

## How the OOM-Killer Deleted my Namespace (and Other Kubernetes Tales)

Laurent Bernaille

## Datadog



Over 400 integrations Over 1,500 employees Over 12,000 customers Runs on millions of hosts Trillions of data points per day 10000s hosts in our infra 10s of k8s clusters with 100-4000 nodes Multi-cloud Very fast growth









- We have shared a few stories already
  - Kubernetes the Very Hard Way
  - 10 Ways to Shoot Yourself in the Foot with Kubernetes
  - Kubernetes DNS horror stories
- The scale of our Kubernetes infra keeps growing
  - New scalability issues
  - But also, much more complex problems





## How The OOM-killer Deleted my Namespace





# "My namespace and everything in it are completely gone"





#### Feb 17 14:30: Namespace and workloads deleted







#### Audit logs

- Delete calls for all resources in the namespace
- No delete call for the namespace resource

Why was the namespace deleted??



## **Deletion call, 4d before**

on CloudNativeCon



#### Audit logs for the namespace



Feb 13th (resource deletion happened on Feb 17th)





#### Engineer did not delete the namespace

### But engineer was migrating to new deployment method



## Spinnaker deploys (v1)



North America 2020







## Helm 3 deploys (v2)









## **Big difference**



North America 2020







## **Big difference**



North America 2020



Α Α А B' В B' В Spinnaker C' С C' С Х Α А HELM В B' B' В C' C' С С



## **Big difference**



North America 2020





v2: removing a resource from the chart deletes it from the cluster



#### 1- Chart refactoring: move namespace out of it

### 2- Deployment: deletes namespace

But why did it take 4 days ?



## **Namespace Controller**



497	// deleteAllContent will use the dynamic client to delete each resource identified in groupVersionResources.	
498	// It returns an estimate of the time remaining before the remaining resources are deleted.	
499	<pre>// If estimate &gt; 0, not all resources are guaranteed to be gone.</pre>	
500	<pre>func (d *namespacedResourcesDeleter) deleteAllContent(</pre>	
501	<pre>namespace string, namespaceDeletedAt metav1.Time) (int64, error) {</pre>	
502	estimate := int64(0)	
503	glog.V(4).Infof("namespace controller – deleteAllContent – namespace: %s", namespace)	
504	resources, err := d.discoverResourcesFn()	
505	<pre>if err != nil {</pre>	
506	return estimate, err	
507	_ }	
508	<pre>// TODO(sttts): get rid of opCache and pass the verbs (especially "deletecollection") down into the deleter</pre>	
509	<pre>deletableResources := discovery.FilteredBy(discovery.SupportsAllVerbs{Verbs: []string{"delete"}}, resources)</pre>	
510	<pre>groupVersionResources, err := discovery.GroupVersionResources(deletableResources)</pre>	Discover all API resources
511	<pre>if err != nil {</pre>	
512	return estimate, err	
513	}	
514	var errs []error	
515	<pre>for gvr := range groupVersionResources {</pre>	
516	<pre>gvrEstimate, err := d.deleteAllContentForGroupVersionResource(gvr, namespace, namespaceDeletedAt)</pre>	
517	<pre>if err != nil {</pre>	Delete all resources in
518	<pre>// If there is an error, hold on to it but proceed with all the remaining</pre>	Delete all resources in
519	// groupVersionResources.	namespace
520	errs = append(errs, err)	
521	}	
522	<pre>if gvrEstimate &gt; estimate {</pre>	
523	estimate = gvrEstimate	
524	}	🔰 @lhernai
525	}	



deletableResources := discovery.FilteredBy(discovery.SupportsAllVerbs{Verbs: []string{"delete"}}, resources)
groupVersionResources, err := discovery.GroupVersionResources(deletableResources)



#### => If discovery fails, namespace content is not deleted

Kubernetes 1.14+: Delete as many resources as possible Kubernetes 1.16+: Add Conditions to namespaces to surface errors



## Namespace Controller logs





#### Controller-manager contained many occurrences of

unable to retrieve the complete list of server APIs: metrics.k8s.io/v1beta1: the server is currently unable to handle the request

