

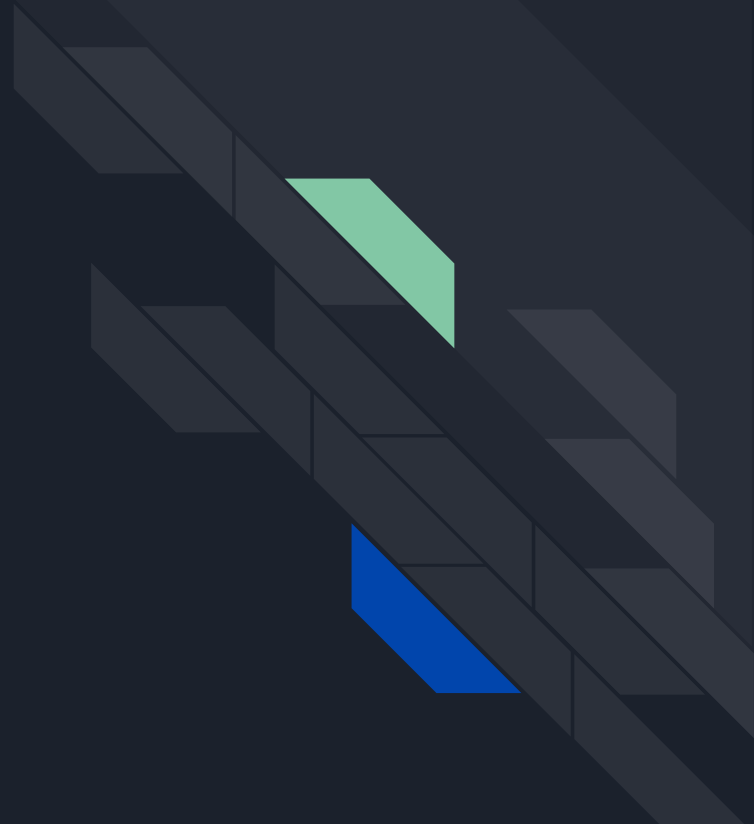
# **Kata Containers**

## Story of a container runtime

Sébastien Boeuf, Software Engineer  
Intel Corporation

# Agenda

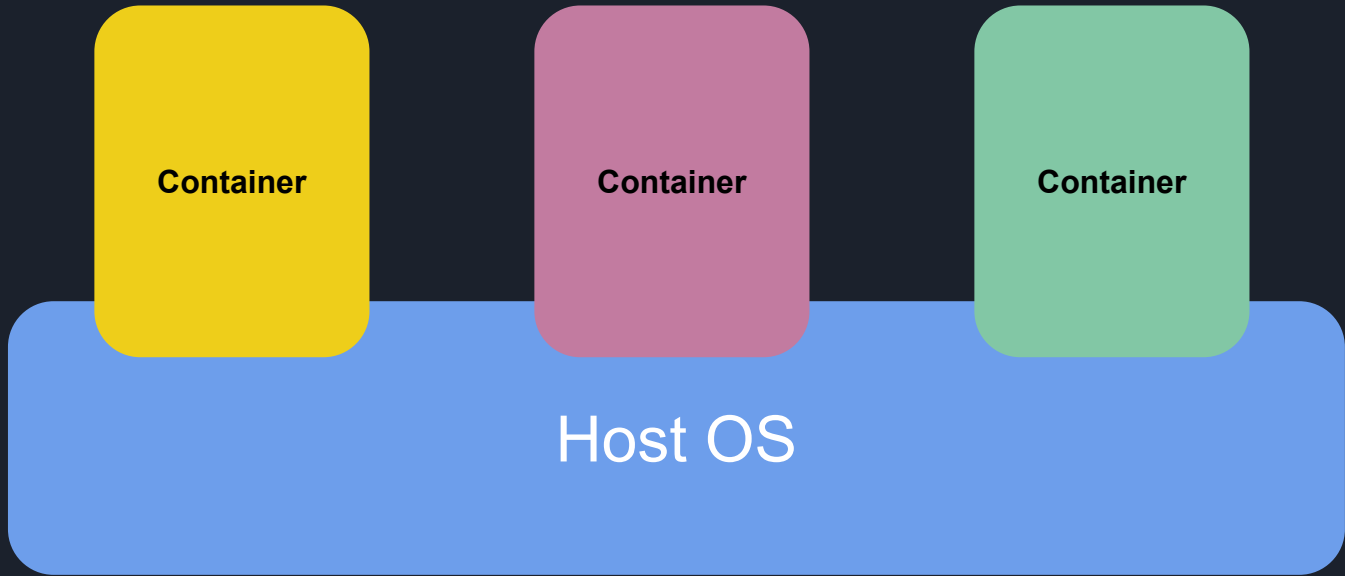
- Why Kata Containers?
- Acceptance
- Community growth
- Ecosystem influence
- Hypervisor flexible



