# Don't Peek Into my Container!

Christophe de Dinechin, Red Hat
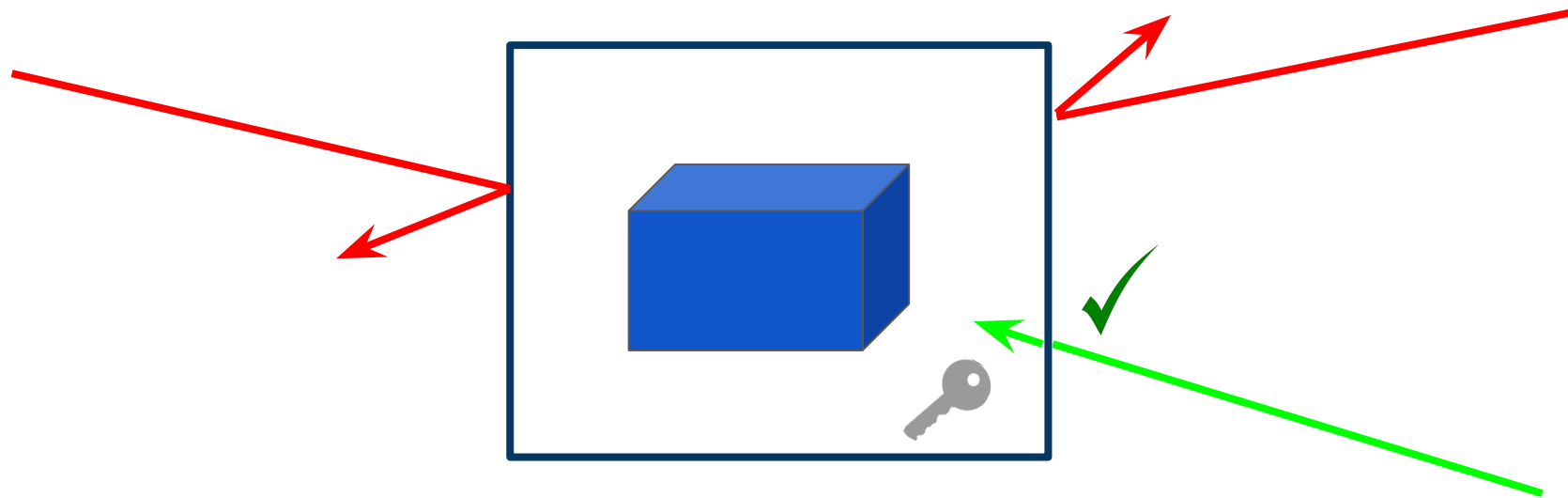
Alice Frosi, Red Hat

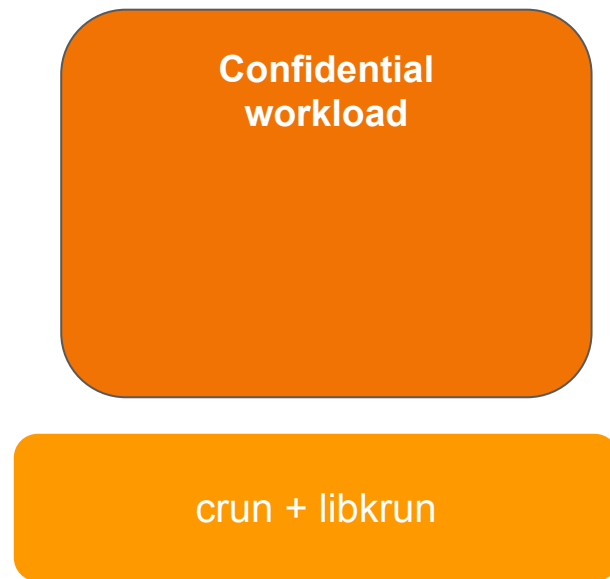Sergio López Pascual, Red Hat

**Red Hat**

# Today's Topics

- ▶ Confidential computing

- ▶ Confidential workloads with k8s and libkrun

- ▶ SEV-enabled libkrun
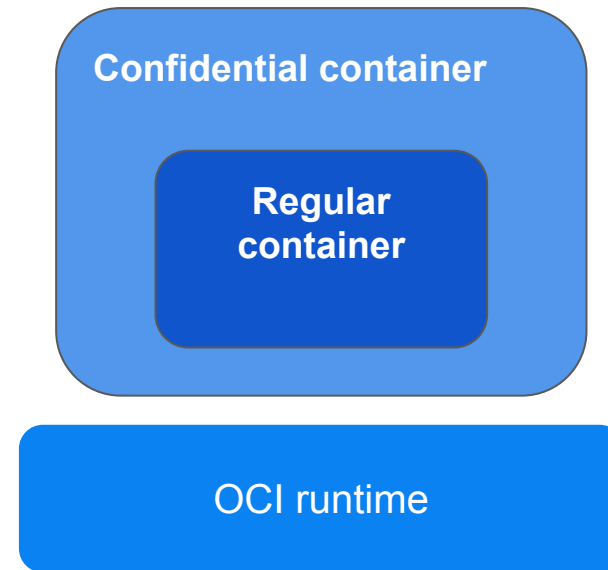
- ▶ From Kata Containers to Confidential Containers

**"Confidential Computing is the protection of data in use by performing computation in a hardware-based Trusted Execution Environment"**

https://confidentialcomputing.io/whitepaper-02-latest/

**Confidential workloads** are transformed containerized workloads into a special form that can be deployed with libkrun and confidential computing technologies

**Confidential containers** are the deployment of a regular containers with an OCI runtime (e.g Kata Containers) and confidential computing technologies
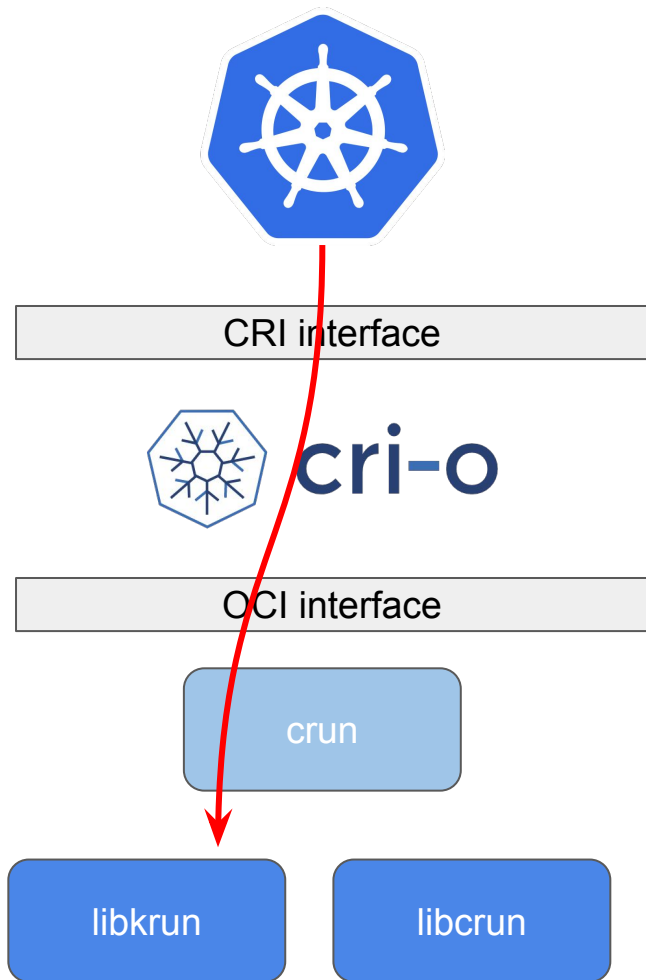
## Confidential workloads

