov 15 14:03:07.162 > {"duration":4000000,"objectRef":{"resource":"events","namespace":"demo","name":"nginx-68c5f9f877-d7qqg.15d757ab19061b3d"},'
Nov 15 14:03:06.897 > {"duration":3000000,"objectRef":{"resource":"events","namespace":"demo","name":"nginx-68c5f9f877-d7qqg.15d757ab09445624"},'
Nov 15 14:03:06.111 > {"duration":4000000,"objectRef":{"resource":"events","namespace":"demo","name":"nginx-68c5f9f877-d7qqg.15d757aada77ac19"},'
Nov 15 14:03:06.100 > {"duration":2000000,"objectRef":{"resource":"secrets","namespace":"demo","name":"default-token-dqmwz","apiVersion":"v1"},'
Nov 15 14:03:05.908 > {"duration":3000000,"objectRef":{"resource":"pods","subresource":"status","namespace":"demo","name":"nginx-68c5f9f877-d7qqg"},'
Nov 15 14:03:05.900 > {"duration":6000000,"objectRef":{"resource":"pods","namespace":"demo","name":"nginx-68c5f9f877-d7qqg","apiVersion":"v1"},'
```

Create events to show progression

MountVolume.Setup succeeded for volume "default-token-dqmwz"

pulling image "nginx"

Successfully pulled image "nginx"

Created container

Started container

Node API calls



KubeCon



CloudNativeCon

North America 2019

```
Nov 15 14:03:07.904 > {"duration":4000000,"objectRef":{"resource":"pods","subresource":"status","namespace":"demo","name":"nginx-68c5f9f877-d7qg"},'
Nov 15 14:03:07.900 > {"duration":2000000,"objectRef":{"resource":"pods","namespace":"demo","name":"nginx-68c5f9f877-d7qqg","apiVersion":"v1"},'
Nov 15 14:03:07.393 > {"duration":5000000,"objectRef":{"resource":"events","namespace":"demo","name":"nginx-68c5f9f877-d7qqg.15d757ab26df72e4"},'
Nov 15 14:03:07.248 > {"duration":4000000,"objectRef":{"resource":"events","namespace":"demo","name":"nginx-68c5f9f877-d7qqg.15d757ab1e1587b1"},'
Nov 15 14:03:07.162 > {"duration":4000000,"objectRef":{"resource":"events","namespace":"demo","name":"nginx-68c5f9f877-d7qqg.15d757ab19061b3d"},'
Nov 15 14:03:06.897 > {"duration":3000000,"objectRef":{"resource":"events","namespace":"demo","name":"nginx-68c5f9f877-d7qqg.15d757ab09445624"},'
Nov 15 14:03:06.111 > {"duration":4000000,"objectRef":{"resource":"events","namespace":"demo","name":"nginx-68c5f9f877-d7qqg.15d757aada77ac19"},'
Nov 15 14:03:06.100 > {"duration":2000000,"objectRef":{"resource":"secrets","namespace":"demo","name":"default-token-dqmzw","apiVersion":"v1"},'
Nov 15 14:03:05.908 > {"duration":3000000,"objectRef":{"resource":"pods","subresource":"status","namespace":"demo","name":"nginx-68c5f9f877-d7qg"},'
Nov 15 14:03:05.900 > {"duration":6000000,"objectRef":{"resource":"pods","namespace":"demo","name":"nginx-68c5f9f877-d7qqg","apiVersion":"v1"},'
```

Finalize pod creation

Get pod

Update container status to “Running”

Takeaways



KubeCon



CloudNativeCon

North America 2019

- A lot of interactions between kube components
- Audit logs give a great understanding of this!
- Events are spiky and get generate a lot of logs
- Events have a 1h default TTL, but stay in audit logs

Let's identify some problems using audit logs

Outline



KubeCon



CloudNativeCon

North America 2019

1. Background: The Kubernetes API
2. Audit Logs
3. Configuring Audit Logs
4. 10000 foot view for a large cluster
5. Understanding Kubernetes Internals
- 6. Troubleshooting examples**



KubeCon



CloudNativeCon

North America 2019

Troubleshooting examples



Troubleshooting



KubeCon



CloudNativeCon

North America 2019

- Understand what happened
 - “Why was a resource deleted?”
