



Dynamic Storage Provisioning of Manila/CephFS Shares on Kubernetes

Róbert Vašek <robert.vasek@codefreax.org>

Ricardo Rocha <ricardo.rocha@cern.ch> @ahcorporto



home.cern

Table of Contents

Introduction

Container Storage Interface

CSI CephFS

Manila shares with Kubernetes

Results, numbers, plots...

We are here!

Introduction

Container Storage Interface

CSI CephFS

Manila shares with Kubernetes

Results, numbers, plots...



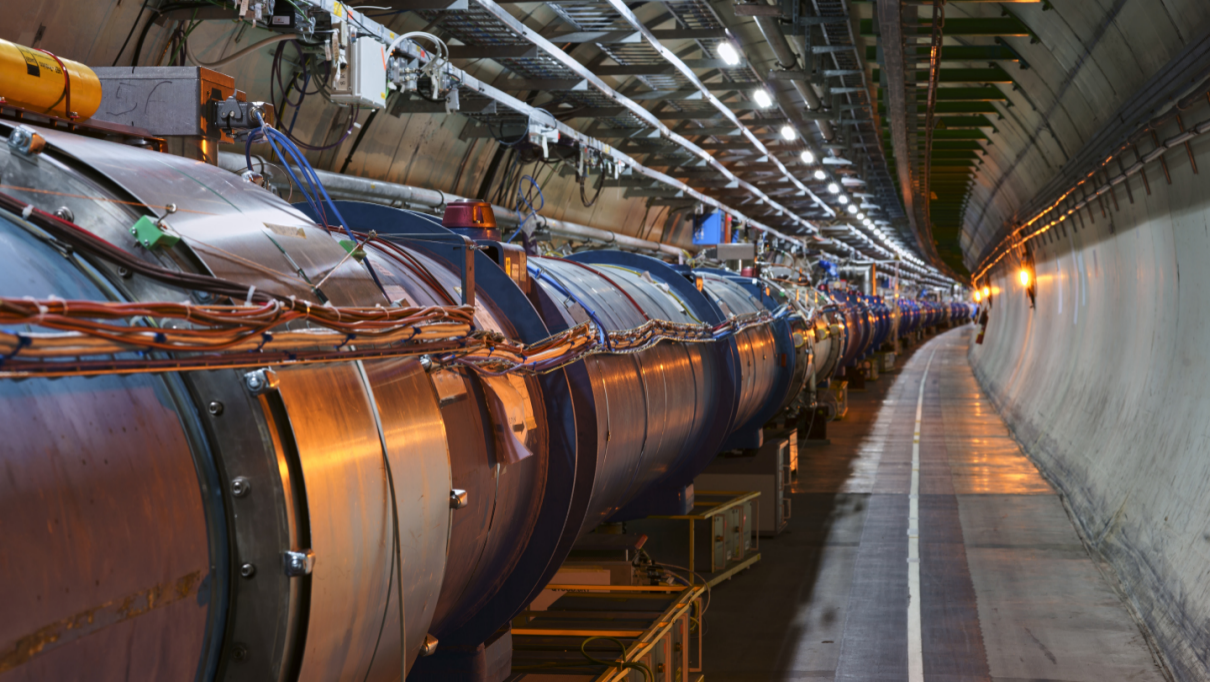
Founded in 1954

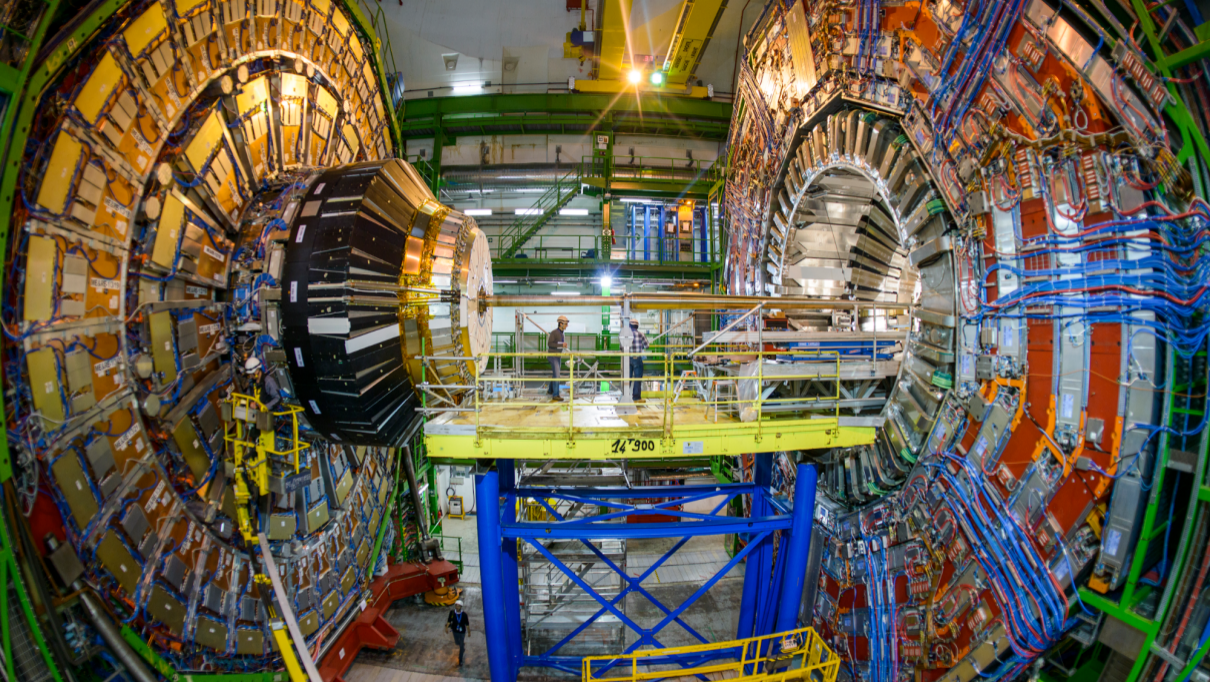
Fundamental Science

What is 96% of the universe made of?

