# How Container Runtimes matter in Kubernetes?

Kunal Kushwaha
NTT OSS Center

# About me

- **Works @ NTT Open Source Software Center**
- **Contributes to containerd and other related projects.**
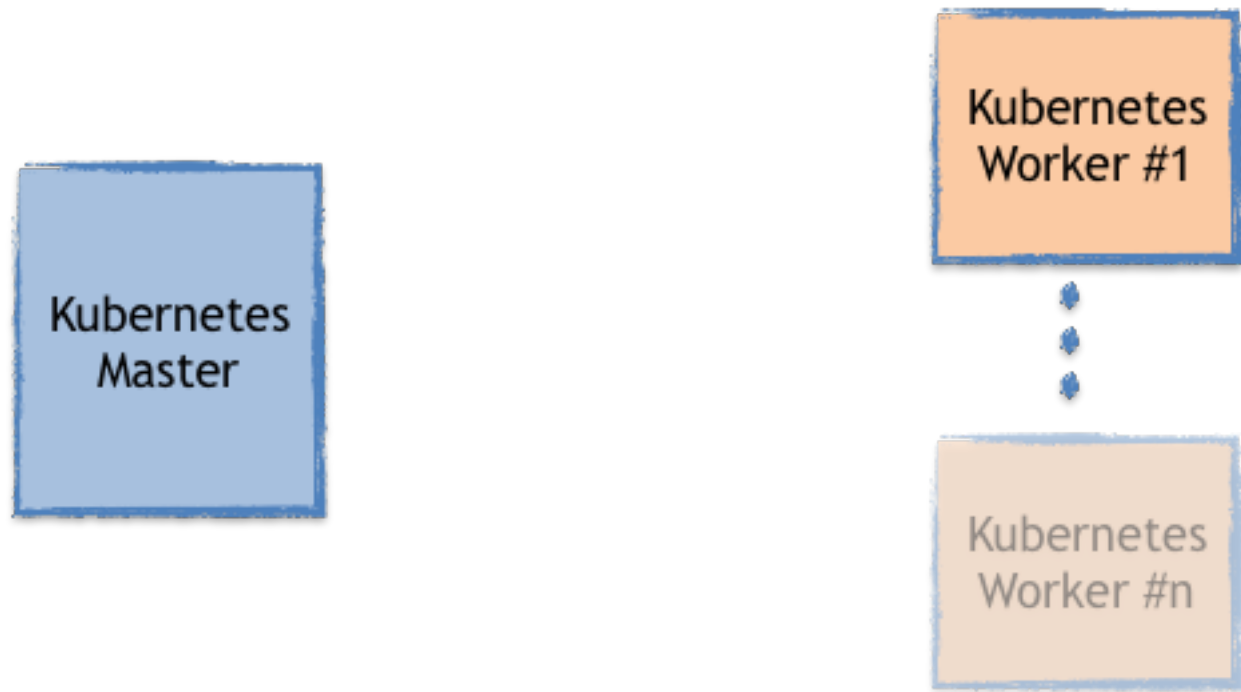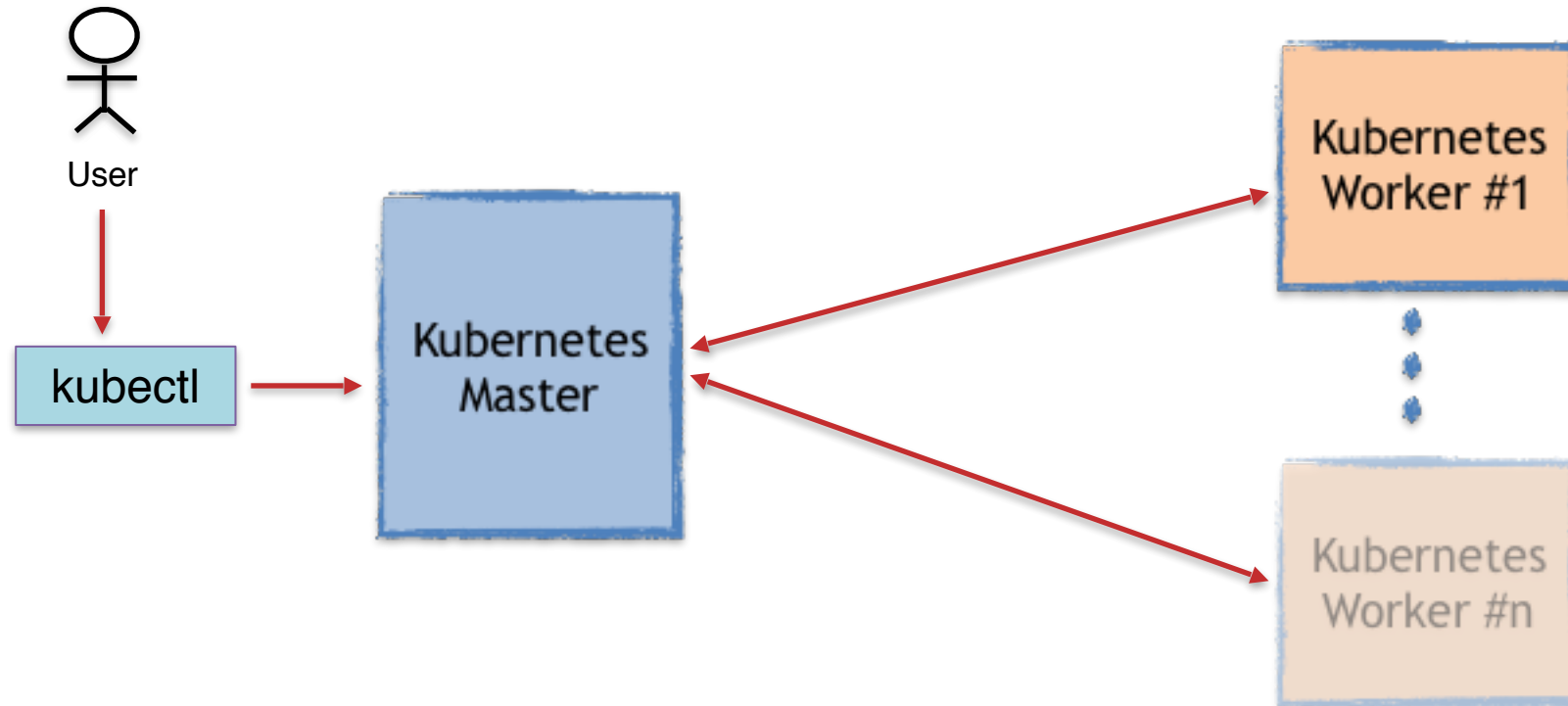- **Docker community leader, Tokyo**

