# So what if I don't want my Persistent Storage to be yet another Bind Mount
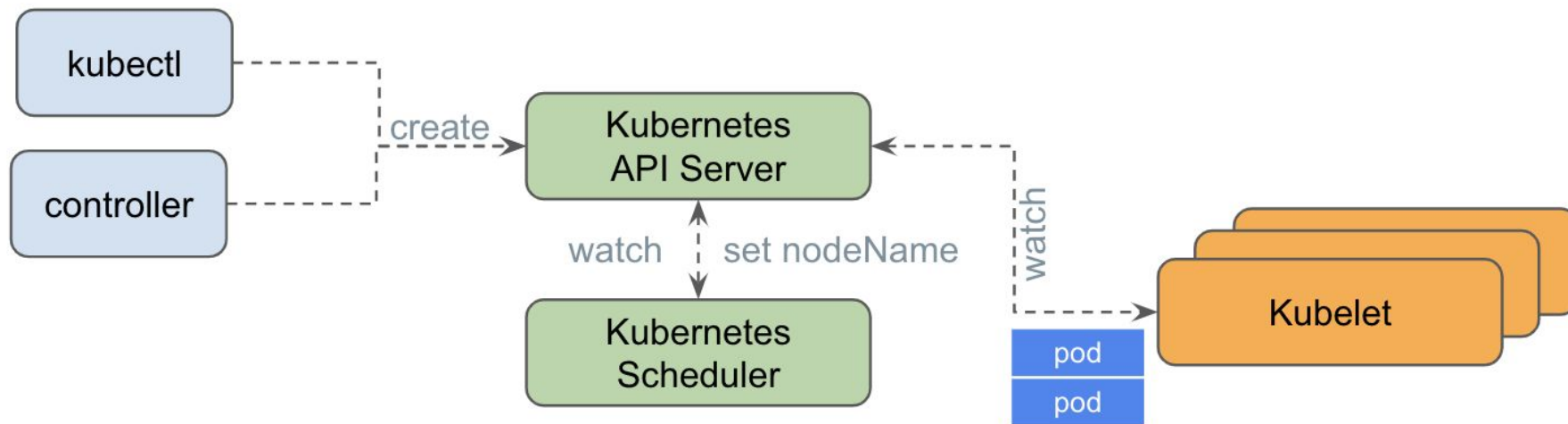
*Deep Debroy and Feng Wang*

# Introduction

- Current Lifecycle Flow for Pods with Persistent Storage
  - bind mounts

- Alternatives to the Standard Flow for Persistent Storage
  - from a MicroVM perspective
  - challenges

- Surfacing "Direct Assignment" Flows to CSI Plugins
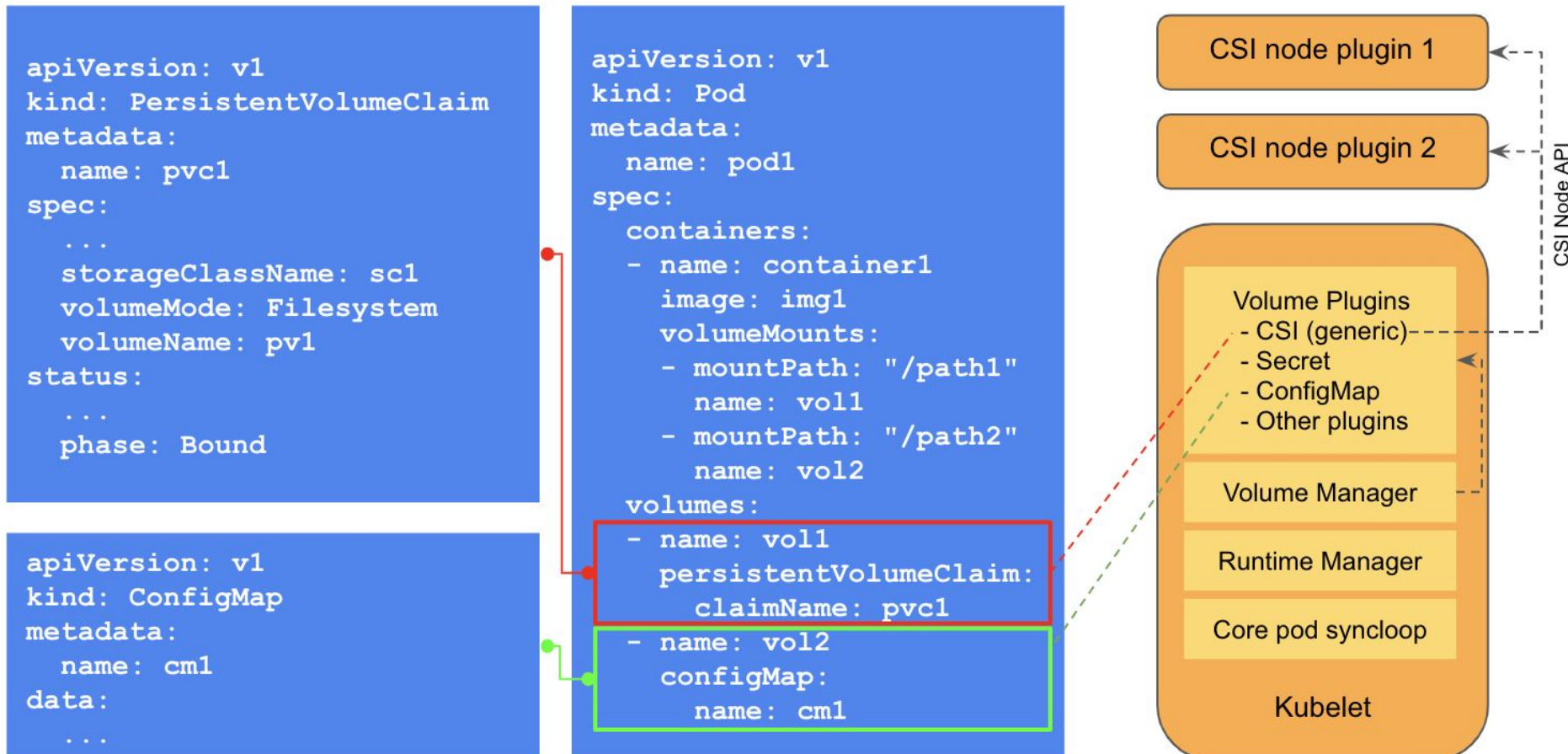  - with Kubernetes and Kata awareness
  - generic support

