



OpenShift Virtualization

Technical Deep Dive

Alfred Bach
PSA Red Hat
Field Partner and Learning Team

09:00 Welcome and Intro

09:15 OpenShift Basic installation on Bare Metal Cluster Layout discussion Compact, Single Node, Regular, HCP

10:30 Installing the OpenShift Virtualisation Operator and configuration

11:00 BREAK

11:15 Setting up Local storage (Single Node and OpenShift Data Foundation) Networking options, installing NMStat and SR-IOV

12:00 Lunch Break

13:00 Backup and Disaster Recovery (OADP and ODF)

13:30 Migration Options (MTV and Ansible Migration)

14:30 Migration Risks

15:00 Q&A

Install OpenShift Container Platform (OCP)

OpenShift 4.16 Supported Providers

Installation Experiences



Outposts



AWS Local Zones



Azure Stack Hub



Alibaba Cloud



Bare Metal



Google Cloud



IBM Cloud



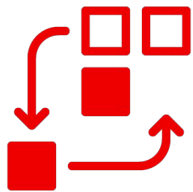
IBM Power Systems

IBM Z

NUTANIX

RED HAT
OPENSTACK
PLATFORM

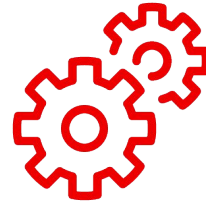
RED HAT
VIRTUALIZATION



Full Stack Automation

Installer Provisioned Infrastructure

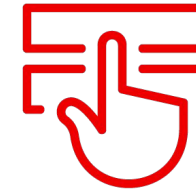
- Auto-provisions infrastructure
- *KS like
- Enables self-service



Pre-existing Infrastructure

User Provisioned Infrastructure

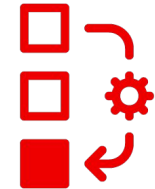
- Bring your own hosts
- You choose infrastructure automation
- Full flexibility
- Integrate ISV solutions



Interactive - Connected

Assisted Installer

- Hosted web-based guided experience
- Agnostic, bare metal, and vSphere only
- ISO Driven



Local - Disconnected

Agent-based Installer

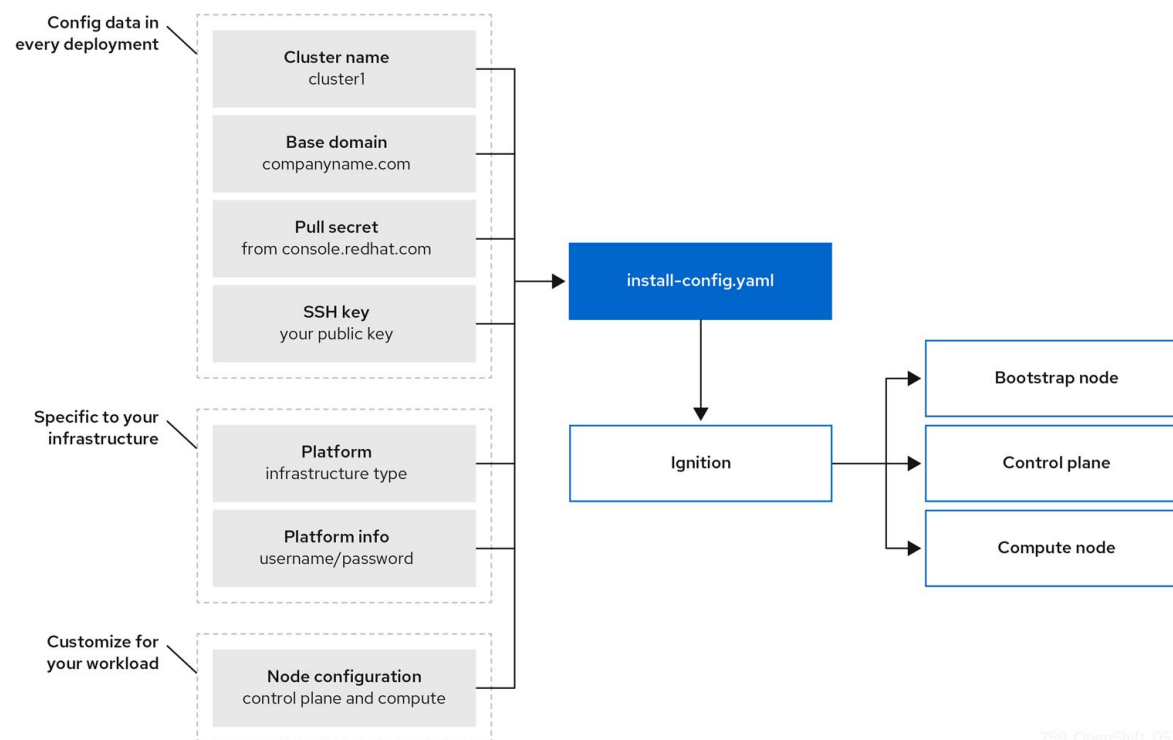
- Disconnected bare metal deployments
- Automated installations via CLI
- ISO driven



OpenShift Bootstrap Process: Self-Managed Kubernetes

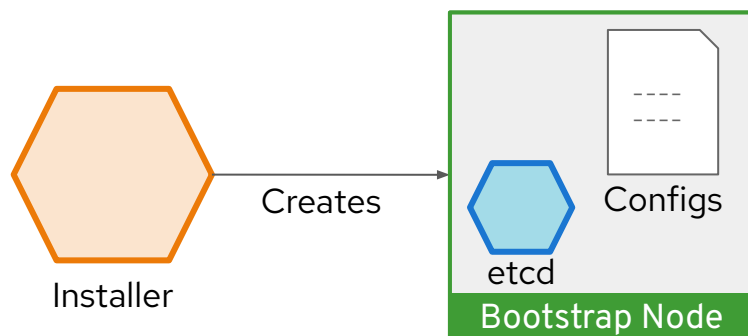
How to boot a self-managed cluster:

- OpenShift 4 is unique in that management extends all the way down to the operating system
- Every machine boots with a configuration that references resources hosted in the cluster it joins enabling cluster to manage itself
- Downside is that every machine looking to join the cluster is waiting on the cluster to be created
- Dependency loop is broken using a bootstrap machine, which acts as a temporary control plane whose sole purpose is bringing up the permanent control plane nodes
- Permanent control plane nodes get booted and join the cluster leveraging the control plane on the bootstrap machine
- Once the pivot to the permanent control plane takes place, the remaining worker nodes can be booted and join the cluster



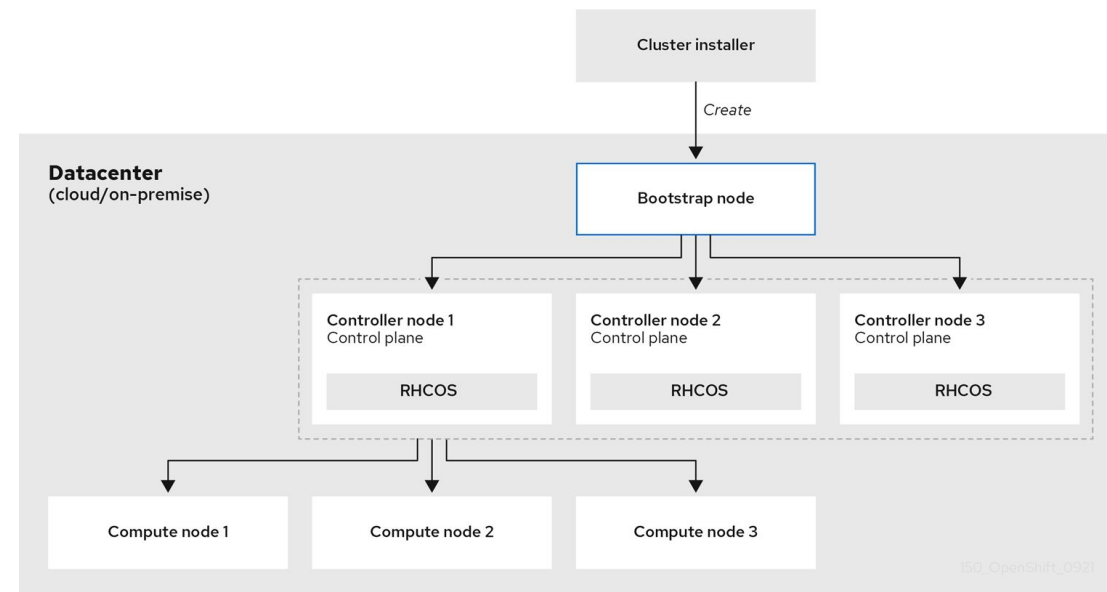
258_OpenShift_0523

OpenShift Bootstrap Process: Step by Step

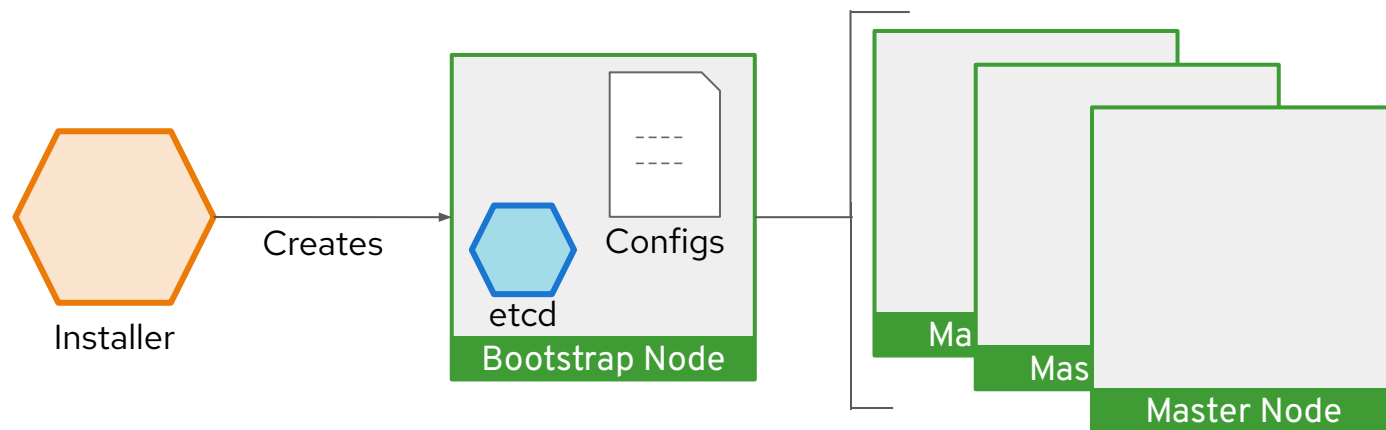


Bootstrapping process step by step:

1. Bootstrap machine boots and starts hosting the remote resources required for master machines to boot. Runs one instance of etcd

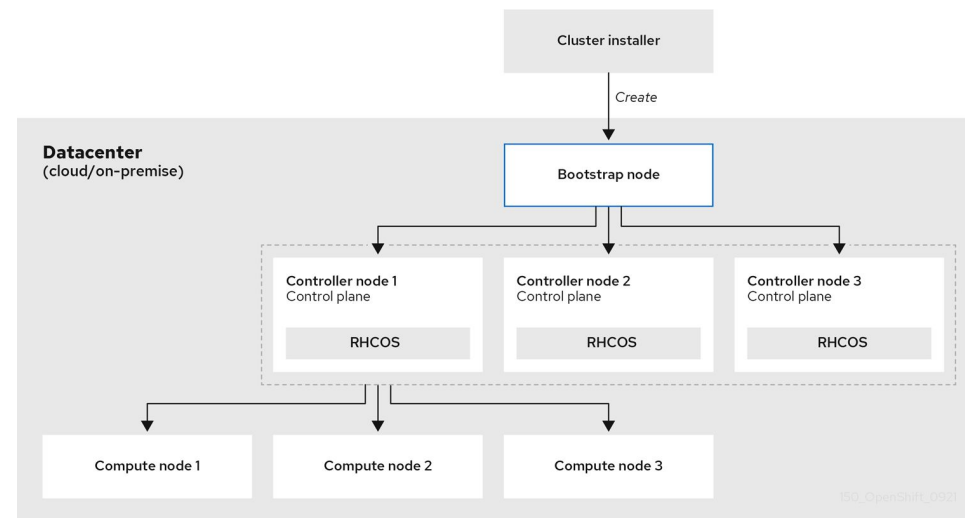


OpenShift Bootstrap Process: Step by Step

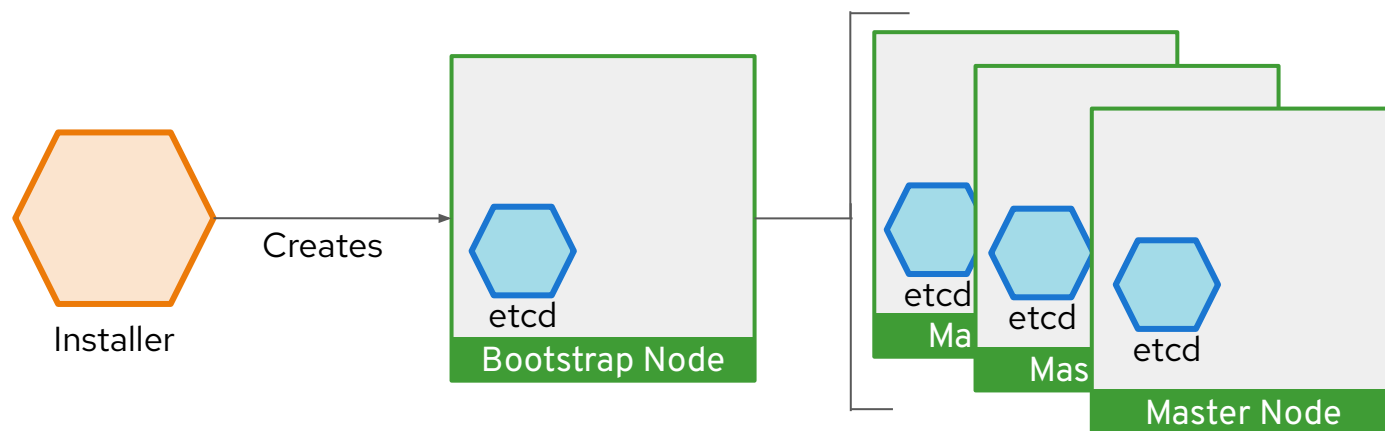


Bootstrapping process step by step:

1. Bootstrap machine boots and starts hosting the remote resources required for master machines to boot. Runs one instance of etcd
2. Master machines fetch the remote resources from the bootstrap machine and finish booting.



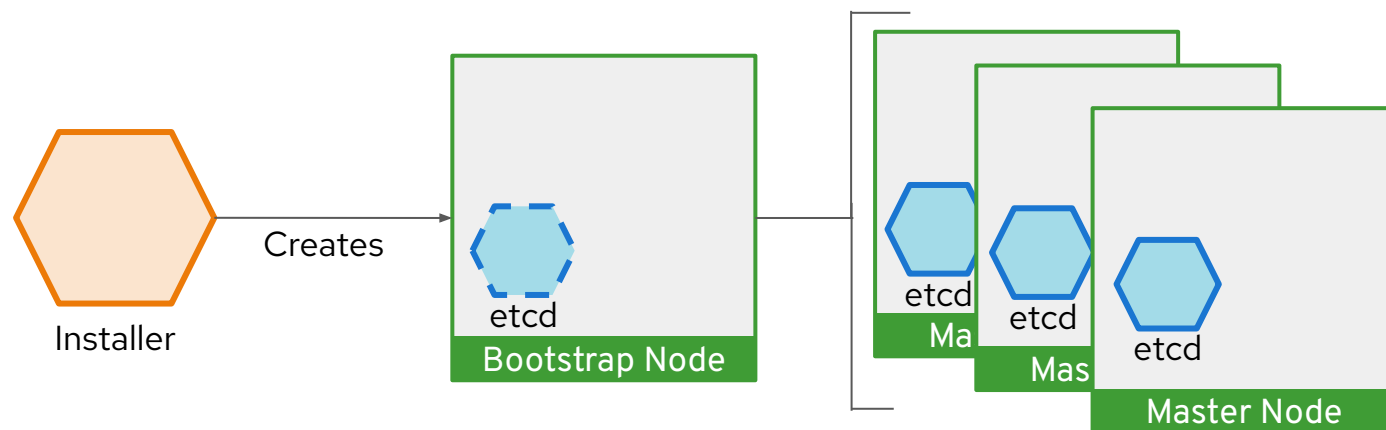
OpenShift Bootstrap Process: Step by Step



Bootstrapping process step by step:

1. Bootstrap machine boots and starts hosting the remote resources required for master machines to boot. Runs one instance of etcd
2. Master machines fetch the remote resources from the bootstrap machine and finish booting.
3. Master machines use the bootstrap node to scale the etcd cluster to 4 total instances.

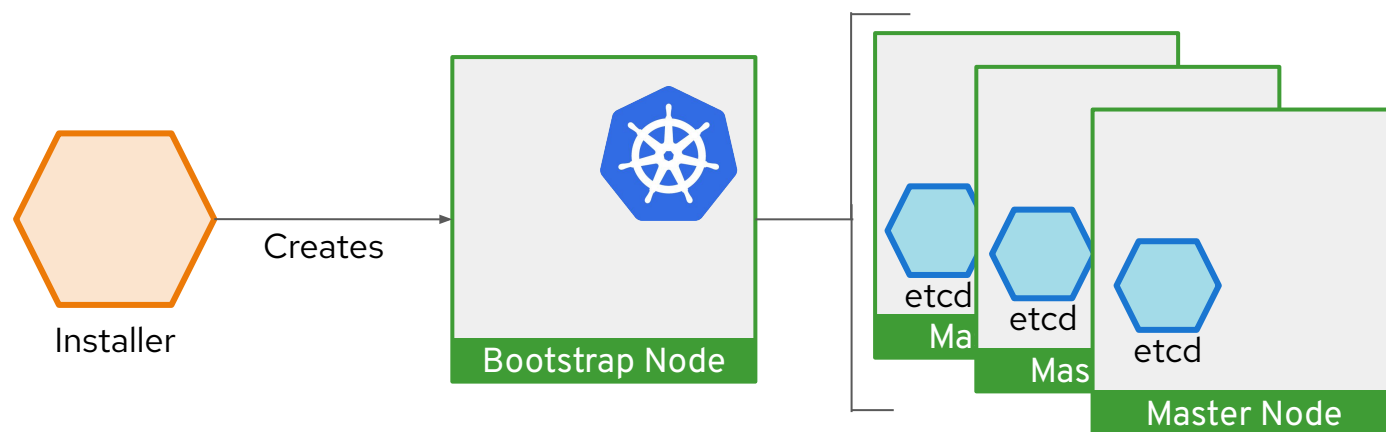
OpenShift Bootstrap Process: Step by Step



Bootstrapping process step by step:

1. Bootstrap machine boots and starts hosting the remote resources required for master machines to boot. Runs one instance of etcd
2. Master machines fetch the remote resources from the bootstrap machine and finish booting.
3. Master machines use the bootstrap node to scale the etcd cluster to 4 total instances.
4. The Etcd operator scales itself down off the bootstrap node, leaving the etcd instance count to 3