@kunalkushwaha

# Agenda

- **Kubernetes Architecture.**
- **What is CRI (Container Runtime Interface)**
- **What is OCI (Open Container Initiative)**
- **CRI & OCI Implementations**
- **Why runtimes affect Kubernetes.**
- **Runtime Benchmarking results**
- **Analyzing for various workloads**
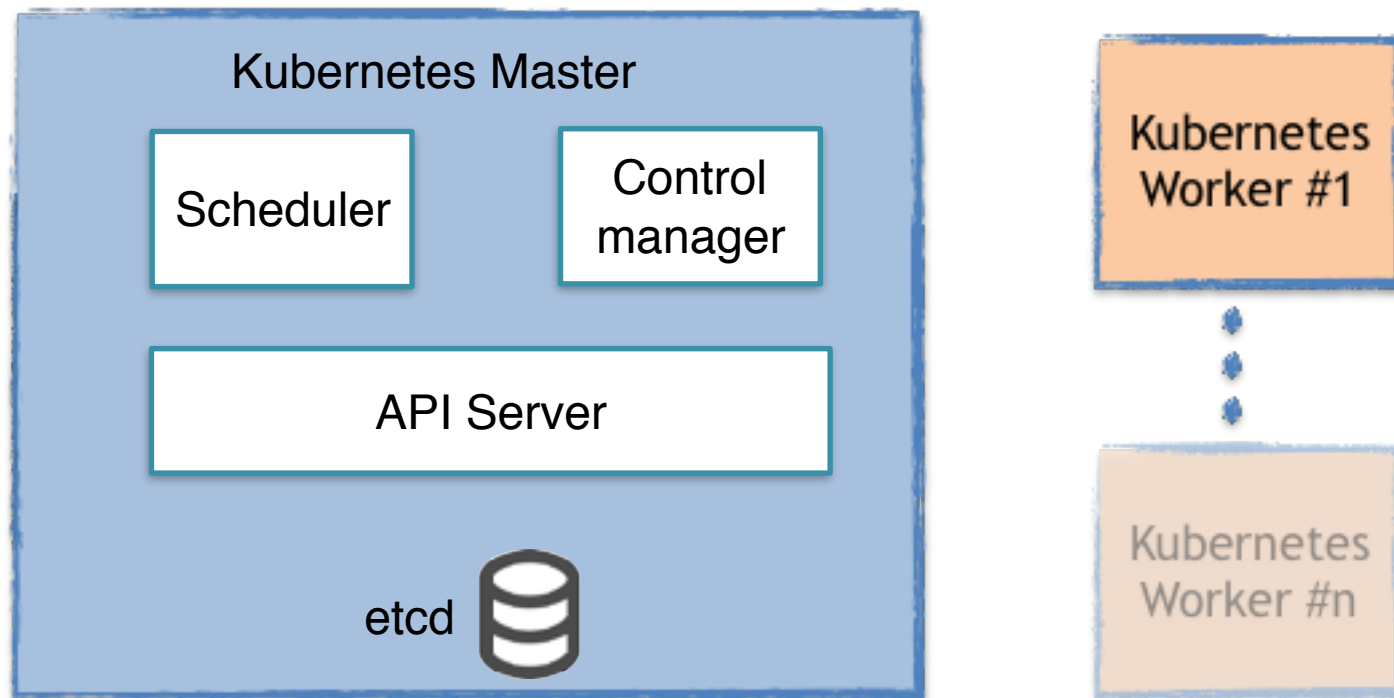- **Summary**