- Debug performance regressions/improve performances
 - “Which application is responsible for so many calls?”
- Also, identify issues by looking at HTTP status codes

Status codes



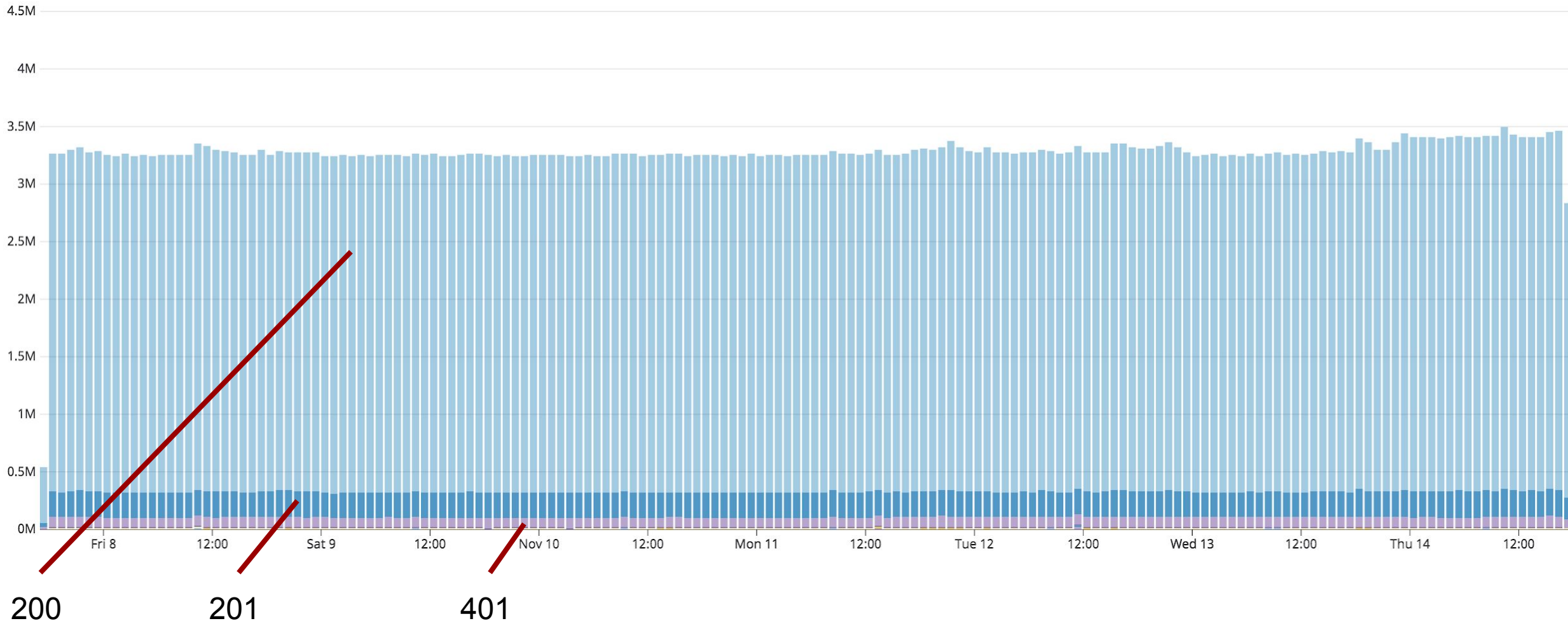
KubeCon



CloudNativeCon

North America 2019