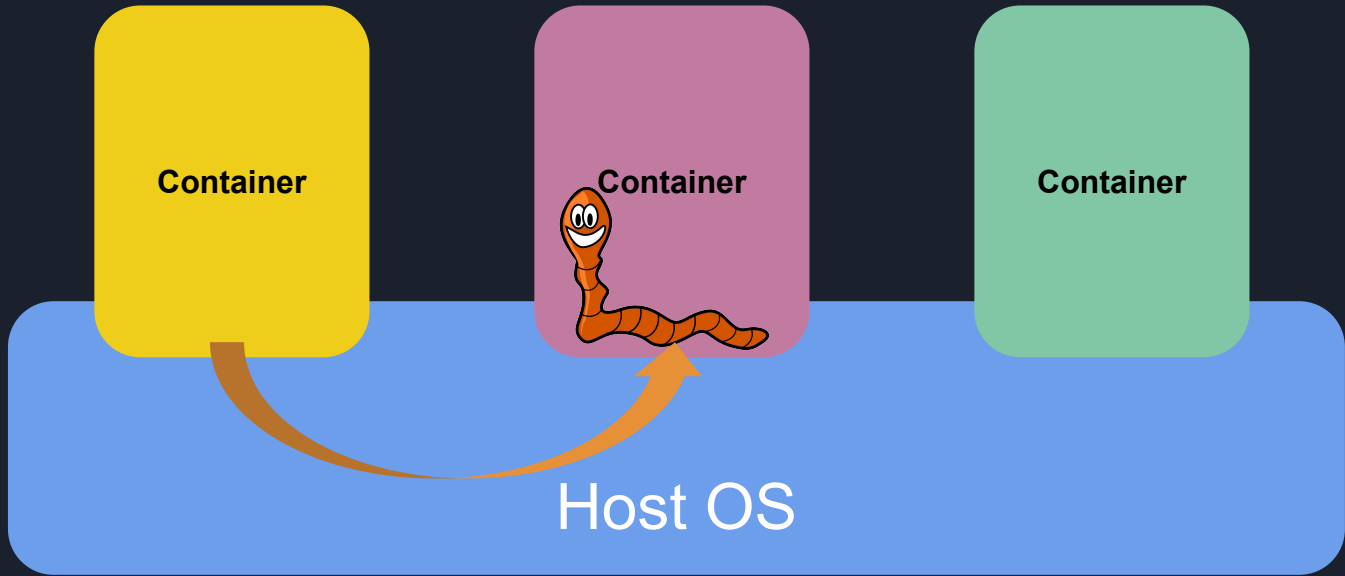




# Containers



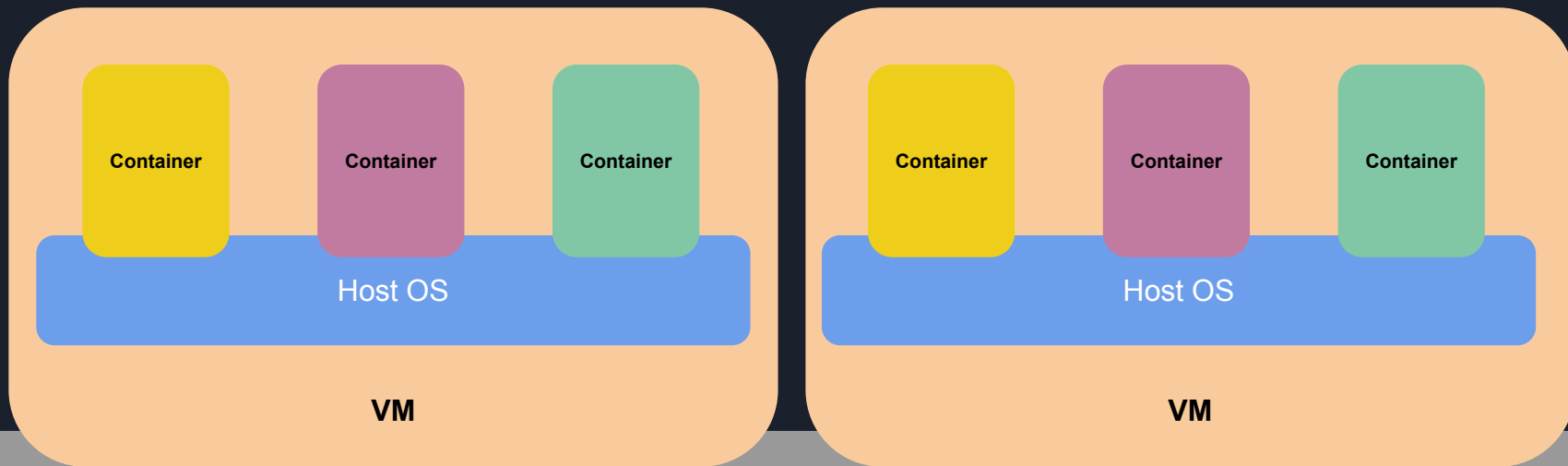
# Security threat







# Manual isolation



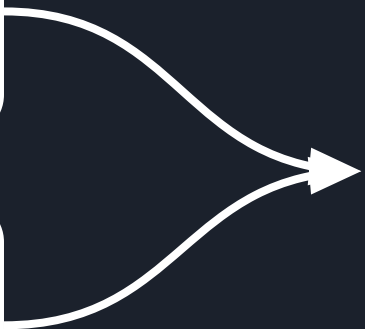
Baremetal server



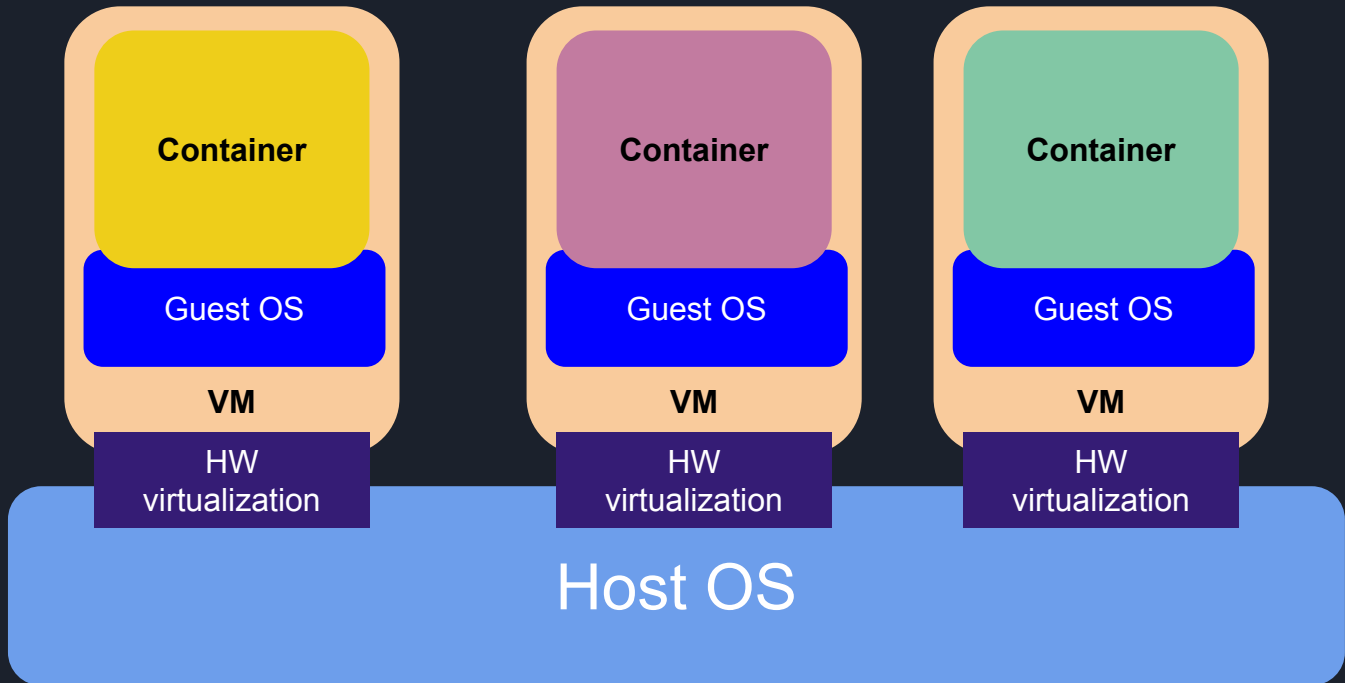


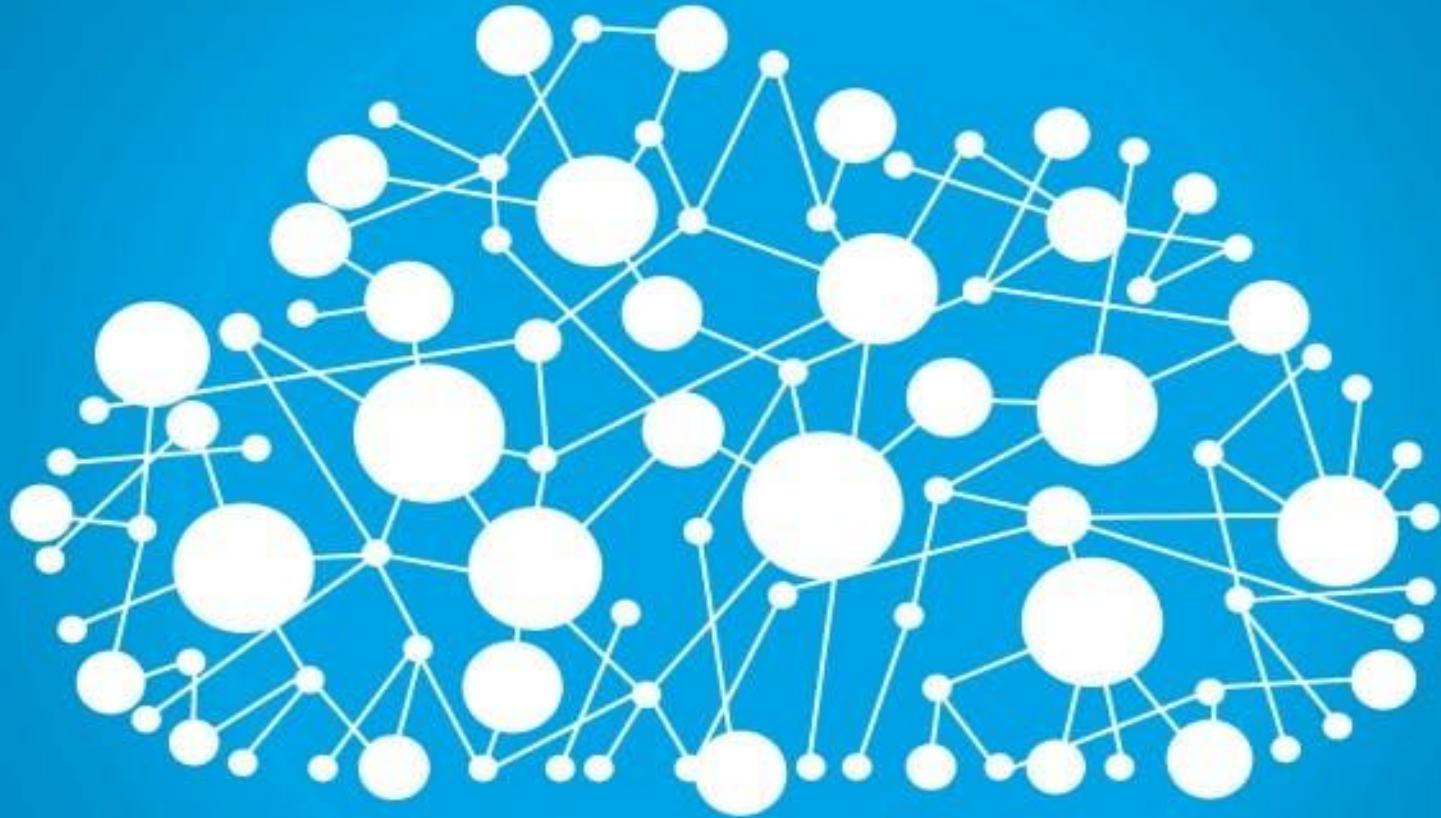