# Kubernetes Architecture

Kubernetes Master

Kubernetes Worker #1

Kubernetes Worker #n

A typical Kubernetes cluster
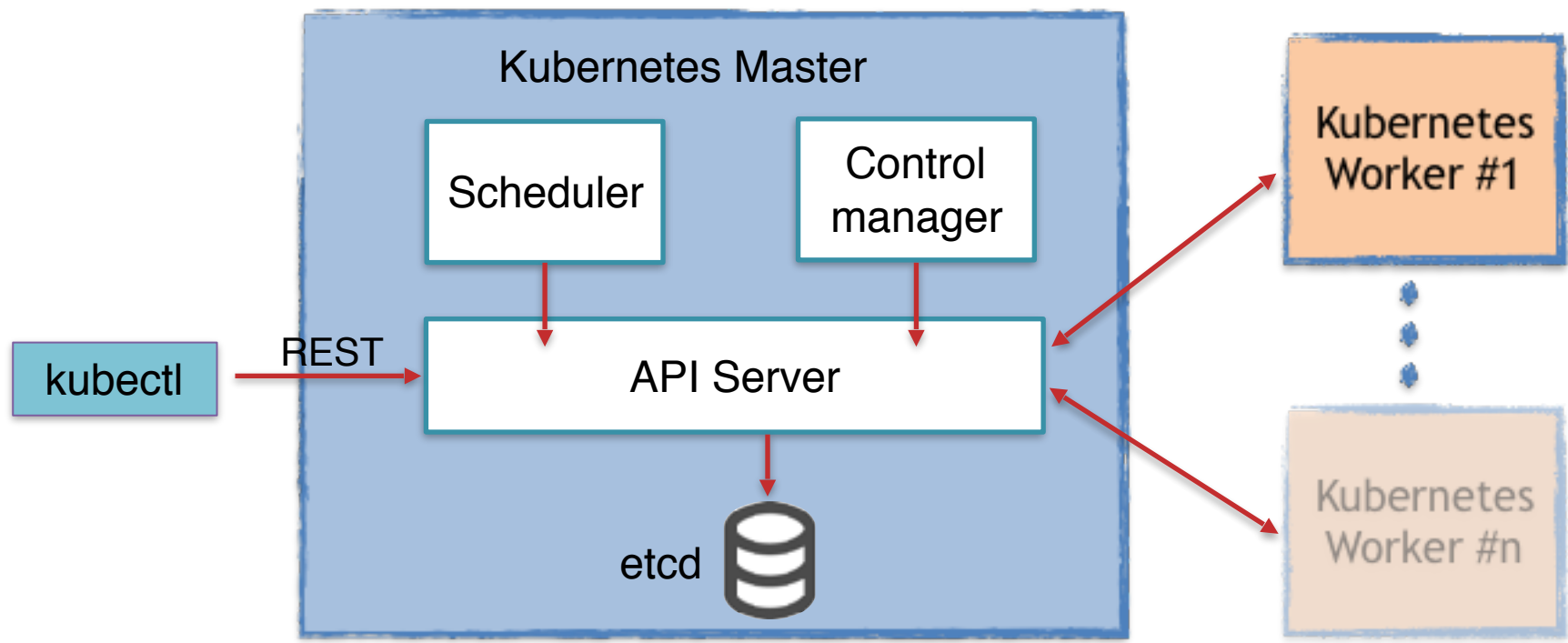
# Kubernetes Cluster Overview



- kubectl is tool for user to interact with k8s cluster.
- Master node interpret the command and if required interact with worker nodes.
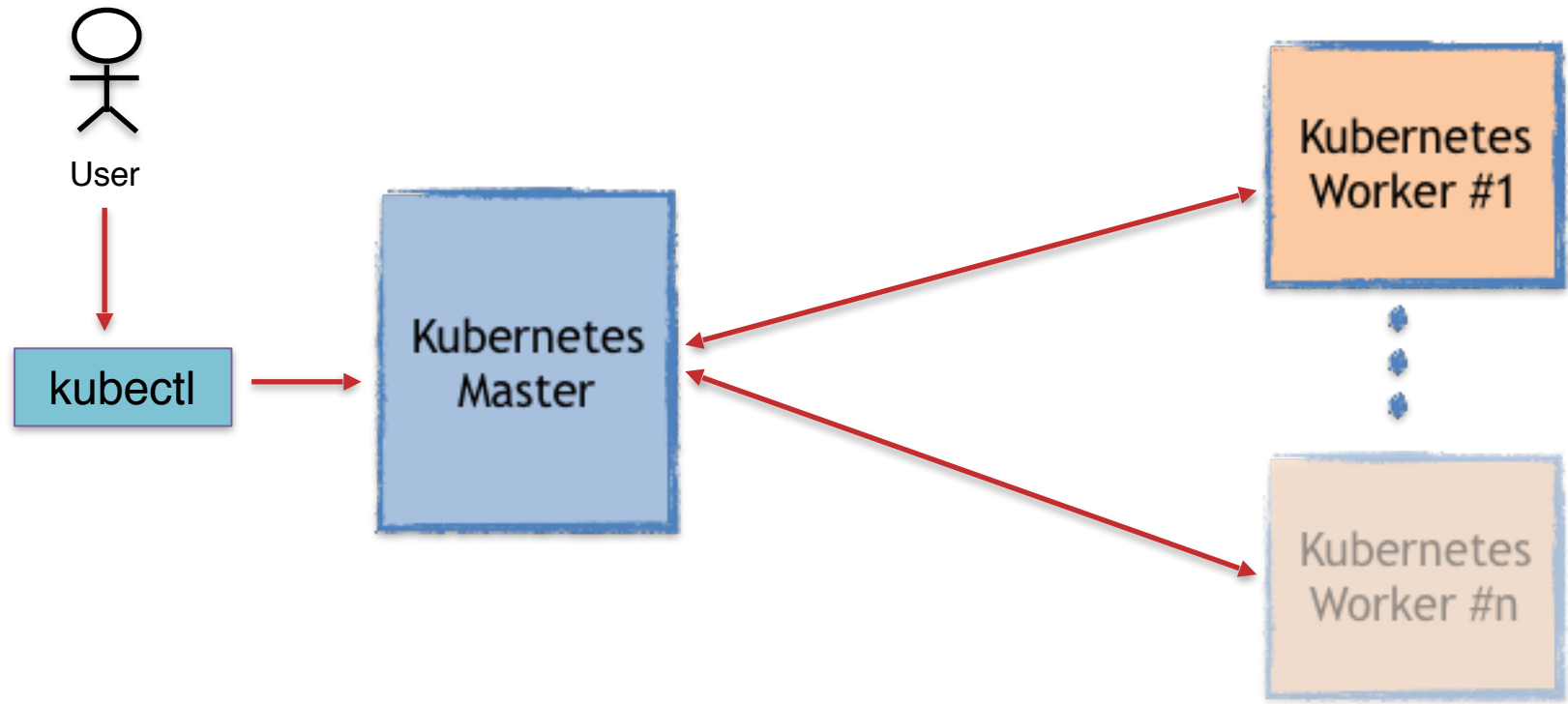
# Master Node Overview