Calls by status code



4xx only



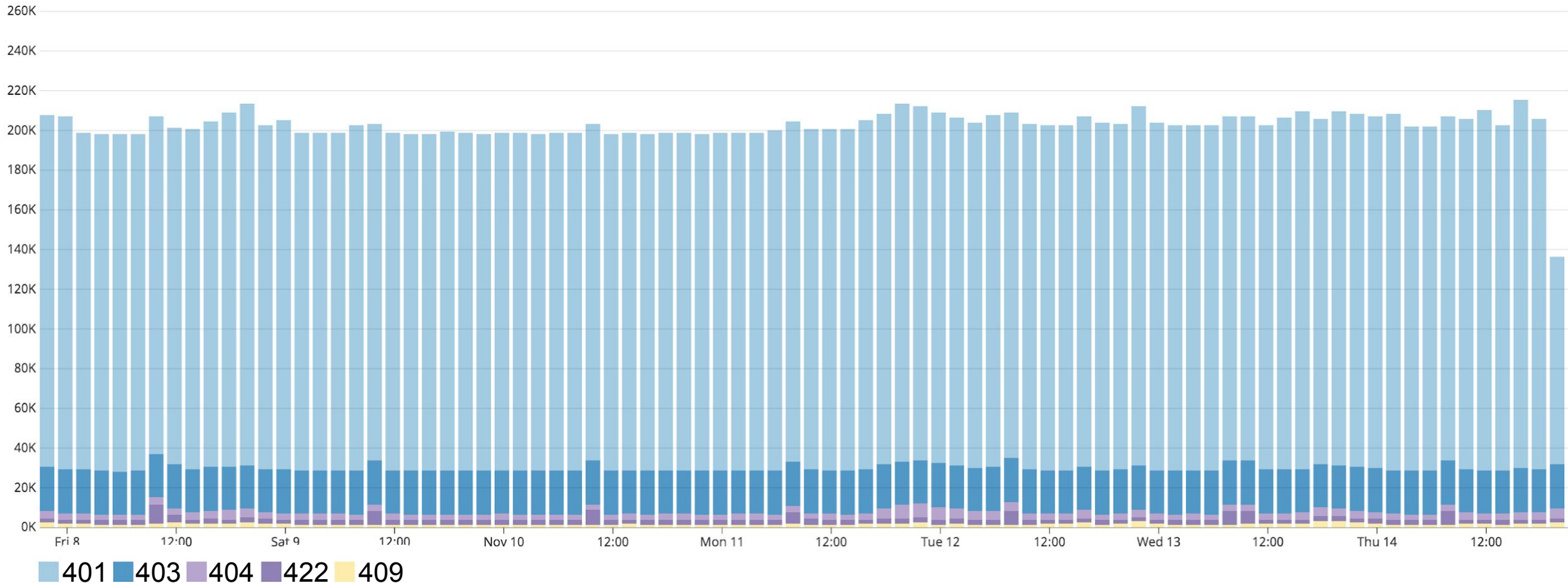
KubeCon



CloudNativeCon

North America 2019

4xx by status code



4xx only



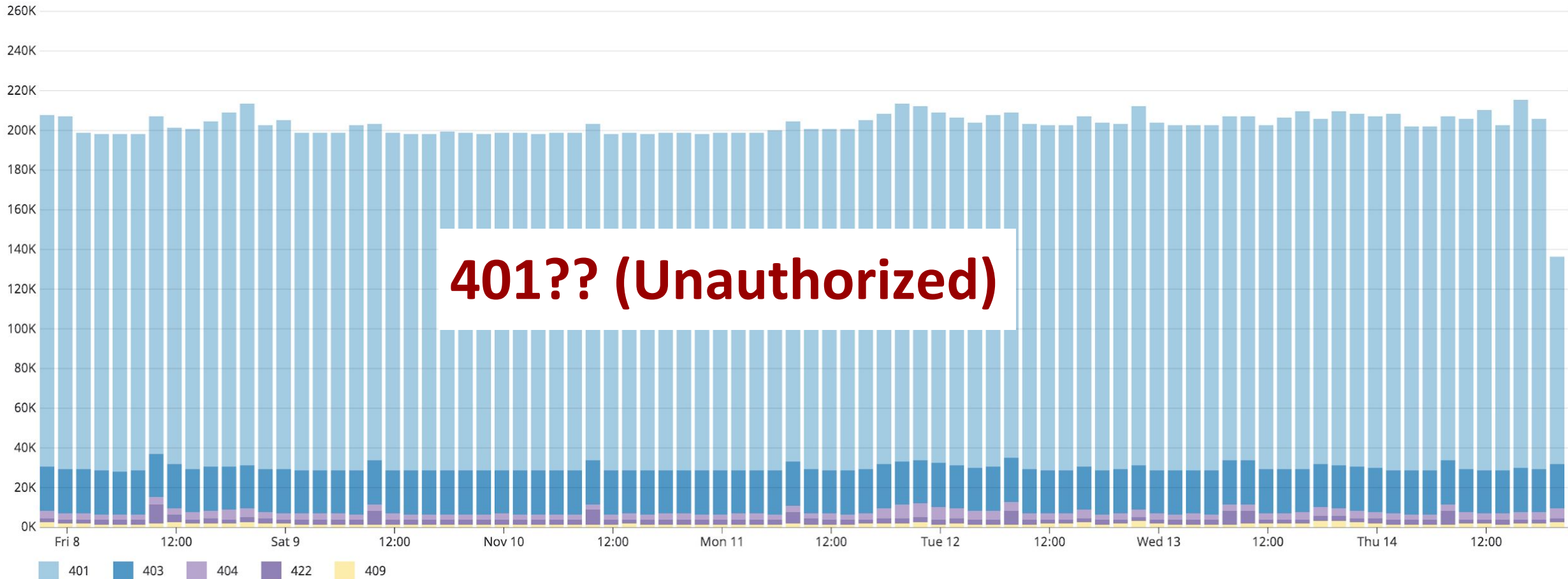
KubeCon



CloudNativeCon

North America 2019

4xx by status code



Analyzing 401s