- ▸ Confidentiality at container level

- ▸ Single container per encrypted VM

- ▸ Deploy a special form of container image with a single layer

- ▸ Simpler architecture and reusing the existing k8s infrastructure

## Confidential containers

- ▸ Confidentiality at pod level

- ▸ Multiple containers per encrypted VM

- ▸ Use encrypted layered container images

- ▸ Part of the infrastructure is moved inside the trusted environment (e.g image offloading)
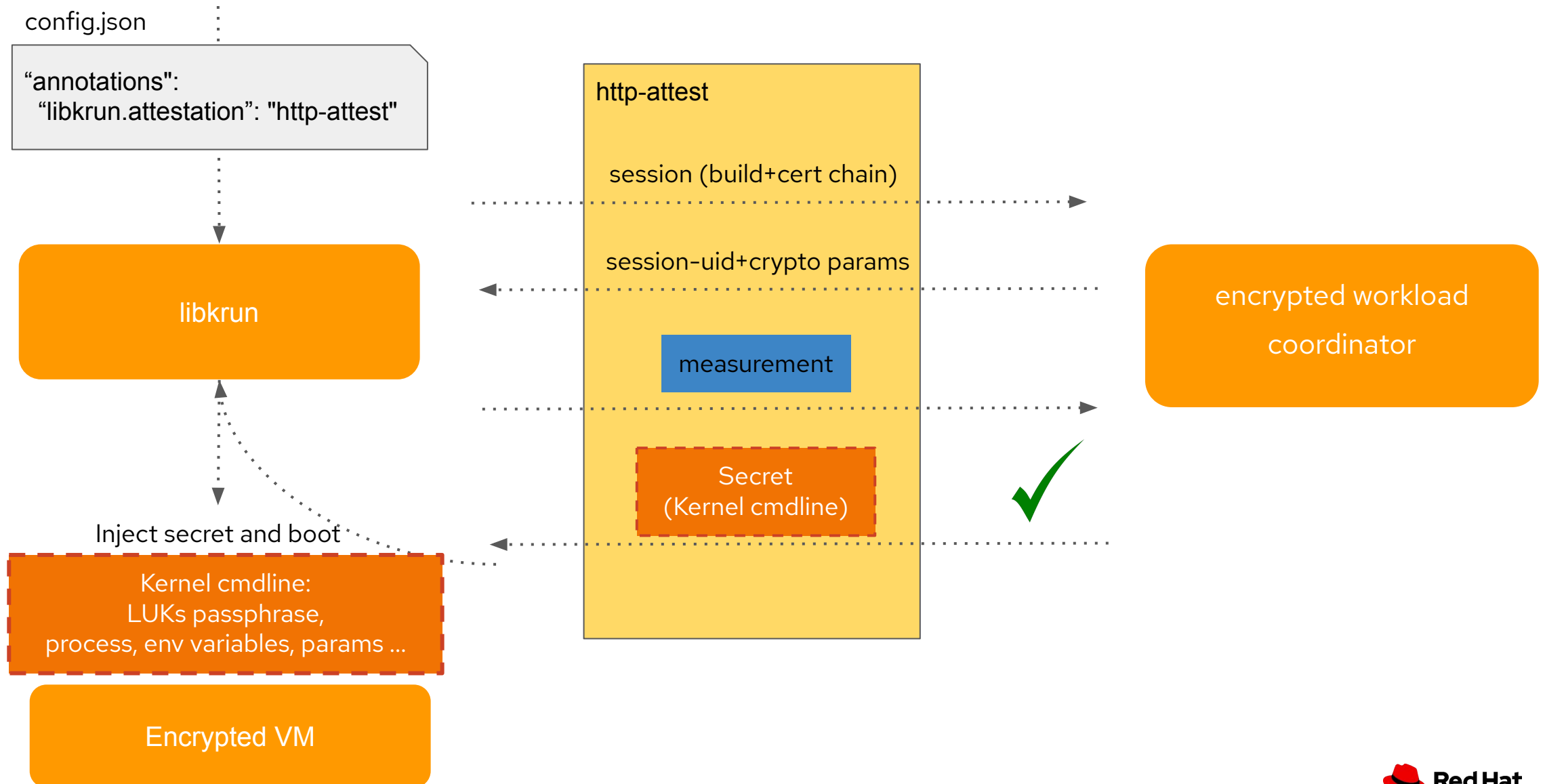
Red Hat

# Confidential workloads with k8s and libkrun

Red Hat

Receive the pod information

Schedule the workload on a node

CRI interface

Pull the container image on the node

Prepare the bundle with the rootfs of the container

Create the config.json with the container information

OCI interface

crun

Create and run the container

libkrun          libcrun
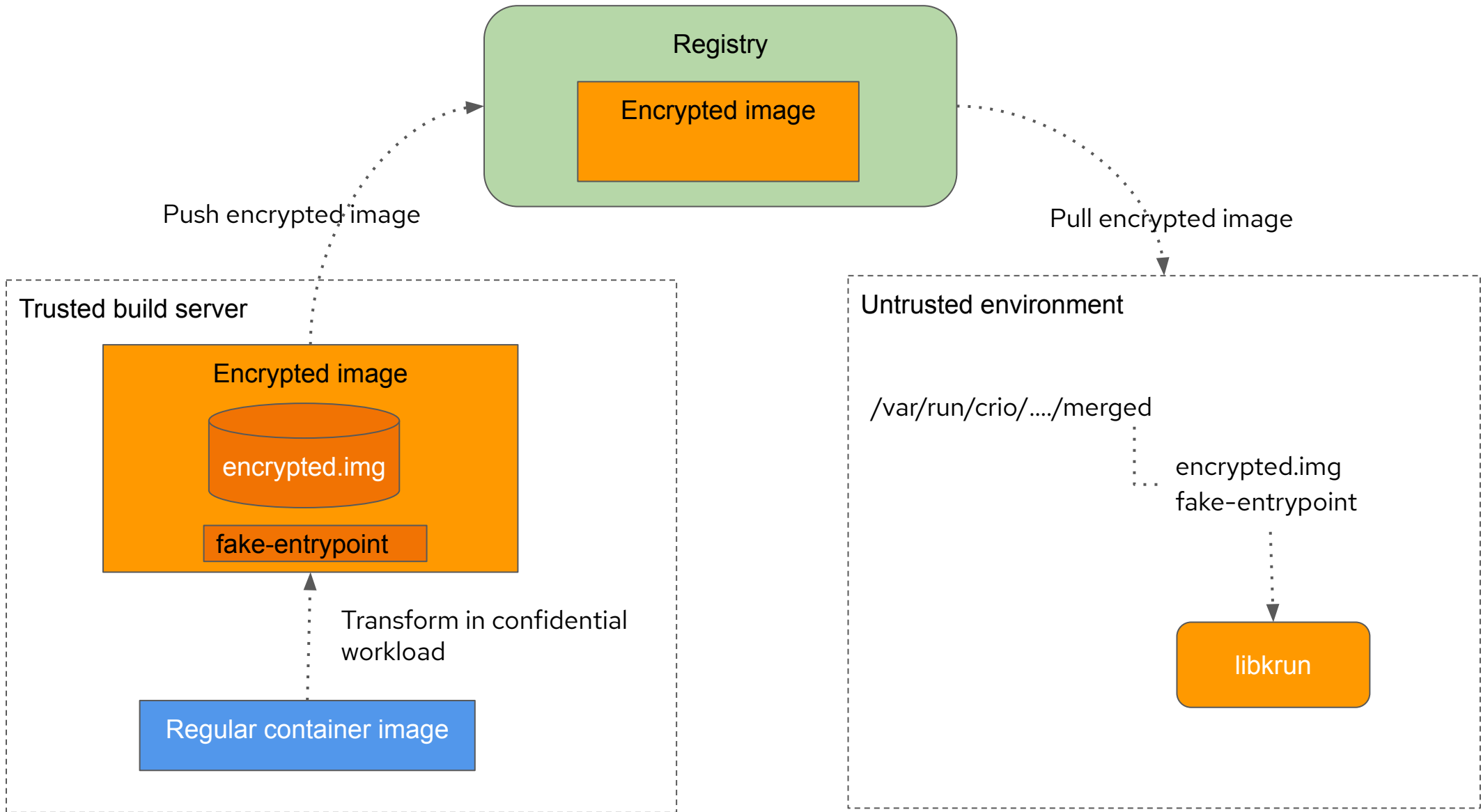
```
apiVersion: v1
kind: Pod
metadata:
  name: pod-krun-sev
  annotations:
    libkrun.attestation: "http-attest"
    run.oci.handler: "krun-sev"
spec:
  containers:
  - image: encrypted/nginx-tls
    name: krun-sev
    command: ["/fake-entrypoint"]
    ports:
    - containerPort: 443
  nodeSelector:
    sev: "true"
```

