# Pods-as-Volumes: Effortlessly Integrating Storage Systems and Middleware into Kubernetes

Alberto Faria

INESC TEC & University of Minho

Ricardo Macedo

João Paulo INESC TEC & University of Minho

INESC TEC & University of Minho INE

#### **Kubernetes**

- A container orchestration system
- Automates many aspects of app deployment
  - **Scheduling** to nodes, **scaling**, handling **failures**, ...
- Declarative objects specify desired state
  - Kubernetes **autonomously** adjusts actual state
- **Pods** are sets of one or more containers
  - Defined by Pod objects
  - Containers of a pod run **together** on a single node



apiVersion: v1

kind: Pod
metadata:

name: my-pod

spec:

#### containers:

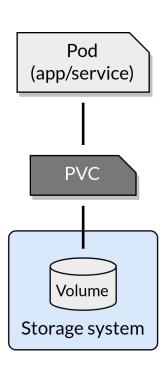
- **name**: container

image: ubuntu:20.04

command: [ sleep, infinity ]

### **Kubernetes: Volumes**

- Can manage **storage** resources through the **volume** abstraction
  - Each volume corresponds to a file system or block device
- User creates a volume by defining a PersistentVolumeClaim (PVC)
  - Describes **desired properties** of the volume (e.g., storage capacity)
- PVC objects also (indirectly) specify a volume provisioner
  - Provisioner implements all volume (de)allocation logic
  - Interacts with underlying **storage system**
- Pod containers can access a volume by mounting its PVC
  - Provisioner coordinates with storage system to expose volume to containers



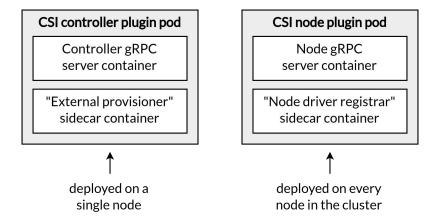
## **Kubernetes: Container Storage Interface (CSI)**

- Kubernetes includes several built-in provisioners
  - CephFS provisioner, NFS provisioner, iSCSI provisioner, ...
- Can build provisioners by implementing the Container Storage Interface (CSI)
  - gRPC interface specification: https://github.com/container-storage-interface/spec
  - Standardizes interaction between **container orchestration** systems and **storage** systems



## A problem: Implementing CSI takes effort

- Implementing CSI fully and correctly can be difficult and time-consuming
  - Must typically build **2** gRPC servers
  - Correct error handling and cleanup is tricky
- Must deploy components across all nodes of the cluster
  - CSI controller plugin pod: gRPC server container, "external provisioner" sidecar container
  - **CSI node plugin** pod: gRPC server container, "node driver registrar" sidecar container



## Another problem: Kubernetes storage is inflexible

- Limited ability to manage storage stacks
- Can only represent volumes and attach them directly to pods
  - Can't define intermediate **layers** to create more complex **stacks**
- Leads to lack of modularity and reimplementation of features
- E.g., no **general** mechanism to enable volume **encryption** 
  - Each storage system/provisioner must implement their own
  - Or each client pod must **embed** encryption logic
- Problem is exacerbated when combining several storage functionalities
  - *E.g.*, applying **compression** *and* **encryption**

## A solution: Pods-as-Volumes (PaV)

- Kubernetes plugin that simplifies the implementation of new volume provisioners
- Allows specifying logic for creating, deleting, exposing volumes as pod templates
  - Pods automatically instantiated from templates when needed
- Easier integration of storage systems into Kubernetes
- Enables straightforward creation of **storage middleware** components
- Open source: https://github.com/albertofaria/pav

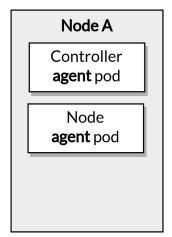
## PavProvisioner objects

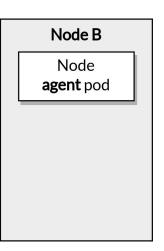
- New object kind: PavProvisioner
- Each PavProvisioner object implements a volume provisioner
- Pod templates embedded in object definition
- Can use PavProvisioner in the same manner as built-in provisioners

```
apiVersion: pav.albertofaria.github.io/v1alpha1
kind: PavProvisioner
metadata: ...
spec:
  provisioningModes: ...
  volumeValidation:
    volumeModes:
    accessModes:
    minCapacity: ...
    maxCapacity: ...
    podTemplate: ...
  volumeCreation:
    volumeHandle: ...
    capacity: ...
    podTemplate: ...
  volumeDeletion:
    podTemplate: ...
  volumeStaging:
    podTemplate: ...
  volumeUnstaging:
    podTemplate: ...
```