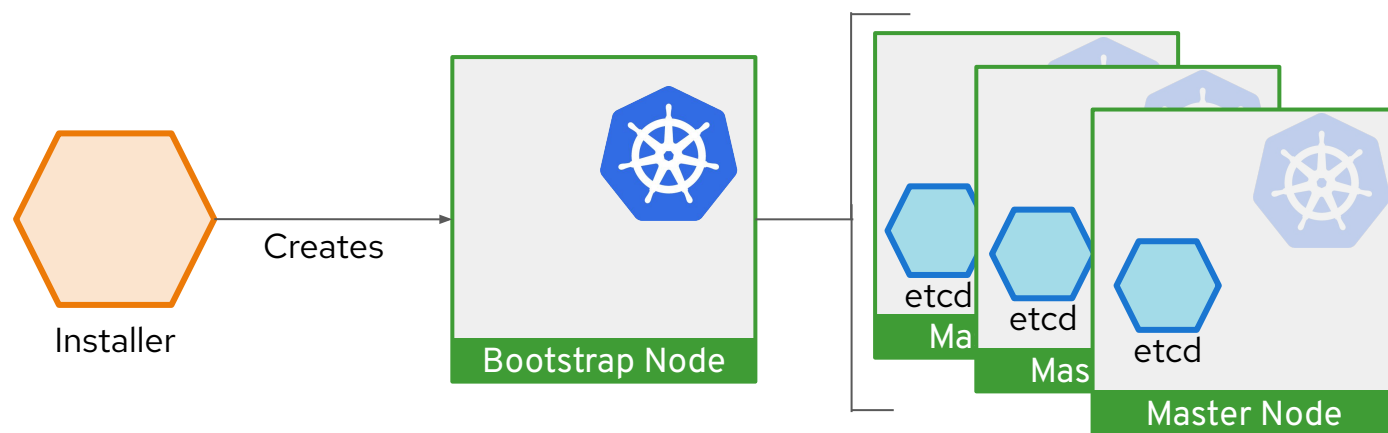
OpenShift Bootstrap Process: Step by Step



Bootstrapping process step by step:

1. Bootstrap machine boots and starts hosting the remote resources required for master machines to boot. Runs one instance of etcd
2. Master machines fetch the remote resources from the bootstrap machine and finish booting.
3. Master machines use the bootstrap node to scale the etcd cluster to 3 instances.
4. The Etcd operator scales itself down off the bootstrap node, then scales back up to 3; all on the Masters
5. Bootstrap node starts a temporary Kubernetes control plane using the newly-created etcd cluster.

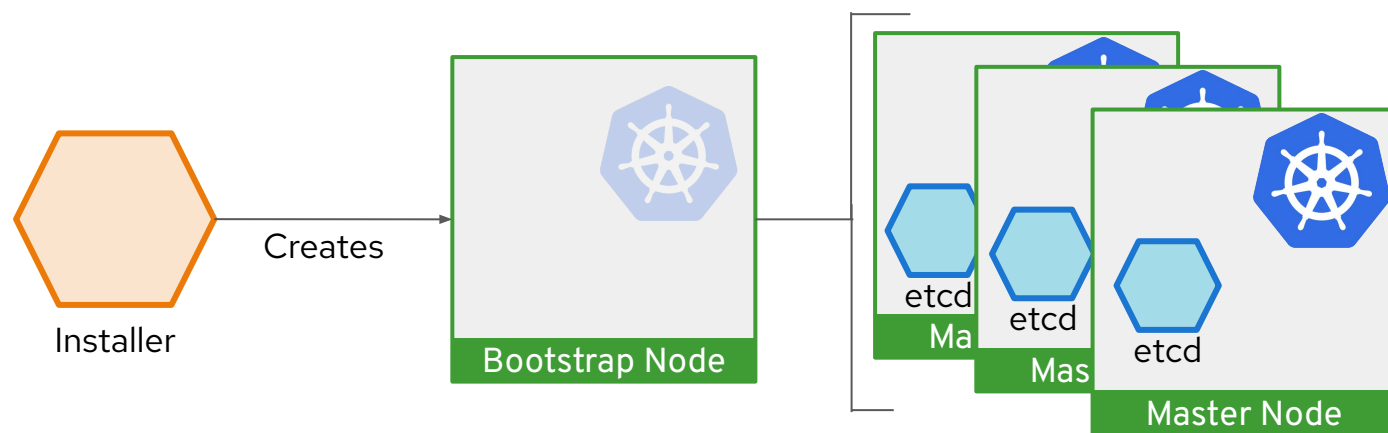
OpenShift Bootstrap Process: Step by Step



Bootstrapping process step by step:

1. Bootstrap machine boots and starts hosting the remote resources required for master machines to boot. Runs one instance of etcd
2. Master machines fetch the remote resources from the bootstrap machine and finish booting.
3. Master machines use the bootstrap node to scale the etcd cluster to 3 instances.
4. The Etcd operator scales itself down off the bootstrap node, then scales back up to 3; all on the Masters
5. Bootstrap node starts a temporary Kubernetes control plane using the newly-created etcd cluster.
6. Temporary control plane schedules the production control plane to the master machines.

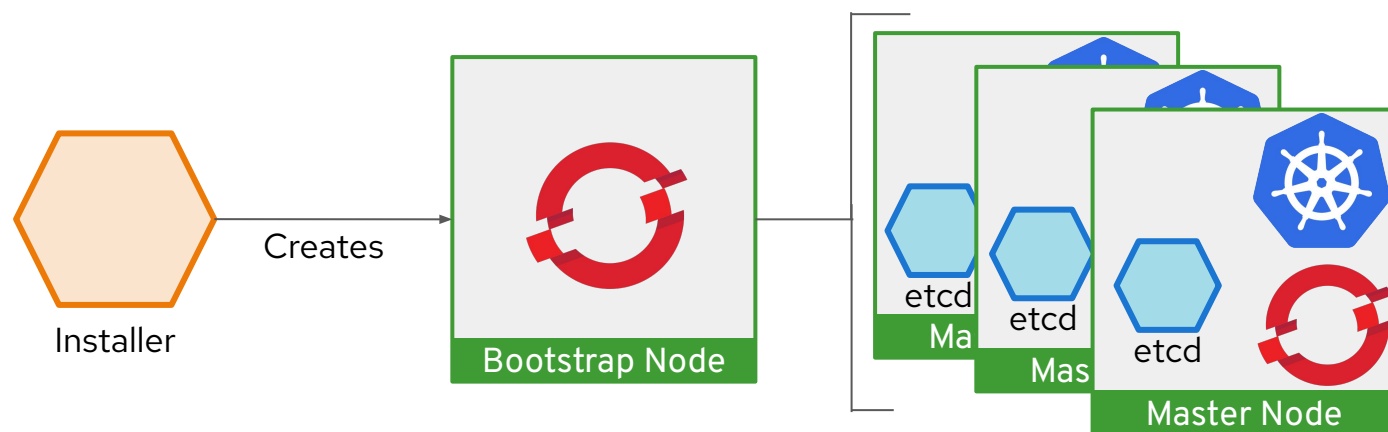
OpenShift Bootstrap Process: Step by Step



Bootstrapping process step by step:

1. Bootstrap machine boots and starts hosting the remote resources required for master machines to boot. Runs one instance of etcd
2. Master machines fetch the remote resources from the bootstrap machine and finish booting.
3. Master machines use the bootstrap node to scale the etcd cluster to 3 instances.
4. The Etcd operator scales itself down off the bootstrap node, then scales back up to 3; all on the Masters
5. Bootstrap node starts a temporary Kubernetes control plane using the newly-created etcd cluster.
6. Temporary control plane schedules the production control plane to the master machines.
7. Temporary control plane shuts down, yielding to the production control plane.

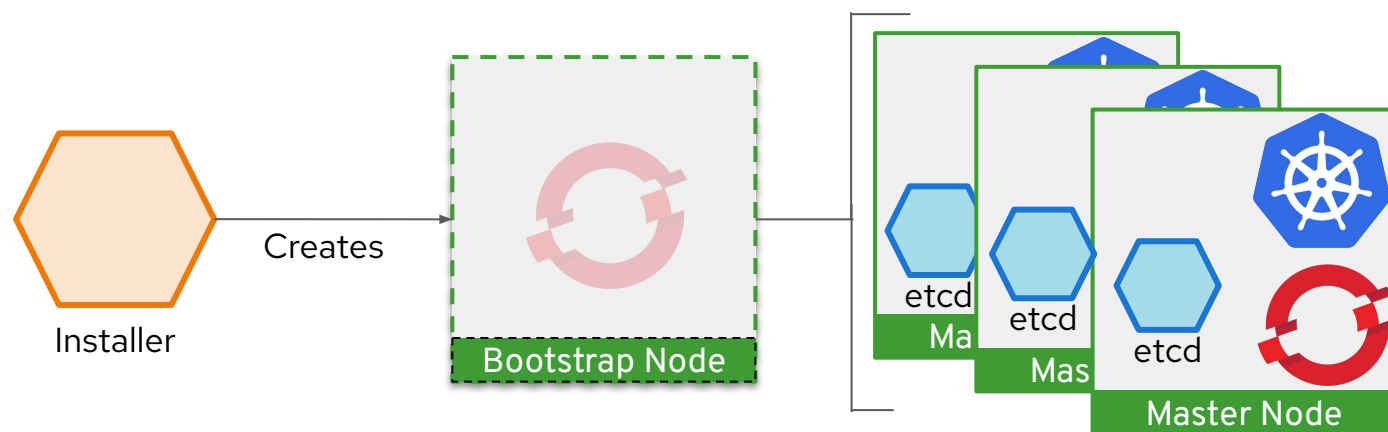
OpenShift Bootstrap Process: Step by Step



Bootstrapping process step by step:

1. Bootstrap machine boots and starts hosting the remote resources required for master machines to boot. Runs one instance of etcd
2. Master machines fetch the remote resources from the bootstrap machine and finish booting.
3. Master machines use the bootstrap node to scale the etcd cluster to 3 instances.
4. The Etcd operator scales itself down off the bootstrap node, then scales back up to 3; all on the Masters
5. Bootstrap node starts a temporary Kubernetes control plane using the newly-created etcd cluster.
6. Temporary control plane schedules the production control plane to the master machines.
7. Temporary control plane shuts down, yielding to the production control plane.
8. Bootstrap node injects OpenShift-specific components into the newly formed control plane.