Red Hat

config.json

"annotations":
  "libkrun.attestation": "http-attest"

libkrun

http-attest

session (build+cert chain)

session–uid+crypto params

measurement

Secret
(Kernel cmdline)

✔

encrypted workload

coordinator

Inject secret and boot

Kernel cmdline:
LUKs passphrase,
process, env variables, params ...

Encrypted VM

Red Hat

# SEV-enabled libkrun for Confidential Workloads

Red Hat

# The need for a Minimal Firmware (I)

Original functionality

**Guest Memory**

libkrun (VMM)

Initial Page Tables

e820

MP table

Linux Zero Page

Kernel Image

initramfs

These tables are written by the VMM and become part of the launch measurement!

Red Hat

# The need for a Minimal Firmware (II)

**Before** starting the VM

**After** starting the VM

SEV-enabled functionality

Guest Memory

libkrun (VMM)

Minimal FW

Kernel Image

initramfs

Initial Page Tables

e820

MP table

Linux Zero Page

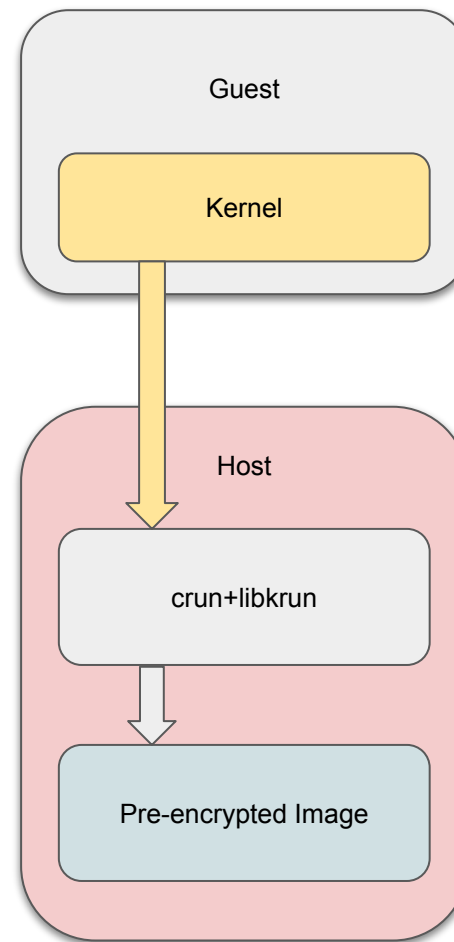These tables are written by the Firmware and are **NOT** part of the launch measurement!

Red Hat

# Replacing virtio-fs with virtio-blk (I)

Regular libkrun
with virtio-fs

SEV-enabled
libkrun with
virtio-blk

Guest

Kernel

Host

crun+libkrun

Image expanded into a
directory on the Host

Guest

Kernel

Host

crun+libkrun

Pre-encrypted Image

# Replacing virtio-fs with virtio-blk (II)

▸ libkrun uses virtio-fs because fits nicely with the container isolation use case.

▸ But, for the Confidential Workloads use case, virtio-fs is not the best solution.

　· Even if we could find an acceptable filesystem-level encryption mechanism, the implementation will leak too much information.

　· The implementation is quite large and complex (is, by far, the largest component in libkrun) compared with virtio-blk, and requires a large number of syscalls, which implies a more permissive seccomp filter.

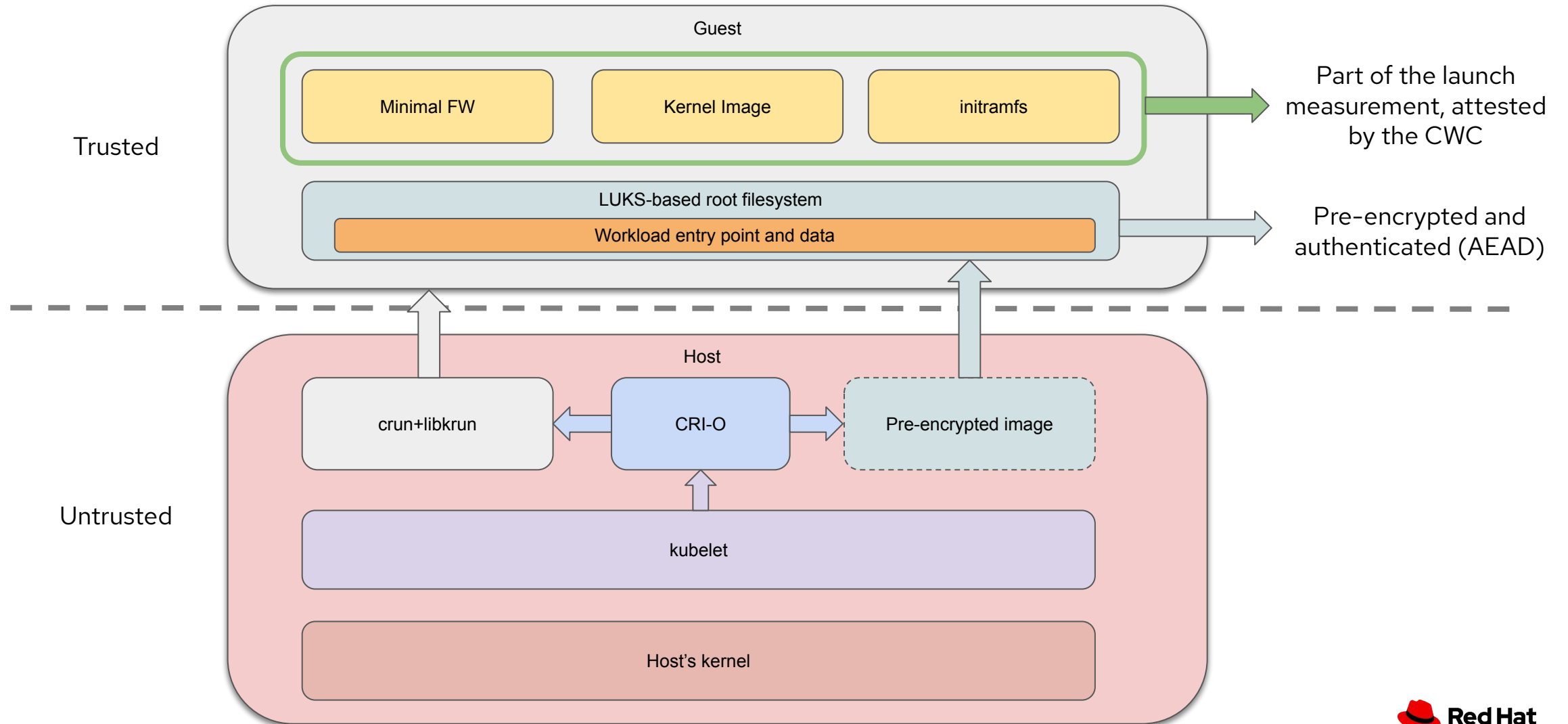　　· Lines of Code: virtio-fs = 7444, virtio-blk = 1325

# Replacing virtio-fs with virtio-blk (III)

▶ Using virtio-blk allows us to easily rely on LUKS.