- Conclusion

# Pod Startup from Kubelet's Perspective

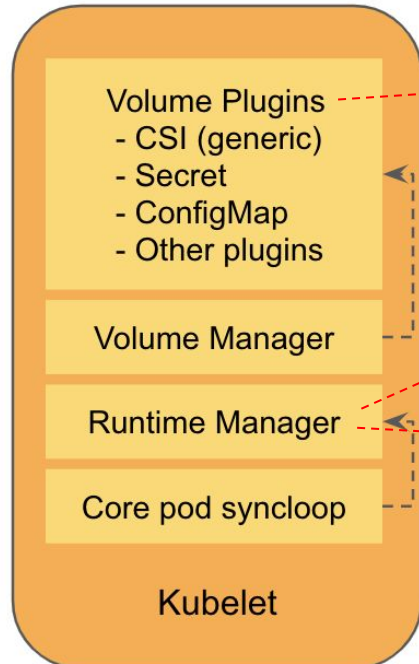- Kubelet gets notified about pods scheduled on the node

- Kubelet mounts the specified volumes in the pod:
  - In-line volumes (configMap, secrets, etc)
  - Persistent Volumes bound to PVCs

# Pod Startup from Kubelet's Perspective
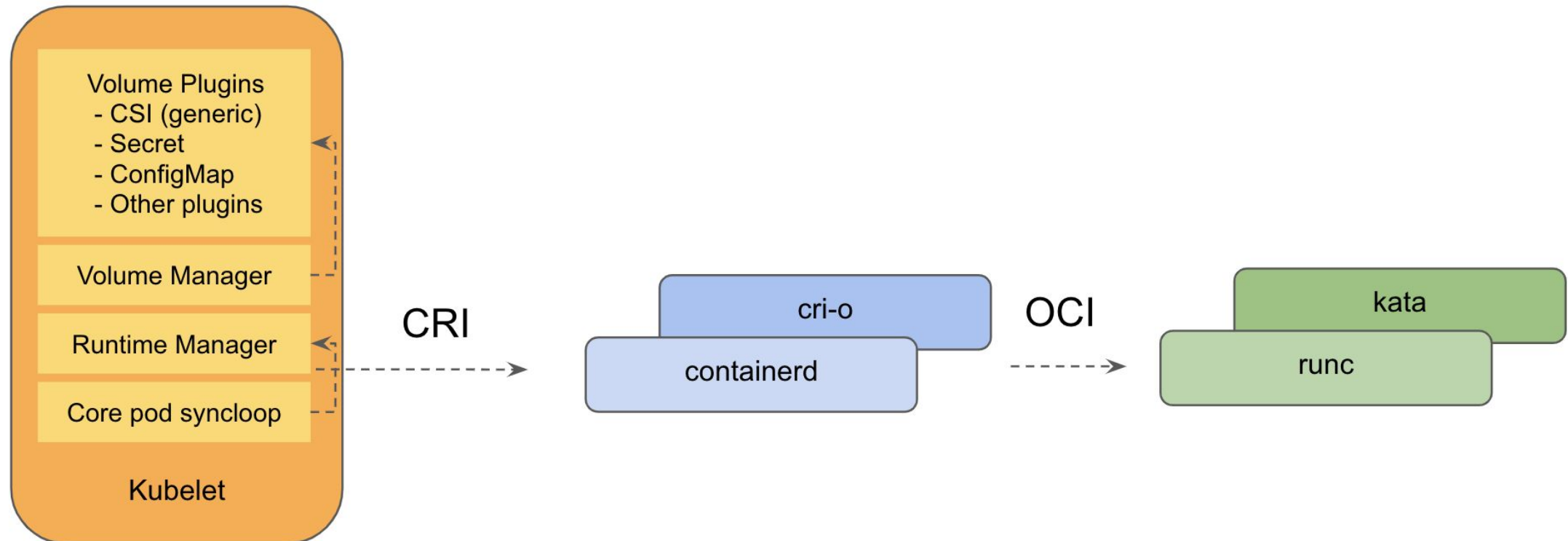
- Kubelet prepares the file system mount

  - Apply FsGroup based on FsGroupChangePolicy

  - Probe and prepare subpaths

  - Determine Security-Enhanced Linux
    (SELinux) labelling



```yaml
apiVersion: v1
kind: Pod
metadata:
  name: pod1
spec:
  securityContext:
    fsGroup: 4059
    fsGroupChangePolicy: "OnRootMismatch"
    seLinuxOptions:
      level: "s0:c123,c456"
  containers:
  - name: container1
    image: image1
    volumeMounts:
    - mountPath: "/path1"
      name: vol1
      subPath: /subpath1
  volumes:
  - name: vol1
    persistentVolumeClaim:
      claimName: pvc1
```
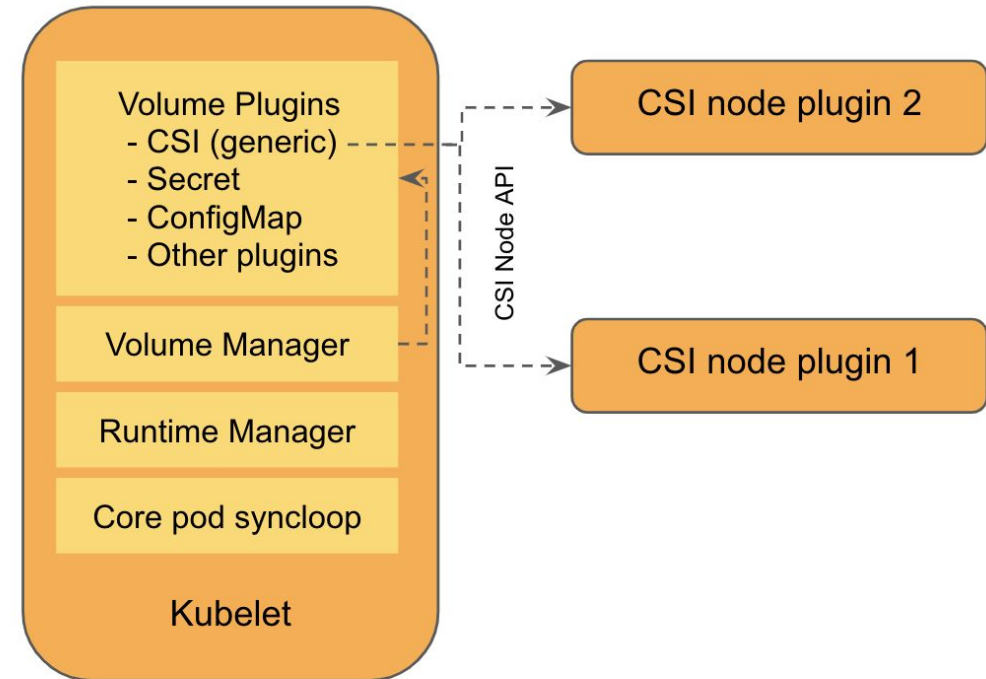
Volume Plugins
- CSI (generic)
- Secret
- ConfigMap
- Other plugins