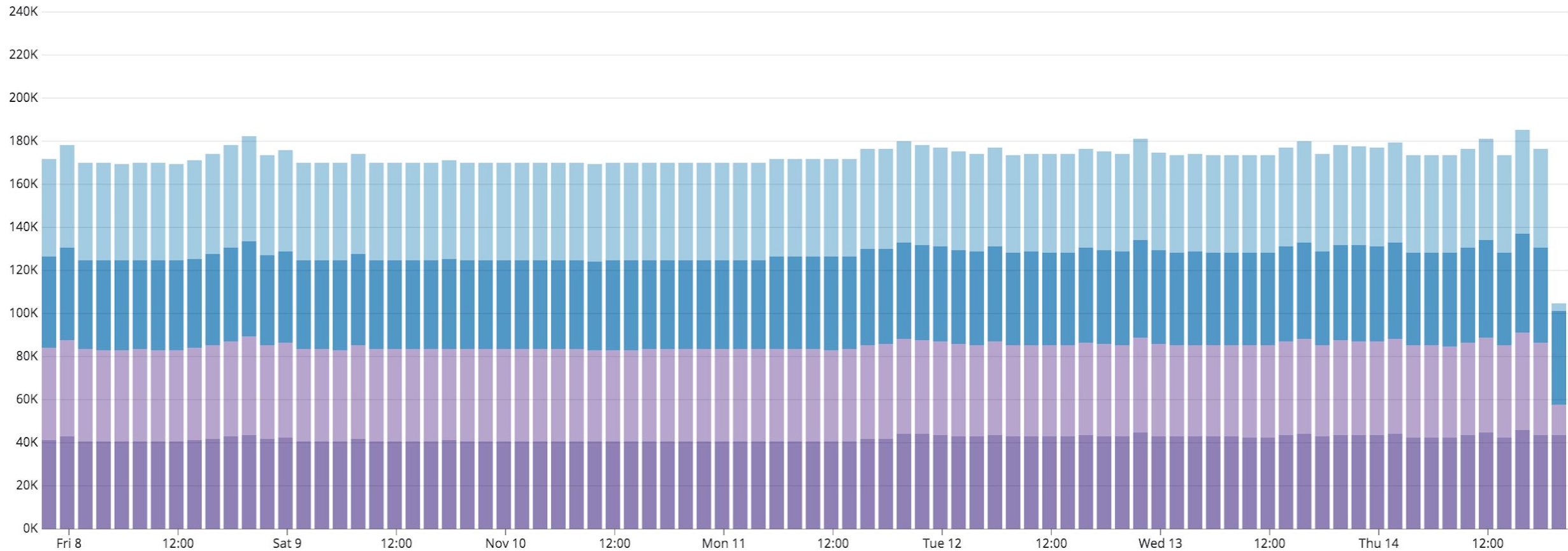
KubeCon



CloudNativeCon

North America 2019

401s by source IP



4 nodes only, turns out they had expired certificates

Analyzing 403s



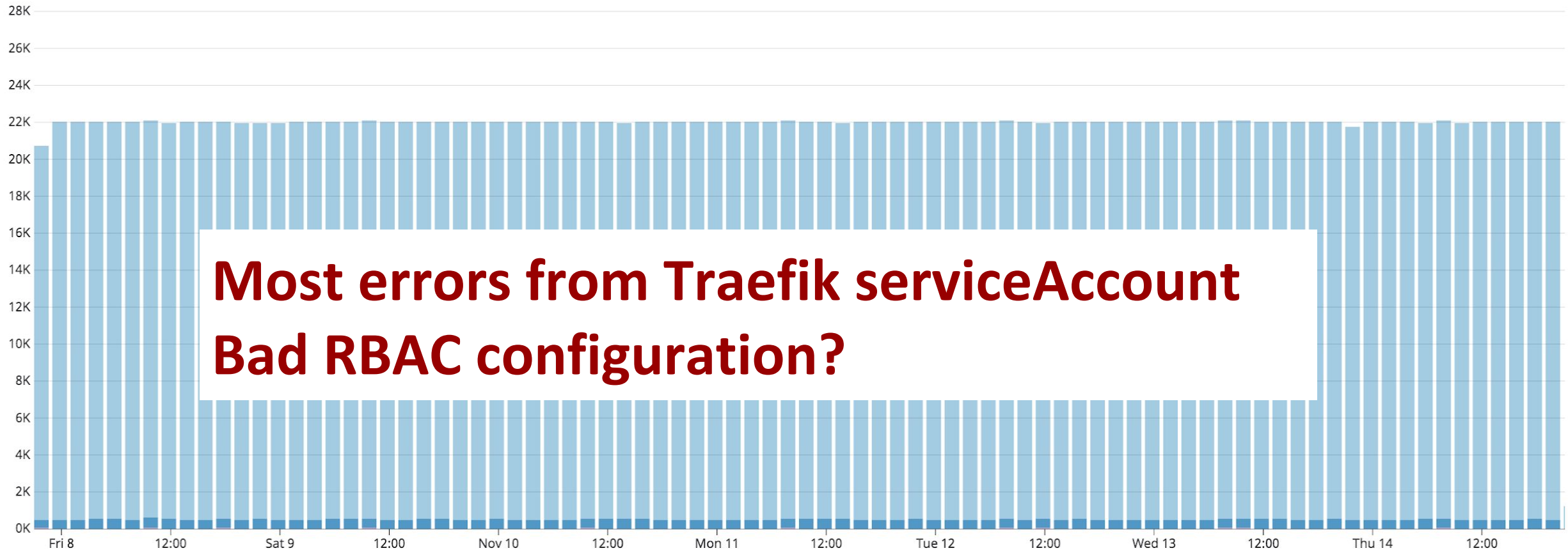
KubeCon



CloudNativeCon

North America 2019