What was the state of matter just after the Big Bang?

Why isn't there anti-matter in the universe?





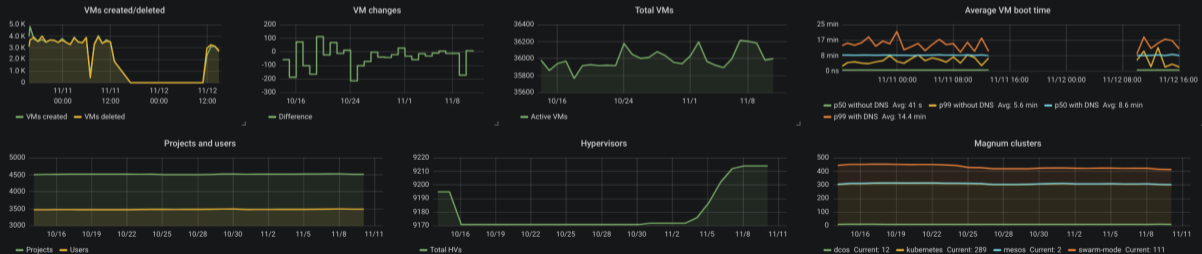
Cloud resources



Openstack services stats

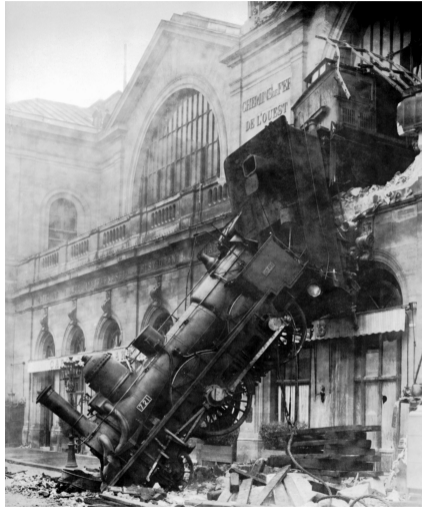


Resource overview by time



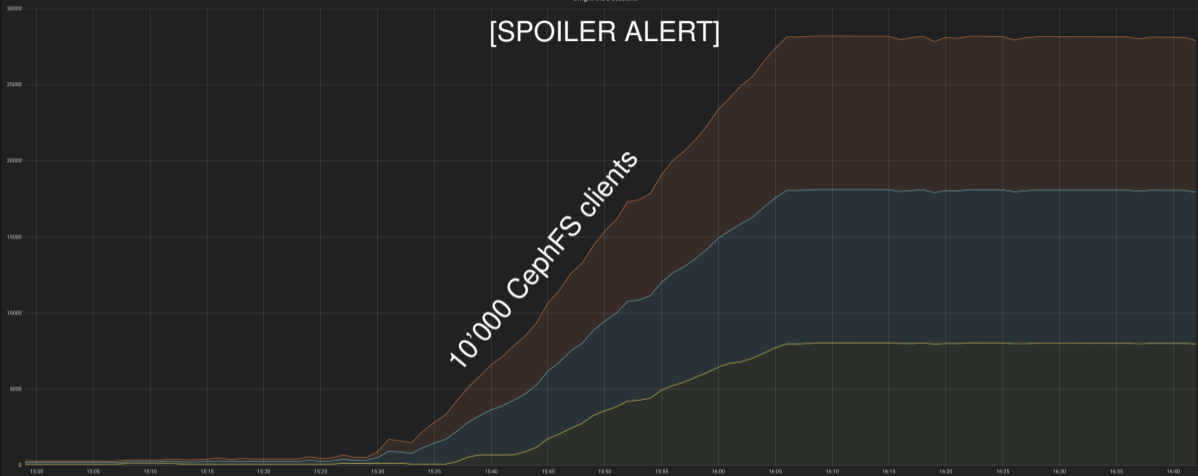
Dynamic Storage Provisioning of Manila/CephFS Shares on Kubernetes

...working title *"From a train wreck to a train ride"*



[SPOILER ALERT]

10'000 CephFS clients



	max	avg	current
ceph:dwright@d0	8042	2925	7876
ceph:dwright@d1	0	0	0
ceph:dwright@d2	10089	5366	5963
p05153020942913	10990	5366	5963

We are here!