Volume Manager

Runtime Manager

Core pod syncloop

Kubelet

# Pod Startup from Kubelet's Perspective

- Kubelet invokes CRI implementation to
  - Create pod sandbox
  - Pull container images
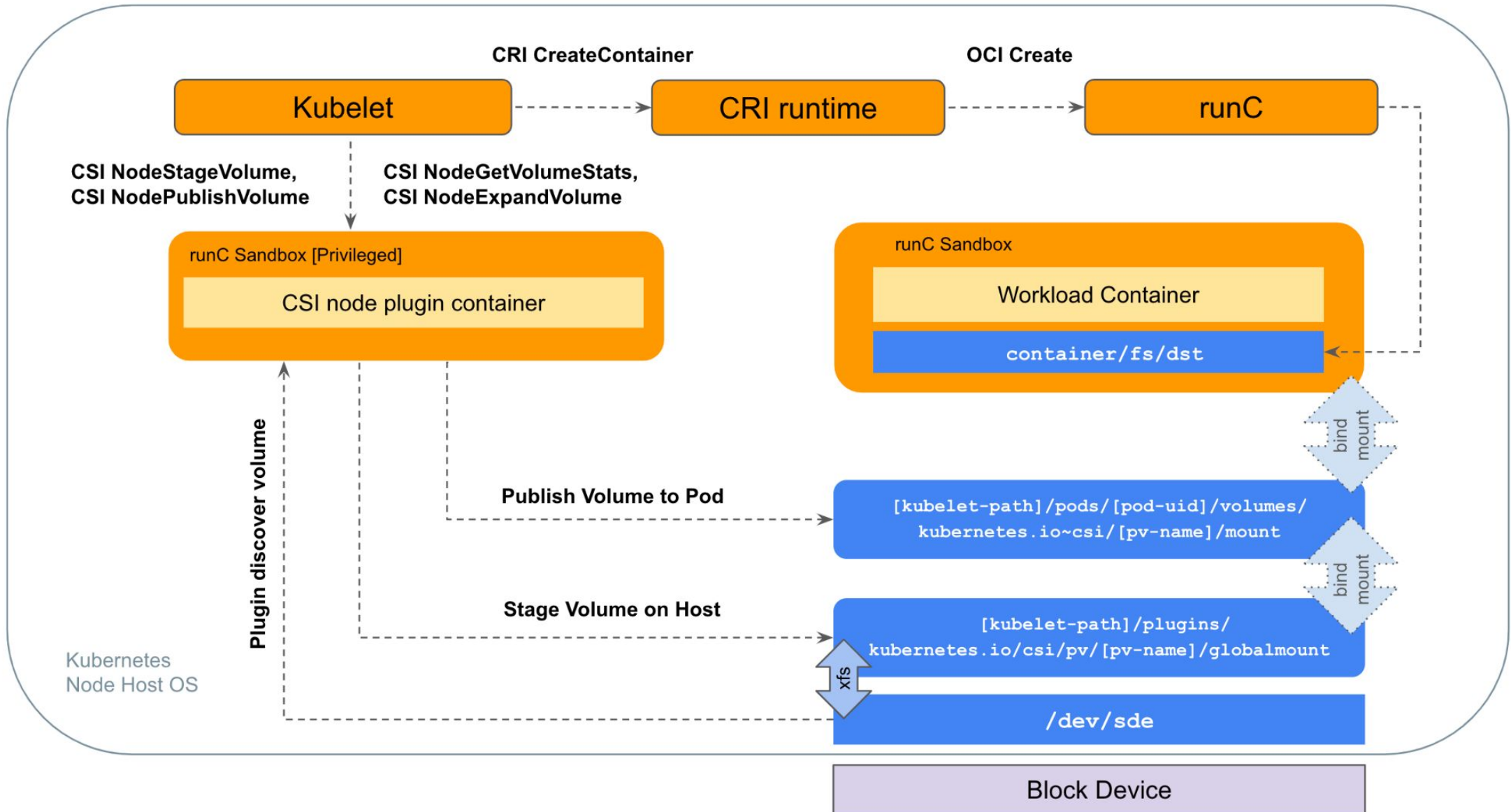  - Create containers