403s by user



Analyzing 403s for this user



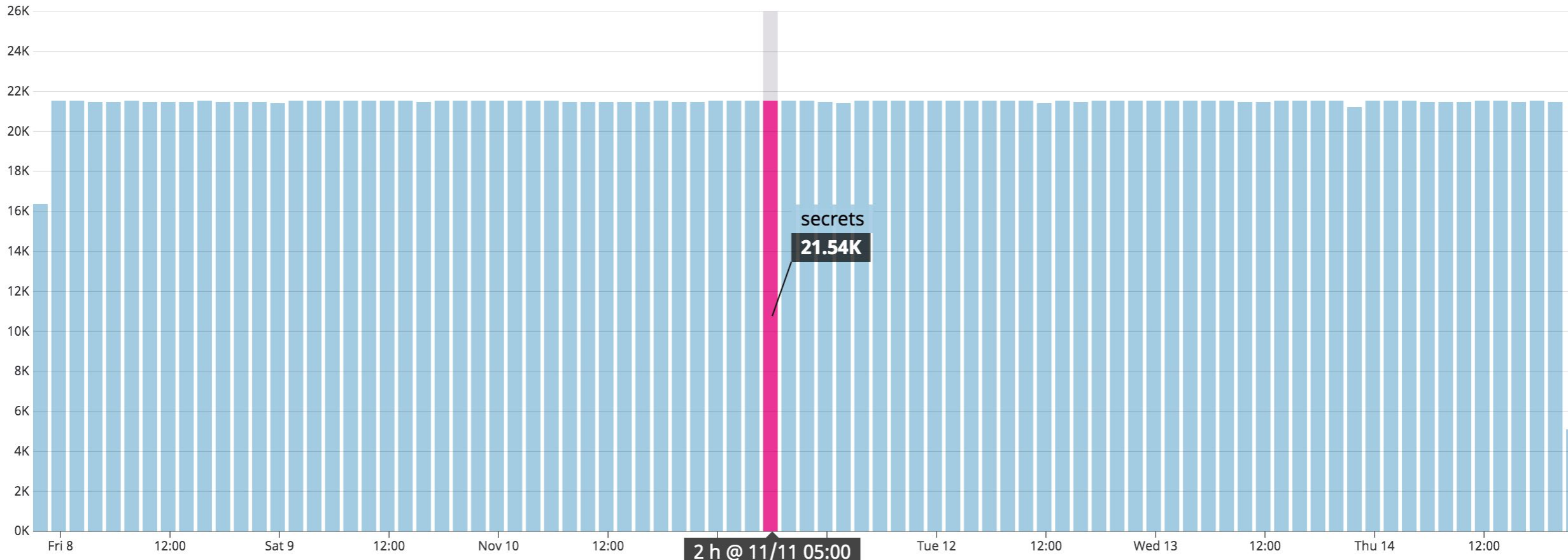
KubeCon



CloudNativeCon

North America 2019

403s for Traefik serviceAccount by resource



We use Traefik without Kubernetes secrets but it still tries to list them

What about 422?