Important components of Kubernetes Master Node

# Master Node Control Flow



- API Server plays a central part for cluster communication
- etcd store all definition of kubernetes resources
- Scheduler and Control Manager push commands for workers via API Server

# Kubernetes Architecture

# Kubernetes Worker Overview



Important components of Kubernetes Worker Node

# Kubernetes Worker Control Flow



- - Kubelet is the primary Node agent. API Server talks to Kubelet.
- - Service Proxy enables user to access applications running on node.
- - Docker running on node is used for creating Pods.

# Kubernetes Worker Control Flow



- Kubelet is the primary Node agent. API Server talks to Kubelet.
- Service Proxy enables user to access applications running on node.
- Docker running on node is used for creating Pods.

# Kubernetes Worker Overview

2014

Kubernetes Master

## Kubernetes Worker

Service Proxy

docker

Kubelet

rkt

Pod

Pod

With alternative container runtimes, Kubelet code gets bloated to support each.

# Container Runtime Interface

Introduced in Kubernetes 1.5 *. (2016)

Interfaces for gRPC service for Runtime & Image Management

Container centric interfaces

Pod containers as Sandbox containers

Current status: v1alpha2

*https://github.com/kubernetes/kubernetes/blob/release-1.5/docs/proposals/container-runtime-interface-v1.md

# Kubelet with CRI



CRI solves supporting various runtime alternatives with no change in Kubelet

# Container Runtime

# What is Container Runtime

Provides core primitives to manage containers on host

Container execution & supervision

Network Interfaces and management

Image management

Manage local storage

e.g. LXC, Docker, rkt

# Open Container Initiative

Container runtime & Image specification

Runtime specs define input to create a container

Multiple platform supported (Linux, Windows, Solaris & VM)

`runc` is default implementation of OCI Runtime Specs

Current Runtime Specs status : v1.0.1

# Gap between Kubelet & OCI runtime

| **Kubelet Requirements for Runtime** | **OCI Runtime** |
|---|---|
| Manage images (pull / push / rm ..) | Do not understand concept of image |
| Talks CRI / gRPC | Input is OCI specs (json and rootfs) |
| Prepare environment to successfully instantiate container. | Consume the rootfs and container config file (json) |
| Prepare network for pod | Attach network as pre-start hook. |

# Runtime in Kubernetes



Apart from OCI, *another* **runtime component** is required

# Runtime in Kubernetes

Kubernetes Worker

High-level Runtime

Kubernetes Master

Kubelet

Container Runtime

OCI Runtime

**CRI**

- High level runtime implement CRI gRPC services
- Take care of all prerequisite to successfully operate OCI runtimes

# Runtime in Kubernetes

Kubernetes Worker

Kubernetes Master

Kubelet

High-level Runtime

Container Runtime

Low-level Runtime

OCI Runtime

**CRI**

**OCI**

- OCI runtime works as low-level runtime
- High-level runtime provides inputs to OCI runtime as per OCI Specs

# CRI Implementations

- **Dockershim**
- **CRI-O**
- **Containerd**
- **Frakti**
- **rktlet**

# Dockershim



- Embedded into Kubelet.
- Dockershim talks to docker, which manage pods.
- Default CRI implementation & enjoy majority in current kubernetes deployments

# CRI-O



- CRI-O reduces the one extra hop from docker.
- CRI-O uses CNI for providing networking to pods.
- Monolithic design (understands CRI and outputs OCI compatible)
- Works with all OCI runtimes.

# containerD



- containerD, with revised scope eliminates the extra hop required by docker.
- Redesigned storage drivers for simplicity and better performance.
- Extensible design, CRI service runs as plugin.
- Uses CNI for networking
- Works with all OCI runtimes.

# Frakti



Kubernetes Worker

Kubelet — Dockershim — Frakti — docker — Pod

CRI

OCI — Hyped runV

VM Pod

- Frakti runtime was designed to support VM based runtime to kubernetes.
- It supports mixed runtimes
    - Linux containers for privilege containers and runV containers for rest
    - Though uses dockershim to use linux containers, result into extra hops
- Also supports Unikernels

# Frakti v2- Coming soon

Kubernetes Worker

Kubelet

**CRI**

containerd

CRI Plugin

Frakti Plugin

runC

Kata containers

Pod

VM Pod

- Frakti v2 will be implemented as runtime plugin for containerD.
- Reduce extra hops and implementation effort too.