# Volume Operations during Pod Lifetime

- Kubelet reports file system stats

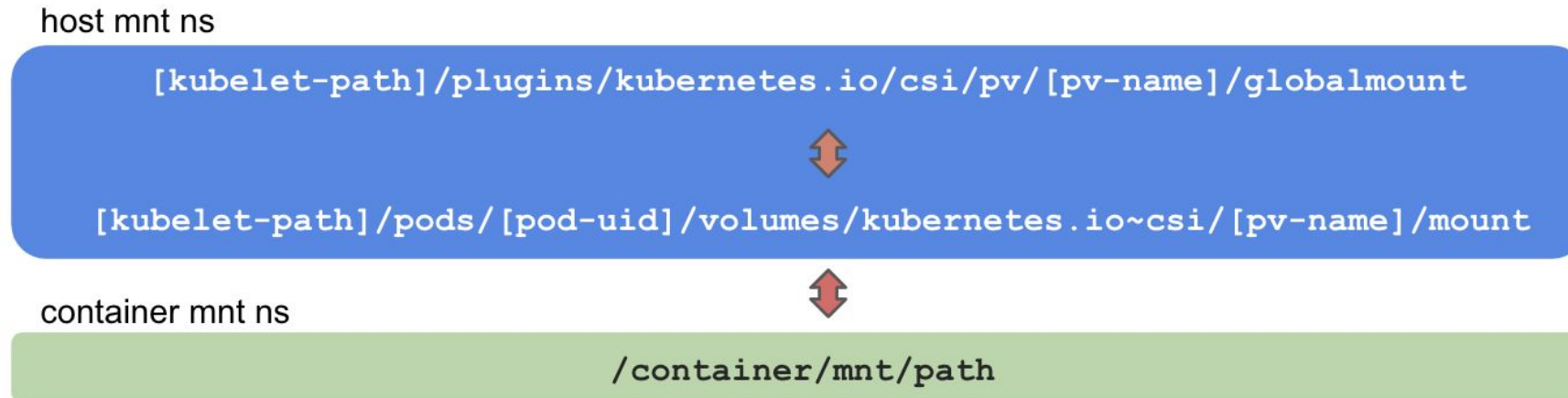- Kubelet resizes the mounted volume

# Regular Path for FS Operations on a PV

- NodeStaging global path → NodePublish target path
  - Created during CSI NodePublish

- NodePublish target path → Path inside container
  - Created during OCI create container

# Current Assumptions in Kubelet

- All volumes have file-systems mounted before container bring-up

  - whenever PV VolumeMode is FileSystem

- Post mount actions *invoked on mounted file system* by Kubelet (and container runtime):

  - fsGroup ownership

  - subpath checks

  - SELinux labelling

**How about a different sequence?**

## How about a different sequence?

mount volumes *after* container bringup *without* using raw block mode
[ use case: microVM environments like Kata ]
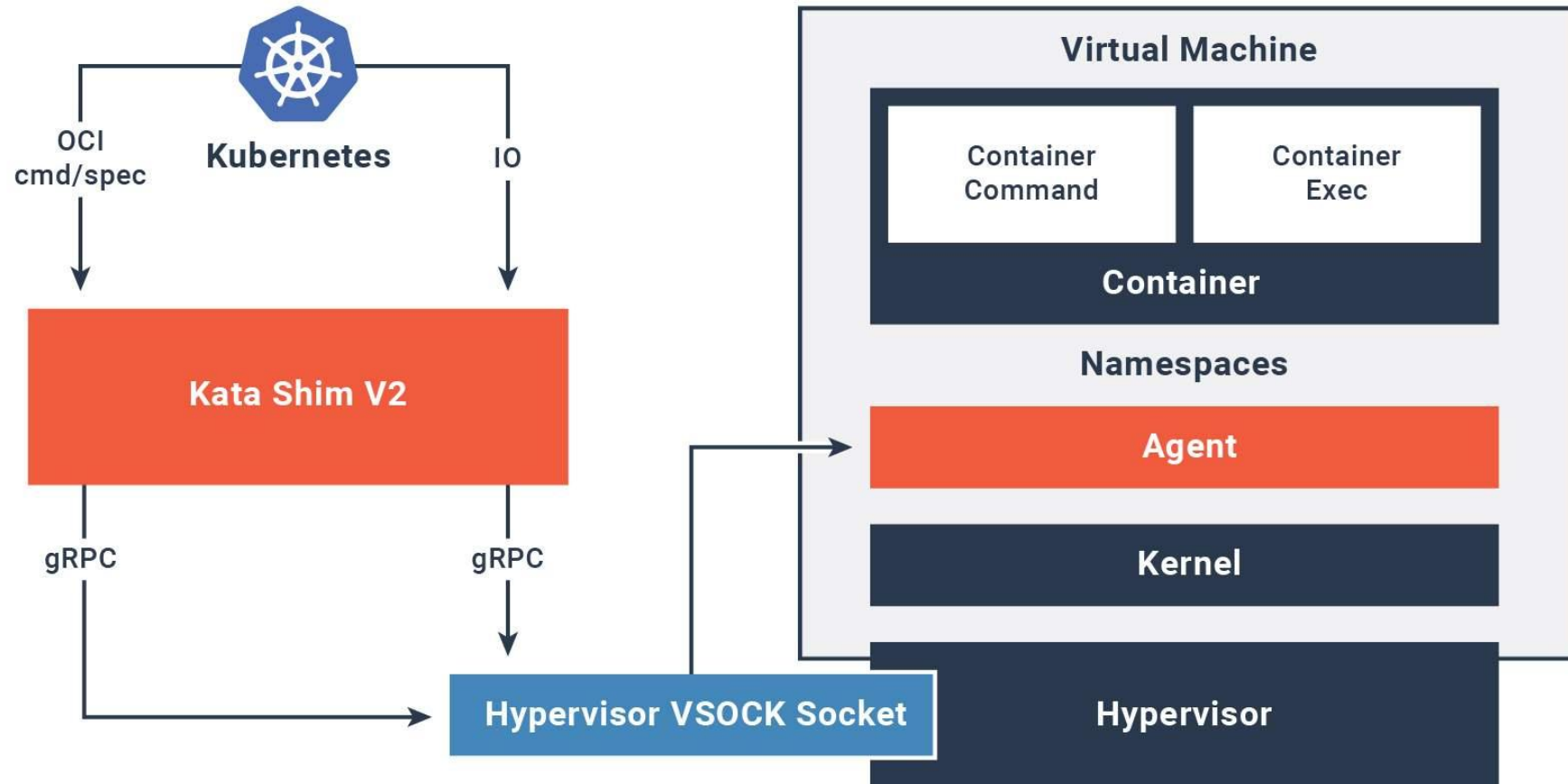
# MicroVM-based Runtime

Kata Containers – A container runtime that runs workloads in a virtual machine.