No context so hard to link with the incident but

- This is the error from the Discovery call
- We have a likely cause: metrics.k8s.io/v1beta1



## metrics-server



- Additional controller providing node/pod metrics
- Registers an ApiService
  - Additional Kubernetes APIs
  - Managed by the external controller

apiVersion: apiregistration.k8s.io/v1
kind: APIService

metadata:

name: v1beta1.metrics.k8s.io

spec:

group: metrics.k8s.io
version: v1beta1

service:

name: metrics-server
namespace: kube-system

Additional APIs (group and version)

Where should the apiserver proxy calls for this APIs



## **Events so far**



- Feb 13th: namespace deletion call
- Feb 13-17th
  - Namespace controller can't discover metrics APIs
  - Namespace content is not deleted
  - $\circ~$  At each sync, controller retries and fails
- Feb 17th
  - **?**
  - Namespace controller is unblocked
  - Everything in the namespace is deleted



## **Events so far**



- Feb 13th: error leads to namespace deletion call
- Feb 13-17th
  - Namespace controller can't discover metrics APIs
  - Namespace content is not deleted
  - $\circ~$  At each sync, controller retries and fails
- Feb 17th
  - OOM-killer kills metrics-server and it is restarted
  - Metrics-server is now working fine
  - Namespace controller is unblocked
  - $\circ~$  Everything in the namespace is deleted





#### Why was metrics-server failing?



## **Metrics-server setup**





#### Pod with two containers

- metrics-server
- sidecar: rotate server certificate, and signal metrics-server to reload

#### Requires shared pid namespace

- shareProcessNamespace (pod spec)
- PodShareProcessNamespace (feature gate on Kubelet and Apiserver)



## **Metrics-server deployment**



Deployed on a pool of nodes dedicated to core controllers ("system")



CloudNativeCon

orth America 2020

## **Metrics-server deployment**



metrics-server

Subset of nodes using an old image without PodShareProcessNamespace

- metrics-server: fine until it was rescheduled on an old node
- At this point, server certificate is no longer reloaded



CloudNativeCon

America 2020

## **Full chain of events**



- Before Feb 13th
  - metrics-server scheduled on a "bad" node
  - controller discovery calls fail
- Feb 13th: error leads to namespace deletion call
- Feb 13-17th
  - Namespace controller can't discover metrics APIs
  - Namespace content is not deleted
  - At each sync, controller retries and fails
- Feb 17th
  - OOM-killer kills metrics-server and it is restarted
  - Metrics-server is now working fine (certificate up to date)
  - Namespace controller is unblocked
  - Everything in the namespace is deleted



## Key take-away



- Apiservice extensions are great but can impact your cluster
  - If the API is not used, they can be down for days without impact
  - But some controllers need to discover all apis and will fail
- In addition, Apiservice extensions security is hard
  - $\circ~$  If you are curious, go and see

PKI the Wrong Way: Simple TLS Mistakes and Surprising Consequences - Tabitha Sable, Datadog





### Virtual

## New node image does not work and it's Weird...





- We build images for our Kubernetes nodes
- A small change pushed
- Nodes become NotReady after a few days
- Change completely unrelated to kubelet / runtime
- We can't promote to Prod



## Investigation



#### Node description

onditions:				
Туре	Status	Reason	Message	
OutOfDisk	False	KubeletHasSufficientDisk	kubelet has sufficient disk space available	
MemoryPressure	False	KubeletHasSufficientMemory	kubelet has sufficient memory available	
DiskPressure	False	KubeletHasNoDiskPressure	kubelet has no disk pressure	
PIDPressure	False	KubeletHasSufficientPID	kubelet has sufficient PID available	
Ready	False	KubeletNotReady	container runtime is down	