# OCI Runtimes

runC

runV

Clear Containers

kata-runtime

gVisor

- Default OCI specs implementation
- Isolation based on Namespace, cgroups, secomp & MAC (AppArmor, SELinux)

# OCI Runtimes

| |
|---|
| runC |
| runV |
| Clear Containers |
| kata-runtime |
| gVisor |

- Default OCI specs implementation
- Isolation based on Namespace, cgroups, secomp & MAC (AppArmor, SELinux)

- OCI compliant VM based runtime
- Uses optimized qemu & KVM.
- A light weight guest kernel is used.

# OCI Runtimes

| | |
|---|---|
| runC | - Default OCI specs implementation<br>- Isolation based on Namespace, cgroups, secomp & MAC (AppArmor, SELinux) |
| runV | - OCI compliant VM based runtime<br>- Uses qemu & KVM.<br>- A light weight guest kernel is used. |
| Clear Containers | - Hardware-virtualized containers using Intel's VT-x<br>- Utilize DAX "direct access" feature of 4.0 kernel |
| kata-runtime | |
| gVisor | |

# OCI Runtimes

| | |
|---|---|
| runC | - Default OCI specs implementation<br>- Isolation based on Namespace, cgroups, secomp & MAC (AppArmor, SELinux) |
| runV | - OCI compliant VM based runtime<br>- Uses qemu & KVM.<br>- A light weight guest kernel is used. |
| Clear Containers | - Hardware-virtualized containers using Intel's VT-x<br>- Utilize DAX "direct access" feature of 4.0 kernel |
| kata-runtime | - Best of runV & cc-containers<br>- 1.0 Release (22nd May, 2018)<br>- Under active development |
| gVisor | |

# OCI Runtimes

| | |
|---|---|
| runC | - Default OCI specs implementation<br>- Isolation based on Namespace, cgroups, secomp & MAC (AppArmor, SELinux) |
| runV | - OCI compliant VM based runtime<br>- Uses qemu & KVM.<br>- A light weight guest kernel is used. |
| Clear Containers | - Hardware-virtualized containers using Intel's VT-x<br>- Utilize DAX "direct access" feature of 4.0 kernel |
| kata-runtime | - Best of runV & cc-containers<br>- 1.0 Release (22nd May, 2018)<br>- Under active development |
| gVisor | - Sandbox based containers<br>- Intercepts application system call acts like kernel.<br>- similar approach as User Mode Linux (UML)<br>- Under active development |

# Final candidates for Evaluation

| High-level Runtime |
|:---:|

Dockershim

CRI-O

containerD

| Low-level Runtime |
|:---:|

runC

Kata containers — runV

Kata containers — clear containers

# Why runtimes affect kubernetes

# Kubernetes Architecture

Kubernetes Master

Kubernetes Worker #1

Service Proxy

Runtime

Kubelet

Kubernetes Worker #n

Service Proxy

Runtime

Kubelet

- Kubernetes offers variety of choices to tune the system

# Kubernetes Architecture



- Kubernetes offers variety of choices to tune the system
- Once rest of components finalized
    - for deployment and management runtime is only variable factor.
    - For application performance only low level runtime matters.

# Performance benchmarking

**Application deployment performance**

- Container operations ( Create, start, stop, remove)

**Application Performance**

- Containerization / Virtualization overhead.

# Performance benchmarking process

- Prerequisite :
  - Pull Sandbox Image
  - Pull Container Image (ubuntu:latest)

## Benchmark Environment

```
Architecture:         x86_64
CPU(s):               8
Core(s) per socket:   4
Model name:           i7-3630QM CPU @ 2.40GHz
Virtualization:       VT-x
Kernel :              linux 4.15
OS :                  Ubuntu
```

### 4 Threads x 50

| Create | Start | Stop | Delete |
|--------|-------|------|--------|

- Create
  - Create & Run PodSandbox
  - Create Application Container

- Start
  - Start Application Container

- Stop
  - Stop Application Container

- Delete
  - Delete Application Container
  - Stop PodSandbox
  - Delete PodSandbox

- Rootfs prepared from Image
- Writable area for container
- CNI plugin invocation for Network

# runC Performance

**Software versions**

```
Containerd : v1.1.0
cri-o      : v1.10.1
Docker     : 18.05.0.ce
Runc       : v1.0
             git #69663f0bd4b
```



Legend: ■ containerd  ■ cri-o  ■ dockershim

**Create**
- containerd: 0.26
- cri-o: **0.73**
- dockershim: 0.64

**Start**
- containerd: 0.17
- cri-o: **0.03**
- dockershim: 0.75

**Stop**
- containerd: 0.19
- cri-o: **0.24**
- dockershim: 0.58

**Delete**
- containerd: 0.27
- cri-o: **0.53**
- dockershim: 0.19

X-axis: Seconds (0, 0.2, 0.4, 0.6, 0.8)