Introduction

Container Storage Interface

CSI CephFS

Manila shares with Kubernetes

Results, numbers, plots...

Container Storage Interface - motivation



Storage

Container Storage Interface - motivation



CO 1



CO 2



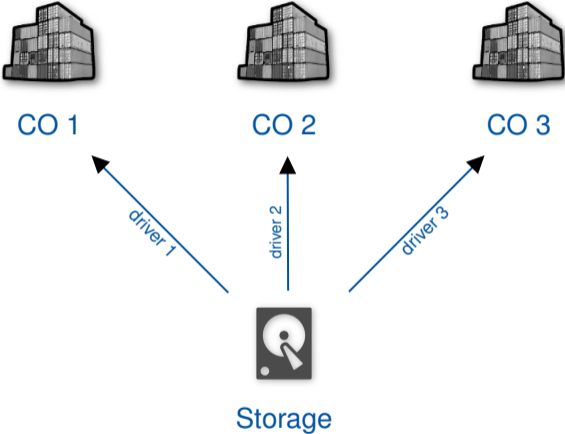
CO 3



Storage

CO - Container Orchestrator

Container Storage Interface - motivation



CO - Container Orchestrator

Container Storage Interface - motivation

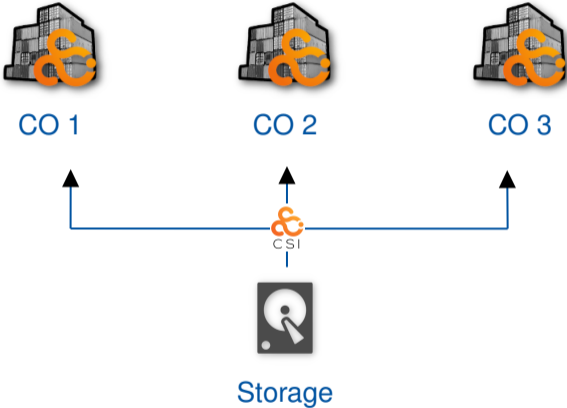
From driver's POV:

- ▶ Lack of standardization
- ▶ Higher development and maintenance costs

From CO's POV:

- ▶ Volume plugin development is tightly coupled with release cycles of the CO
- ▶ Bugs in volume plugins can crash critical components
- ▶ Volume plugins get full privileges
- ▶ Difficult dependency management

Container Storage Interface - motivation



CO - Container Orchestrator

Container Storage Interface

Overview

- ▶ Industry standard for cluster-wide storage plugins
- ▶ Collaboration of communities incl. Kubernetes, Mesos, Docker and Cloud Foundry
- ▶ Defines the protocol between a CO and a plugin
- ▶ Plugins are CO-agnostic
- ▶ Write once – use everywhere, *just works*TM



CONTAINER
STORAGE
INTERFACE

Container Storage Interface

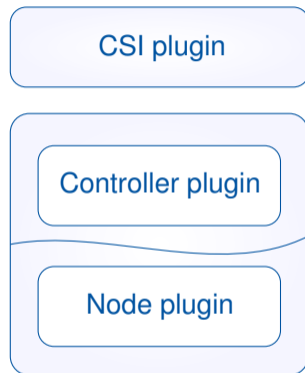
- ▶ First alpha released in Dec 2017
- ▶ Working implementation in Kubernetes 1.9 already, a lot of changes since then, some of those were breaking
- ▶ Other COs soon to follow

December 2017	v0.1.0
March 2018	v0.2.0
June 2018	v0.3.0
	⋮
just today	v1.0.0-rc2
end of Nov 2018	v1.0.0

CSI Services

CSI RPC services (endpoints):

- ▶ *Identity service*: allows a CO to query for plugin's capabilities, health probes and other metadata. Must be implemented by both controller and node plugins, you'll see why in a bit.
- ▶ *Controller service*: creates, deletes, lists volumes and their snapshots.
- ▶ *Node service*: (un)stages, (un)publishes volumes on a node.



CSI Architecture



CSI RPCs quick overview

*Controller Service**

- ▶ CreateVolume
- ▶ DeleteVolume
- ▶ ControllerPublishVol*
- ▶ ControllerGetCaps
- ▶ ...

Node Service

- ▶ NodeStageVolume*
- ▶ NodePublishVolume
- ▶ NodeGetCapabilities
- ▶ ...

