

RED HAT
SUMMIT

BOSTON, MA
JUNE 23-26, 2015

High-Performance OpenStack for Science and Data Analytics

Steven Carter
Solutions Architect, Cisco Systems
26 June 2015

#redhat #rhsummit



Agenda

- Overview
- Networking
- Storage
- Hardware
- Schedulers
- Use Cases

HPC on OpenStack?

Similarities:

- Compatible Workloads:
 - HPC workloads among the first stateless workloads
- Achieve Economies of Scope
 - Major vendors working on innovation in the space

Differences:

- Science requires data!
 - To/From the WAN
 - To/From Central Storage
 - Between Instances
- Science requires more Compute/Memory
 - Less/No Oversubscription
 - GPUs

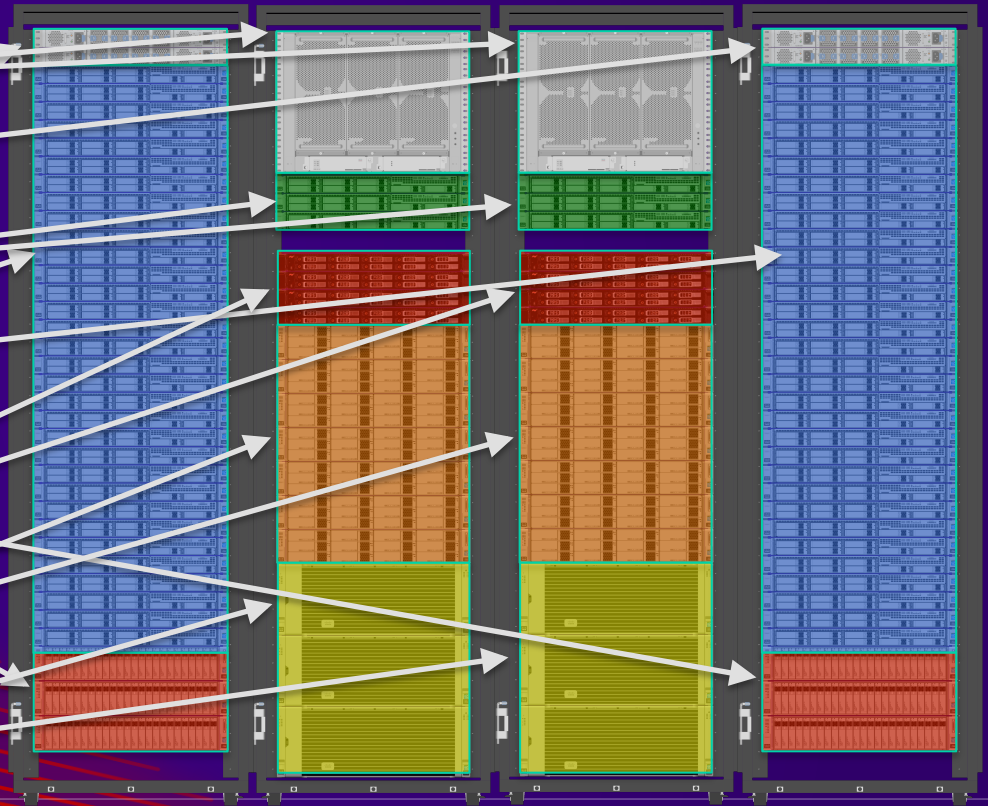
System Overview

- Purpose Built
- Scalable Architecture
- Easily extensible by adding Scalable Units
- Network architected for Data Movement, Low Latency
- Hardware optimized for Cycles/Memory



System Overview

- Fabric Spine
- Fabric Leaf
- Control Nodes
- Compute Nodes
- Local Block Storage
- Shared SSD
- Shared File
- Shared Object