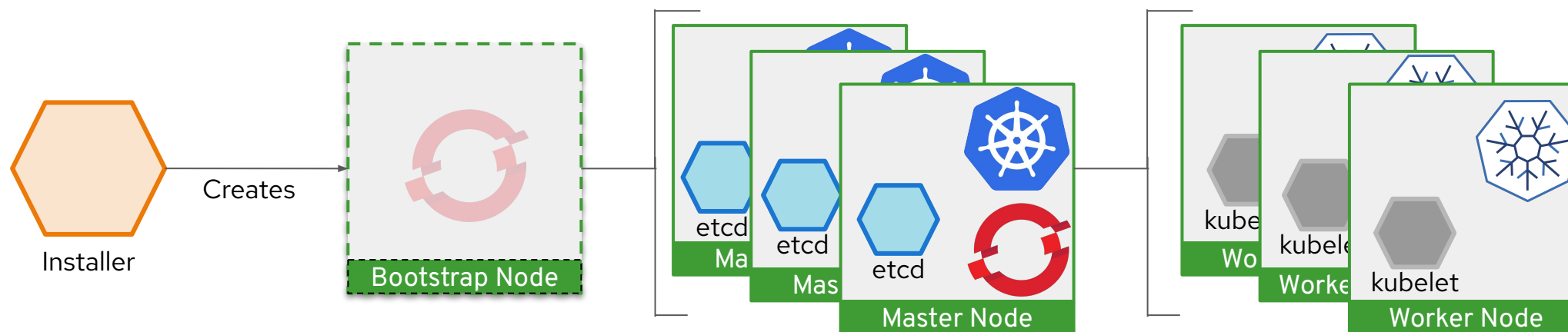
OpenShift Bootstrap Process: Step by Step



Bootstrapping process step by step:

1. Bootstrap machine boots and starts hosting the remote resources required for master machines to boot. Runs one instance of etcd
2. Master machines fetch the remote resources from the bootstrap machine and finish booting.
3. Master machines use the bootstrap node to scale the etcd cluster to 3 instances.
4. The Etcd operator scales itself down off the bootstrap node, then scales back up to 3; all on the Masters
5. Bootstrap node starts a temporary Kubernetes control plane using the newly-created etcd cluster.
6. Temporary control plane schedules the production control plane to the master machines.
7. Temporary control plane shuts down, yielding to the production control plane.
8. Bootstrap node injects OpenShift-specific components into the newly formed control plane.
9. Installer then tears down the bootstrap node or if user-provisioned, this needs to be performed by the administrator.

OpenShift Bootstrap Process: Step by Step



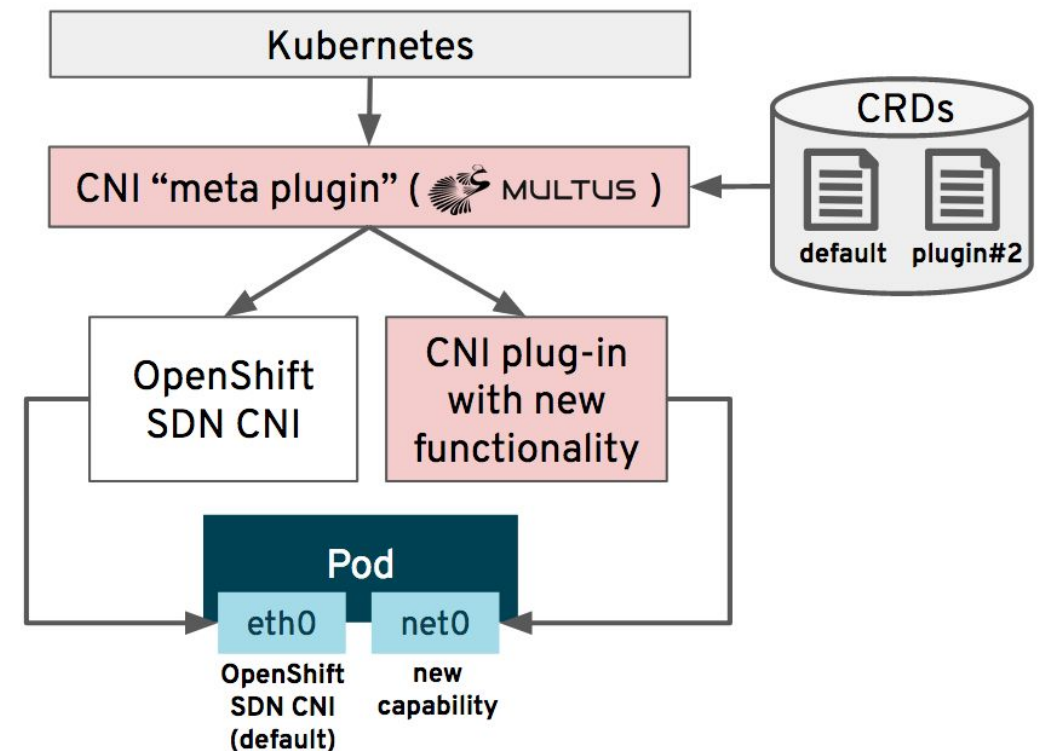
Bootstrapping process step by step:

1. Bootstrap machine boots and starts hosting the remote resources required for master machines to boot. Runs one instance of etcd
2. Master machines fetch the remote resources from the bootstrap machine and finish booting.
3. Master machines use the bootstrap node to scale the etcd cluster to 3 instances.
4. The Etcd operator scales itself down off the bootstrap node, then scales back up to 3; all on the Masters
5. Bootstrap node starts a temporary Kubernetes control plane using the newly-created etcd cluster.
6. Temporary control plane schedules the production control plane to the master machines.
7. Temporary control plane shuts down, yielding to the production control plane.
8. Bootstrap node injects OpenShift-specific components into the newly formed control plane.
9. Installer then tears down the bootstrap node or if user-provisioned, this needs to be performed by the administrator.
10. Worker machines fetch remote resources from masters and finish booting.

Network

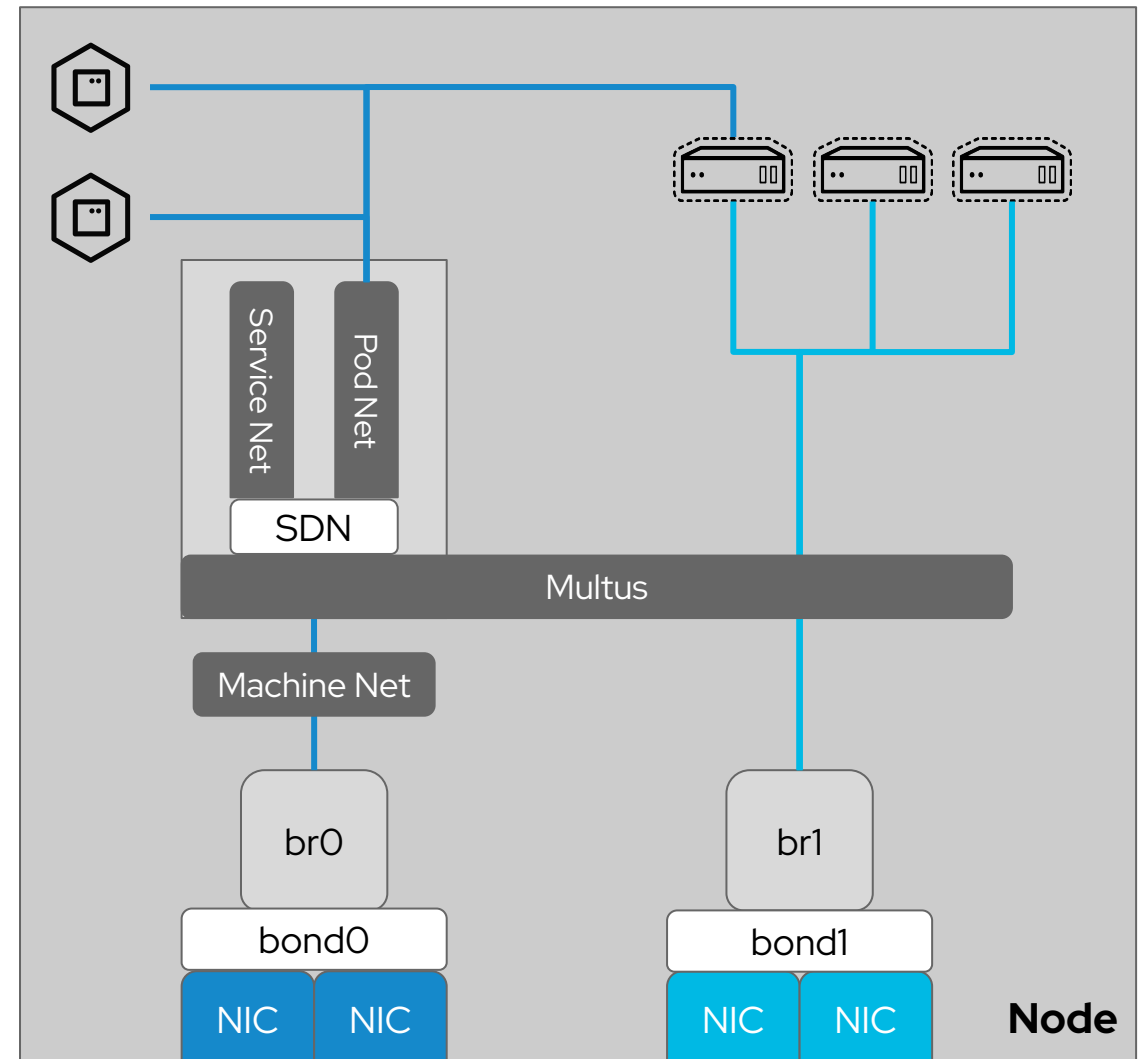
Virtual Machine Networking

- Virtual machines optionally connect to the standard pod network
 - OpenShift SDN, OVNKubernetes
 - Partners, such as Calico, are also supported
- Additional network interfaces accessible via Multus:
 - Bridge, SR-IOV, OVN secondary networks
 - VLAN and other networks can be created at the host level using nmstate
- When using at least one interface on the default SDN, Service, Route, and Ingress configuration applies to VM pods the same as others