- **Secure** – container workload runs on a dedicated kernel

- **Lightweight** – faster startup than a full-blown virtual machine

- **Compatible** – compliant to OCI container format and containerd/crio shim interface.
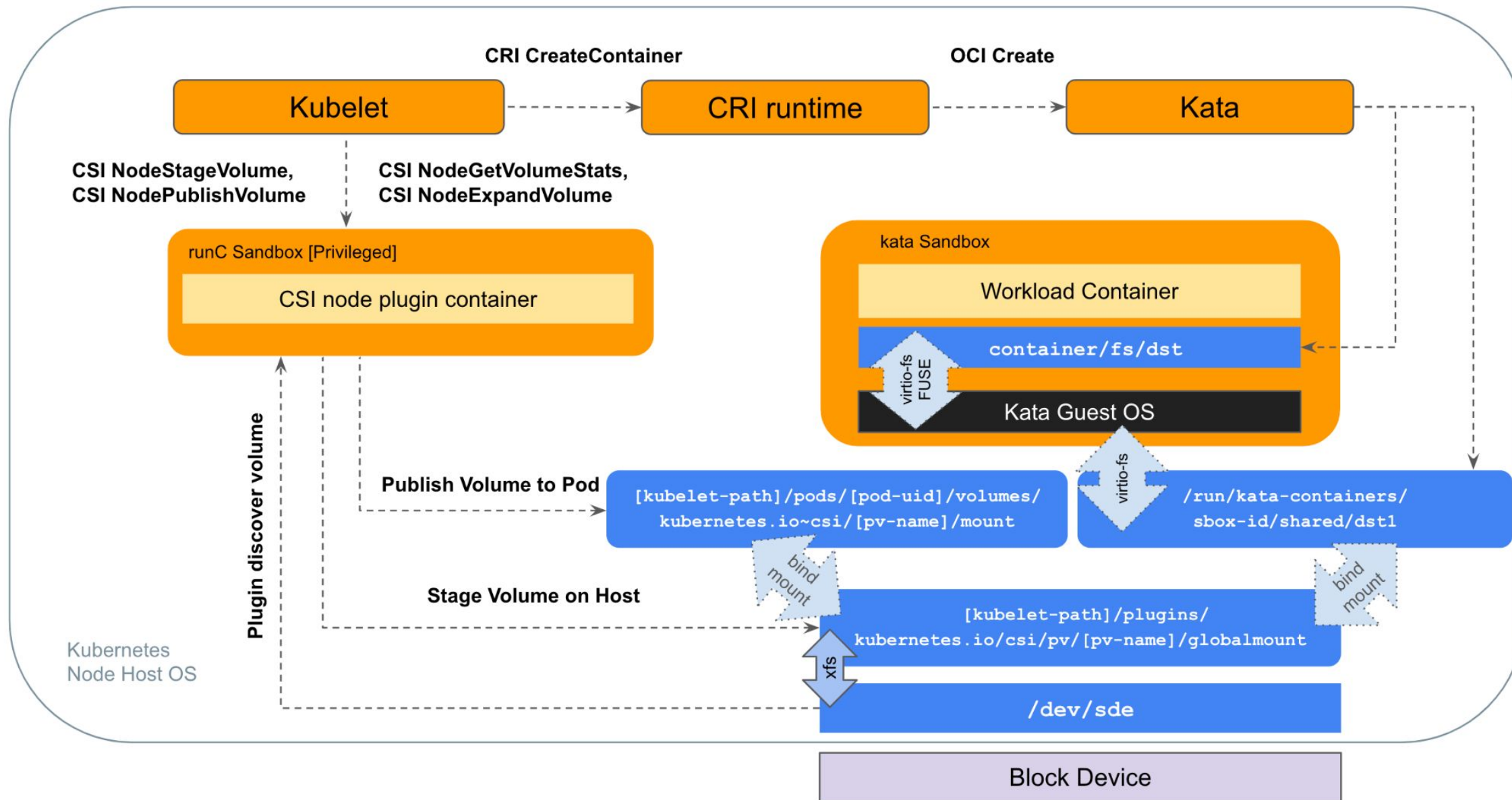
# Kata Containers architecture