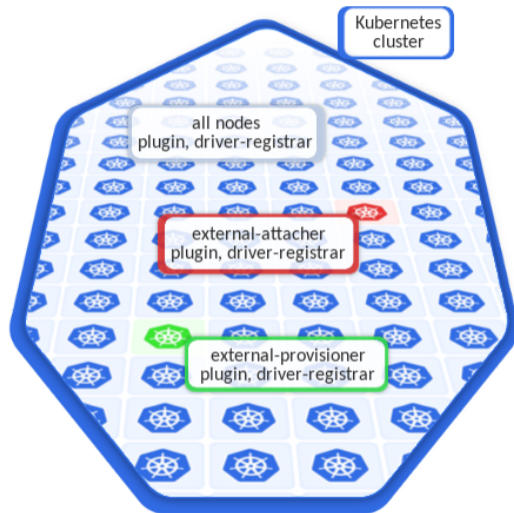
Identity Service

- ▶ GetPluginInfo
- ▶ GetPluginCapabilities
- ▶ ...

* optional

CSI in Kubernetes

- ▶ In-tree CSI volume plugin in kubelet
Node(Un)StageVolume
Node(Un)PublishVolume
- ▶ Side-car containers
 - ▶ **driver-registrar**
plugin discovery,
registers the driver with kubelet
 - ▶ **external-provisioner**
CreateVolume
DeleteVolume
 - ▶ **external-attacher**
ControllerPublishVolume
ControllerUnpublishVolume



We are here!

Introduction

Container Storage Interface

CSI CephFS

Manila shares with Kubernetes

Results, numbers, plots...

CSI CephFS overview

`github.com/ceph/ceph-csi`

- ▶ Provides an interface between a CSI-enabled Container Orchestrator and the Ceph cluster
- ▶ Provisions and mounts CephFS volumes
- ▶ Supports both the kernel CephFS client and the CephFS FUSE driver



CSI CephFS overview

Compared to Kubernetes in-tree CephFS volume plugin

- ▶ In-tree volume plugins to be eventually migrated to CSI
- ▶ Decoupled from Kubernetes
- ▶ Ability to choose between mounting tools
- ▶ Planned support for volume expansion, snapshots

We are here!

Introduction

Container Storage Interface

CSI CephFS

Manila shares with Kubernetes

Results, numbers, plots...

Manila external provisioner for Kubernetes overview

github.com/kubernetes/cloud-provider-openstack

- ▶ Provisions new Manila shares, fetches existing ones
- ▶ Maps them to Kubernetes *PersistentVolume* objects
- ▶ Currently supports CephFS shares only (both in-tree CephFS plugin and csi-cephfs)
- ▶ Supports authentication using both user credentials as well as trustees
 - ▶ Magnum → Kubernetes + manila-provisioner StorageClass + trustee secrets = Manila support out-of-the-box
- ▶ The future is in CSI



We are here!

Introduction

Container Storage Interface

CSI CephFS

Manila shares with Kubernetes

Results, numbers, plots...

Benchmarks

Goals

1. Verify the CSI CephFS implementation for common use cases
2. Verify the Manila Provisioner implementation
3. Test CSI CephFS driver behavior on a heavy loaded cluster

Client

1. Kubernetes v1.12.1, csi-cephfs 0.3.1

Clusters

1. *Dwight*: 3x24 HDD OSDs, 3 MDS, Ceph Luminous Bluestore
2. *Jim*: 300 SSD OSD, 2 MDS, Ceph Luminous Bluestore, hyper-converged

Benchmarks

Methodology

1. Provision s CephFS shares using *manila-provisioner*
2. Create a Deployment with r replicas, sized so we get one pod per node
3. Mount s provisioned shares into each pod using csi-cephfs (fuse)
4. Measure time taken for all pods to become *Running*, MDS sessions, hcr/s

Tests

1. *idle*: do nothing
2. *busy*: unpack a large archive (linux kernel)

Parameters

1. $s = 100, r = 100$; 10'000 idle clients
2. $s = 10, r = 100$; 1'000 busy clients

Idle benchmark - attempt #1

Our very first test of csi-cephfs with concurrent workloads

Preparation

- ▶ 10 CephFS shares
- ▶ 100 replicas
- ▶ The goal is to have 1'000 idle clients running

Idle benchmark - attempt #1

Our very first test of csi-cephfs with concurrent workloads

Preparation

- ▶ 10 CephFS shares
- ▶ 100 replicas
- ▶ The goal is to have 1'000 idle clients running

Outcome

- ▶ :(

Idle benchmark - attempt #1

```
Aug 30 09:51:45 cci-cephfs-scale-003-n3tf4nqlzisk-minion-56.cern.ch runc[3255]:  
↳ E0830 09:51:45.410380      3270 csi_attacher.go:137] kubernetes.io/csi:  
↳ attacher.WaitForAttach failed for volume  
↳ [pvc-c5848f32-ac39-11e8-bbfb-02163e01b7c5] (will continue to try):  
↳ volumeattachments.storage.k8s.io  
↳ "csi-4f2dbe5cb257e7d7b172c4a1e6a1d26bfff82dabeb91e441c527d46f368f1615" is  
↳ forbidden: User "system:node:cci-cephfs-scale-003-n3tf4nqlzisk-minion-56" cannot  
↳ get volumeattachments.storage.k8s.io at the cluster scope: no path found to  
↳ object
```

Idle benchmark - attempt #1

```
User "system:node:NODE_NAME" cannot get volumeattachments at the  
cluster scope: no path found to object
```