- LUSK2 has the ability to combine dm-crypt and dm-integrity.

- Provides both confidentiality and integrity protection (Authenticated Encryption with Additional Data, AEAD).

- Protects against all known attacks, except data replay, which would require specialized hardware storage.

- Reference: <u>Practical Cryptographic Data Integrity Protection with Full Disk Encryption Extended Version</u> by Milan Brož, Mikuláš Patočka and Vashek Matyáš.

# The need for an initramfs

▶ The binary we use in libkrun to set up the environment inside the guest is bundled in the integrated virtio-fs server.

- Without virtio-fs, we needed an alternative.

▶ We've incorporated a simple initramfs.

- Includes a variant of the binary to set up the environment, a static version of cryptsetup, and some support directories and device nodes.
- For SEV-SNP and TDX cases, it'll likely include a small attestation client.
- Opens the LUKS device (using the injected secret) and continues doing the usual environment adjustments before executing the workload entry point.
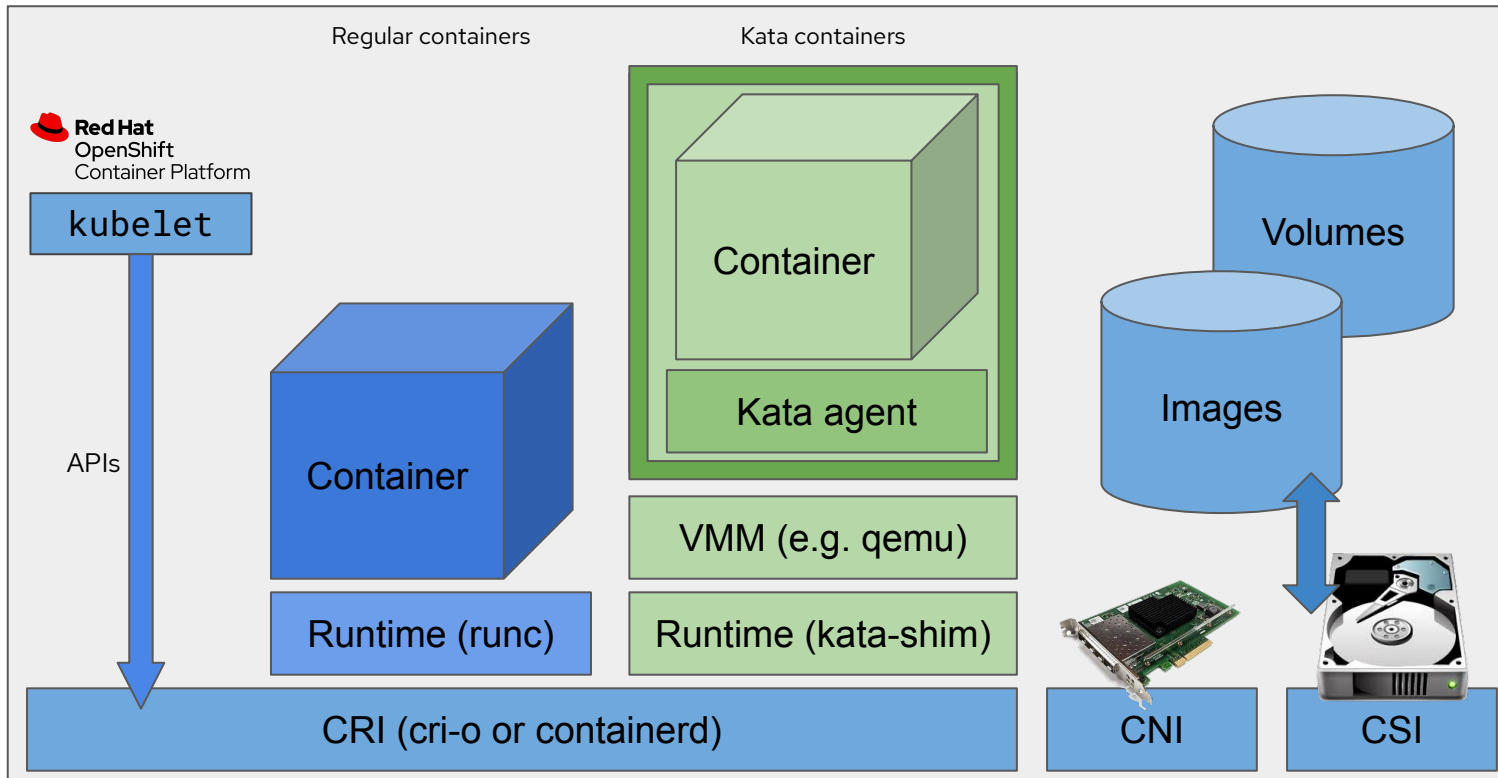
Red Hat

# The Big Picture

**Trusted**

**Guest**

| Minimal FW | Kernel Image | initramfs |
|---|---|---|

Part of the launch measurement, attested by the CWC

LUKS-based root filesystem

Workload entry point and data

Pre-encrypted and authenticated (AEAD)

**Untrusted**

**Host**

crun+libkrun ← CRI-O → Pre-encrypted image

kubelet

Host's kernel

17

# From Kata containers to Confidential containers

Red Hat

# Kata Containers overview

▸ **Run Containers** described the usual way (e.g. same yaml file, images, storage, networking...)

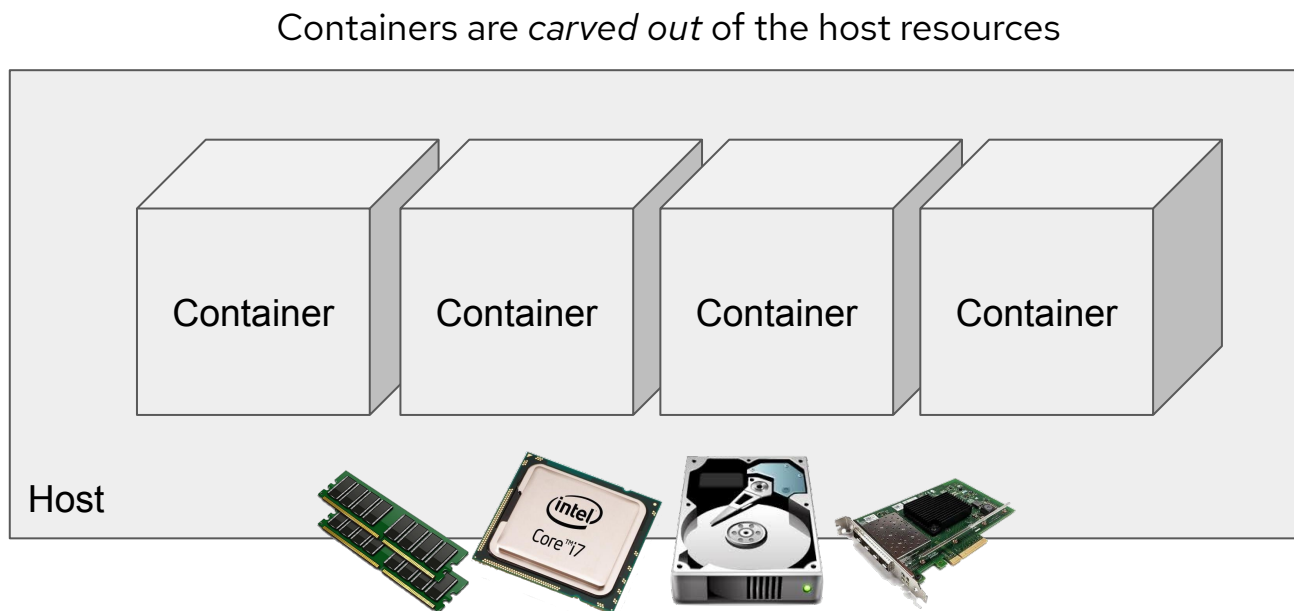▸ **... in Virtual Machines** with their own independent kernel and very little user space