## **Runtime is down?**



North America 2020



Let's look at containerd

\$ crictl pods		
POD ID	STATE	NAME
681bb3aec7fba	Ready	local-volume-provisioner-jtwz:
53f122f351da2	Ready	datadog-agent-qctt2-81jgp
cledb4f19713c	Ready	localusers-z26n9
01cf5cafd6bc9	Ready	node-monitoring-ws2xf
39f01bdaade86	Ready	kube2iam-zzdt6
3fcf520680a63	Ready	kube-proxy-sr7tn
5ecd0966634f1	Ready	node-local-dns-m2zbp

NAMESPACE local-volume-provisioner datadog-agent prod-platform node-monitoring kube2iam kube-system coredns

#### Looks completely fine







#### Something is definitely wrong

remote\_runtime.go:434] Status from runtime service failed: rpc error: code = DeadlineExceeded desc =
context deadline exceeded

kubelet.go:2130] Container runtime sanity check failed: rpc error: code = DeadlineExceeded desc =
context deadline exceeded

kubelet.go:1803] skipping pod synchronization - [container runtime is down]

#### "Status from runtime service failed" ➤ Is there a CRI call for this?

// Status returns the status of the runtime.
rpc Status(StatusRequest) returns (StatusResponse) {}



## **Testing with crictl**



- crictl has no "status" subcommand
- but "crictl info" looks close, let's check

// Info sends a StatusRequest to the server, and parses the returned StatusResponse.
func Info(cliContext \*cli.Context, client pb.RuntimeServiceClient) error {
 request := &pb.StatusRequest{Verbose: !cliContext.Bool("quiet")}
 logrus.Debugf("StatusRequest: %v", request)
 r, err := client.Status(context.Background(), request)



## **Crictl info**



#### \$ crictl info ^C

- Hangs indefinitely
- But what does Status do?
- Almost nothing, except this

44	<pre>// Check the status of the cni initialization</pre>
45	<pre>if err := c.netPlugin.Status(); err != nil {</pre>
46	<pre>networkCondition.Status = false</pre>
47	<pre>networkCondition.Reason = networkNotReadyReason</pre>
48	<pre>networkCondition.Message = fmt.Sprintf("Network plugin returns error: %v", err)</pre>
49	}



## **CNI** status



- // Status returns the status of CNI initialization. 124
- func (c \*libcni) Status() error { 125
- 126 c.RLock()
- defer c.RUnlock() 127
- if len(c.networks) < c.networkCount { 128
  - return ErrCNINotInitialized
- } 130
- return nil 131
- 132 }

129

- Really not much here right?
- But we have calls that hang and a lock
- Could this be related?



## **Containerd goroutine dump**

#### Blocked goroutines?

<pre>containerd[737]:</pre>	goroutine	107298383	[semacquire,	2 minutes]:
<pre>containerd[737]:</pre>	goroutine	107290532	[semacquire,	4 minutes]:
<pre>containerd[737]:</pre>	goroutine	107282673	[semacquire,	6 minutes]:
<pre>containerd[737]:</pre>	goroutine	107274360	[semacquire,	8 minutes]:
containerd[737]:	goroutine	107266554	[semacquire,	10 minutes]:

#### The runtime error from the kubelet also happens every 2mn

#### Goroutine details match exactly the codepath we looked at

containerd[737]: goroutine 107298383 [semacquire, 2 minutes]: ... containerd[737]: .../containerd/vendor/github.com/containerd/go-cni.(\*libcni).Status() containerd[737]: .../containerd/vendor/github.com/containerd/go-cni/cni.go:126 containerd[737]: .../containerd/vendor/github.com/containerd/cri/pkg/server.(\*criService).Status(...) containerd[737]: .../containerd/containerd/vendor/github.com/containerd/cri/pkg/server.(\*criService).Status(...)