**Performance difference due to high level runtime**

**Low-level runtime (runC) is constant in all**

**cri-o and docker share same graph driver design, could be reason for high create time.**

**containerD perform better in almost all case.**

39

# Latency with runC

**Less is better**



| | containerd | cri-o | dockershim |
|---|---|---|---|

**Time to start**
- containerd: 0.43
- cri-o: 0.76
- dockershim: 1.39

**Time to stop**
- containerd: 0.46
- cri-o: 0.77
- dockershim: 0.77

X axis: 0, 0.35, 0.7, 1.05, 1.4 — **Seconds**

Time before application start running in runC container

Time before resources are released after application stops

**cri-o & containerD both perform better than docker**

**In performance, containerD performs better than cri-o**

# Kata-runtime Performance

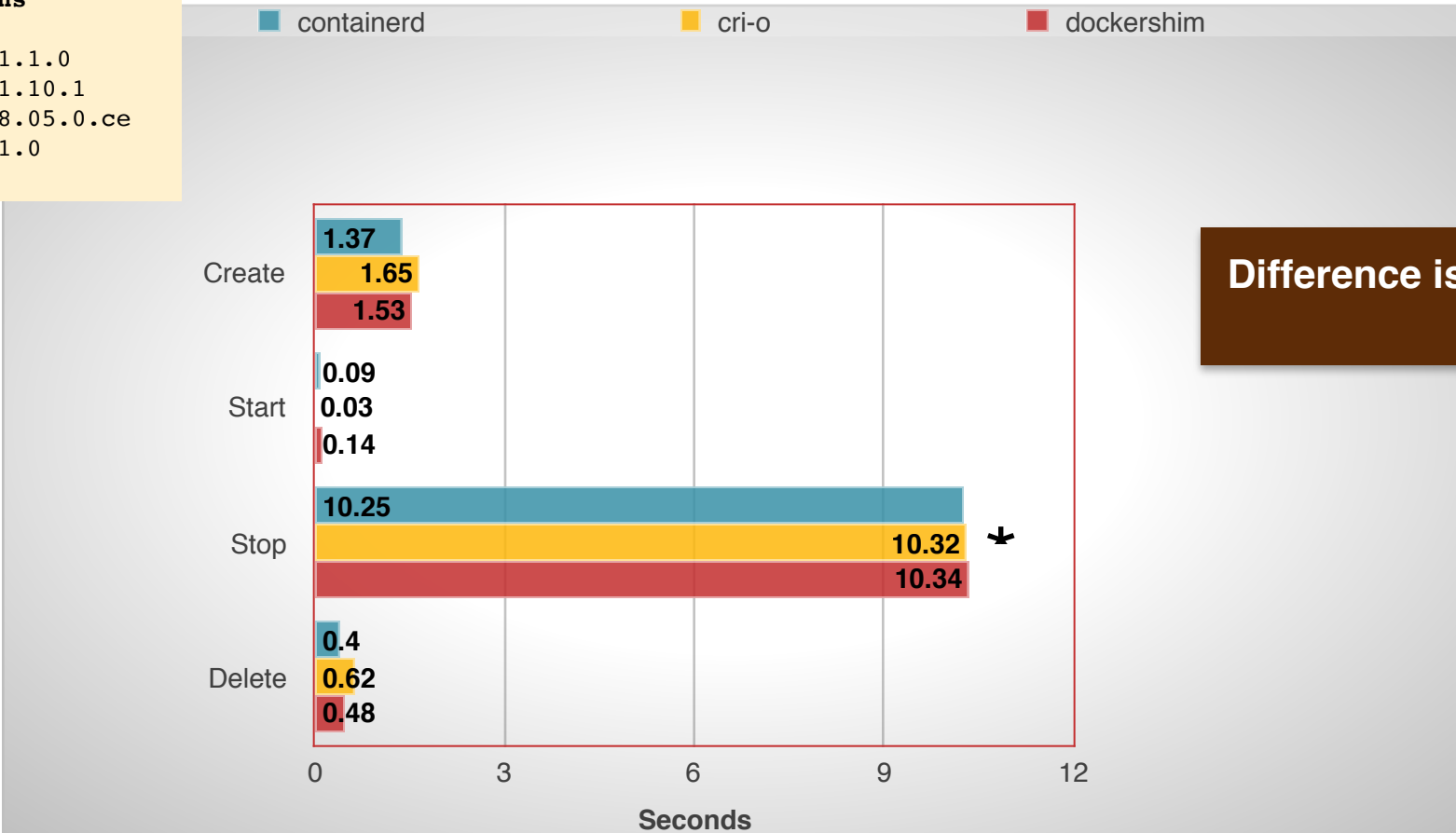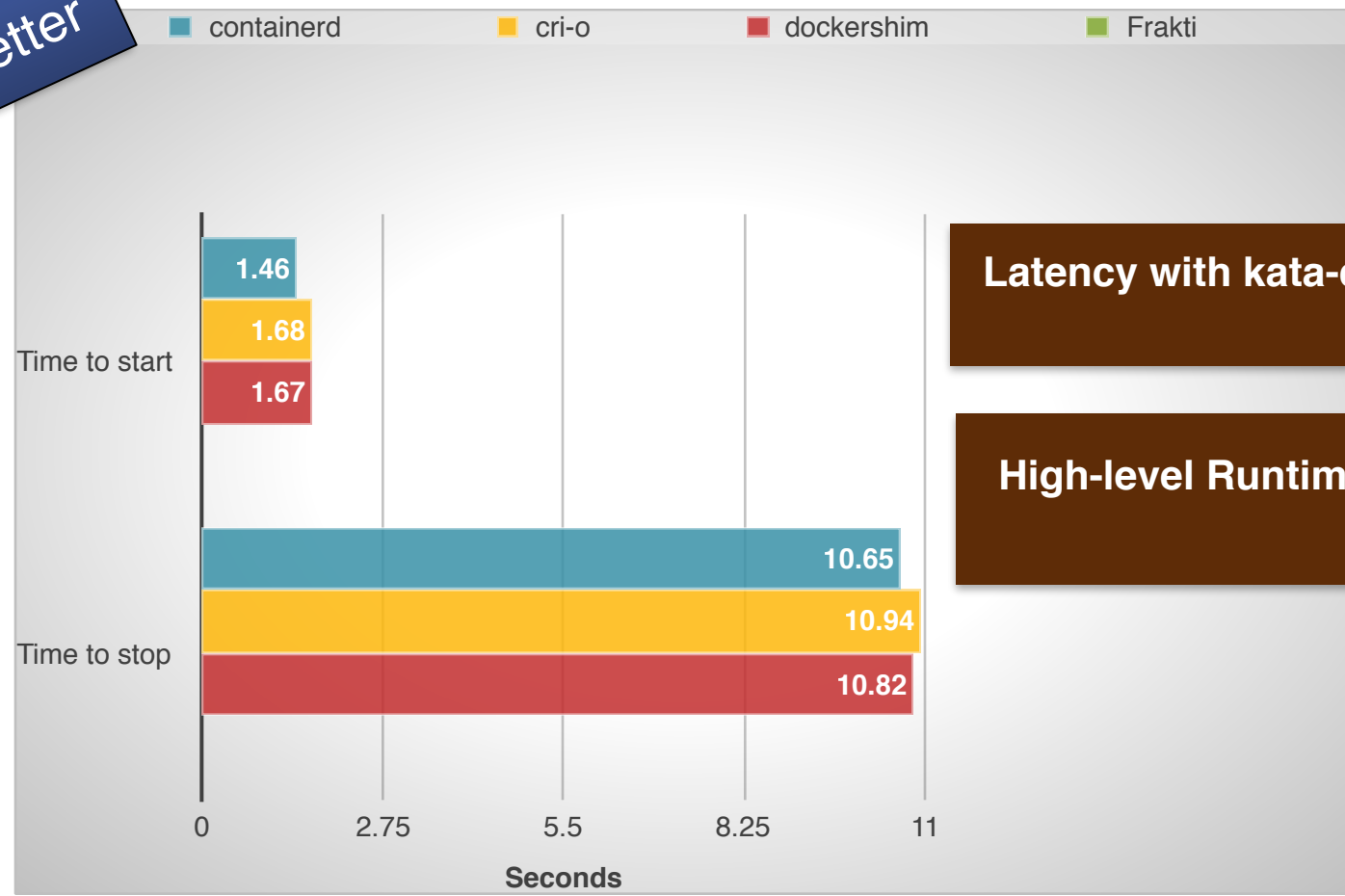Difference is mainly due to high level runtime performance.