The ecosystem of **containers**
The sandboxing of **virtualization**

**CRI**: Container *Runtime* Interface
**CNI**: Container *Networking* Interface
**CSI**: Container *Storage* Interface

19

# Problem statement:
# Can we trust the host?

Red Hat

▸ **Containers** run on a host, often managed by a third party, like a cloud provider

▸ **Sandboxing** goes only one way, protecting the host from containers, not the other way round

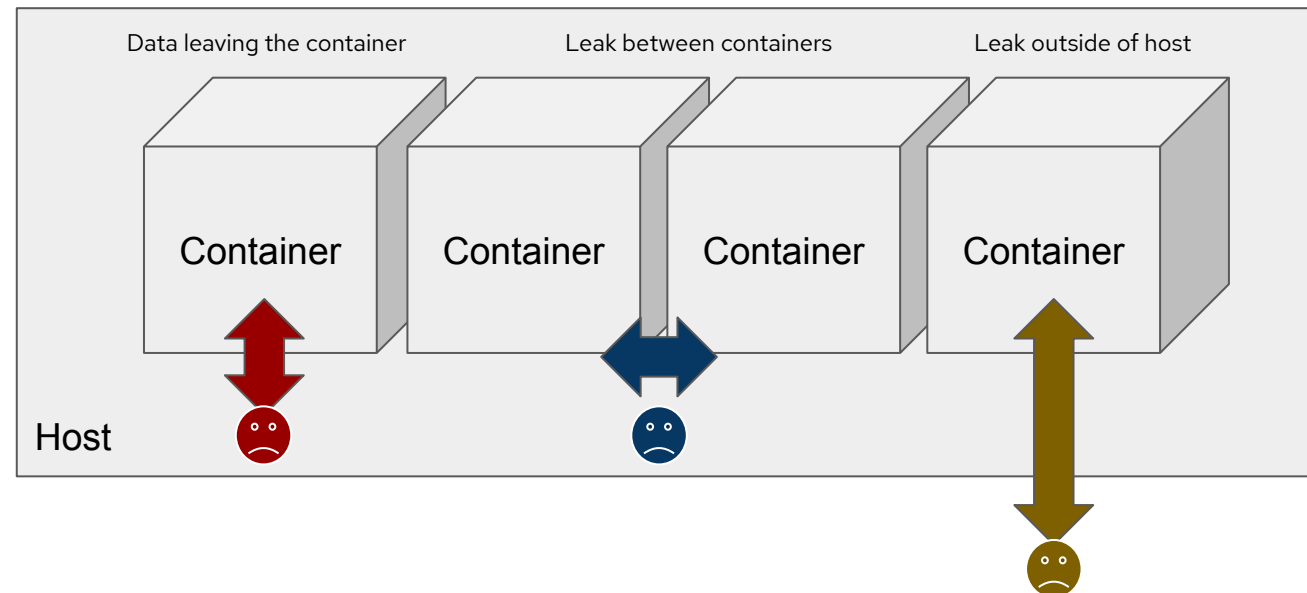▸ **Resources** belong to the host, which owns them and has free access

Containers are *carved out* of the host resources



| Container | Container | Container | Container |

Host

What do you need to do if you start considering the host as **hostile**?

21

▸ **Data exposure** of information held in the container is possible

▸ **Multiple tenants** may not want to share the same host because of confidentiality risk

▸ **Legal concerns** may preclude the use of containers if you cannot guarantee confidentiality
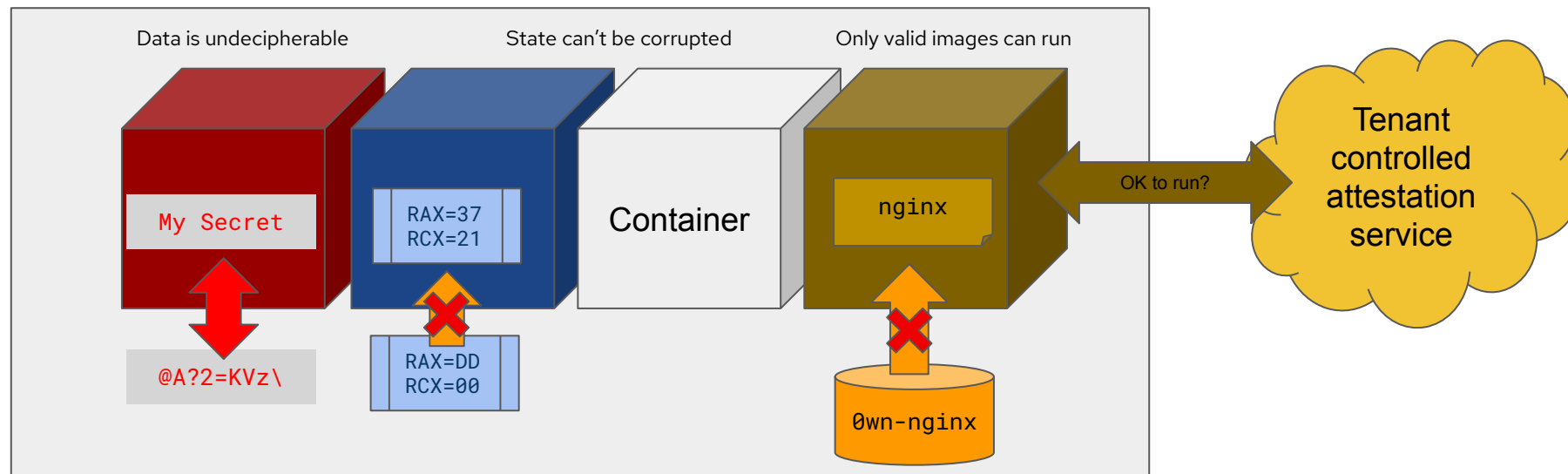
There is a potential for unwanted data leaks