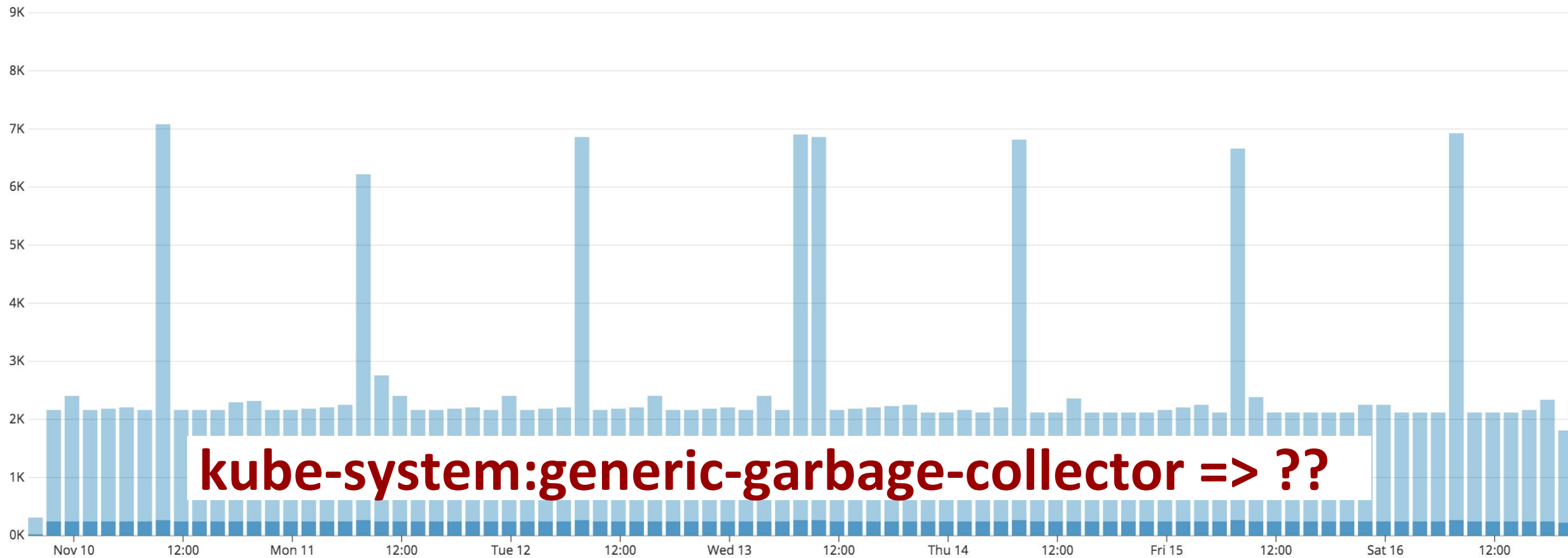
KubeCon



CloudNativeCon

North America 2019

422 by user



What is failing?