# Legacy

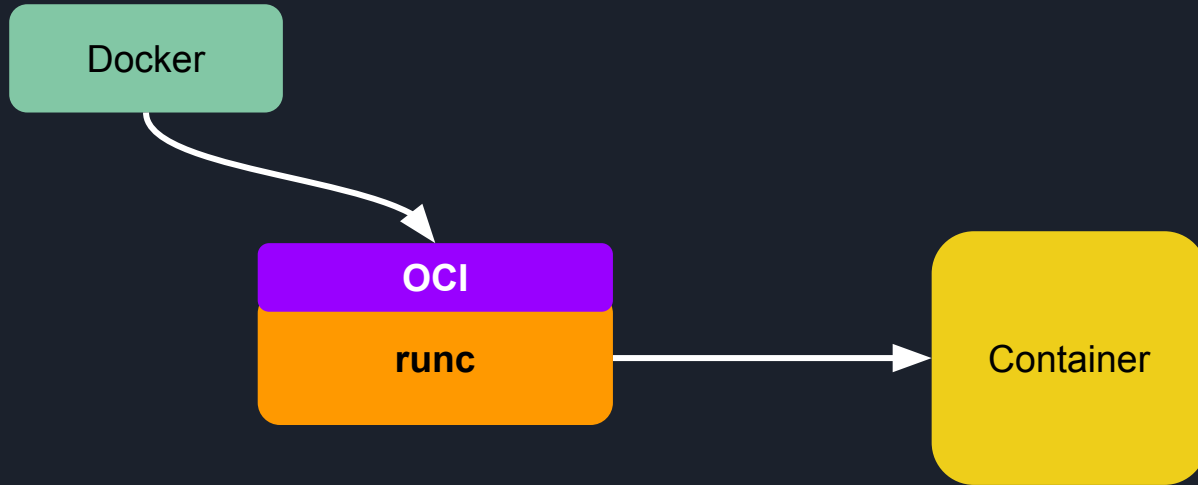


# Kata Containers

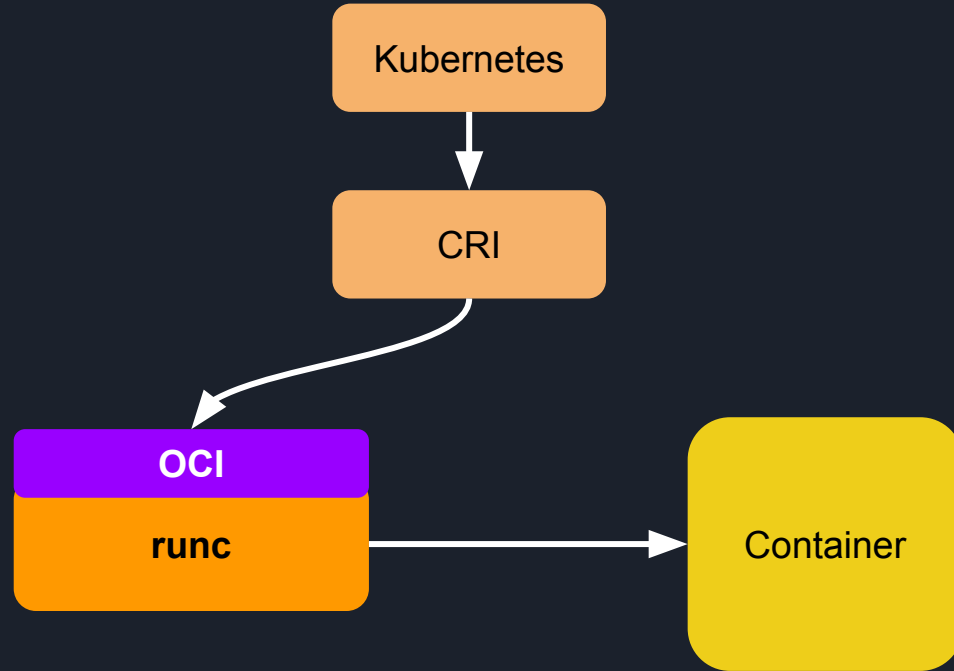




# Container ecosystem

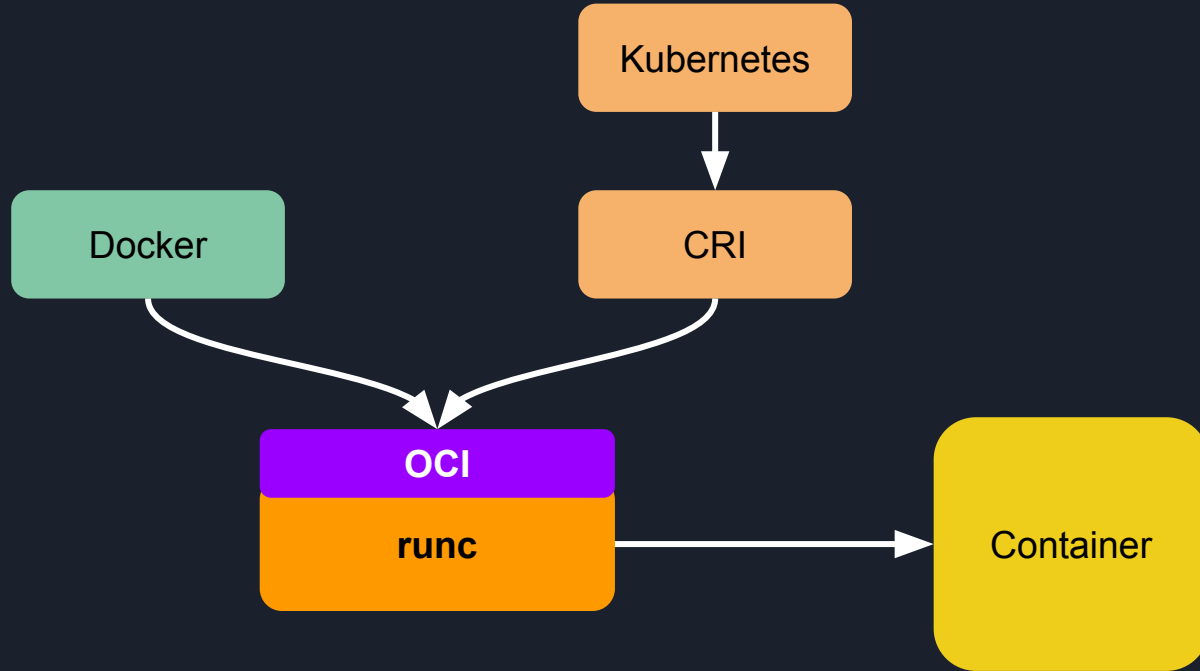


# Container ecosystem

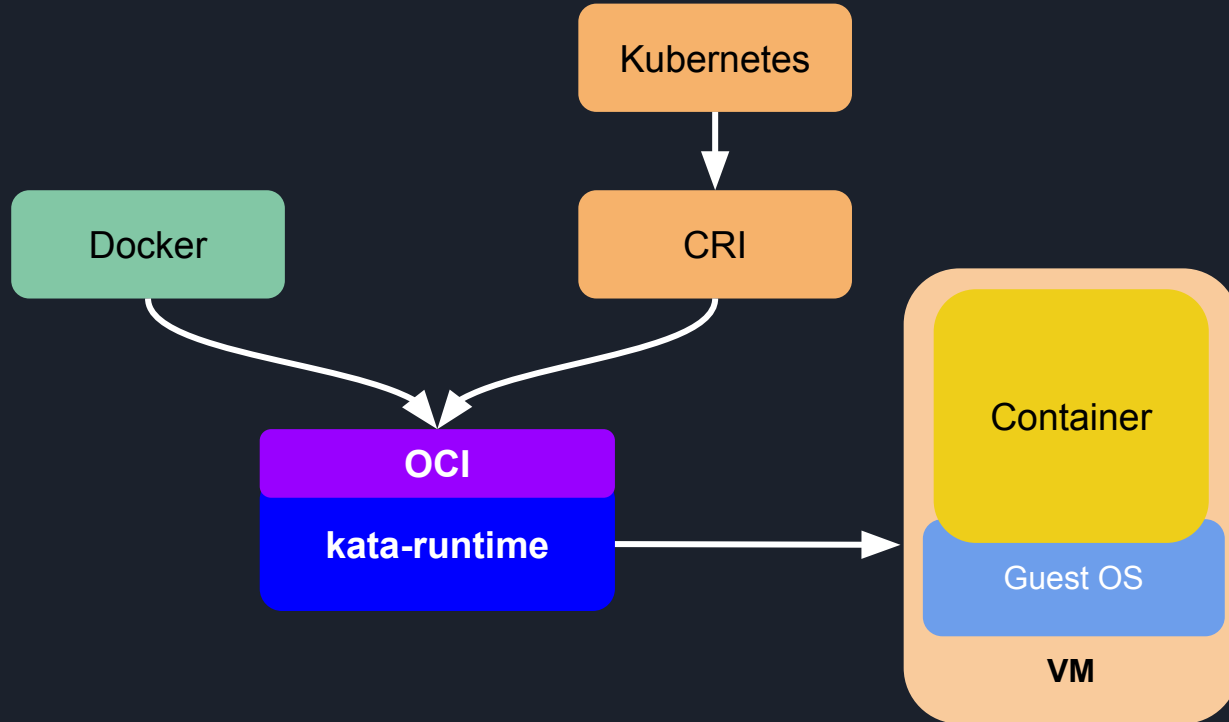