# Enabling technology:
# Confidential Computing

Red Hat

# Confidential Computing: more than encryption...

▸ **Memory encryption** prevents the host from getting data out of guest memory

▸ **Integrity protection** offers guarantees about guest state corruption

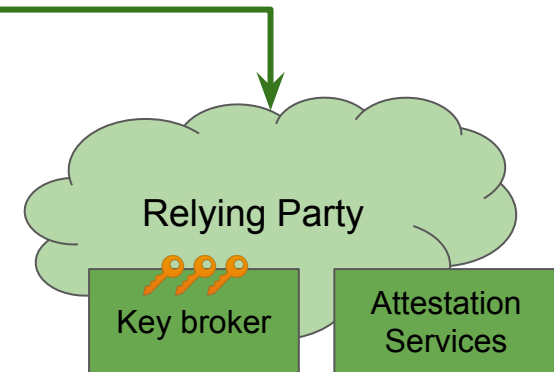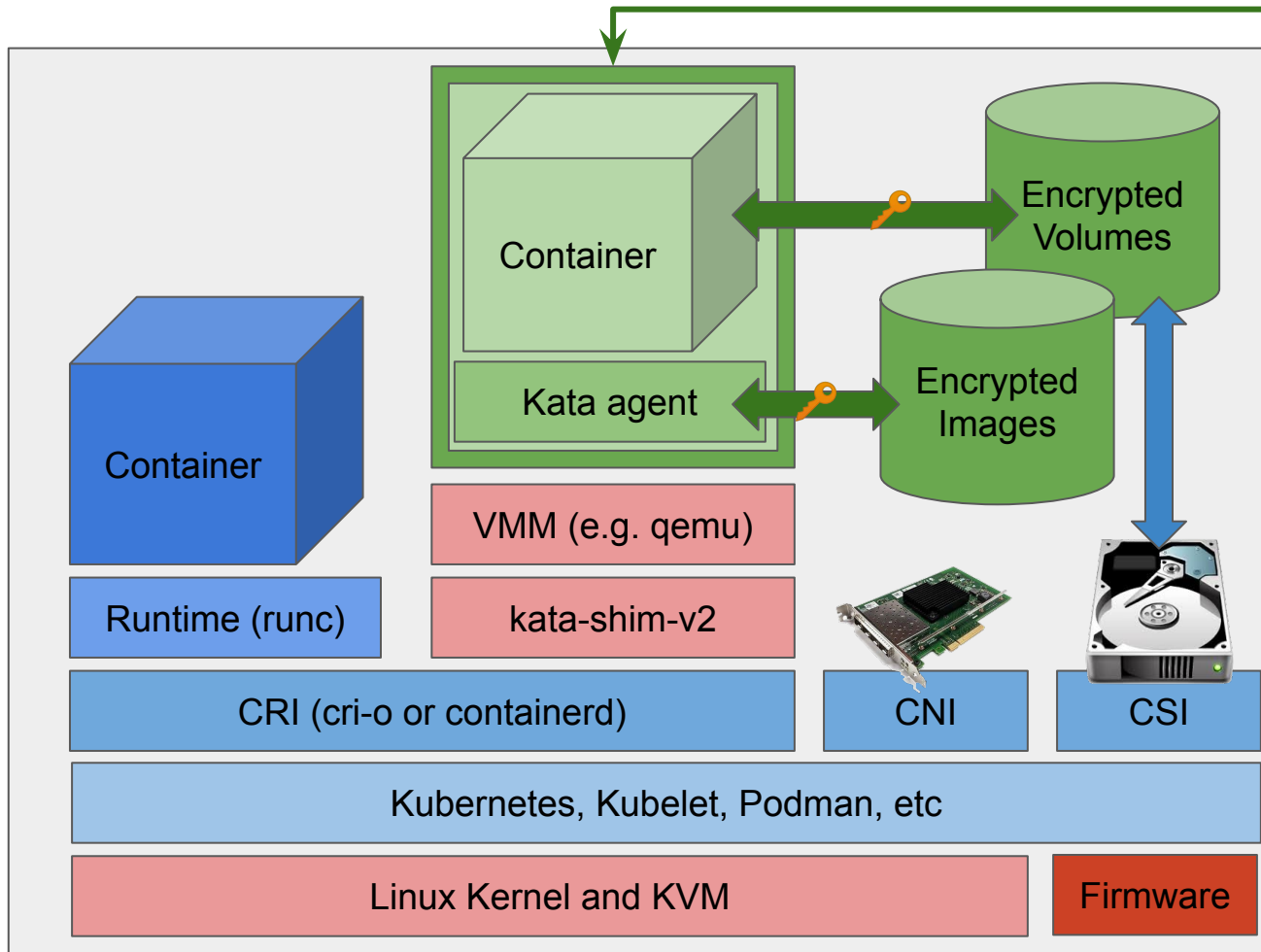▸ **Attestation** lets the guest owner (tenant) validate what runs in the guest

# Many vendor-specific technologies

- **AMD** offers Secure Encrypted Virtualization (SEV)
  - SEV-ES adds Encrypted State (e.g. CPU register file)
  - SEV-SNP adds Secure Nested Pages (integrity protection for memory and more)
- **Intel** offers Trusted Domain Extensions (TDX)
- **IBM S390** offers Secure Execution (SE)
- **Power** offers Protected Execution Facility (PEF)
- **Arm** announced Confidential Computing Architecture (CCA)

- All these technologies are based on **virtualization**
- Each of these technologies works in a slightly[1] different way. There be zombies 🧟

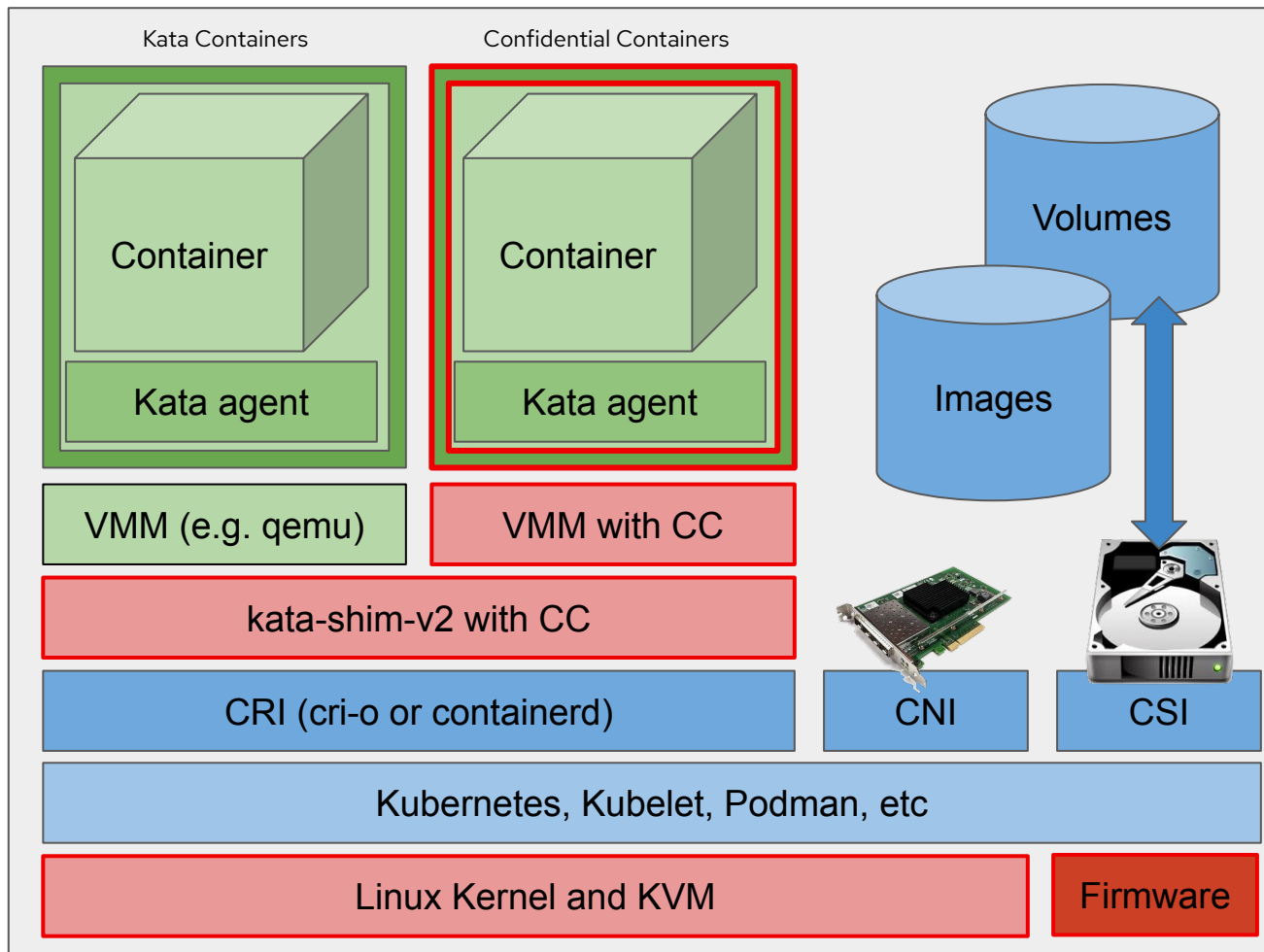[1] For a slightly understated definition of "slightly"