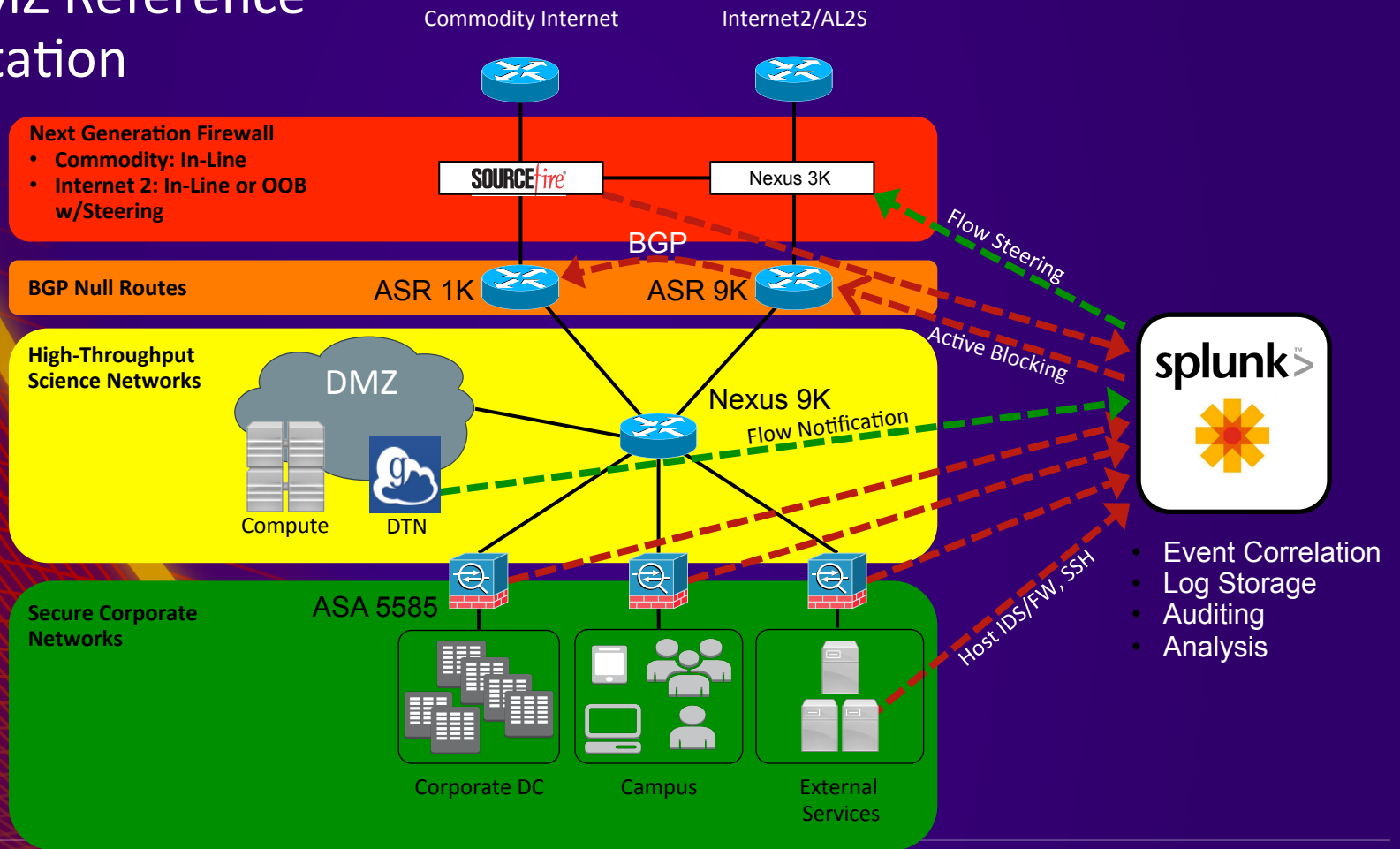


Networking: Science DMZ

#redhat #rhsummit



Science DMZ Reference Implementation



Science DMZ Components



splunk™

- Most informed view of the status of the network, servers, and apps.
- Provides event correlation by severity, rate, and between events
- Provides for auditing and reporting capabilities
- Leverage existing skill by writing logic in Splunk search language



- Provides single northbound REST API
- Application specific logic in between
- Hardware specific southbound APIs



globus

- Software-as-a-Service (SaaS) solution to manage transfers
- GridFTP extensions provides parallelism (i.e., the use of multiple socket connections between pairs of data movers), restart markers, and data channel security.
- GridFTP control plane provides the source and destination information for the flows it sets up
- Effectively authenticates flows before they bypass security

OpenFlow Data Flow Steering

Base setup depending on mode:

Out-Of-Band IDS:

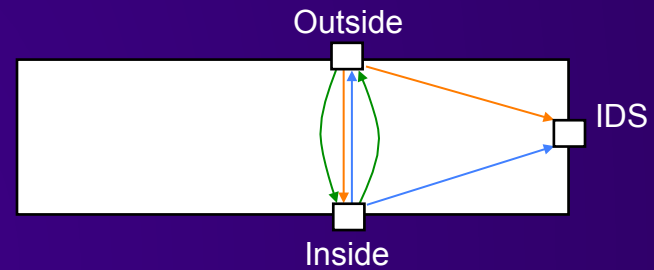
- `<priority>100</priority>`
- `<in-port>54</in-port>`
- `<output-node-connector>52</output-node-connector>`
- `<output-node-connector>25</output-node-connector>`

In-Band Firewall/IPS:

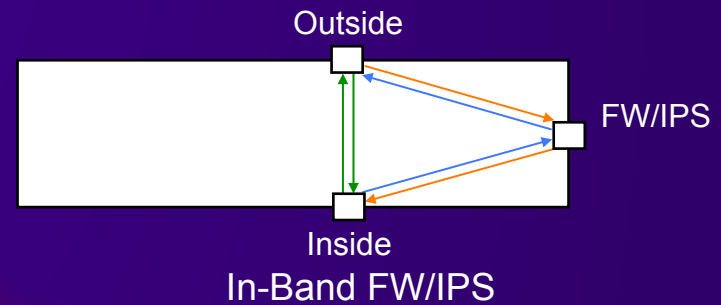
- `<priority>100</priority>`
- `<in-port>54</in-port>`
- `<output-node-connector>25</output-node-connector>`
- `<in-port>25</in-port>`
- `<output-node-connector>52</output-node-connector>`

Bypass operation the same for both modes

- `<priority>200</priority>`
- `<in-port>54</in-port>`
- `<output-node-connector>52</output-node-connector>`