### **Architecture**

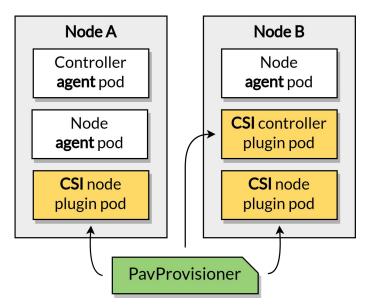
- Two main components: agent and CSI plugin
- The **agent** is deployed once per Kubernetes **cluster** 
  - Controller agent, deployed as a single pod for the entire cluster
  - Node agent, deployed on all nodes





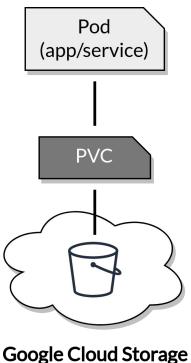
## **Architecture (cont.)**

- For each PavProvisioner, a CSI plugin instance is deployed
  - **CSI controller plugin**, deployed as a single pod for the entire cluster
  - **CSI node plugin**, deployed on all nodes
- Kubernetes **interacts** with CSI plugins
  - CSI plugins **delegate** most work to the agent



## **Use cases: Google Cloud Storage integration**

- Google's **object storage** cloud service
- Users create **buckets**, which contain **objects** 
  - **Libraries** in many languages to access buckets
  - Buckets can also be **mounted** locally as a **file system**
- No **built-in** Kubernetes integration
- We create one with PaV
  - Each volume is a bucket
  - Bucket is **mounted** into client pods as a **file system**



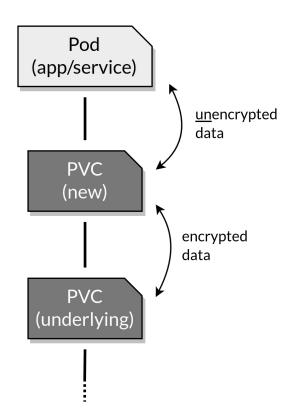
## **Use cases: Google Cloud Storage integration (cont.)**

- **62 lines** of YAML
- No coding required
- Relies on existing gsutil and gcsfuse tools

```
1 apiVersion: pav.albertofaria.github.io/v1alpha1
2 kind: PavProvisioner
3 metadata: { name: gcs-provisioner }
4 spec:
    provisioningModes: [ Dynamic, Static ]
                                                                          volumeStaging:
    volumeCreation:
      capacity: 1Ei
                                                                     38
                                                                            podTemplate:
                                                                              metadata: { namespace: "{{params.secretNamespace}}" }
      podTemplate:
         metadata: { namespace: "{{params.secretNamespace}}" }
                                                                     40
                                                                                restartPolicy: Never
         spec: &gsutil-pod-spec
                                                                     41
                                                                                containers:
11
           restartPolicy: Never
                                                                     43
                                                                                  - name: gcsfuse
12
           containers:
             - &gsutil-container
                                                                                    image: albertofaria/gcsfuse: 0.36.0
                                                                                    command: [ /bin/bash, -c ]
               name: gsutil
               image: albertofaria/gsutil:5.2
                                                                     46
15
                                                                                    args:
               command: [ gsutil, -o,
                                                                     48
                                                                                        mkdir /pav/volume &&
                 Credentials:gs_service_key_file=/secret/key,
                 -o. "GSUtil:default_project_id={{
                                                                     49
                                                                                        gcsfuse -o=allow_other \
18
                                                                                          --key-file=/secret/key --dir-mode=777 \
19
                 params.projectId }}" ]
                                                                     51
                                                                                          --temp-dir=/temp --file-mode=666 \
20
               args: [ mb, -b, "on",
                                                                     52
                                                                                          --stat-cache-ttl=0 --type-cache-ttl=0 \
21
                 -1, "{{ params.location or 'US' }}",
                 "gs://{{ defaultVolumeHandle }}" ]
                                                                     53
                                                                                          {{ volumeHandle|tobash }} /pav/volume &&
                                                                     54
                                                                                        touch /pav/ready &&
23
               volumeMounts:
                                                                     55
24
                 - { name: secret, mountPath: /secret }
                                                                                         sleep infinity
25
                                                                     56
                                                                                    securityContext: { privileged: true }
           volumes:
26
             - &secret-volume
                                                                     57
                                                                                    volumeMounts:
                                                                                      - { name: secret, mountPath: /secret }
27
               name: secret
                                                                     58
28
               secret: { secretName: "{{params.secretName}}" }
                                                                     59
                                                                                      - { name: temp, mountPath: /temp }
29
    volumeDeletion:
                                                                     60
                                                                                volumes:
       podTemplate:
                                                                     61
                                                                                  - *secret-volume
31
         metadata: { namespace: "{{params.secretNamespace}}" }
                                                                                  - { name: temp, emptyDir: {} }
32
         spec:
33
           <<: *gsutil-pod-spec
34
           containers:
35
             - <<: *gsutil-container
                                                                                                                                12
36
               args: [ rm, -r, "gs://{{ volumeHandle }}" ]
```