*
- Bug in Stop logic, while invoked through CRI
  - Takes < 2 seconds, if done directly through docker or containerD

# Latency with Kata

# kata vs runV vs clear-containers



**Software versions**

```
Containerd : v1.1.0
Docker     : 18.05.0.ce
Frakti     : v1.10.0
runV       : v1.0.0
```

Legend: ■ kata + containerd  ■ cc-containers + containerd  ■ runV + frakti

| | Seconds |
|---|---|
| **Create** | 1.37 / 2.67 / 0.6 |
| **Start** | 0.09 / 0.35 / 0.69 |
| **Stop** | 10.25 / 1.49 / 0.54 |
| **Delete** | 0.4 / 0.81 / 0.32 |

**Stop function of cc-containers & runV looks normal. Hence fix required for kata containers.**

**Kata containers performance is in-between runV and cc-runtime.**

# Latency with VM based runtimes

**Less is better**



Legend: ■ kata + containerd    ■ cc-container + containerd    ■ runV + frakti    ■

**Time to start:**
- kata + containerd: 2.27
- cc-container + containerd: 3.02
- runV + frakti: 1.29

**Time to stop:**
- kata + containerd: 10.65
- cc-container + containerd: 2.3
- runV + frakti: 0.86

X-axis: Seconds (0, 2.75, 5.5, 8.25, 11)