- Pod, service, and machine network are configured by OpenShift automatically
 - Use kernel parameters (dracut) for configuration at install - **bond0** in the example to the right
- Use the NMstate Operator to configure additional host network interfaces
 - **bond1** and **br1** in the example to the right
- VMs and Pods connect to one or more networks simultaneously

The following slides show an example of how this setup is configured



GUI-based host network configuration

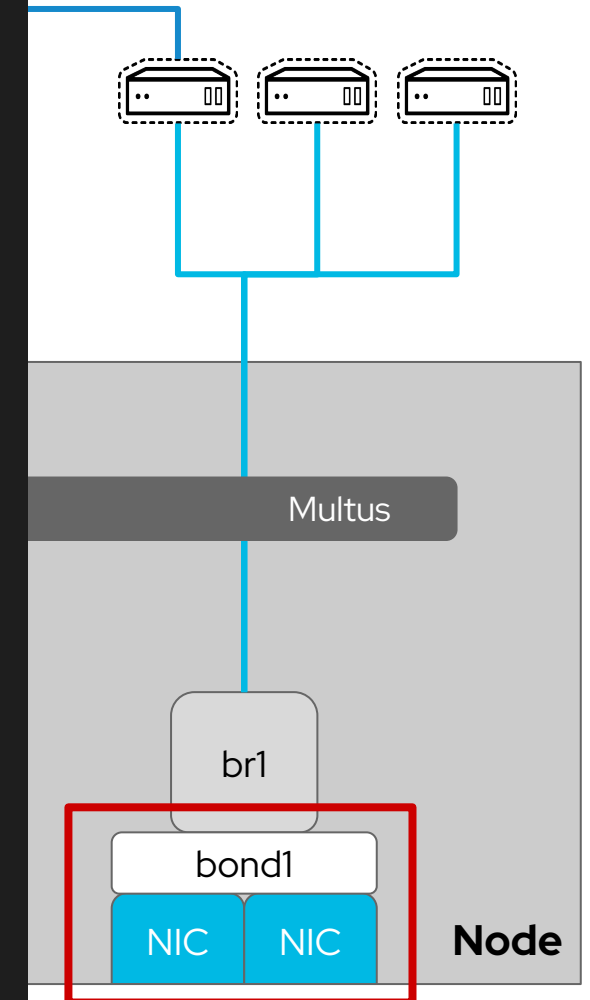
- Apply NMstate configuration using a form in the OpenShift admin console
- Create and configure
 - Ethernet interface IP (static, DHCP)
 - Bonds - mode 1-6 bonds with options, including IP configuration
 - Linux bridge configuration utilizing ethernet and/or bonds for “uplinks”
- Specify node selectors to have configuration automatically applied to matching nodes

The screenshot displays the Red Hat OpenShift Admin Console interface. On the left is a sidebar menu with categories like Administrator, Home, Operators, Workloads, Virtualization, Migration, Networking, Services, Routes, Ingresses, NetworkPolicies, NetworkAttachmentDefinitions, NodeNetworkConfigurationPolicy (selected), NodeNetworkState, Storage, Builds, Observe, Compute, User Management, and Administration. The main panel is titled 'Create NodeNetworkConfigurationPolicy' with an 'Edit YAML' link. It contains a description of NMstate configuration, a checkbox for applying the policy to specific node subsets (checked), a 'Policy name' field with 'worker-bond0-br1', a 'Description' field, and a 'Policy Interface(s)' section. This section lists 'Bonding bond0' and 'Bridge br1'. The 'Bridge br1' configuration includes fields for 'Interface name' (br1), 'Network state' (Up), 'Type' (Bridge), 'IP configuration' (IPv4), 'Port' (bond0), and an 'Enable STP' checkbox. A 'Create' button is at the bottom.

Host bond configuration

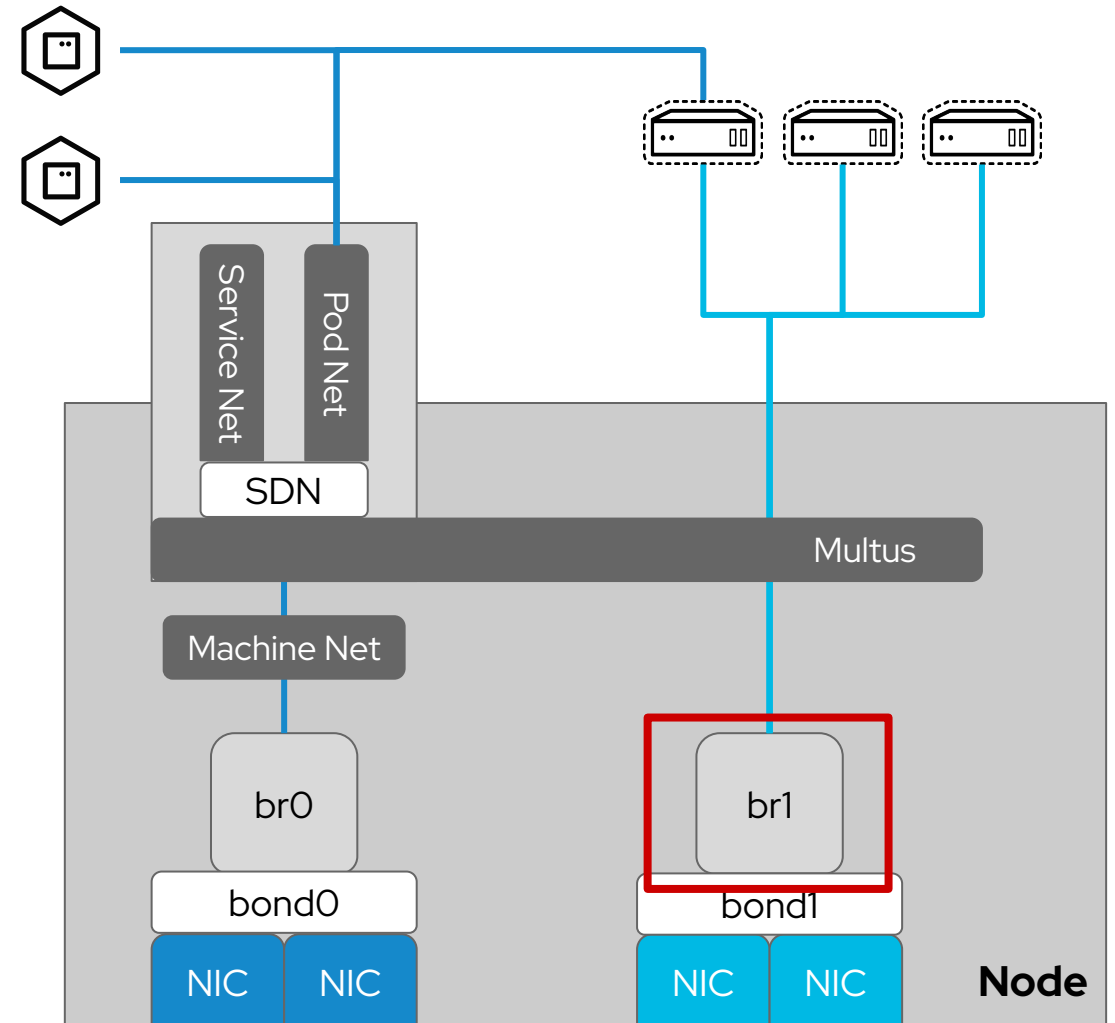
- NodeNetworkConfigurationPolicy (NNCP)
 - Nmstate operator CRD
 - Configure host network using declarative language
- Applies to all nodes specified in the nodeSelector, including newly added nodes automatically
- Update or add new NNCPs for additional host configs

```
1  apiVersion: nmstate.io/v1alpha1
2  kind: NodeNetworkConfigurationPolicy
3  metadata:
4    name: worker-bond1
5  spec:
6    nodeSelector:
7      node-role.kubernetes.io/worker: ""
8    desiredState:
9      interfaces:
10       - name: bond1
11         type: bond
12         state: up
13         ipv4:
14           enabled: false
15         link-aggregation:
16           mode: balance-alb
17           options:
18             miimon: '100'
19           slaves:
20             - eth2
21             - eth3
22           mtu: 1450
```