## **Seems CNI related**



#### Let's bypass containerd/kubelet

\$ ip netns add debug
\$ CNI\_PATH=/opt/cni/bin cnitool add pod-network /var/run/netns/debug

#### Works great and we even have connectivity

```
$ ip netns exec debug ping 8.8.8.8
PING 8.8.8.8 (8.8.8.8) 56(84) bytes of data.
64 bytes from 8.8.8.8: icmp_seq=1 ttl=111 time=2.41 ms
64 bytes from 8.8.8.8: icmp_seq=2 ttl=111 time=1.70 ms
64 bytes from 8.8.8.8: icmp_seq=3 ttl=111 time=1.76 ms
64 bytes from 8.8.8.8: icmp_seq=4 ttl=111 time=1.72 ms
```





\$ CNI\_PATH=/opt/cni/bin cnitool del pod-network /var/run/netns/debug
^c

- hangs indefinitely
- and we can even find the original process still running

ps aux | grep cni root 367 0.0 0.1 410280 8004 ? Sl 15:37 0:00 /opt/cni/bin/cni-ipvlan-vpc-k8s-unnumbered-ptp







- On this cluster we use the Lyft CNI plugin
- The delete calls fails for cni-ipvlan-vpc-k8s-unnumbered-ptp
- After a few tests, we tracked the problem to this call

vethPeerIndex, \_ = netlink.VethPeerIndex(&netlink.Veth{LinkAttrs: \*vethIface.Attrs()})

• Weird, this is in the netlink library!





#### link, veth: fix stack corruption from retrieving peer index #498



<b>a</b>				
	borkmann commented on Nov 6, 2019	Contributor 😳 •••		
F	For 4.20 and newer kernels VethPeerIndex() causes a stack corruption as			
the kernel is copying more data to golang user space than originally				

#### Many thanks to Daniel Borkmann!







- Packer configured to use the latest Ubuntu 18.04
- The default kernel changed from 4.15 to 5.4
- Netlink library had a bug for kernels 4.20+
- Bug only happened on pod deletion







#### Containerd

- Status calls trigger config file reload, which acquires a RWLock on the config mutex
- CNI Add/Del calls acquire a RLock on the mutex
- If CNI Add/Del hangs, Status calls block
- Does not impact containerd 1.4+ (separate fsnotify.Watcher to reload config)
- cni-ipvlan-vpc-k8s ("Lyft" CNI plugin)
  - Bug in Netlink library with kernel 4.20+
  - Does not impact cni-ipvlan-vpc-k8s 0.6.2+ (different function used)





Nodes are not abstract compute resources

- Kernel
- Distribution
- Hardware

Error messages can be misleading

- For CNI errors, runtime usually reports a clear error
- Here we had nothing but a timeout







## How a single app can overwhelm the control plane





- Users report connectivity issue to a cluster
- Apiservers are not doing well





## What's with the Apiservers?



#### Apiservers can't reach etcd

**KubeCon** 

CloudNativeCon

orth America 2020

#### Apiservers crash/restart











#### Etcd memory usage is spiking

#### Etcd is oom-killed





- Cluster size has not significantly changed
- The control plane has not been updated
- So it is very likely related to api calls



## **Apiserver requests**

KubeCon CloudNativeCon



@lbernail





#### **Understanding Apiserver caching**





## Why are list calls expensive?





#### What happens for GET/LIST calls?



- Resources have versions (ResourceVersion, RV)
- For GET/LIST with RV=X

If cachedVersion >= X, return cachedVersion

else wait up to 3s, if cachedVersion>=X return cachedVersion

else Error: "Too large resource version"

- ➢ RV=X: "at least as fresh as X"
- RV=0: current version in cache





#### What about GET/LIST without a resourceVersion?



- If RV is not set, apiserver performs a Quorum read against etcd (for consistency)
- This is the behavior of
  - $\circ$  kubectl get
  - client.CoreV1().Pods("").List() with default options (client-go)





time curl 'https://cluster.dog/api/v1/pods'
real Om4.631s
time curl 'https://cluster.dog/api/v1/pods?resourceVersion=0'
real Om1.816s

- Test on a large cluster with more than 30k pods
- Using table view ("application/json; as=Table; v=v1beta1;g=meta.k8s.io, application/json")
  - Only ~25MB of data to minimize transfer time (full JSON: ~1GB)



time curl 'https://cluster.dog/api/v1/pods?labelSelector=app=A'

real 0m3.658s

time curl 'https://cluster.dog/api/v1/pods?labelSelector=app=A&resourceVersion=0'
real 0m0.079s