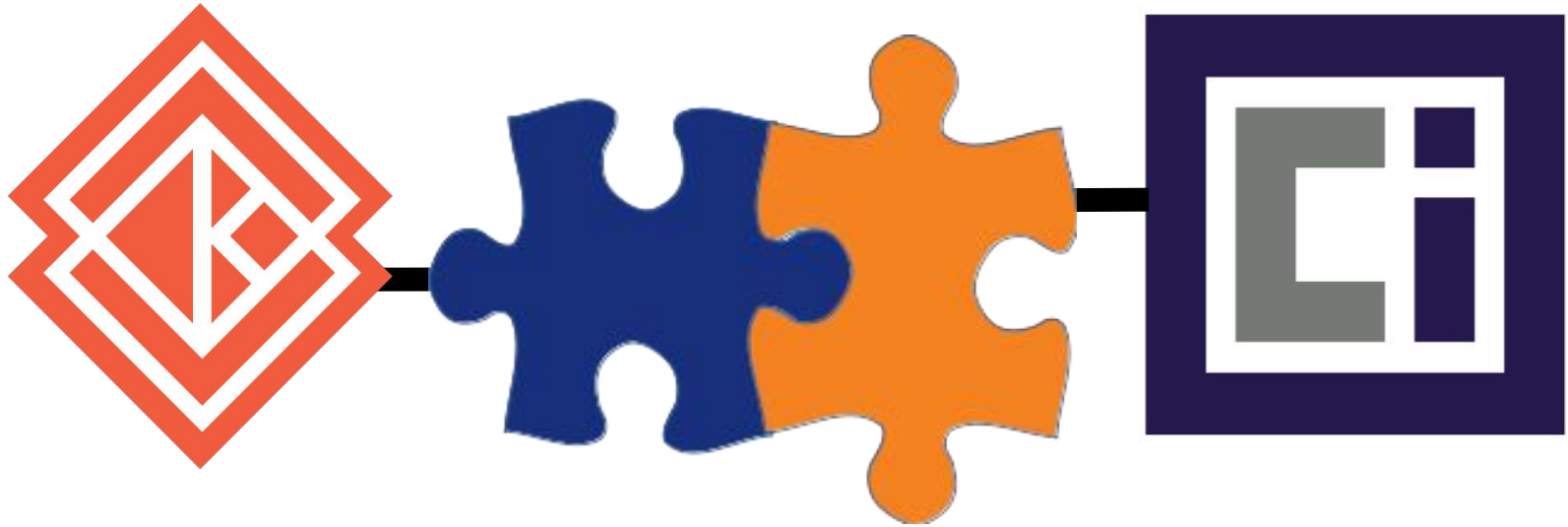


# Container ecosystem



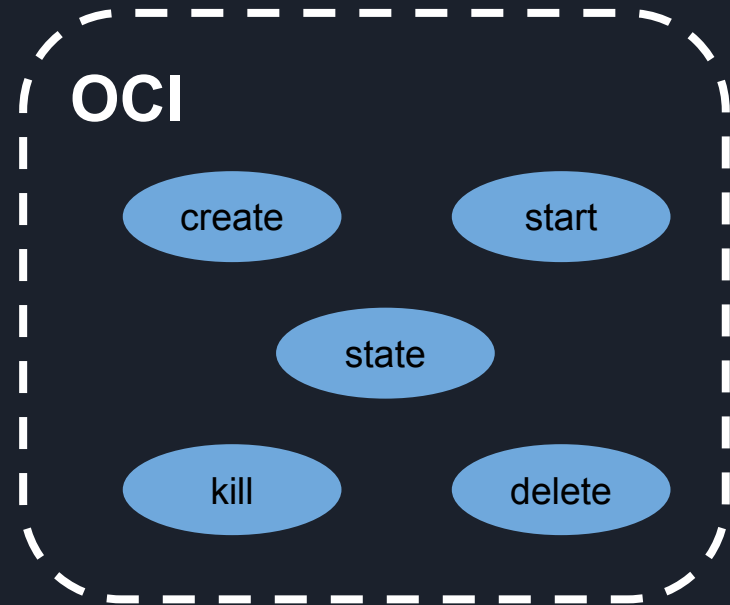
# Seamless integration







# OCI compatible





# OCI compatible

**runc**

exec

list

pause

resume

run

update

**OCI**

create

start

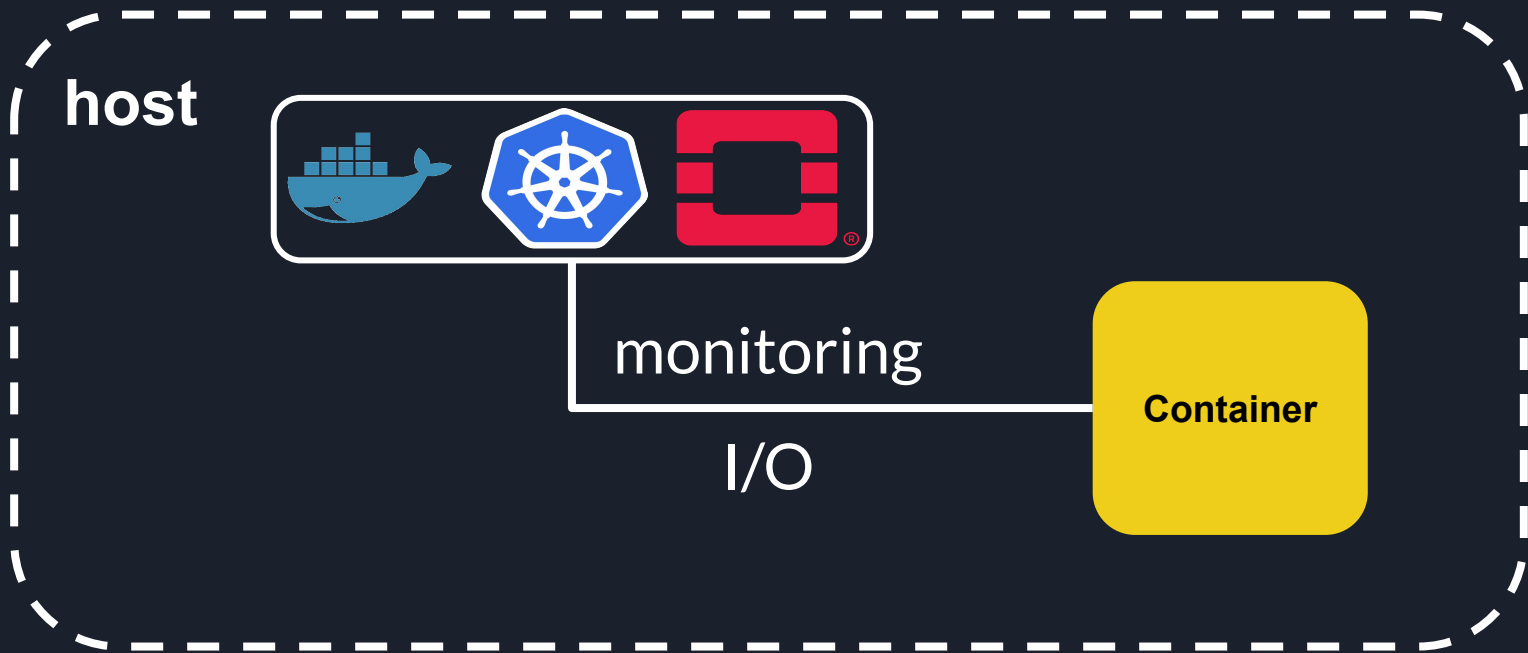