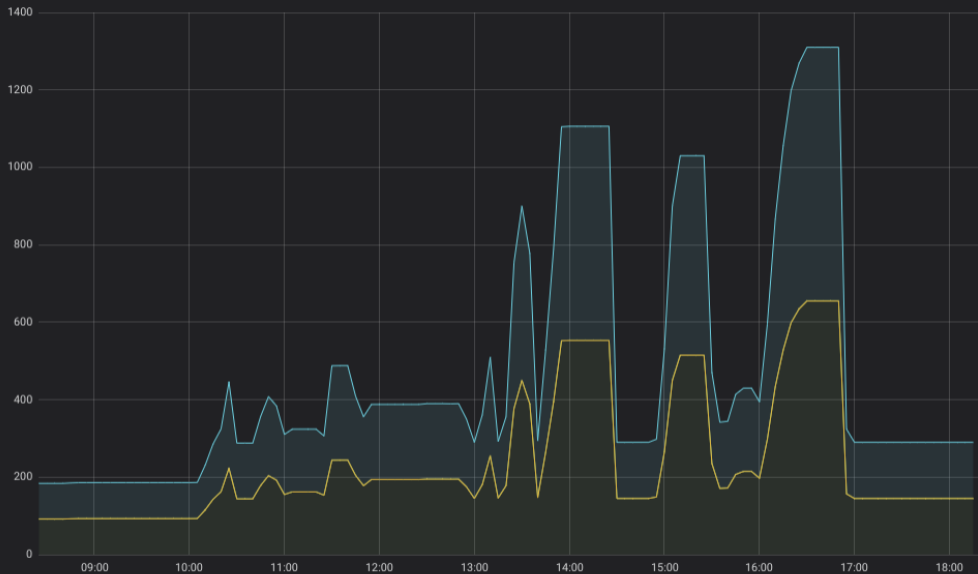
- ▶ Provisioning of shares worked just fine
- ▶ Some pods survived
- ▶ Others that reported this error would never recover

Idle benchmark - attempt #2

A script that:

- ▶ Scales the deployment in small increments
- ▶ Kills pods that take too long to create (got stuck in the *VolumeAttachment* error)

dwight: MDS Sessions



— cephdwightmds0
— cephdwightmds2

	max	avg	current
cephdwightmds0	655	238	145
cephdwightmds2	655	238	145

Idle benchmark - attempt #2

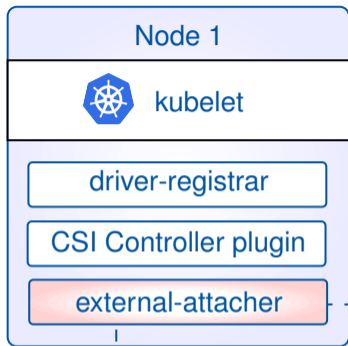
Outcome

- ▶ We've managed to get up to 655 concurrent clients (could be even more)
- ▶ Slow and ugly but somehow working

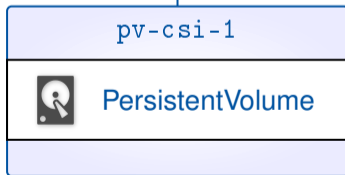
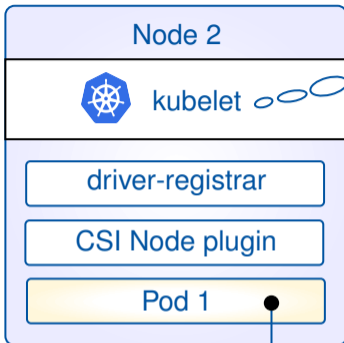
Idle benchmark - attempt #3

"Third time's the charm"?

- ▶ Kubernetes 1.12, driver-registrar 0.4 released
- ▶ Kubelet plugin registration of CSI drivers
- ▶ *CSISkipAttach*
 - ▶ Skips the creation of *VolumeAttachment* objects
 - ▶ Volumes are marked as attached immediately