# Separate Trust Realms: Platform, Tenant and Host



Container

Runtime (runc)

Container

Kata agent

VMM (e.g. qemu)

kata-shim-v2

Encrypted Volumes

Encrypted Images

CRI (cri-o or containerd)

CNI

CSI

Kubernetes, Kubelet, Podman, etc

Linux Kernel and KVM

Firmware

Relying Party

Key broker

Attestation Services

**Relevant Trust Realms**

**Trusted Platform**: Offers confidentiality guarantees using hardware-level cryptographic enforcement.

**Host**: Offers and manages the resources used to run the container (CPU, memory, I/O, etc)

**Tenancy**: Confidential area carved out of the host, but not visible nor accessible to it.
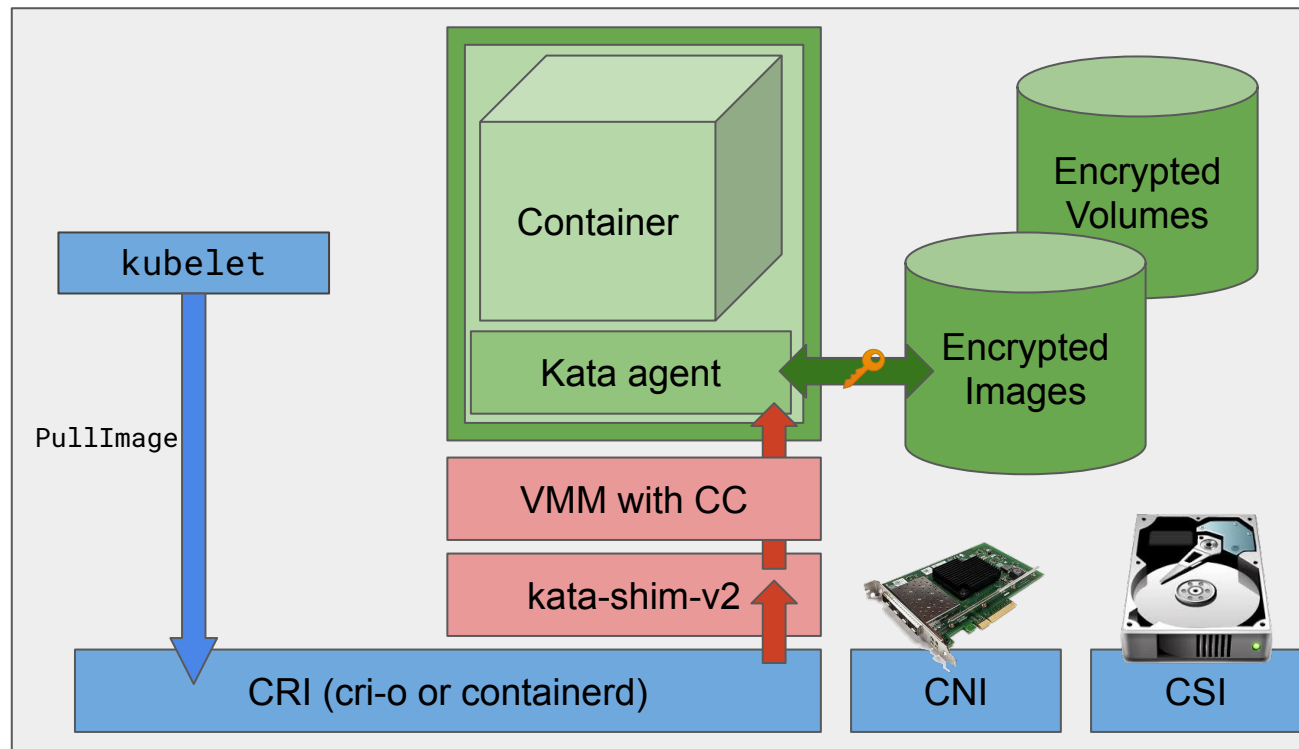
# Enabling Confidential Computing for Kata Pods



Kata Containers

Confidential Containers

Container

Kata agent

Container

Kata agent

Volumes

Images

VMM (e.g. qemu)

VMM with CC

kata-shim-v2 with CC

CRI (cri-o or containerd)

CNI

CSI

Kubernetes, Kubelet, Podman, etc

Linux Kernel and KVM

Firmware

**Impacted components**:

**Kata runtime**: Pass right options to VMM

**VMM**: Enable encryption, etc, when setting up VM

**Kernel**: Low-level hardware support, e.g. SEV, TDX

**Firmware**: Special services, e.g. page validation

**Hardware**: Encryption in memory controller

The Kata development is done for most platforms

# Securing Image Download

▸ **Pull Image** from *inside* the guest instead of pulling it from the host

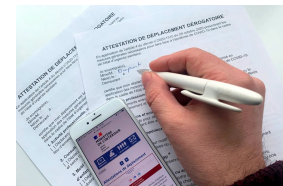▸ **Store Images** on an encrypted volume, where only guest has decryption keys