Diagram source: https://katacontainers.io/learn/

# Regular Path to Mount a PV in Kata

# Performance and Security Trade-offs

- **Performance**

    virtio-fs has worse performance compared to virtio-blk/virtio-scsi
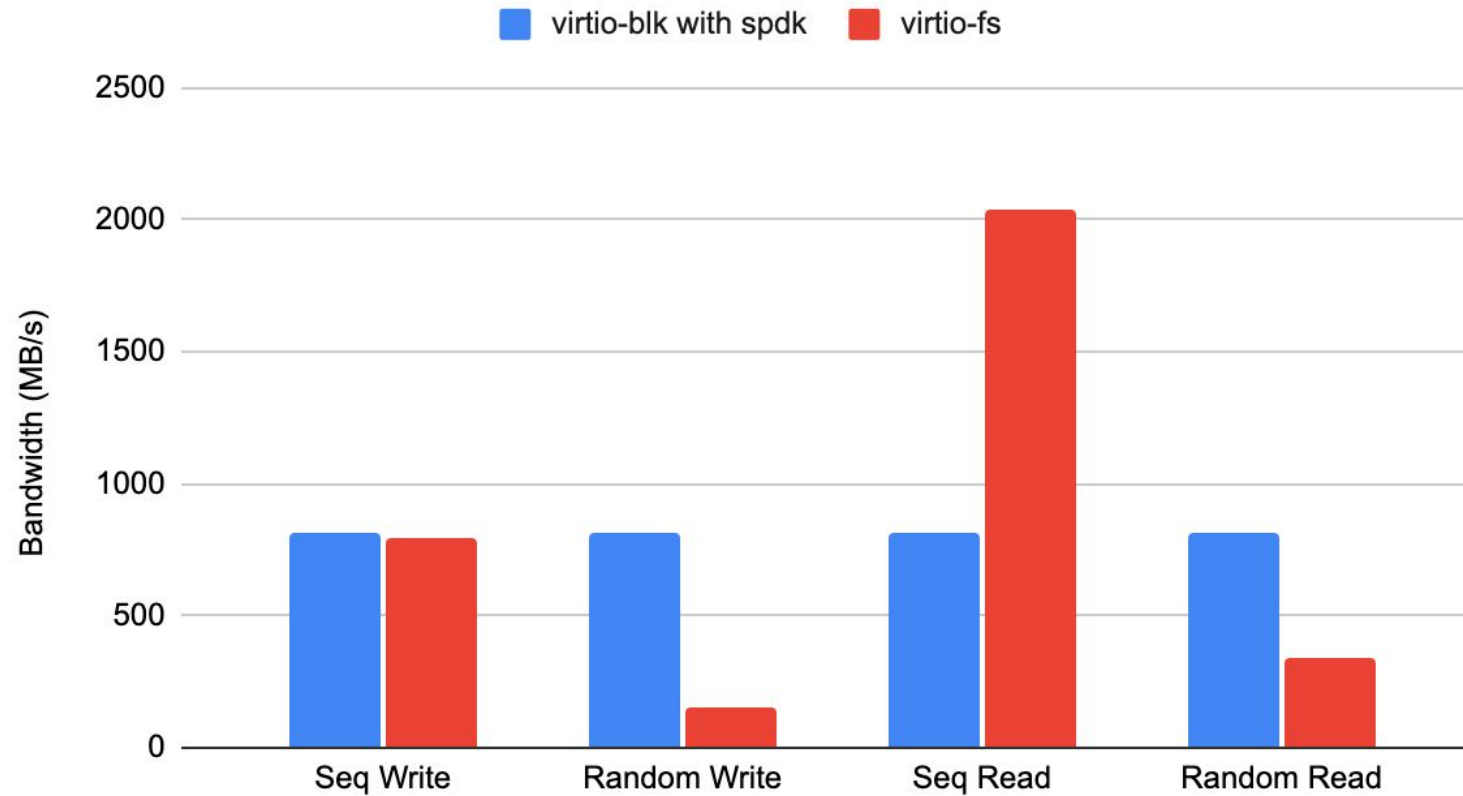
- **Isolation**