- Call with RV="" is slightly faster (less data to send)
- Call with RV=0 is much much faster
   > Filtering is performed on cached data
- When RV="", all pods are still retrieved from etcd and then filtered on apiservers



## Why not filter in etcd?

KubeCon CloudNativeCon Virtual

etcd key structure
/registry/{resource type}/{namespace)/{resource name}

So we can ask etcd for:

- a specific resource
- all resources of a type in a namespace
- all resources of a type
- no other filtering / indexing

curl 'https://cluster.dog/api/v1/pods?labelSelector=app=A'

real 0m3.658s <== get all pods (30k) in <u>etcd</u>, filter on apiserver

curl 'https://cluster.dog/api/v1/namespaces/datadog/pods?labelSelector=app=A'

real 0m0.188s <== get pods from datadog namespace (1000) in <u>etcd</u>, filter on apiserver

curl 'https://cluster.dog/api/v1/namespaces/datadog/pods?labelSelector=app=A&resourceVersion=0'
real 0m0.058s <== get pods from datadog namespace in <u>apiserver cache</u>, filter



## Informers instead of List



orth America 2020



#### How do informers work?



- Informers are much better because
  - They maintain a local cache, updated on changes
  - They start with a LIST using RV=0 to get data from the cache
- Edge case
  - Disconnections: can trigger a new LIST with RV="" (to avoid getting back in time)
  - Kubernetes 1.20 will have EfficientWatchResumption (alpha)

https://github.com/kubernetes/enhancements/tree/master/keps/sig-api-machinery/1904-efficient-watch-resumption







- LIST calls go to etcd by default and can have a huge impact
- LIST calls with label filters still retrieve everything from etcd
- Avoid LIST, use Informers
- kubectl get uses LIST with RV=""
- kubectl get could allow setting RV=0
  - Much faster: better user experience
  - Much better for etcd and apiservers
  - Trade-off: small inconsistency window

Many, many thanks to Wojtek-t for his help getting all this right



- We know the problem comes from LIST calls
- What application is issuing these calls?







@lbernail

#### Cumulated query time by user for list calls



A single user accounts for 2+ days of query time over 20 minutes!





#### Cumulated query time by user for list/get calls over a week

VERB	USER USERNAME	SUM:DURATION
		2.19wk
list		7.05day
	system:serviceaccount:nodegroup-controller:nodegroup-cor	2.84day
	system:serviceaccount:kube-system:node-controller	46.68hr
	@datadoghq.com	11.1hr
	system:serviceaccount:datadog-agent:datadog-agent-cluste	7.7hr
	system:serviceaccount:kube2iam:kube2iam	4.96hr
get		3.34day
	adatadoghq.com	8.99hr
	kube-scheduler	4.25hr
	@datadoghq.com	2.59hr
	@datadoghq.com	2.46hr
	@datadoghq.com	2.18hr



## **Nodegroup controller?**



- An in-house controller to manage pools of nodes
- Used extensively for 2 years
- But recent upgrade deployed: deletion protection
  - Check if pods are running on pool of nodes
  - $\circ~$  Deny nodegroup deletion if it is the case





On nodegroup delete

- List all nodes in this nodegroup based on labels
   => Some groups have 100+ nodes
- 2. List all pods on each node filtering on bound node
  => List all pods (30k) for node
  => Performed in parallel for all nodes in the group
  => The bigger the nodegroup, the worse the impact

#### All LIST calls here retrieve full node/pod list from etcd



## What we learned



- List calls are very dangerous
  - The volume of data can be very large
  - Filtering happens on the apiserver
  - Use Informers (whenever possible)
- Audit logs are extremely useful
   Who did what when?
  - Which users are responsible for processing time



## Conclusion







- Apiservice extensions are powerful but can harm your cluster
- Nodes are not abstract compute resources
- For large (1000+ nodes) clusters
  - Understanding Apiserver caching is important
  - Client interactions with Apiservers are critical
  - Avoid LIST calls





## Thank you

We're hiring! https://www.datadoghq.com/careers/

laurent@datadoghq.com @lbernail @lbernail