Host bridge configuration

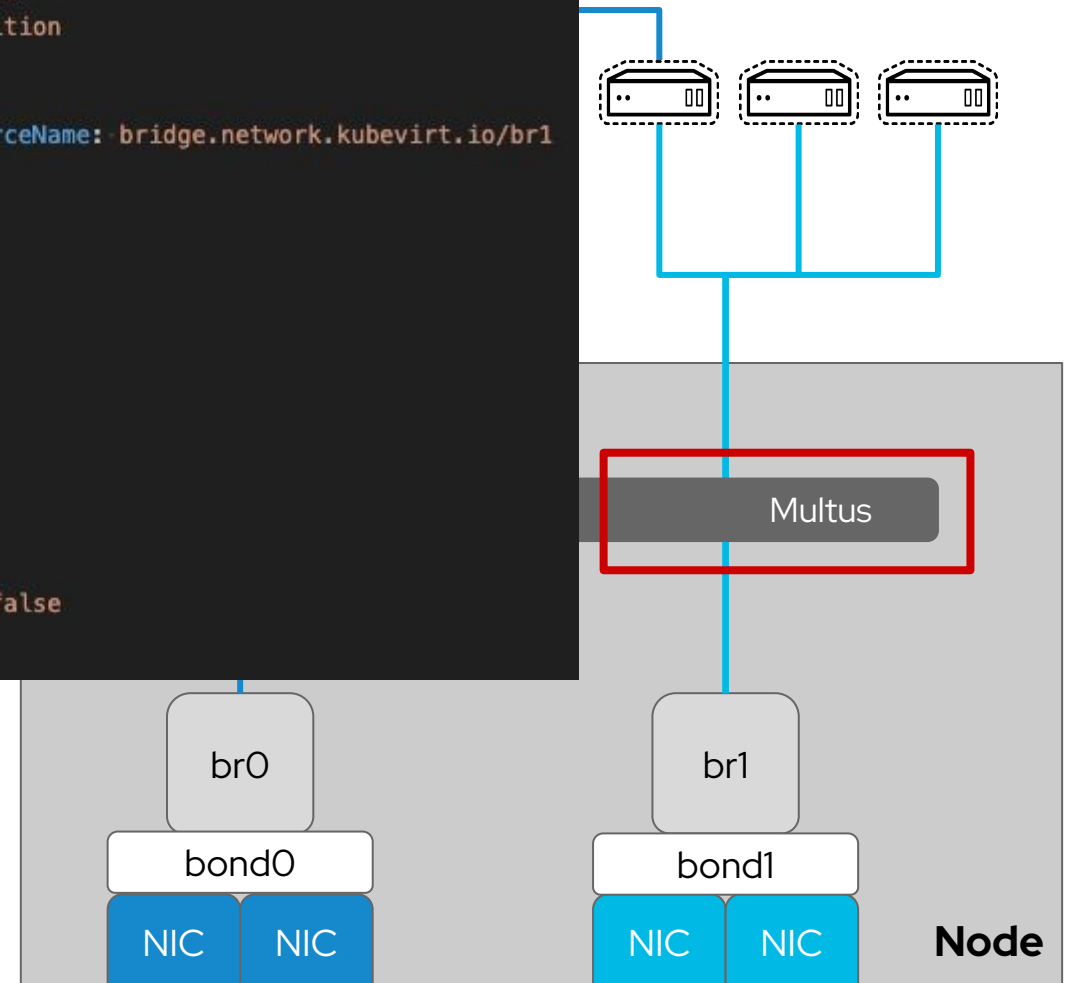
```
1  apiVersion: nmstate.io/v1alpha1
2  kind: NodeNetworkConfigurationPolicy
3  metadata:
4    name: worker-bond1-br1
5  spec:
6    nodeSelector:
7      node-role.kubernetes.io/worker: ""
8    desiredState:
9      interfaces:
10       - name: br1
11         description: br1 with bond1
12         type: linux-bridge
13         state: up
14         ipv4:
15           enabled: false
16         bridge:
17           options:
18             stp:
19               enabled: false
20           port:
21             - name: bond1
```



Network Attachment Definition configuration

- net-attach-def configures multus to allow the VM to access an underlying resource
 - Optionally define VLAN tags
- Limited to the namespace it's created in
 - Except the default namespace, which is available to all

```
1  apiVersion: k8s.cni.cncf.io/v1
2  kind: NetworkAttachmentDefinition
3  metadata:
4    annotations:
5      k8s.v1.cni.cncf.io/resourceName: bridge.network.kubevirt.io/br1
6    name: vlan-93
7    namespace: default
8  spec:
9    config: >-
10     -{
11       "name": "vlan-93"
12       "type": "cnv-bridge",
13       "cniVersion": "0.3.1",
14       "bridge": "br1",
15       "vlan": 93,
16       "macspoofchk": true,
17       "ipam": {},
18       "preserveDefaultVlan": false
19     }
20
```



Host network configuration status

- Use the admin console to view the NodeNetworkState (OpenShift 4.14+)
- Detailed configuration information for host networking including
 - IP and MAC addresses
 - Bond configuration
 - Bridge configuration
- Review and troubleshoot host network configuration

The screenshot displays the Red Hat OpenShift Admin Console interface. The left-hand navigation menu is expanded, and 'NodeNetworkState' is selected, indicated by a red '1'. The main panel shows the 'NodeNetworkState' for a specific node. At the top, there's a search bar and a filter dropdown. Below this, a table lists network interfaces. The 'ethernet' section is expanded, showing a list of interfaces with their IP addresses and MAC addresses. The 'linux-bridge' section is also expanded, showing a list of interfaces. Red callouts '2' and '3' highlight the 'ethernet' and 'linux-bridge' sections respectively.

Name	IP address	Ports	MAC address
ethernet			
enp1s0 ↑	172.22.0.71/24	-	DE:AD:BE:EF:00:04
enp2s0 ↑	192.168.123.104/24	-	52:54:00:00:00:04
enp3s0 ↑	-	-	52:54:00:00:01:04
linux-bridge			
br-flat ↑	-	1	52:54:00:00:01:04

Connecting Pods to networks

- Multus uses CNI network definitions in the NetworkAttachmentDefinition to allow access
 - `net-attach-def` are namespaced
 - Pods cannot connect to a `net-attach-def` in a different namespace
- `cnv-bridge` and `cnv-tuning` types are used to enable VM specific functions
 - MAC address customization
 - MTU and promiscuous mode
 - `sysctls`, if needed
- Pod connections are defined using an annotation
 - Pods can have many connections to many networks

```
1  apiVersion: k8s.cni.cncf.io/v1
2  kind: NetworkAttachmentDefinition
3  metadata:
4    name: br1-public
5    annotations:
6      k8s.v1.cni.cncf.io/resourceName: bridge.network.kubevirt.io/br1
7  spec:
8    config: '{
9      "cniVersion": "0.3.1",
10     "name": "br1-public",
11     "plugins": [
12       {
13         "type": "cnv-bridge",
14         "bridge": "br1"
15       },
16       {
17         "type": "cnv-tuning"
18       }
19     ]
20   }'
```

```
1  kind: Pod
2  apiVersion: v1
3  metadata:
4    name: application-pod
5    annotations:
6      k8s.v1.cni.cncf.io/networks: bond1-br1
```


Connecting VMs to networks

- Virtual machine interfaces describe NICs attached to the VM
 - `spec.domain.devices.interfaces`
 - Model: virtio, e1000, pcnet, rtl8139, etc.
 - Type: masquerade, bridge
 - MAC address: customize the MAC
- The networks definition describes the connection type
 - `spec.networks`
 - Pod = default SDN
 - Multus = secondary network using Multus
- ***Using the GUI makes this easier*** and removes the need to edit / manage connections in YAML

```
1  apiVersion: kubevirt.io/v1alpha3
2  kind: VirtualMachine
3    name: demo-vm
4  spec:
5    template:
6      spec:
7        domain:
8          devices:
9            interfaces:
10              - bridge: {}
11                model: virtio
12                name: nic-0
13          hostname: demo-vm
14          networks:
15            - multus:
16              networkName: bond1-br1
17              name: nic-0
```

Storage

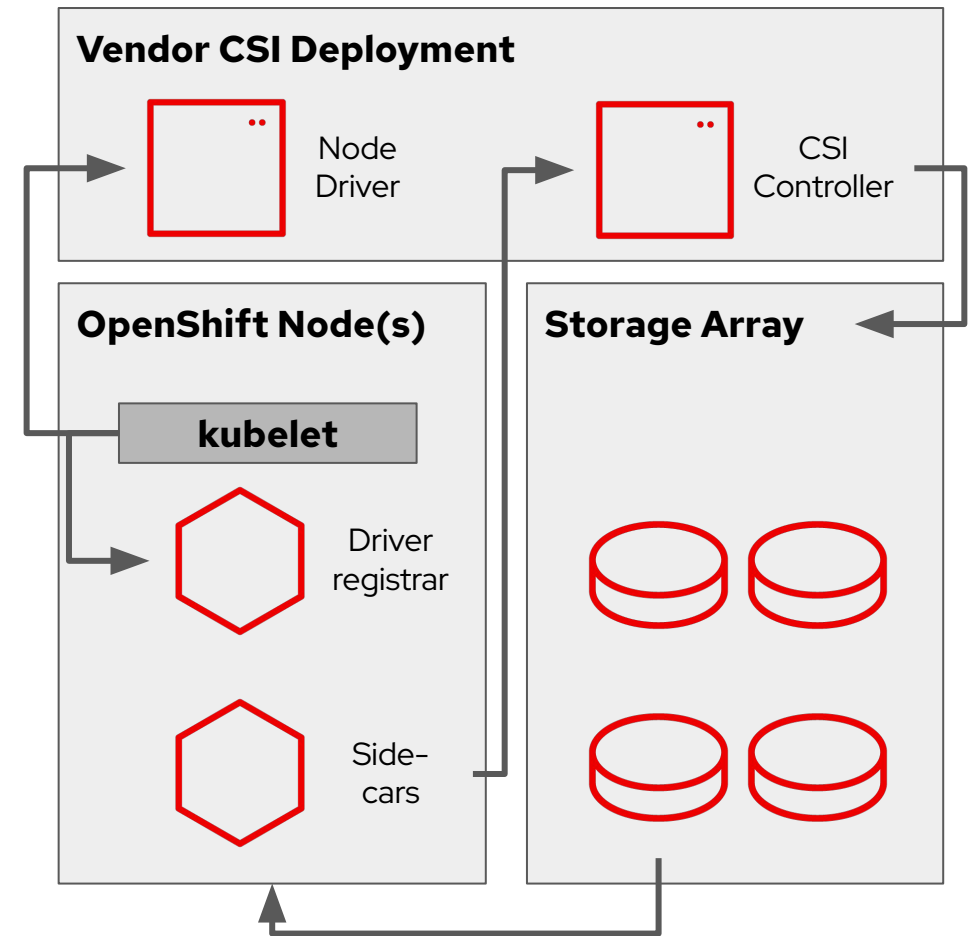
Virtual Machine Storage

- OpenShift Virtualization uses the Kubernetes PersistentVolume (PV) paradigm
- PVs can be backed by
 - CSI drivers, including partners and ODF
 - Local storage using host path provisioner / logical volume Operator
- Use dynamically or statically provisioned PVs
- RWX is **required** for live migration
- Disks are attached using VirtIO or SCSI controllers
 - Connection order specified in the VM definition
- Boot order customized via VM definition