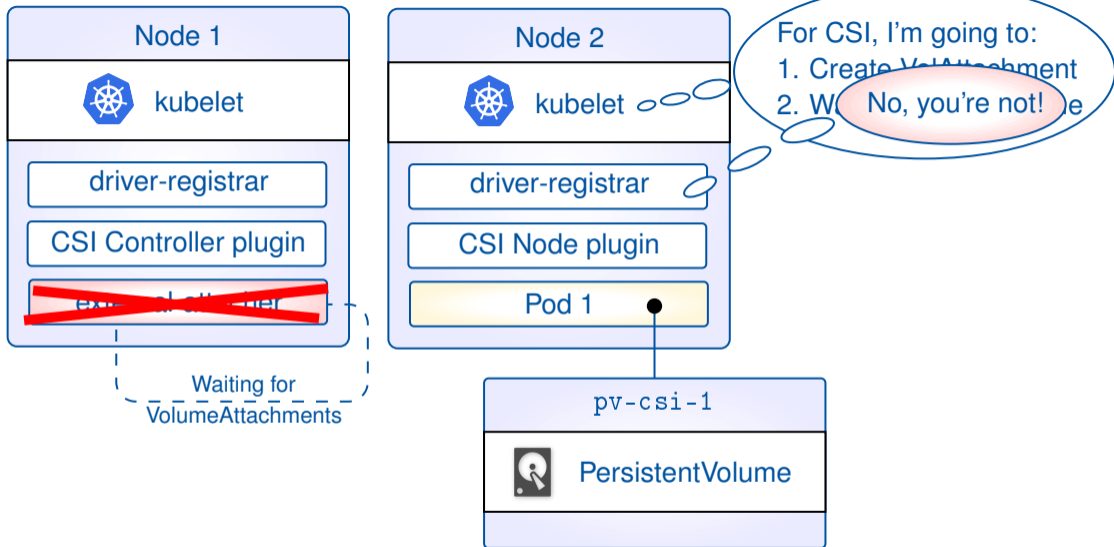


Waiting for
VolumeAttachments



For CSI, I'm going to:
1. Create VolAttachment
2. Wait for Attached=true

* external-provisioner omitted from image



* external-provisioner omitted from image

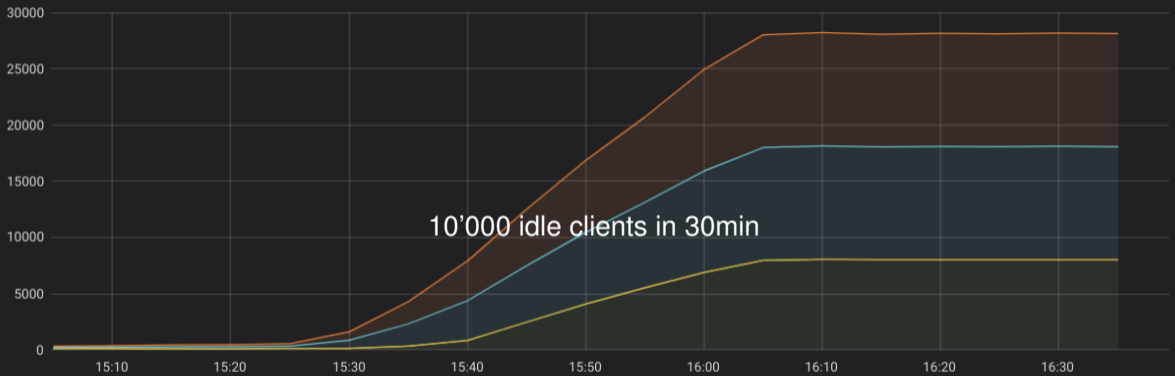
Idle benchmark - attempt #3

- ▶ We resumed our tests using the new versions of Kubernetes and driver-registrar
- ▶ Parameters: 100 CephFS Shares * 100 replicas = 10'000 idle clients
- ▶ Gradual, gentle scale up

Idle benchmark - attempt #3

```
15:29:28 : scaling to 5 replicas  
15:32:33 : scaling to 10 replicas  
...  
15:57:18 : scaling to 80 replicas  
15:58:51 : scaling to 85 replicas  
16:00:29 : scaling to 90 replicas  
16:02:26 : scaling to 95 replicas  
16:04:16 : scaling to 100 replicas
```