state

kill

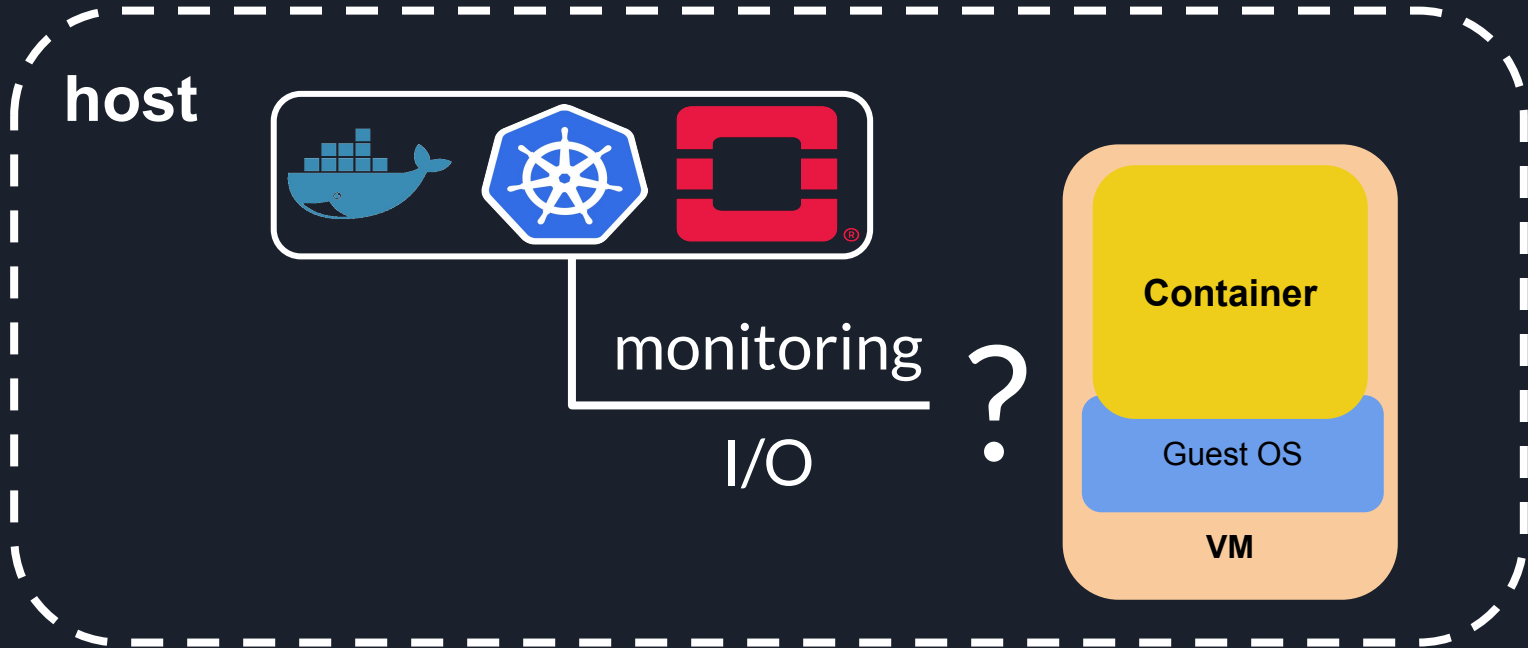
delete



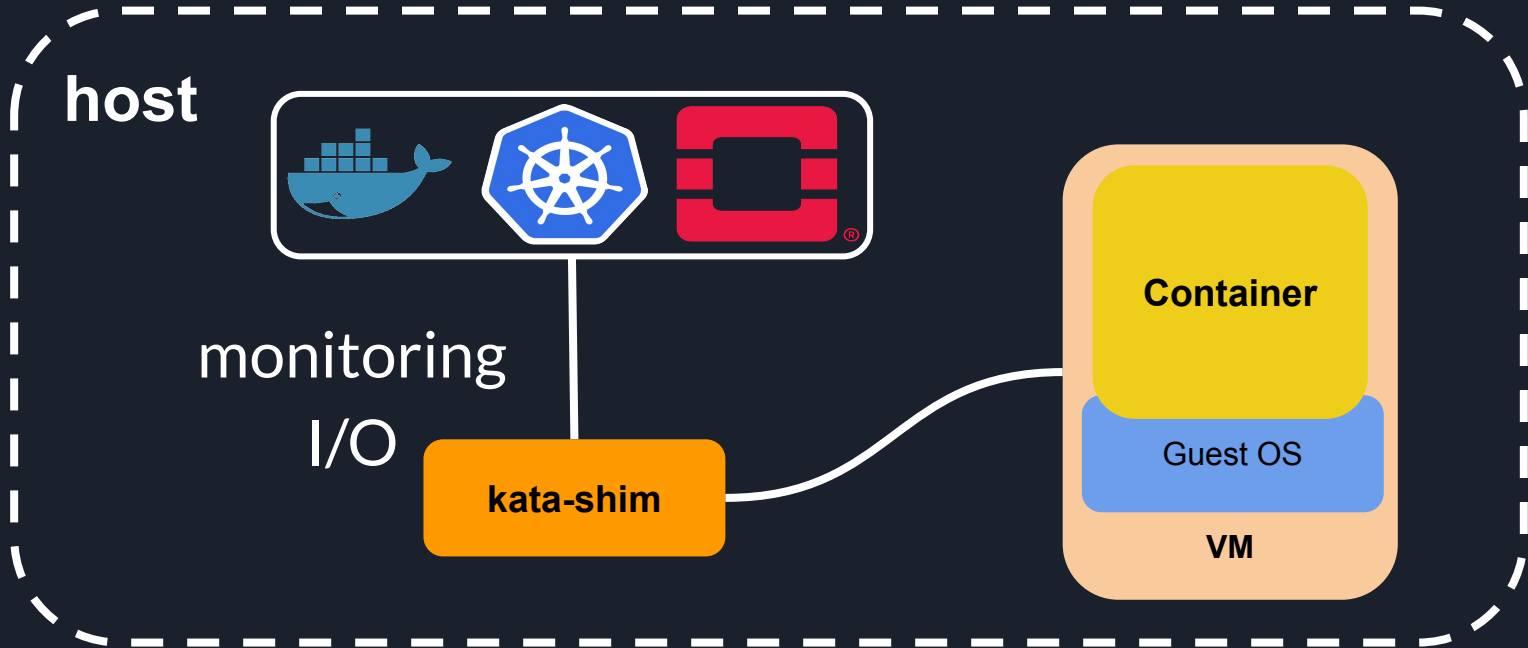
# OCI compatible



# OCI compatible



# OCI compatible







# Community growth

## Additional architectures

- aarch64 (**ARM**)
- ppc64 and s390 (**IBM**)