PersistentVolumeClaim Details	
Name rhel-rootdisk	Status ✔ Bound
Namespace NS default	Capacity 20Gi
Labels app=containerized-data-importer	Access Modes ReadWriteMany
Annotations 12 Annotations ✎	Volume Mode Filesystem
Label Selector No selector	Storage Class SC managed-nfs-storage
Created At Jul 8, 4:18 pm	Persistent Volume PV pvc-a1aac411-2e46-495a-897e-cf3bc2442199
Owner DV rhel-rootdisk	

VM disks in PVCs

- VM disks on FileSystem mode PVCs are created as thin provisioned raw images by default, but can be configured
- Block mode PVCs are attached directly to the VM
 - Many partners support RWX with block mode PVCs
- CSI features, e.g. snapshot and clone, offload operations to the storage device
 - Use DataVolumes to clone VM disks to automatically select the optimal method to clone the disk
 - Use VM details interface for (powered off) VM snaps
- Resize VM disks using standard PVC methods
- Hot add is not supported



DataVolumes

- VM disks can be imported from multiple sources using DataVolumes, e.g. an HTTP(S) or S3 URL for a QCOW2 or raw disk image, optionally compressed
- VM disks are cloned / copied from existing PVCs
- DataVolumes are created as distinct objects or as a part of the VM definition as a `dataVolumeTemplate`
- DataVolumes use the `ContainerizedDataImporter` to connect, download, and prepare the disk image
- DataVolumes create PVCs based on defaults defined in the profile
 - Access mode, snapshot method

```
1  dataVolumeTemplates:
2    - apiVersion: cdi.kubevirt.io/v1alpha1
3      kind: DataVolume
4      metadata:
5        creationTimestamp: null
6        name: vm-rootdisk
7      spec:
8        pvc:
9          accessModes:
10           - ReadWriteMany
11          resources:
12            requests:
13              storage: 20Gi
14          storageClassName: my-storage-class
15          volumeMode: Filesystem
16        source:
17          http:
18            url: 'http://web.server/disk-image.qcow2'
```

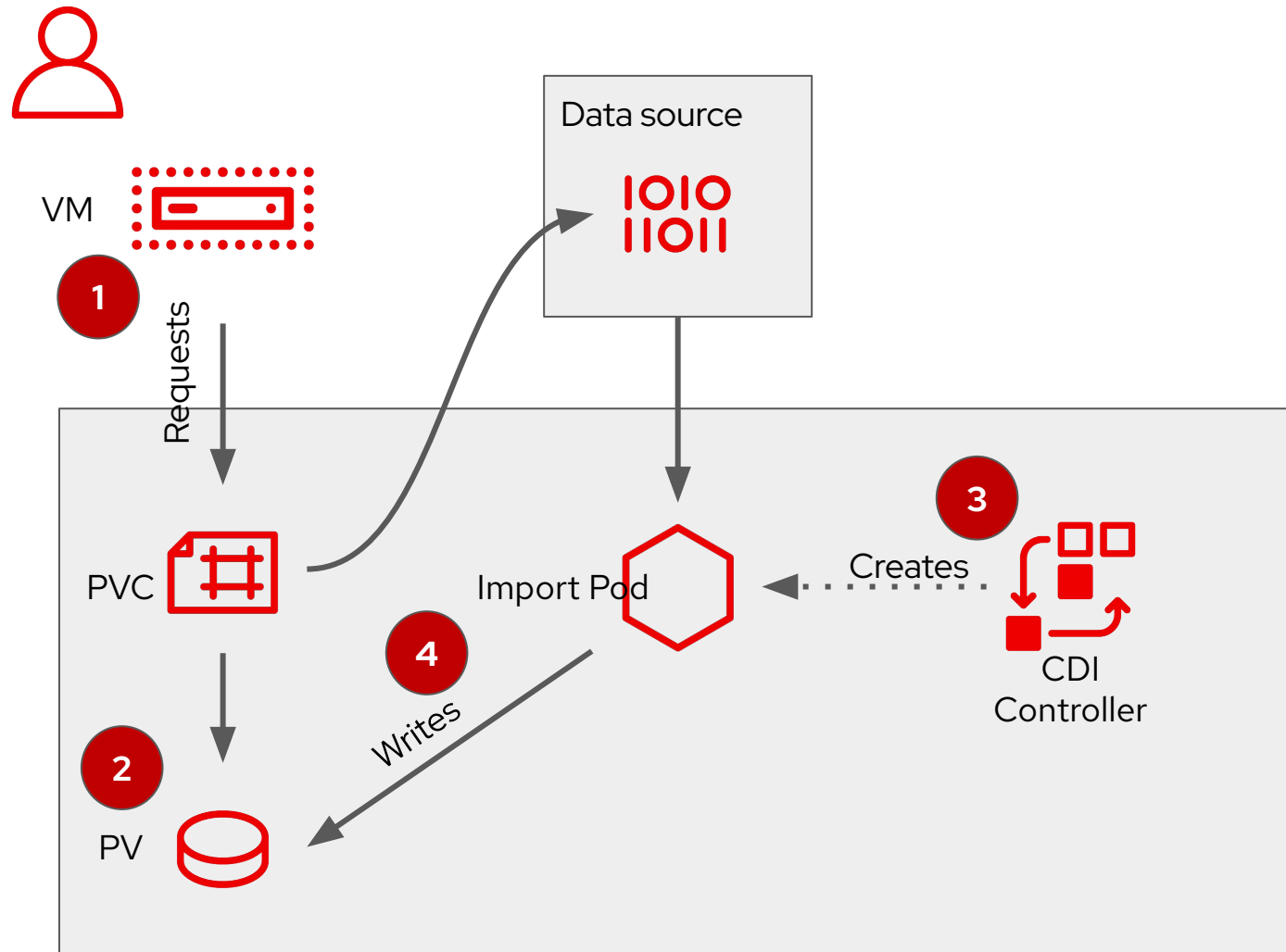
Storage Profiles

- Provide default settings and properties for StorageClasses used by DataVolumes
- Created automatically for every StorageClass
- Preconfigured values for many storage providers, administrator can modify and customize
- DataVolume definitions only need to specify StorageClass, without knowledge of underlying details
 - spec.storage doesn't require fields other than the size and StorageClass

```
1  apiVersion: cdi.kubevirt.io/v1beta1
2  kind: StorageProfile
3  metadata:
4    name: hostpath-provisioner
5  spec:
6    claimPropertySets:
7      - accessModes:
8          - ReadWriteOnce
9        volumeMode:
10         Filesystem
```

```
1  apiVersion: cdi.kubevirt.io/v1alpha1
2  kind: DataVolume
3  metadata:
4    name: rhel8-vm-disk
5  spec:
6    storage:
7      resources:
8        requests:
9          storage: 30Gi
10     storageClassName: hostpath-provisioner
11    source:
12      pvc:
13        name: rhel8-cloud-image
14        namespace: os-images
```

Containerized Data Importer



1. The user creates a virtual machine with a DataVolume
2. The StorageClass is used to satisfy the PVC request
3. The CDI controller creates an importer pod, which mounts the PVC and retrieves the disk image. The image could be sourced from S3, HTTP, or other accessible locations
4. After completing the import, the import pod is destroyed and the PVC is available for the VM

Ephemeral Virtual Machine Disks

- VMs booted via PXE or using a container image can be “diskless”
 - PVCs may be attached and mounted as secondary devices for application data persistence
- VMs based on container images use the standard copy-on-write graph storage for OS disk R/W
 - Consider and account for capacity and IOPS during RHCOS node sizing if using this type
- An `emptyDisk` may be used to add additional ephemeral capacity for the VM

```
1 spec:
2   domain:
3     disks:
4       - bootOrder: 1
5         disk:
6           bus: virtio
7           name: rootdisk
8   volumes:
9     - containerDisk:
10       image: registry.lab.1an:5000/fedora:31
11       name: rootdisk
```


Helper disks

- OpenShift Virtualization attaches disks to VMs for injecting data
 - Cloud-Init
 - ConfigMap
 - Secrets
 - ServiceAccount
- These disks are read-only and can be mounted by the OS to access the data within

```
1 spec:
2   domain:
3     devices:
4       - disk:
5         bus: virtio
6         name: cloudinitdisk
7     volumes:
8       - cloudInitNoCloud:
9         userData: |-
10            #cloud-config
11            password: redhat
12            chpasswd: { expire: False }
13         name: cloudinitdisk
```

Name ↑	Source ↑	Size ↑	Interface ↑	Storage Class ↑	
cloudinitdisk	Other	-	VirtIO	-	⋮

Comparing with traditional virtualization platforms

Live Migration

- Live migration moves a virtual machine from one node to another in the OpenShift cluster
- Can be triggered via GUI, CLI, API, or automatically
- **RWX storage is required**
- Live migration is cancellable by deleting the API object
- Default maximum of five (5) simultaneous live migrations
 - Maximum of two (2) outbound migrations per node, 64MiB/s throughput each
- Uses the SDN by default, customizable to a dedicated network after install