Out-Of-Band IDS



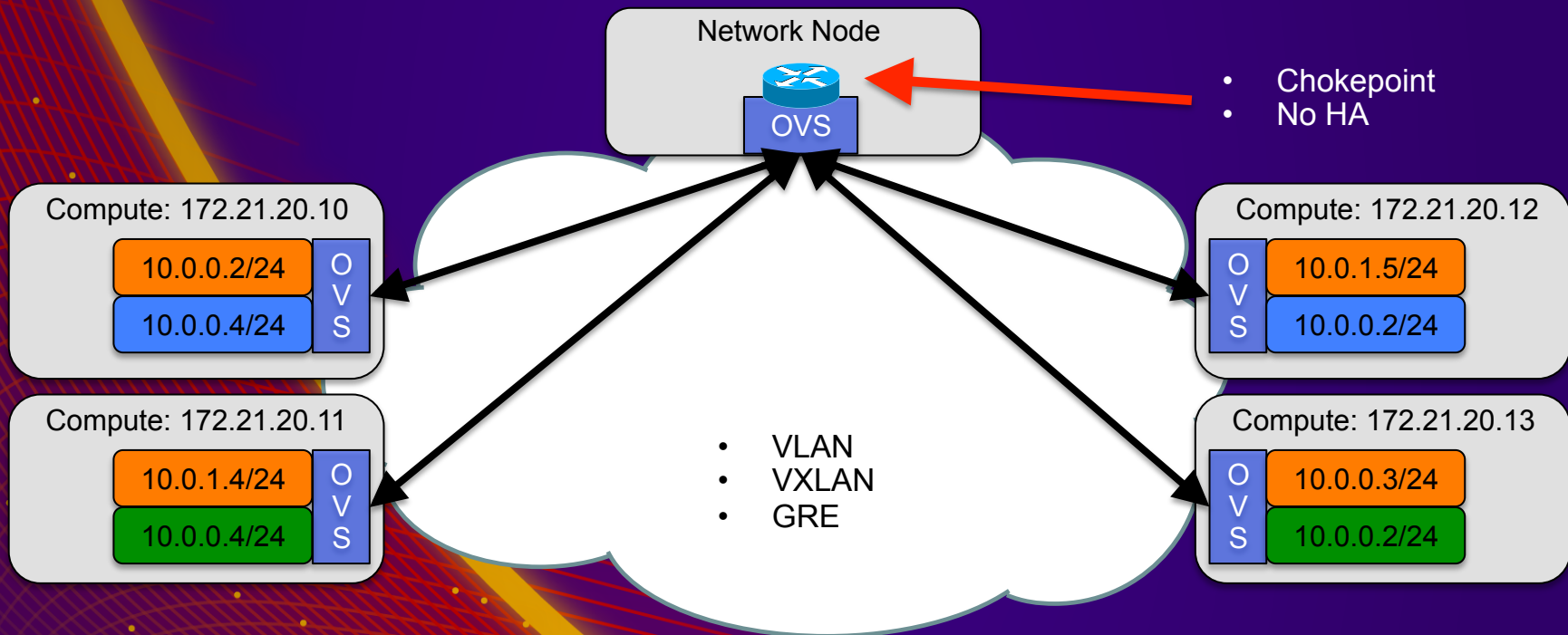
In-Band FW/IPS

Networking: Fabric

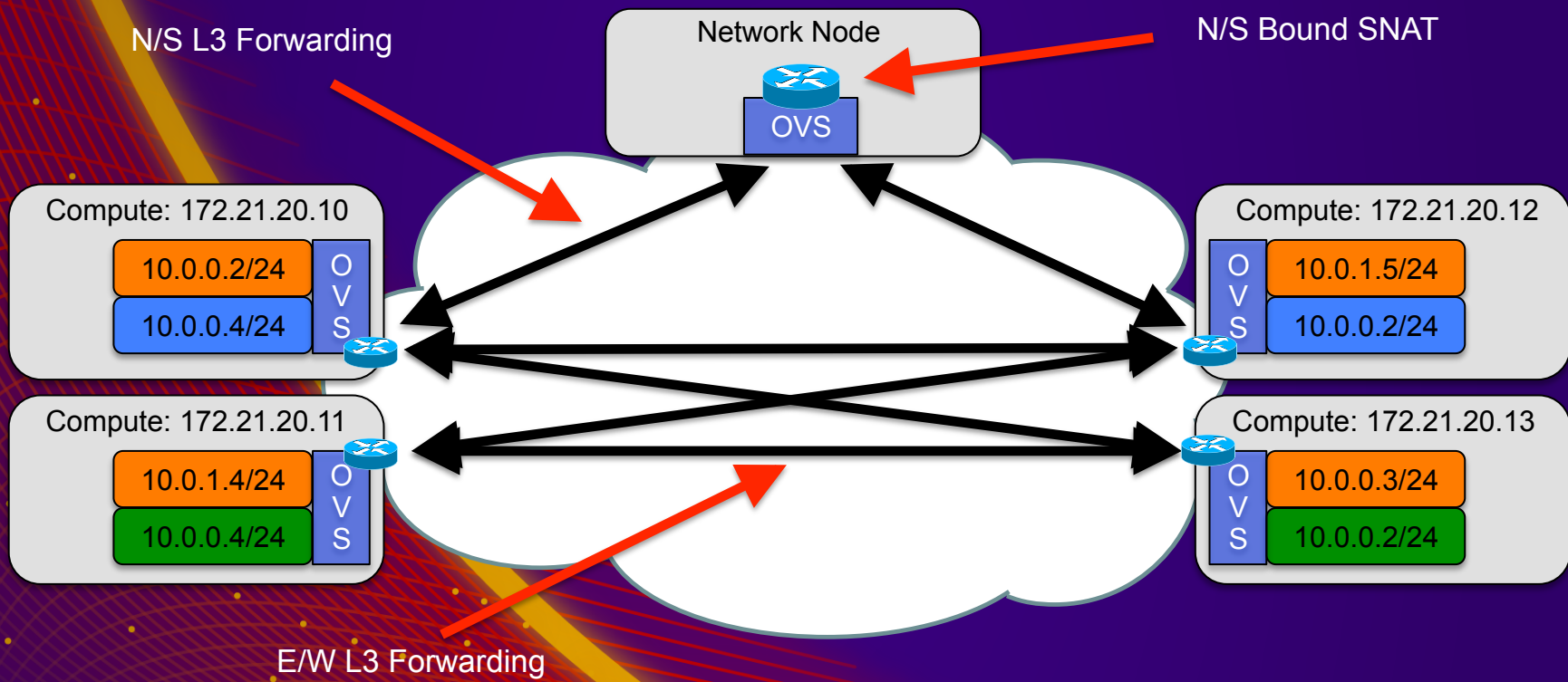
#redhat #rhsummit



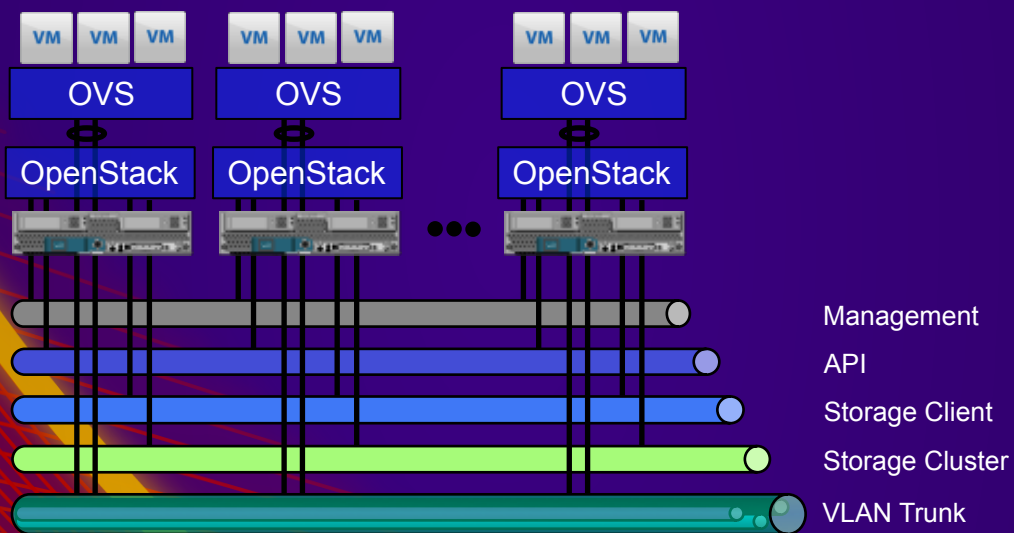
Traditional Neutron Router



Distributed Virtual Routing



Neutron Provider Networks



Overlays

VLAN

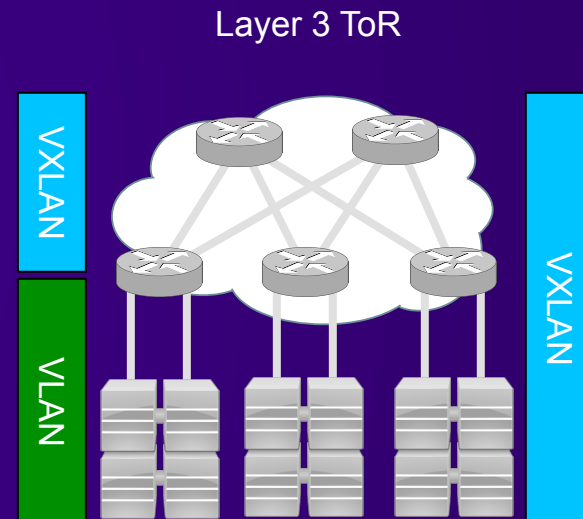
- Performing
- Pervasive
- 12 bits of tenant networks of tenants
- Layer 2 ToR

VXLAN

- Less Performing
- Less Pervasive
- Control plane coalescing around standard EVPN
- 24 bits of tenant networks (twice a much, right!)
- Rapidly evolving
- Layer 3 ToR

GRE

- Blech



Storage

#redhat #rhsummit



Storage Access from Instance

Local to Node:

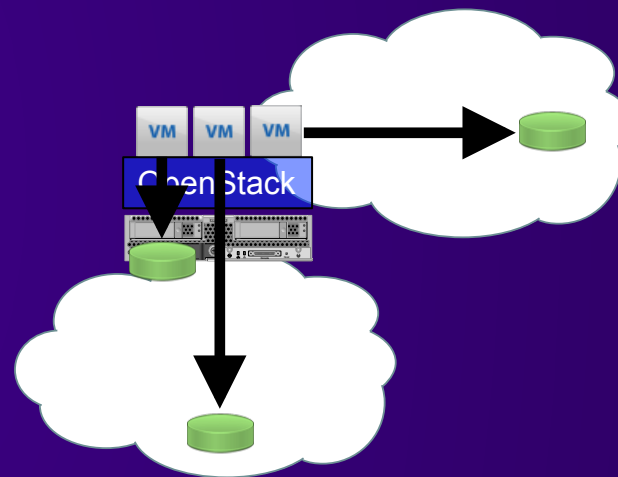
- Local Disks

Remote through OpenStack:

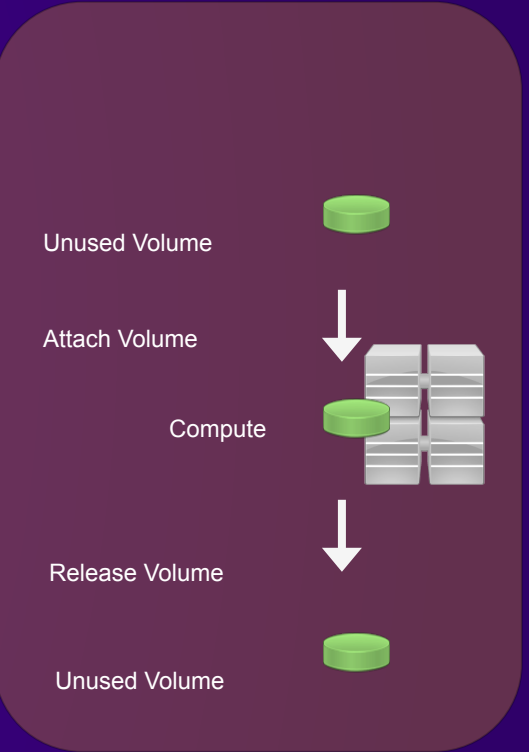
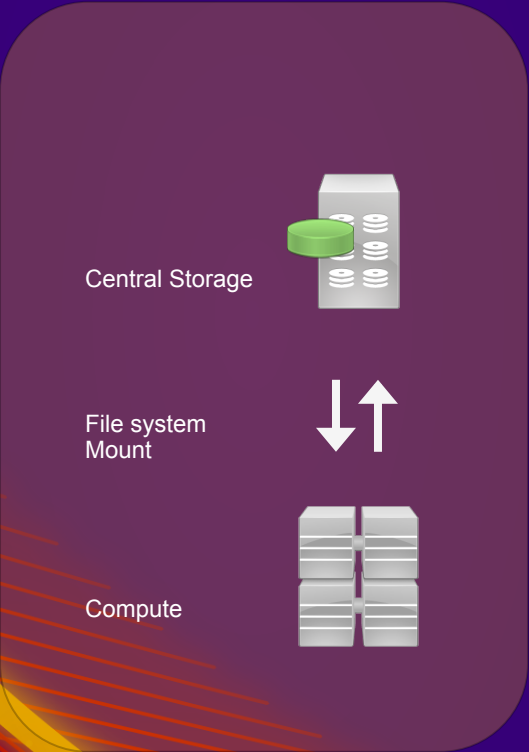
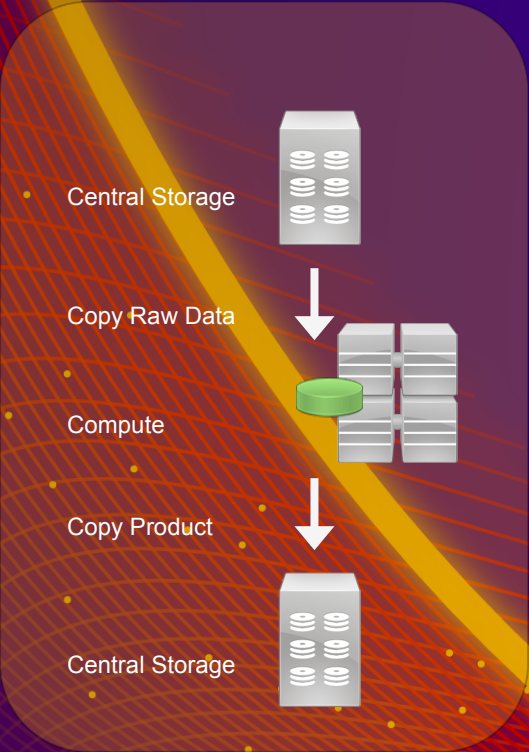
- Cinder
- Manila