## Enhanced stability and production ready

- **Huawei**
- **Baidu**
- **Alibaba**



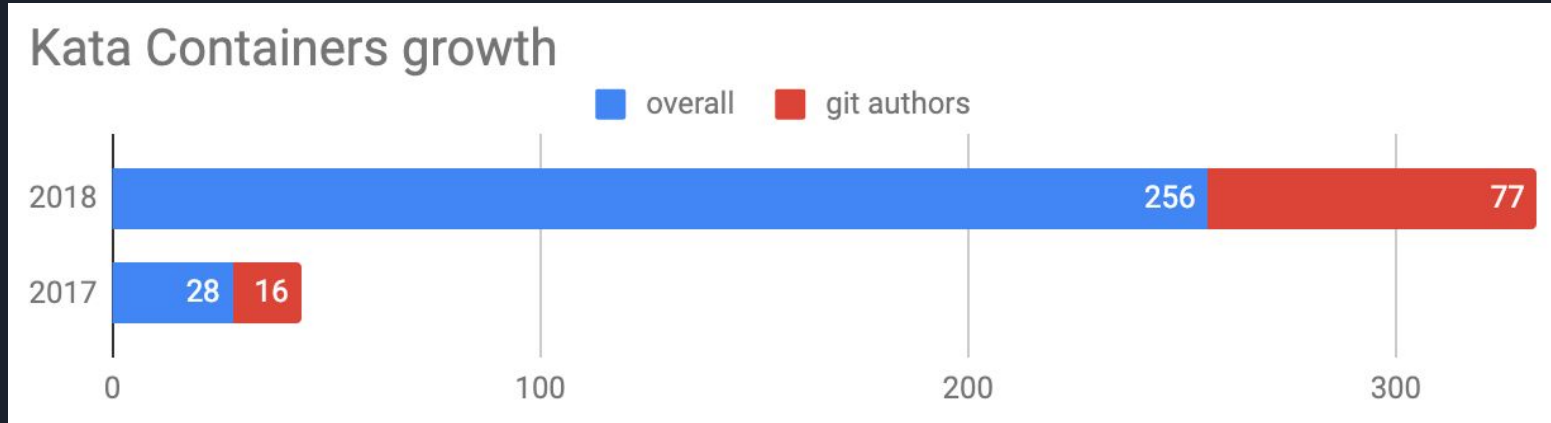


# Community growth

## CI resources

- Vexxhost (**Vexxhost**)
- Azure (**Microsoft**)
- AWS (**Amazon**)
- GCE (**Google**)

# Community growth



2000 pull requests / 100 contributors





# Extend OCI

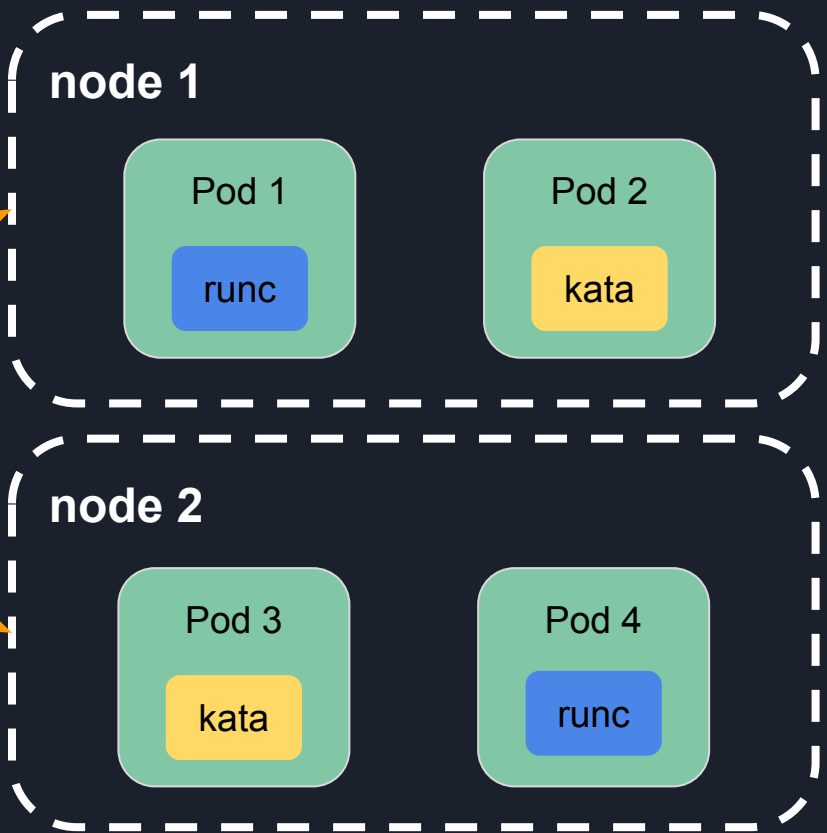
```
// VM contains information for virtual-machine-based containers.
type VM struct {
    // Hypervisor specifies hypervisor-related configuration for virtual-machine-based containers.
    Hypervisor VMHypervisor `json:"hypervisor,omitempty"`
    // Kernel specifies kernel-related configuration for virtual-machine-based containers.
    Kernel VMKernel `json:"kernel"`
    // Image specifies guest image related configuration for virtual-machine-based containers.
    Image VMImage `json:"image,omitempty"`
}
```



# RuntimeClass

```
message RunPodSandboxRequest {  
    // Configuration for creating a PodSandbox.  
    PodSandboxConfig config = 1;  
    // Named runtime configuration to use for this PodSandbox.  
    // If the runtime handler is unknown, this request should be rejected. An  
    // empty string should select the default handler, equivalent to the  
    // behavior before this feature was added.  
    // See https://git.k8s.io/enhancements/keps/sig-node/runtime-class.md  
    string runtime_handler = 2;  
}
```

# RuntimeClass





# Pod overhead

Introduce a `Pod.Spec.Resources` field on the pod to specify the pods overhead.

```
Pod {  
  Spec PodSpec {  
    // overhead is the resource overhead consumed by the Pod, not including  
    // container resource usage. Users should leave this field unset.  
    // +optional  
    Overhead *ResourceRequirements  
  }  
}
```

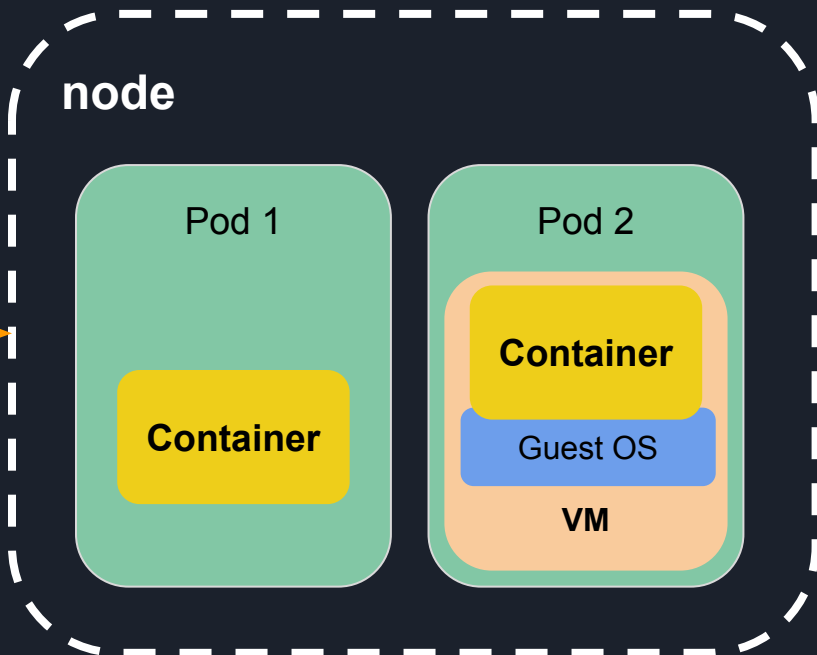
For scheduling, the pod resource requests are added to the container resource requests.