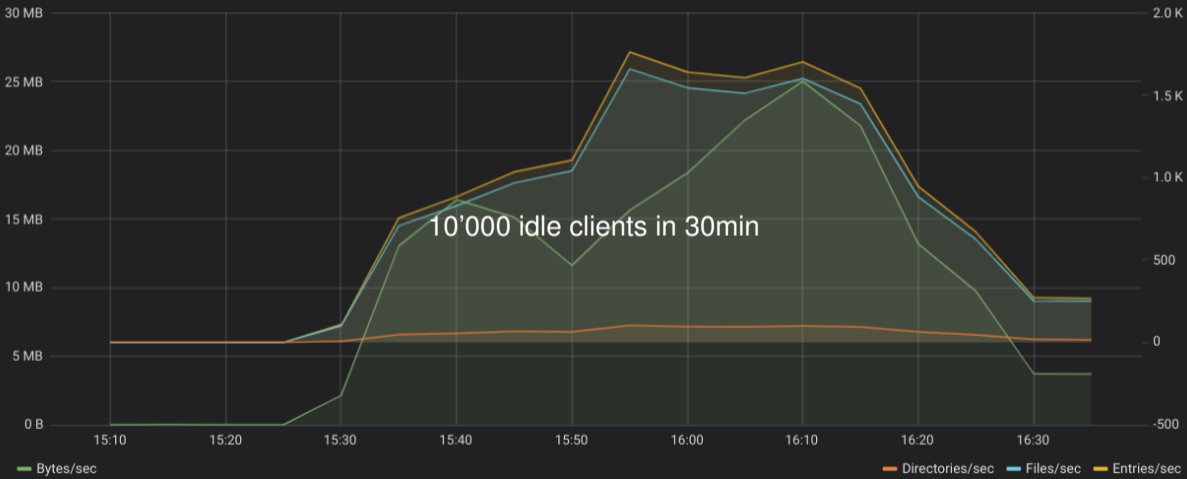
dwight: MDS Sessions



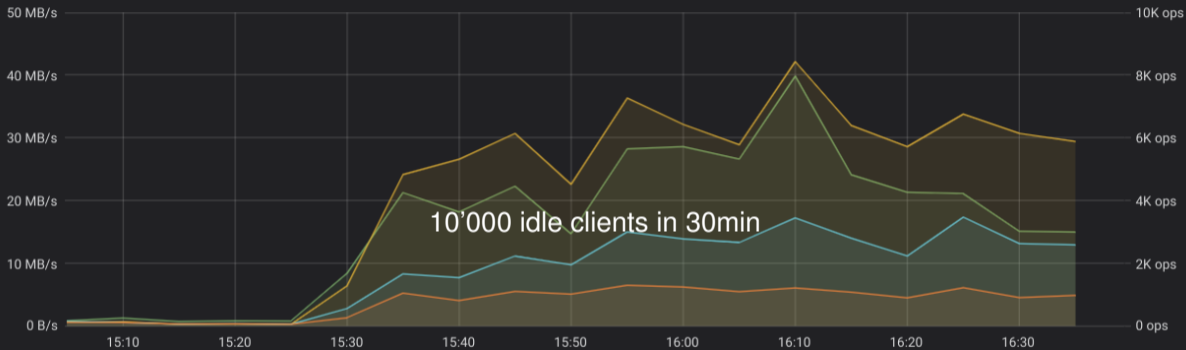
- cephdwightmds0
- cephdwightmds2
- p05153026942913

	max	avg	current
cephdwightmds0	8042	4039	8011
cephdwightmds2	10082	5551	10056
p05153026942913	10082	5552	10056

dwight: Namespace Rates



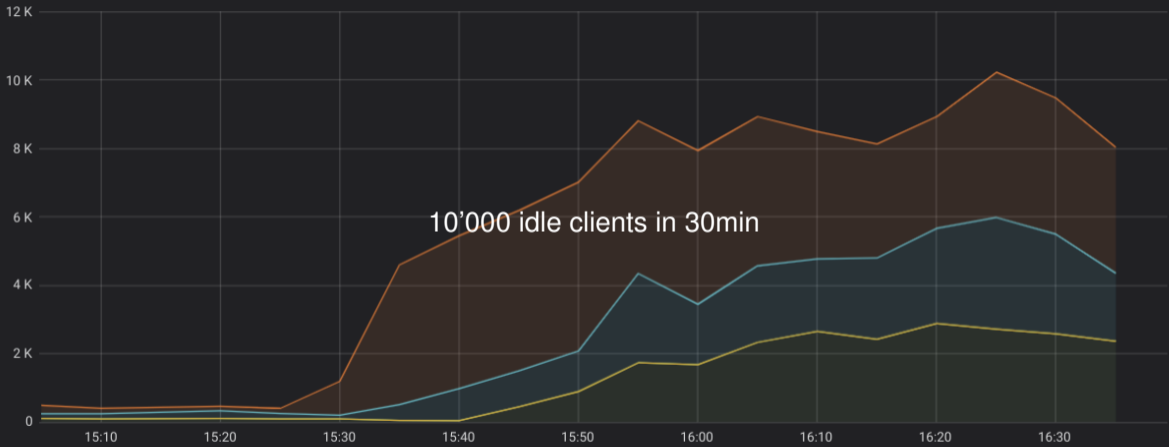
dwight: CephFS Pool IO



10,000 idle clients in 30min

	max	avg	current
Data Bps	39.8 MB/s	16.3 MB/s	15.0 MB/s
Data IOPS (right-y)	8K ops	4K ops	6K ops
Metadata Bps	17.4 MB/s	9.0 MB/s	13.0 MB/s
Metadata IOPS (right-y)	1K ops	770 ops	980 ops

dwight: MDS Client Requests (hcr/s)



10'000 idle clients in 30min

cephdwightmds0 cephdwightmds2 p05153026942913

Idle benchmark - attempt #3

Some bits still to be understood

```
mds.cephdwightmds2 mds.1 redacted:6800/2246956767 2238 : cluster [WRN] evicting
```

```
↪ unresponsive client cci-cephfs-scale-001-3vimalfm74dd-minion-76.cern.ch
```

```
↪ (674953366), after 304.623890 seconds
```

```
mds.cephdwightmds0 mds.2 redacted:6800/2942371108 2167 : cluster [WRN] evicting
```

```
↪ unresponsive client cci-cephfs-scale-001-3vimalfm74dd-minion-76.cern.ch:
```

```
↪ (674953366), after 304.277996 seconds
```