Remote, direct from Instance:

- Swift
- NFS
- Gluster
- Lustre



HPC Data Workflows



Local Temporal Storage

Central Storage: Swift, Possibly File

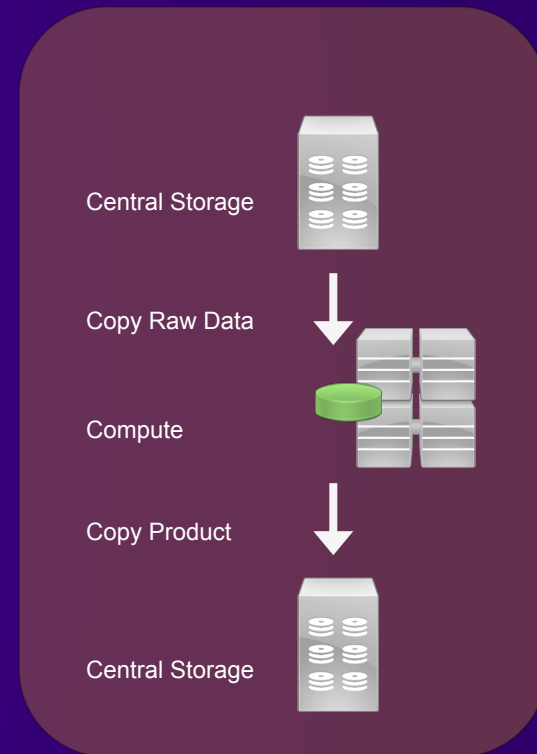
Instance Storage: Local, Cinder

Pros:

- Favors Object Storage – No mounting needed in instance
- Most Portable
- Most Conducive to SSD Usage
- Easiest on Shared Storage

Cons:

- Requires movement on/off compute nodes



Central Persistent Storage

Central Storage: File (NFS, Lustre, Gluster, GPFS, etc)

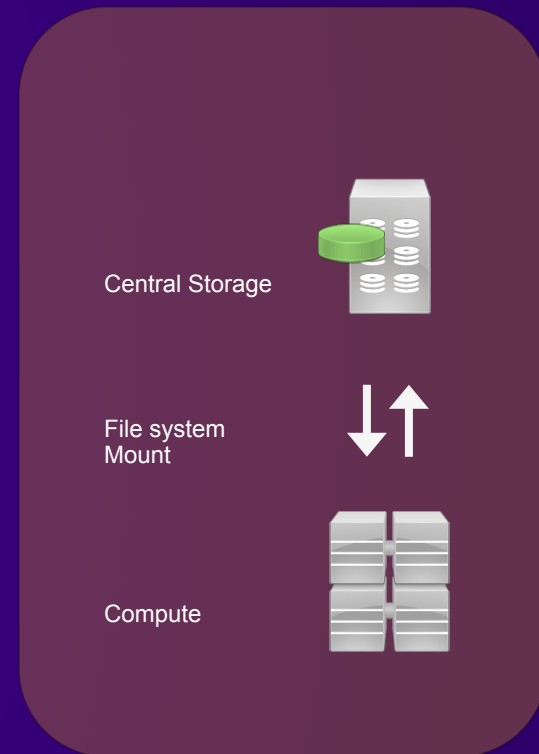
Instance Storage: Minimal

Pros:

- No copies

Cons:

- Hardest on Shared Storage
- Requires good network access for tenants
- Need mount on Instance



Persistent Volume

Central Storage: Cinder (Ceph, NFS, NetApp, etc.)