# Pod overhead

**pod1.yaml**  
cpus: 2  
mem: 256M

**pod2.yaml**  
cpus: 2  
mem: 256M  
  
Overhead:  
- cpus: 1  
- mem: 128M





# Shim v2

CRI



**containerd**  
or  
**CRI-O**

# Shim v2

CRI

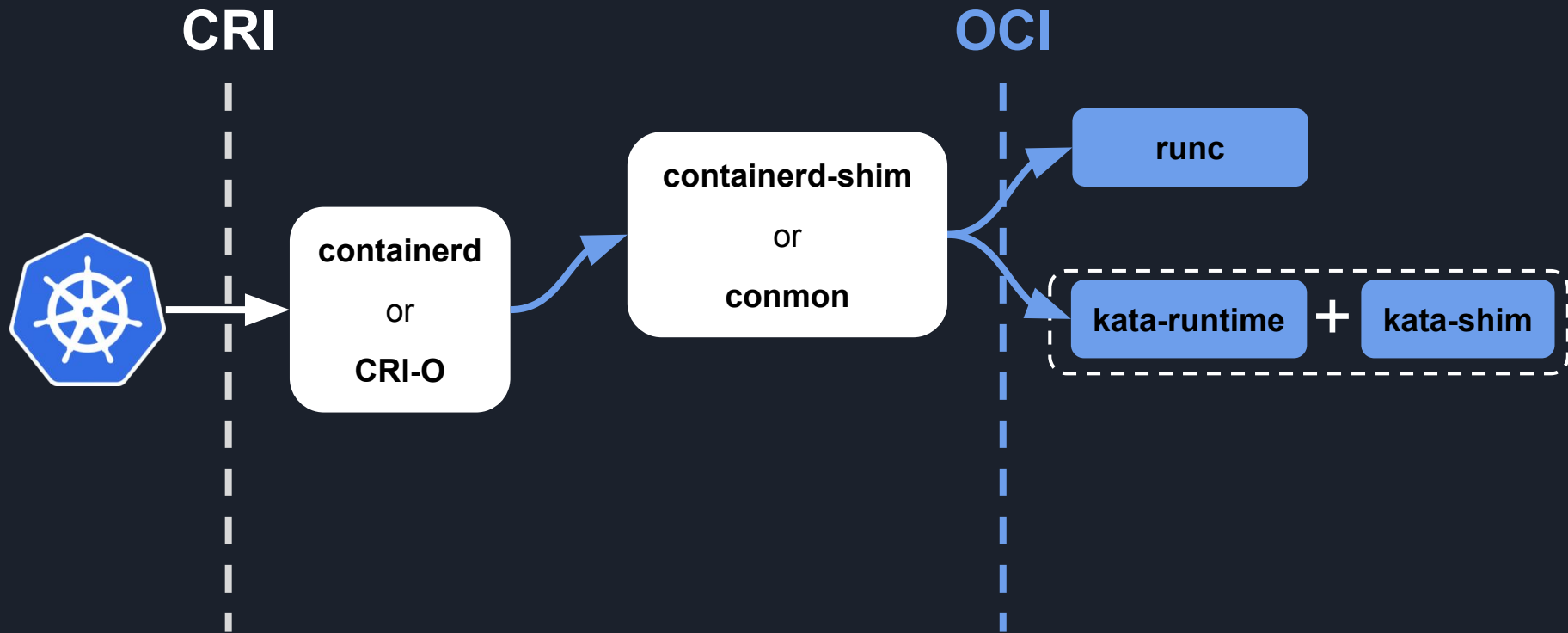


**containerd**  
or  
**CRI-O**

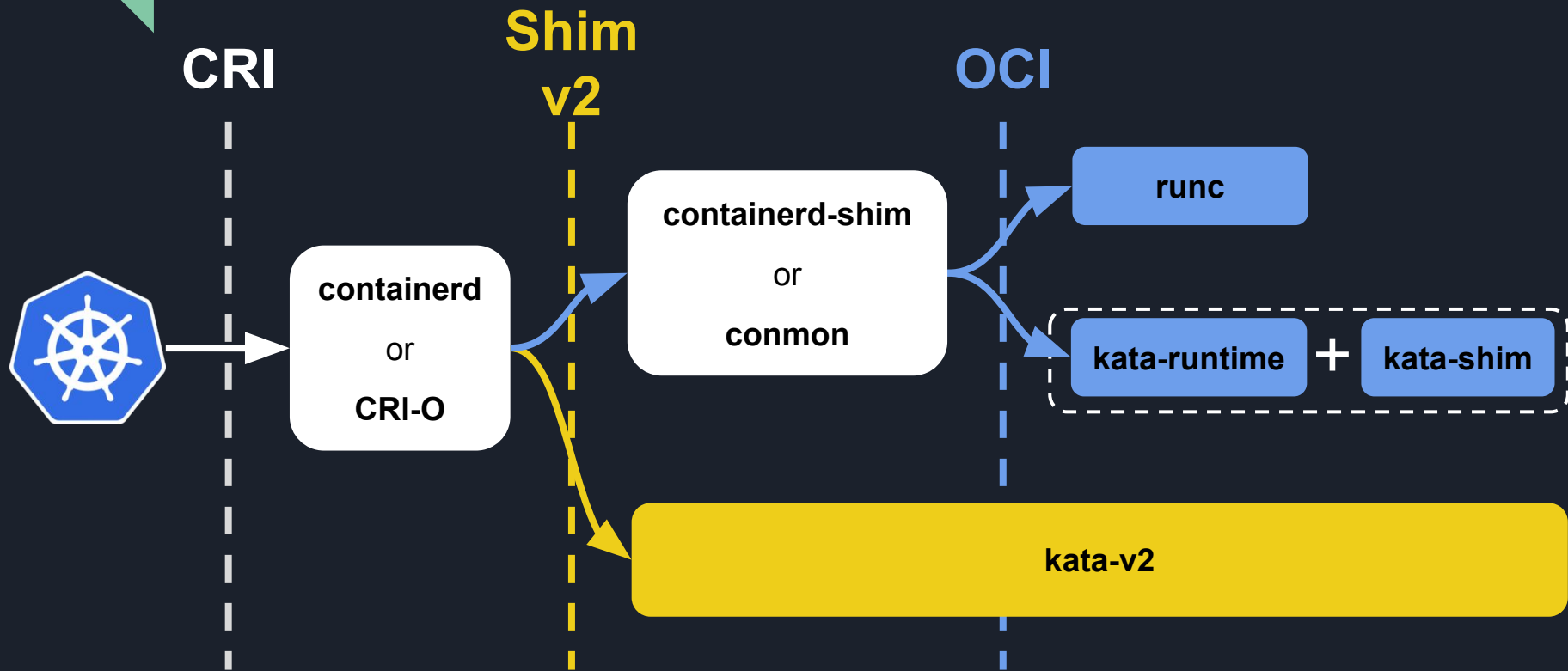


**containerd-shim**  
or  
**conmon**

# Shim v2



# Shim v2



# Shim v2

```
service Task {  
  rpc State(StateRequest) returns (StateResponse);  
  rpc Create(CreateTaskRequest) returns (CreateTaskResponse);  
  rpc Start(StartRequest) returns (StartResponse);  
  rpc Delete(DeleteRequest) returns (DeleteResponse);  
  rpc Pids(PidsRequest) returns (PidsResponse);  
  rpc Pause(PauseRequest) returns (google.protobuf.Empty);  
  rpc Resume(ResumeRequest) returns (google.protobuf.Empty);  
  rpc Checkpoint(CheckpointTaskRequest) returns (google.protobuf.Empty);  
  rpc Kill(KillRequest) returns (google.protobuf.Empty);  
  rpc Exec(ExecProcessRequest) returns (google.protobuf.Empty);  
  rpc ResizePty(ResizePtyRequest) returns (google.protobuf.Empty);  
  rpc CloseIO(CloseIORequest) returns (google.protobuf.Empty);  
  rpc Update(UpdateTaskRequest) returns (google.protobuf.Empty);  
  rpc Wait(WaitRequest) returns (WaitResponse);  
  rpc Stats(StatsRequest) returns (StatsResponse);  
  rpc Connect(ConnectRequest) returns (ConnectResponse);  
  rpc Shutdown(ShutdownRequest) returns (google.protobuf.Empty);  
}
```

wait

stats

resizePty

**No host PID  
assumption!**

**k8s pod  
scaling!**



# Shared filesystem

## Virtio-9p

- Not fully POSIX compliant  $\Rightarrow$  Workload functional issues
- Not performant
- Production should use virtio-blk  $\Rightarrow$  devicemapper