Idle benchmark - attempt #3

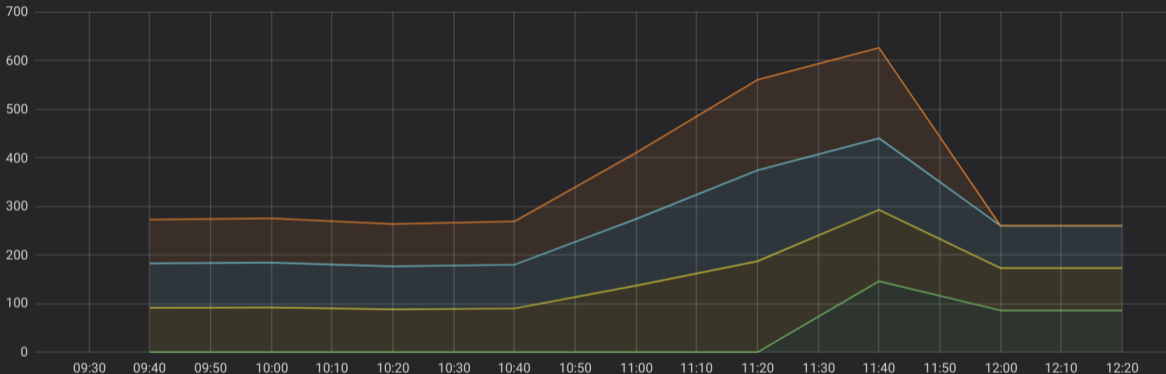
Some bits still to be understood

<code>cci-cephfs-scale-001-3vimalfm74dd-minion-13</code>	Ready	<none>	8d	v1.12.1
<code>cci-cephfs-scale-001-3vimalfm74dd-minion-14</code>	NotReady	<none>	8d	v1.12.1

Busy benchmark - attempt #1

- ▶ Each client extracting the linux kernel
- ▶ Slowly ramp up until it breaks

dwight: MDS Sessions



	max	avg	current
cephdwightmds0	146	35	86
cephdwightmds1	187	112	87
cephdwightmds2	187	112	87
p05153026942913	186	96	0



Dan Van Der Ster 11:44 AM

i think you crashed something
mds's are flapping for some reason
HEALTH_OK now....



Ricardo Brito Da Rocha 11:47 AM

nice!
my day is done then 😊



Robert Vasek 1:47 PM

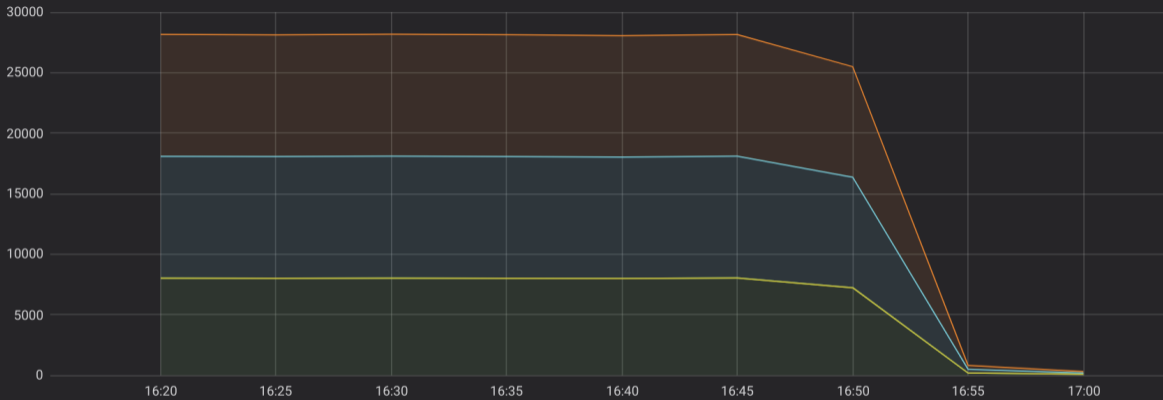
very nice!

Bonus benchmark - client deletion

- ▶ Stop 10'000 clients and delete their shares, simultaneously
- ▶ Kubernetes and driver OK, CephFS OK, Manila share daemon needed a kick

	...		pvc-08d...		1		CEPHFS		deleting		False		Geneva CephFS Testing		nova	
	...		pvc-7a9...		1		CEPHFS		deleting		False		Geneva CephFS Testing		nova	

dwight: MDS Sessions



cephdwightmds0	max	avg	current
cephdwightmds2	8038	6181	93
p05153026942913	10071	7763	109
	10071	7764	109

Conclusion & Next Steps

To recap:

- ▶ Standardized storage interface for Container Orchestrators with CSI
- ▶ Works nicely in Kubernetes, others soon to follow
- ▶ manila-provisioner + csi-cephfs handle large concurrency and scaling well
- ▶ Already in production at CERN

Next Steps:

- ▶ Add support for volume expansion and snapshots
- ▶ Make the Manila Provisioner a CSI plugin

Questions?

- ▶ Robert Vasek <robert.vasek@codefreax.org>
- ▶ Ricardo Rocha <ricardo.rocha@cern.ch> @ahcorporto
- ▶ CSI CephFS: <https://github.com/ceph/ceph-csi>
- ▶ Manila Provisioner: <https://github.com/kubernetes/cloud-provider-openstack>