**runV performs for container operations is best in VM containers.**

**Kata is still in active development**

# Performance Overhead – Low-level runtimes

**More is better**

**Less is better**

**I/O Throughput**

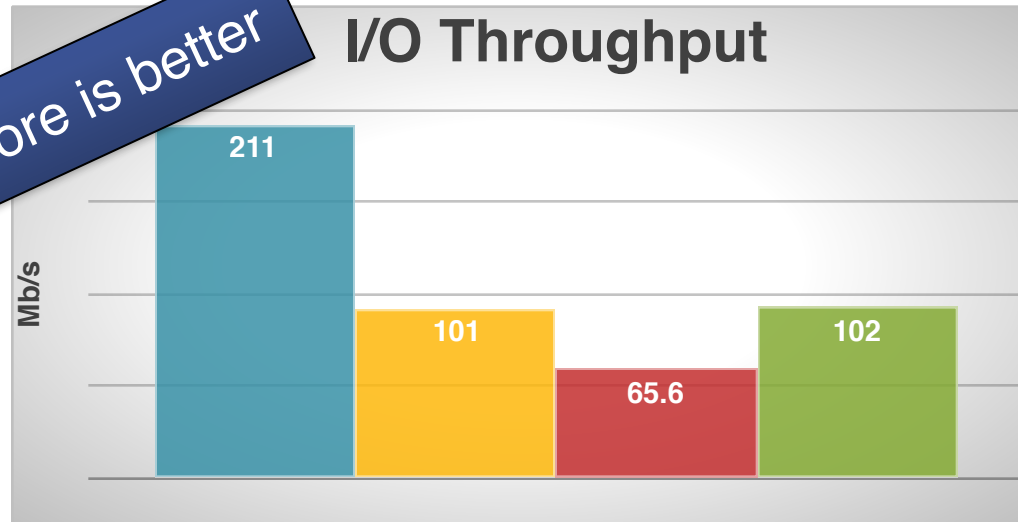| Runtime | Mb/s |
|---|---|
| runC | 211 |
| kata-containers | 101 |
| runV | 65.6 |
| clear containers | 102 |

**Average System Load**

| Runtime | CPU Load |
|---|---|
| runC | 1.62 |
| kata-containers | 3.17 |
| runV | 3.17 |
| clear containers | 3.91 |

■ runC ■ kata-containers ■ runV ■ clear containers

**Runtime performance overhead affect application running inside container.**

**runC perform best in both IO throughput and average CPU load.**

**kata-containers perform best among VM containers.**

# Workloads

# Serverless

- **Host functions instead of applications?**
  - Functions as service
  - e.g. AWS Lambda
- **Ideal Platform**
  - Low latency
  - High parallelism i.e. high density.
  - Low on resources (CPU, Memory)

# Serverless platform

| | containerd + runC | cri-o + runC | Frakti + runV | Any + kata-containers |
|---|---|---|---|---|
| **Latency** | Best | Better | Good | Average |
| Cold start | Best | Better | Better | Average |
| Warm start | Better | Best | Average | Good |
| **Density** | Best | Good | Average | Average |
| **Security** | Good (namespace + seccomp + SELinux) | Good (namespace + seccomp + SELinux) | Best (VM based) | Best (VM Based) |
| **Stability** **Support Cycle** | Stable (defined support cycle for each release) ✓ | Stable/Best with Openshift (Not defined) | Stable (managed by hyper.sh) (not defined) | Under Active development (Not defined) |

# Peak hour demand / Micro Services

- **Mostly applications are of type Micro services.**
- **Ideally immutable**
- **Quick scale up and scale down.**
- **Ideal Platform**
  - Low latency for start application and free resources.
  - Better utilize the host system.

# Mean Time To Recover (MTTR) - DevOps

- **Short Lived containers**
- **Frequent updates**
- **Fast recovery is important.**
- **Low on resources**

# Micro-services / MTTR

| | containerd + runC | cri-o + runC | Frakti + runV | Any + kata-containers |
|---|---|---|---|---|
| **Latency** | Best | Better | Good | Average |
| **Density** | Best | Better | Average | Good |
| **Security** | Good (namespace + seccomp + SELinux) | Good (namespace + seccomp + SELinux) | Best (VM based) | Best (VM Based) |
| **Stability**<br><br>**Support Cycle** | Stable<br><br>(defined support cycle for each release)<br><br>✔ | Stable/Best with Openshift<br><br>(Not defined) | Stable<br><br>(managed by hyper.sh) (not defined) | Under Active development |

# Long running containers

- **Migrated application.**
- **Stateful containers.**
- **Hard to scale containers.**
- **Requirements**
  - Stability
  - Security
  - Performance
  - Migration

# Long running containers

| | containerd + runC | cri-o + runC | Frakti + runV | Any + kata-containers |
|---|---|---|---|---|
| **Stability**<br><br>**Support Cycle** | Best<br><br>(defined support cycle for each release) | Stable/Best with Openshift<br><br>(Not defined) | Good<br><br>(managed by hyper.sh)<br>(not defined) | Under Active development<br><br>(not defined) |
| **Security** | Good<br>(namespace + seccomp + SELinux) | Good<br>(namespace + seccomp + SELinux) | Best<br>(VM based) | Best<br>(VM Based) |
| **Performance Overhead** | Best | Best | Average | Better |
| **Migration** | Required | Required | Required | Required |
| **Governance** | CNCF + OCI ✔ | Kubernetes Incubator + OCI | Kubernetes + hypersh | OpenStack Foundation ✔ |

# Summary

- CRI and OCI enable more choices for container runtimes.

- For Cloud Native workloads, Linux containers based runtimes suite better.

- High level runtime performance do not matter much for long running containers, So low level runtime performance & capabilities become focus.

- VM based runtimes are promising, but still need some time to reach flexibility and usability  as Linux containers runtime.

- Migration of monolithic applications / high security applications to modern platform like kubernetes will get boost with VM based runtimes.

# Few more OCI runtimes

- **Runtime getting ready for OCI complaint**
  - rkt - container runtime from CoreOS
    - https://github.com/rkt/rkt
    - https://github.com/rkt/rkt/issues/3368
  - gVisor - Sandbox based containerization
    - https://github.com/google/gvisor
  - railcar – linux containers in implementation in rust
    - https://github.com/oracle/railcar
    - slow development
  - crun – linux containers in implementation in C
    - https://github.com/giuseppe/crun
    - Fully featured but lack clarity on maintenance and support.

**Thank You**