# Shared filesystem

Redhat developed replacement for virtio-9p ⇒ **virtio-fs**

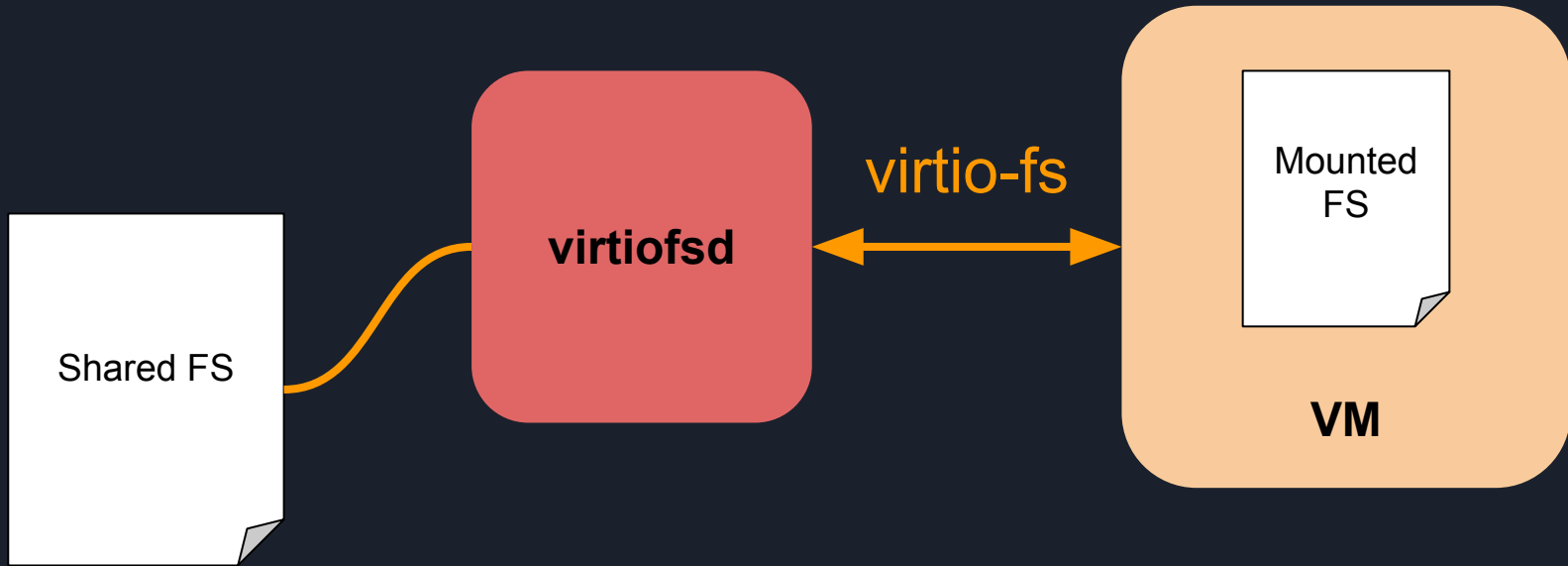
- Fully POSIX compliant ⇒ Solve workload functional issues
- As performant as virtio-blk (with DAX optimization)
- Overlay back into the picture for production

# Shared filesystem

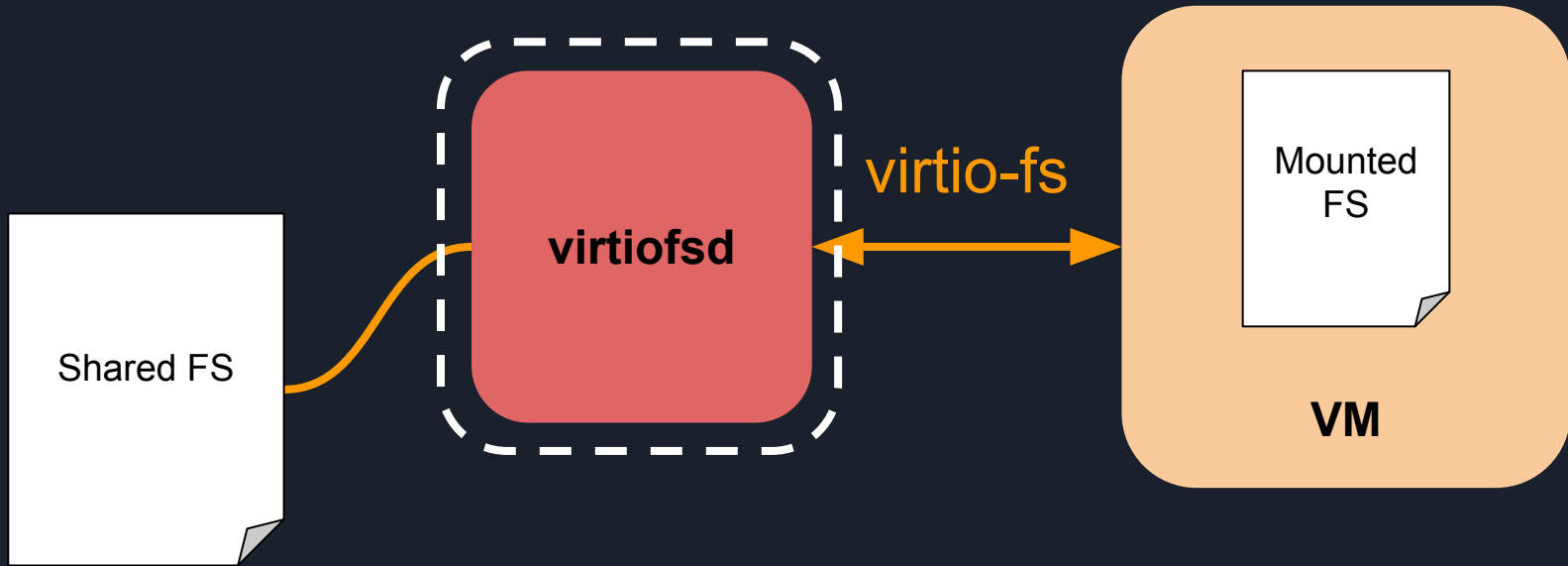




# Shared filesystem



# Shared filesystem







# QEMU/NEMU

- Swiss army knife hypervisor  $\Rightarrow$  Default for Kata
  - Type 2 (KVM)
  - Multi-purpose
  - Extensive device model (virtio-gpu, virtio-crypto, ...)
  - Direct Device Assignment (VFIO)
- Wide codebase in C  $\Rightarrow$  Potential attack surface
- NEMU reduces the attack surface



# Firecracker

- Lightweight hypervisor
  - Type 2 (KVM)
  - Narrow focus: container workloads and FaaS
  - Reduced device model
- Small codebase in Rust  $\Rightarrow$  Highly secure



# ACRN (in progress)

- Lightweight hypervisor
  - Type 1
  - Focus on Automotive and IoT
  - Industry standard **FuSa** (Functional Safety)
- Small codebase in C  $\Rightarrow$  Highly secure



← PICK YOUR OWN



# Takeaways







# Join the fun!

Sources: <https://github.com/kata-containers/runtime>

Get started:

<https://github.com/kata-containers/documentation/blob/master/Developer-Guide.md>

Slack: [katacontainers.slack.com](https://katacontainers.slack.com)

IRC: [#kata-dev@freenode](https://freenode.net/rooms/#kata-dev)

Mailing list: [kata-dev@lists.katacontainers.io](mailto:kata-dev@lists.katacontainers.io)



Thank you