Migration Reason	vSphere	OpenShift Virtualization
Resource contention	DRS	Pod eviction policy, pod descheduler
Node maintenance	Maintenance mode	Maintenance mode, node drain

Automated live migration

- OpenShift / Kubernetes triggers Pod rebalance actions based on multiple factors
 - Soft / hard eviction policies
 - Pod descheduler
 - Pod disruption policy
 - Node resource contention resulting in evictions
 - Pods are `Burstable` QoS class by default
 - All memory is requested in Pod definition, only CPU overhead is requested
- Pod rebalance applies to VM pods equally
- VMs will behave according to the eviction strategy
 - `LiveMigrate` - use live migration to move the VM to a different node
 - No definition - terminate the VM if the node is drained or Pod evicted

VM scheduling

- VM scheduling follows pod scheduling rules
 - Node selectors
 - Taints / tolerations
 - Pod and node affinity / anti-affinity
- Kubernetes scheduler takes into account many additional factors
 - Resource load balancing - requests and reservations
 - Large / Huge page support for VM memory
 - Use scheduler profiles to provide additional hints (for all Pods)
- Resources are managed by Kubernetes
 - CPU and RAM requests less than limit - Burstable QoS by default
 - K8s QoS policy determines scheduling priority: BestEffort class is evicted before Burstable class, which is evicted before Guaranteed class

Node Resource Management

- VM density is determined by multiple factors controlled at the cluster, OpenShift Virtualization, Pod, and VM levels
- Pod QoS policy
 - Burstable (limit > request) allows more overcommit, but may lead to more frequent migrations
 - Guaranteed (limit = request) allows less overcommitment, but may have less physical resource utilization on the hosts
- Cluster Resource Override Operator provides global overcommit policy, can be customized per project for additional control
- Pods request full amount of VM memory and approx. 10% of VM CPU
 - VM pods request a small amount of additional memory, used for libvirt/QEMU overhead
 - Administrator can set this to be overcommitted

High availability

- Node failure is detected by Kubernetes and results in the Pods from the lost node being rescheduled to the surviving nodes
 - Use machine health checks, node health checks, and the Node Self Remediation Operator to expedite detection and recovery
- VMs are not scheduled to nodes which have not had a heartbeat from `virt-handler`, regardless of Kubernetes node state
- Additional monitoring may trigger automated action to force stop the VM pods, resulting in rescheduling
 - May take up to 5 minutes for `virt-handler` and/or Kubernetes to detect failure
 - Liveness and Readiness probes may be configured for VM-hosted applications
 - Machine health checks can decrease failure detection time

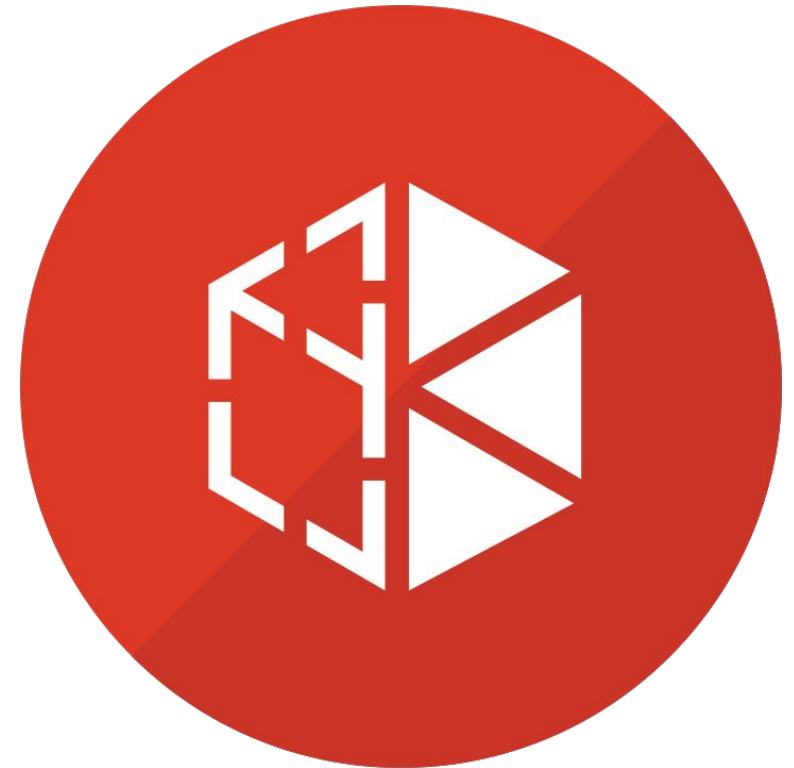
Terminology comparison

Feature	RHV	OpenShift Virtualization	vSphere
Where VM disks are stored	Storage Domain	PVC	datastore
Policy based storage	None	StorageClass	SPBM
Non-disruptive VM migration	Live migration	Live migration	vMotion
Non-disruptive VM storage migration	Storage live migration	N/A	Storage vMotion
Active resource balancing	Cluster scheduling policy	Pod eviction policy, descheduler	Dynamic Resource Scheduling (DRS)
Physical network configuration	Host network config (via nmstate w/4.4)	nmstate Operator, Multus	vSwitch / DvSwitch
Overlay network configuration	OVN	OCP SDN (OpenShiftSDN, OVNKubernetes, and partners), Multus	NSX-T
Host / VM metrics	Data warehouse + Grafana (RHV 4.4)	OpenShift Metrics, health checks	vCenter, vROps

Runtime awareness

Deploy and configure

- OpenShift Virtualization is deployed as an Operator utilizing multiple CRDs, ConfigMaps, etc. for primary configuration
- Many aspects are controlled by native Kubernetes functionality
 - Scheduling
 - Overcommitment
 - High availability
- Utilize standard Kubernetes / OpenShift practices for applying and managing configuration



Compute configuration

- VM nodes should be physical with CPU virtualization technology enabled in the BIOS
 - Nested virtualization *works*, but ***is not supported***
 - Emulation *works*, but ***is not supported*** (and is extremely slow)
- Node labeler detects CPU type and labels nodes for compatibility and scheduling
- Configure overcommitment using native OpenShift functionality - Cluster Resource Override Operator
 - Optionally, customize the project template so that non-VM pods are not overcommitted
 - Customize projects hosting VMs for overcommit policy
- Apply Quota and LimitRange controls to projects with VMs to manage resource consumption
- VM definitions default to all memory “reserved” via a request, but only a small amount of CPU
 - CPU and memory request/limit values are modified in the VM definition

Network configuration

- Apply traditional network architecture decision framework to OpenShift Virtualization
 - Resiliency, isolation, throughput, etc. determined by combination of application, management, storage, migration, and console traffic
 - Most clusters are not VM only, include non-VM traffic when planning
- Node interface on the `MachineNetwork.cidr` is used for “primary” communication, including SDN
 - This interface should be both resilient and high throughput
 - Used for migration and console traffic
 - Configure this interface at install time using kernel parameters, reinstall node if configuration changes
- Additional interfaces, whether single or bonded, may be used for traffic isolation, e.g. storage and VM traffic
 - Configure using nmstate Operator, apply configuration to nodes using selectors on NNCP

Storage configuration

- Shared storage is not required, but *very highly encouraged*
 - **Live migration depends on RWX PVCs**
- Create shared storage from local resources using OpenShift Container Storage
 - RWX file and block devices for live migration
- No preference for storage protocol, use what works best for the application(s)
- Storage backing PVs should provide adequate performance for VM workload
 - Monitor latency from within VM, monitor throughput from OpenShift
- For IP storage (NFS, iSCSI), consider using dedicated network interfaces
 - Will be used for all PVs, not just VM PVs
- Certified CSI drivers are recommended
 - Many non-certified CSI provisioners work, but do not have same level of OpenShift testing
- Local storage may be utilized via the Host Path Provisioner

Deploying a VM operating system

Creating virtual machines can be accomplished in multiple ways, each offering different options and capabilities

- Start by answering the question “Do I want to manage my VM like a container or a traditional VM?”
- Deploying the OS persistently, i.e. “I want to manage like a traditional VM”
 - Methods:
 - Import a disk with the OS already installed (e.g. cloud image) from a URL or S3 endpoint using a DataVolume, or via CLI using virtctl
 - Clone from an existing PVC or VM template
 - Install to a PVC using an ISO
 - VM state will remain through reboots and, when using RWX PVCs, can be live migrated
- Deploying the OS non-persistently, i.e. “I want to manage like a container”
 - Methods:
 - Diskless, via PXE
 - Container image, from a registry
 - VM has no state, power off will result in disk reset. No live migration.
- Import disks deployed from a container image using CDI to make them persistent

Deploying an application

Once the operating system is installed, the application can be deployed and configured several ways

- The application is pre-installed with the OS
 - This is helpful when deploying from container image or PXE as all components can be managed and treated like other container images
- The application is installed to a container image
 - Allows the application to be mounted to the VM using a secondary disk. Decouples OS and app lifecycle. When used with a VM that has a persistently deployed OS this breaks live migration
- The application is installed after OS is installed to a persistent disk
 - cloud-init - perform configuration operations on first boot, including OS customization and app deployment
 - SSH/Console - connect and administer the OS just like any other VM
 - Ansible or other automation - An extension of the SSH/console method, just automated

Additional resources

More information

- OpenShift Virtualization content kit: <https://red.ht/virtkit>
- Documentation:
 - OpenShift Virtualization: <https://docs.openshift.com>
 - KubeVirt: <https://kubevirt.io>
- Demos and video resources:
https://www.youtube.com/playlist?list=PLaR6Rq6Z4lqeQeTosfoFzTyE_QmWZW6n
- Labs and workshops: <https://demo.redhat.com>

The LAB Starts at 13:30 !!

<https://demo.redhat.com/workshop/pmfaa6>

Password:

D54ha4s2

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