    virtio-fs requires file system to be mounted in the host

- **Other Gaps**

    virtio-fs may not surface native file system features

# Performance Comparison of virtio fs vs blk
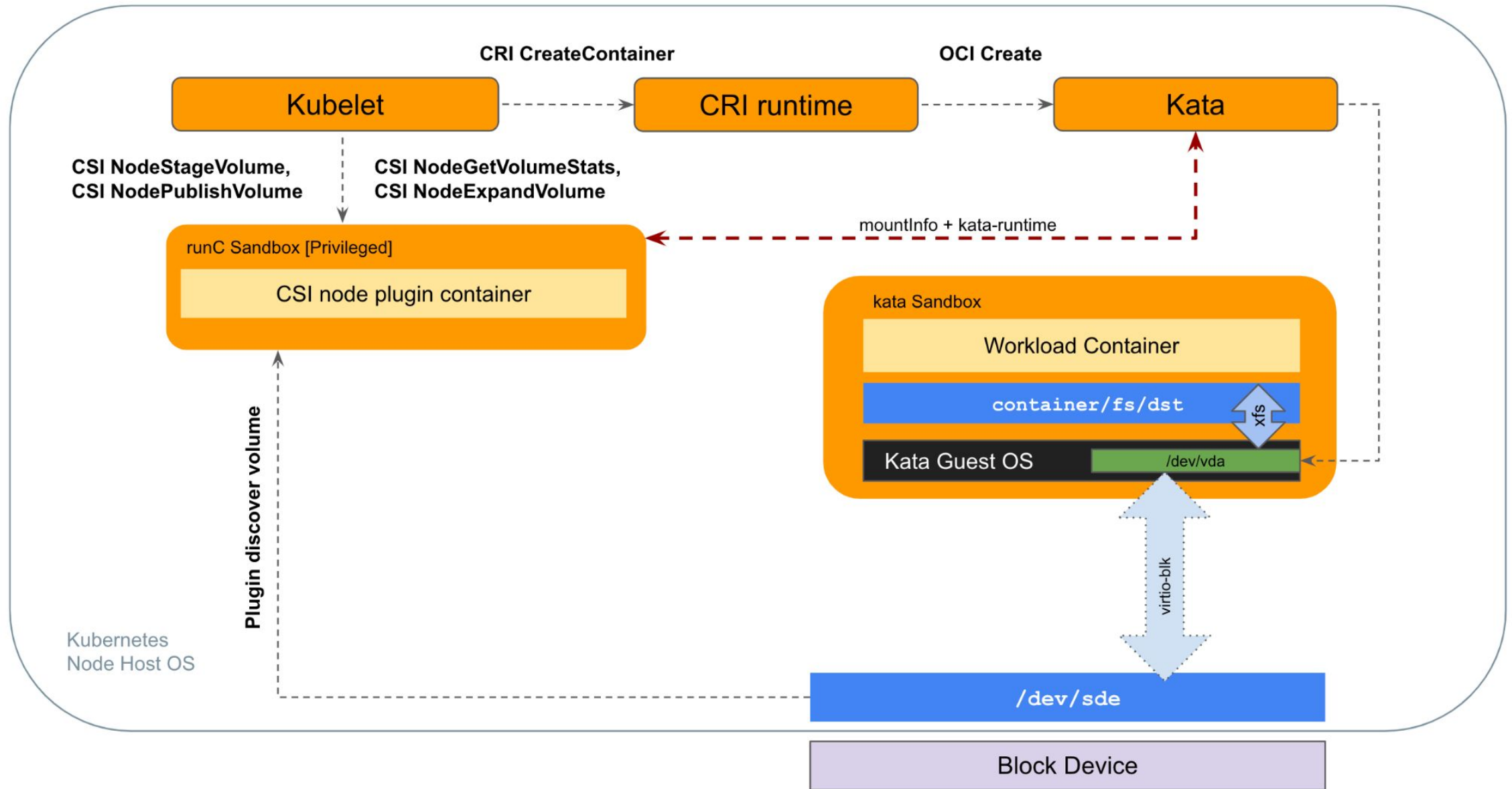


fio bandwidth test

test details: https://gist.github.com/fengwang666/a624c11e26fb7f9035af00fad6cea467
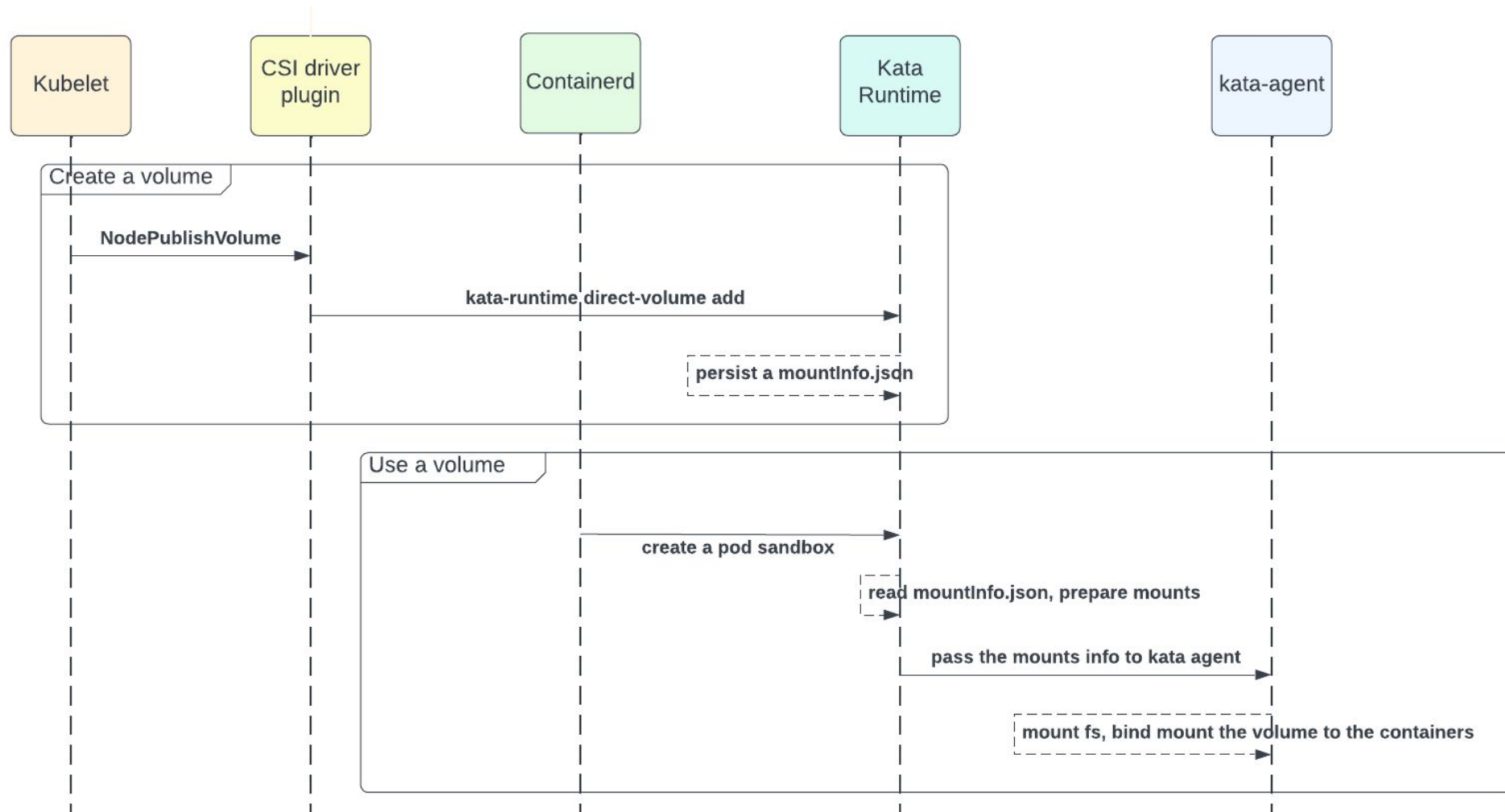
## Direct Assigned Storage