## Use cases: Transparent encryption middleware

- No **general** mechanism for **volume encryption** in Kubernetes
  - Each storage system/provisioner must implement **their own**
  - Or each client pod must embed encryption logic
- We create an **encryption middleware** with PaV
  - **Transparent encryption** for block volumes
- User creates new PVC referencing underlying PVC
  - Client pods mount new PVC



## Use cases: Transparent encryption middleware (cont.)

- **70 lines** of YAML
- No coding required
- Relies on existing cryptsetup utility
- Pod templates mount the underlying PVC

```
1 apiVersion: pav.albertofaria.github.io/v1alpha1
 2 kind: PavProvisioner
 3 metadata: { name: crvpt-provisioner }
    provisioningModes: [ Dynamic ]
                                                                          volumeDeletion:
    volumeValidation: { volumeModes: [ Block ] }
     volumeCreation:
                                                                            podTemplate:
                                                                     43
                                                                              metadata: { namespace: "{{pvc.metadata.namespace}}" }
       podTemplate:
         metadata: { namespace: "{{pvc.metadata.namespace}}" }
                                                                     44
9
                                                                     45
                                                                                <<: *crvptsetup-pod-spec
10
         spec: &cryptsetup-pod-spec
11
           restartPolicy: Never
                                                                     46
                                                                                containers:
                                                                     47
                                                                                   - <<: *crvptsetup-container
12
           containers:
13
                                                                     48
                                                                                    args: [ cryptsetup -q erase /volume ]
             - &cryptsetup-container
                                                                          volumeStaging:
14
               name: cryptsetup
                                                                     50
                                                                            podTemplate:
15
                image: albertofaria/cryptsetup:2.4.1
                                                                              metadata: { namespace: "{{pvc.metadata.namespace}}" }
16
               command: [ /bin/bash. -c ]
                                                                     51
                                                                     52
17
               args:
                                                                     53
                                                                                <<: *cryptsetup-pod-spec
18
                                                                     54
19
                   set -o errexit -o pipefail
                                                                                containers:
20
                                                                     55
                                                                                   - <<: *cryptsetup-container
                   cryptsetup -q luksFormat \
21
                      /volume /secret/passphrase
                                                                     56
                                                                                    args:
22
                   size="$( blockdev --getsize64 /volume )"
                                                                     57
                                                                                         dev={{ volumeHandle|tobash }} &&
23
                   offset="$( cryptsetup luksDump \
                                                                     58
                     /volume --dump-json-metadata |
                                                                     59
                                                                                         cryptsetup open /volume "$dev" \
24
                                                                                           --key-file /secret/passphrase &&
                     jq '.segments."0".offset' | tr -d '"' )"
                                                                     60
25
26
                   echo "$(( size - offset ))" > /pav/capacity
                                                                     61
                                                                                        cp -p "/dev/mapper/$dev" /pav/volume
               securityContext: { privileged: true }
                                                                          volumeUnstaging:
27
28
               volumeMounts:
                                                                     63
                                                                            podTemplate:
                                                                              metadata: { namespace: "{{pvc.metadata.namespace}}" }
29
                 - { name: secret, mountPath: /secret }
                                                                     64
30
               volumeDevices:
                                                                     65
31
                 - { name: underlying, mountPath: /volume }
                                                                     66
                                                                                <<: *crvptsetup-pod-spec
                                                                     67
32
           volumes:
                                                                                containers:
33
             - name: secret
                                                                     68
                                                                                   - <<: *crvptsetup-container
                                                                     69
34
                                                                                    args: [ "cryptsetup close
               secret:
                                                                                      {{ volumeHandle|tobash }} || (( $? == 4 ))" ]
35
                 secretName: "{{ pvc.metadata.annotations[
                                                                     70
36
                    'crvpt/secretName'] }}"
37
             - name: underlying
38
               persistentVolumeClaim:
                                                                                                                                14
                 claimName: "{{ pvc.metadata.annotations[
39
                   'crypt/underlyingClaimName'] }}"
```

## **Summary**

- Implementing new provisioners using CSI is difficult and time-consuming
- Kubernetes' ability to manage storage stacks is limited
  - Insufficient expressiveness to represent and manage storage layers
- We propose Pods-as-Volumes (PaV)
  - **Easily** implement new volume **provisioners**
  - Specify **logic** underlying volume lifecycle and behavior as **pod templates**
- Simplifies integration of storage systems into Kubernetes
- Enables straightforward creation of composable storage middleware components

## PaV is open source!