Pass `PullImage` API to guest

The Kubelet delegates the `PullImage` operation to the `ImageService` in the CRI.
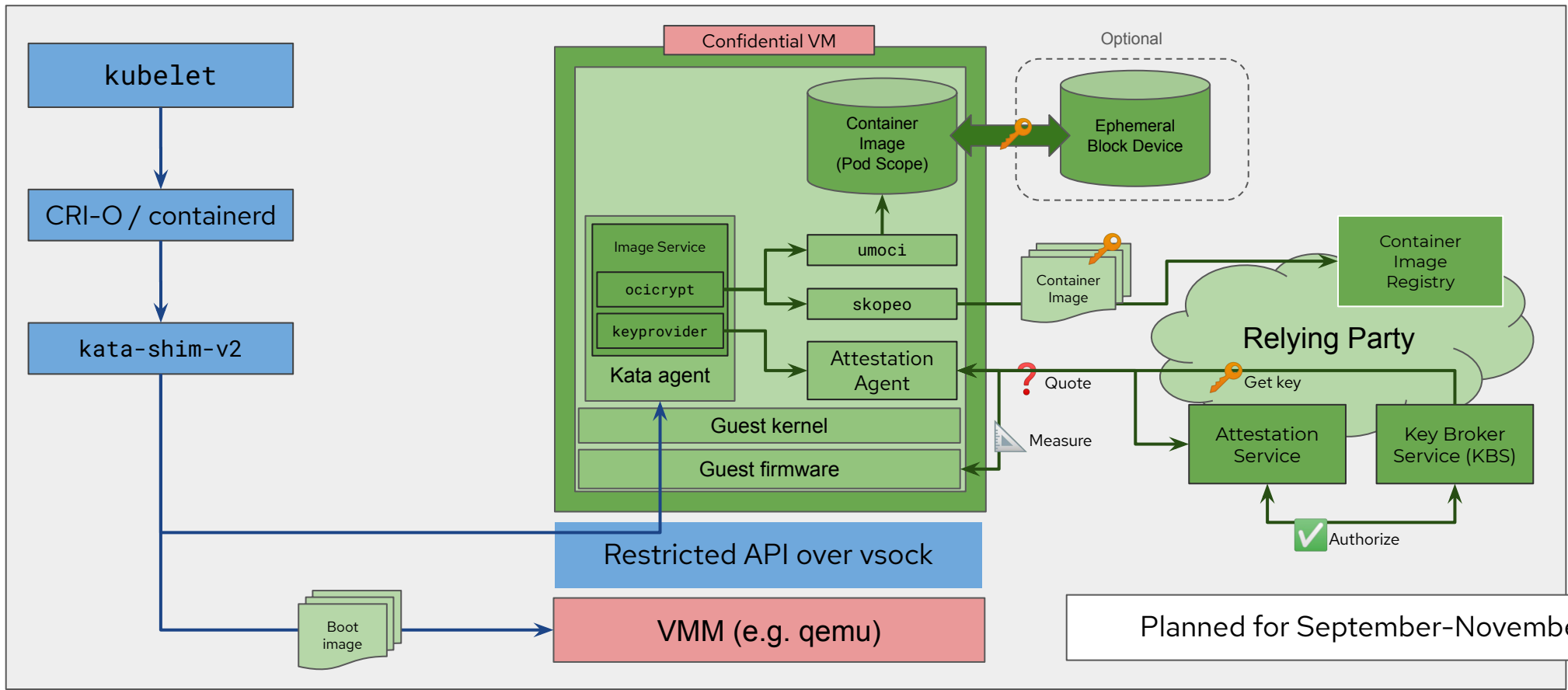
Today, that API does not exist between CRI and `kata-shim-v2`, since container images are currently pulled on the host.
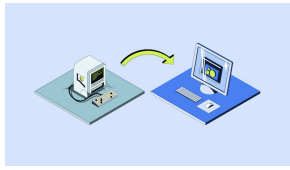
This API situation is relatively typical of the sort of issues we run into for this project.

Also, for initial prototyping, the key has to be pulled out of some magic hat.

28

# Attestation process



Confidential VM

Optional

Container Image (Pod Scope)

Ephemeral Block Device

kubelet

CRI-O / containerd

kata-shim-v2

Image Service

ocicrypt

keyprovider

Kata agent

umoci

skopeo

Attestation Agent

Guest kernel

Guest firmware

Container Image

Container Image Registry

Relying Party

Quote

Measure

Get key

Attestation Service

Key Broker Service (KBS)

Authorize

Restricted API over vsock

Boot image
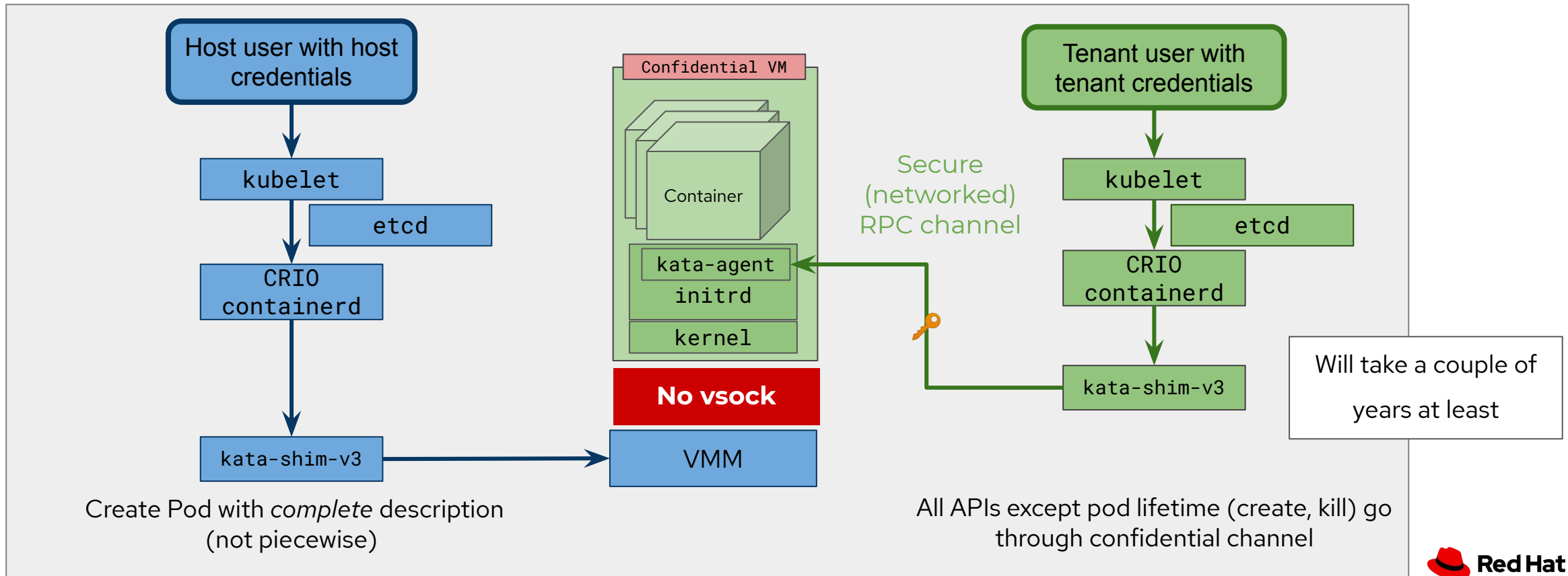
VMM (e.g. qemu)

Planned for September-November 2021

# How do you configure your virtual machine?

▶ **Hot Plugging** is currently used to add memory, CPU or devices to the pod

- The Pod APIs do not give us the information about container sizes

- Resources are dynamically added at container creation time

- This adds a lot of complexity to the runtime, and inefficiency (e.g. fat page tables)

▶ **Integrity** is hard to guarantee if you can change the configuration at runtime

- Memory hot-plugging or ballooning mechanisms conflict with encryption / validation

- Devices, notably pass-through PCI devices with DMA, are also problematic

▶ **Immutable Pods** are fully defined ahead of time, before booting the virtual machine

- This requires many changes in the existing Kubernetes APIs

- Existing APIs may put things "in the wrong place", e.g. send logs to the *host*.

- This will simplify and optimize the non-confidential case, e.g. remove hot-plugging

# The need for a shadow control plane

▸ **Tenants** need their own *isolated* administrative realm (logs, container metrics, …)

▸ **Hosts** manage physical resources (pod creation/destruction, raw disks, H/W metrics…)



Host user with host credentials

kubelet

etcd

CRIO containerd

kata-shim-v3

Create Pod with *complete* description (not piecewise)

Confidential VM

Container

kata-agent
initrd
kernel

**No vsock**

VMM

Secure (networked) RPC channel

Tenant user with tenant credentials

kubelet

etcd

CRIO containerd

kata-shim-v3

Will take a couple of years at least

All APIs except pod lifetime (create, kill) go through confidential channel

31

# Thank you

Red Hat is the world's leading provider of enterprise open source software solutions. Award-winning support, training, and consulting services make Red Hat a trusted adviser to the Fortune 500.

linkedin.com/company/red-hat

youtube.com/user/RedHatVideos

facebook.com/redhatinc

twitter.com/RedHat

**Red Hat**