CSI plugin delegates PV mount and preparation to container runtime
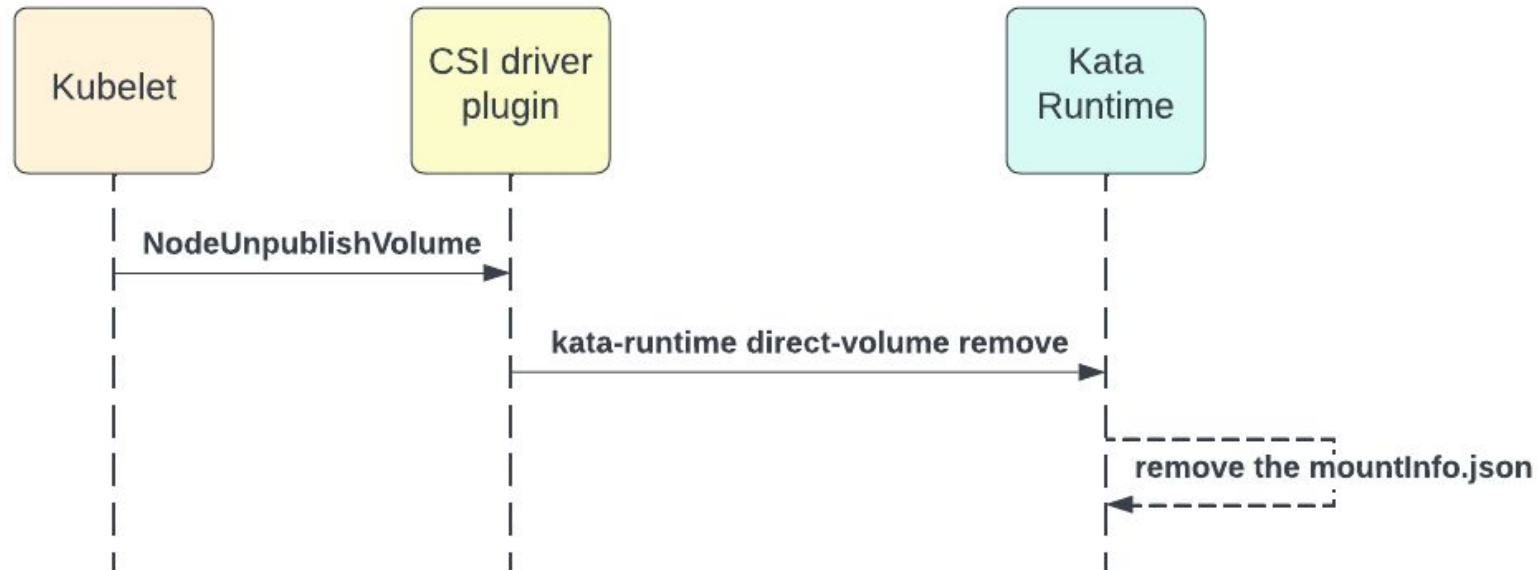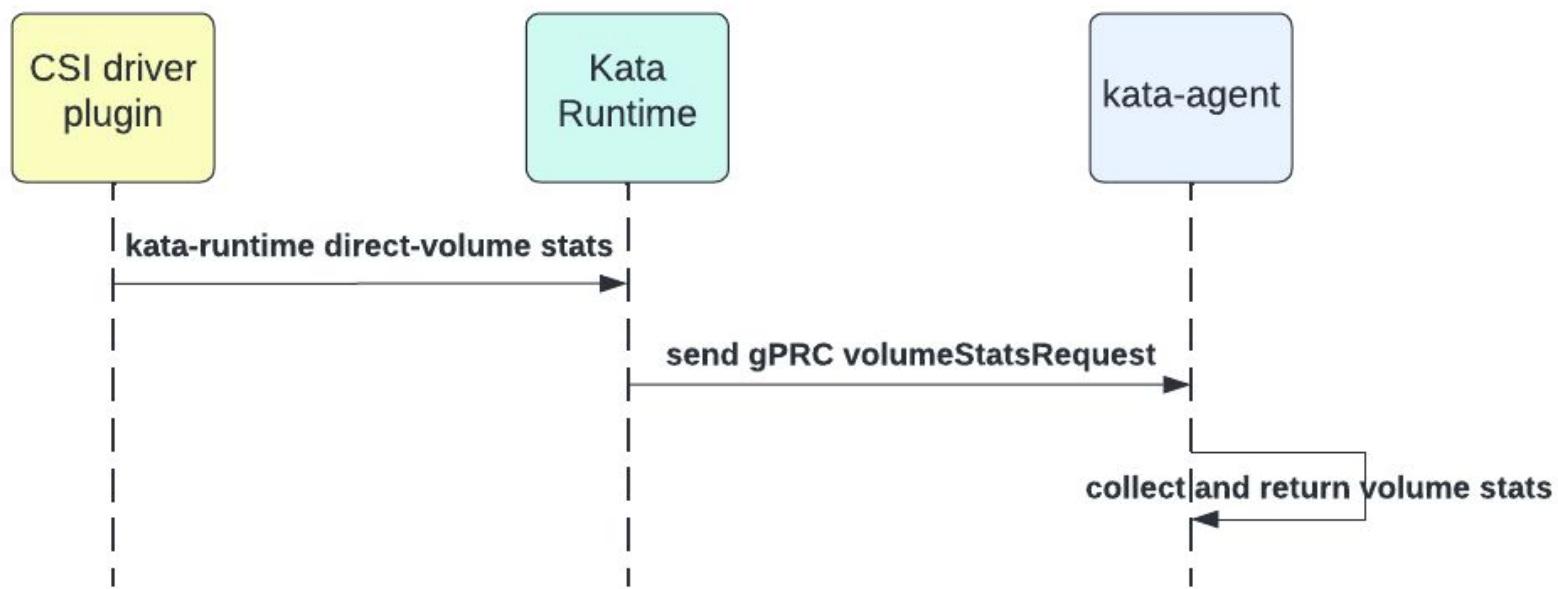
# Direct Assign PV to Container Runtime

# Direct Assigned PV: Mount Volume

# Direct Assigned PV: Unmount Volume

# Direct Assigned PV: Volume Stats

# Direct Assigned PV: Resize Volume

Container Runtime specific delegation logic in CSI Plugin

Demo

Current: Container Runtime specific delegation logic in CSI Plugin

- No changes in Kubelet, CRI, OCI and CSI specs

- Most post-mount configurations work

- CSI Plugin needs to lookup runtimeClass of pod

- Subpath support not possible

# Implementing Direct Assigned Storage

## Future: Container Runtime agnostic delegation logic in CSI Plugin

- Enhancements in Kubelet, RuntimeClass and CSI spec

- Kubelet will match capabilities of CSI plugin and runtime

- CSI plugins can use a common proxy to delegate operations

- All post-mount configurations will be supported



KEP-2857: Runtime Assisted Mounting of Persistent Volumes

Open  **ddebroy** wants to merge 1 commit into `kubernetes:master` from `ddebroy:runtime1`

Future: Container Runtime agnostic delegation logic in CSI Plugin

- Scoped to ReadWriteOncePod access modes for safety

- Should not be used for mount scenarios that require a secret

## KEP-2857: Runtime Assisted Mounting of Persistent Volumes

🔀 Open   **ddebroy** wants to merge 1 commit into `kubernetes:master` from `ddebroy:runtime1`

# Takeaways and Next Steps

- Explored alternatives to the standard mount flow for persistent storage
    - from a microVM perspective

- Ways to delegate mount and post mount configuration to container runtime
    - for block based PVs
    - avoid bind mounts and file-system projections

- Please get involved!
    - KEP-2857 (sig-storage) [#sig-storage in kubernetes.slack.com]
    - Kata community [katacontainers.slack.com]

# Questions

Please scan the QR Code above to
leave feedback on this session