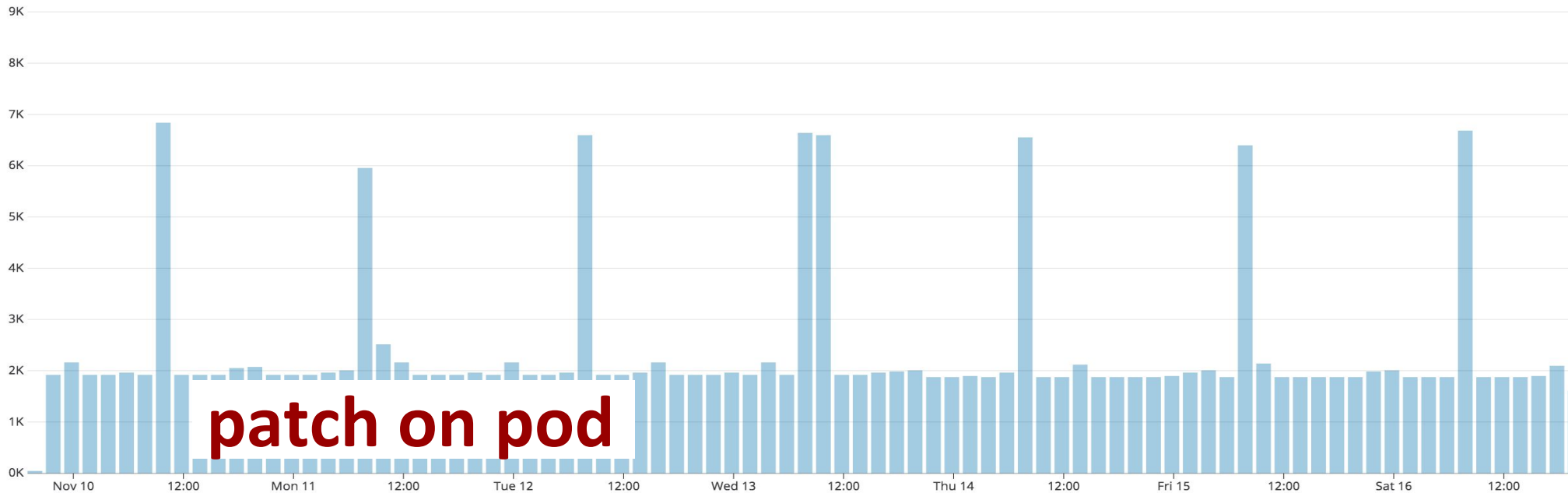
KubeCon



CloudNativeCon

North America 2019

generic-garbage-collector by verb / resource



```
requestObject {
```

```
  metadata {
```

```
    ownerReferences [
```

```
      {"uid":"b430d461-b876-11e9-82b5-42010a50e236", "$patch":"delete"}
```

```
]
```

What is happening?



KubeCon



CloudNativeCon

North America 2019

- Pods are in “Evicted” status and must be kept
- Controlling RS has been deleted and pods should be orphaned
- Garbage collector fails to orphan them (remove ownerRef)
- ReplicaSet has been Terminating for 2 months...

- Root cause: mutating webhook
 - We modify the pod spec at creation
 - Mutating webhook is registered on pods for CREATE/UPDATE
 - We modify immutable fields
 - Garbage collector patch triggers this...

Takeaways



KubeCon



CloudNativeCon

North America 2019

- Looking at calls triggering HTTP errors help find issues
 - Misconfigured RBAC (403)
 - Applications doing calls they shouldn't (403)
 - Expired certificates (401)
 - Other weird things



KubeCon



CloudNativeCon

North America 2019

Conclusion



Conclusion



KubeCon



CloudNativeCon


North America 2019

- Audit logs can be incredibly valuable
 - Low-level understanding of Kubernetes
 - Detection of misconfigurations
 - Troubleshooting of issues
 - Identify performance issues
- Taking advantage of them require some effort
 - Policies are not easy to get right
 - They are verbose and require a tool to analyze them

We're hiring!

Visit our Kubecon booth
or <https://www.datadoghq.com/careers/>

Or contact us directly:

 **@lbernail** laurent@datadoghq.com
@roboll_ roboll@datadoghq.com



KubeCon



CloudNativeCon

North America 2019

Thank you