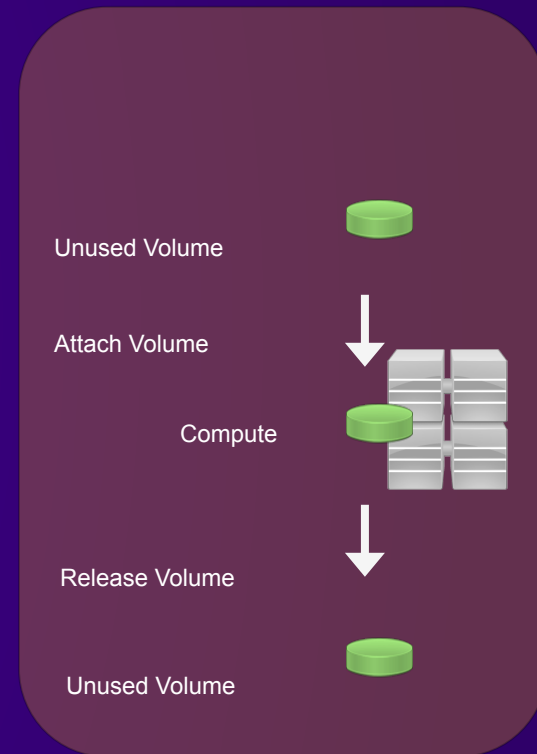
Instance Storage: Minimal

Pros:

- Generally Filesystem
- No copies

Cons:

- Limited ccessability to data (single node access to volumes)



Hardware

Platforms

B200

- 2 sockets, mid-memory, limited storage, dense footprint

M4308 Cartridges

- 1 socket, mid-memory, shared storage, dense footprint

C220

- 2 sockets, mid-memory, limited storage

C240

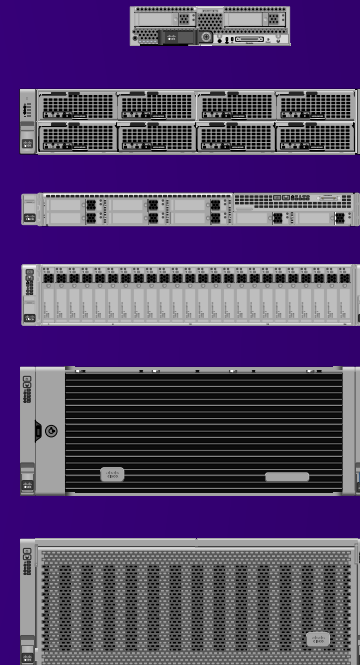
- 2 sockets, mid- to large-memory, storage, expansion

C460

- 4 sockets, large memory, storage, large expansion

C3160

- 2 sockets, mid-memory, large storage



Compute

Storage

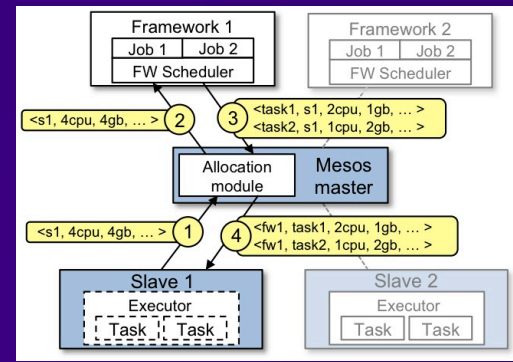
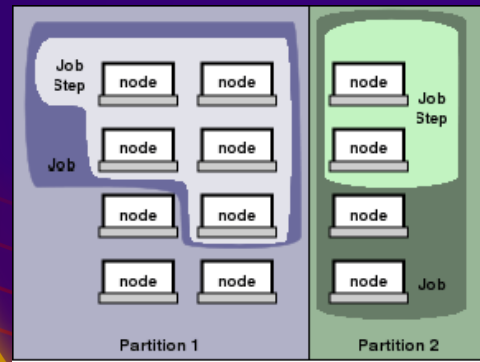
Containers

- More about packaging than security
- Great for portability
- Compelling when running directly on hardware (e.g. Iron)
- Gives better access to hardware for things like GPU and low-latency interconnect
- Need to consider noisy neighbor problems
- Works well with cartridges because it is a 1:1 mapping from container to physical

Scheduling



Traditional HPC



- Traditional batch scheduler
- Coarse grained scheduling
- Allocates entire node to job
- Virtual nodes added to queue
- Optimized for job completion time

- Cloud oriented “Distributed Systems Kernel”
- Fine grained scheduling
- Abstracts CPU, memory, storage, and other resources
- Optimized for throughput and utilization

Extension to Public Clouds

#redhat #rhsummit



Extension to Public Clouds

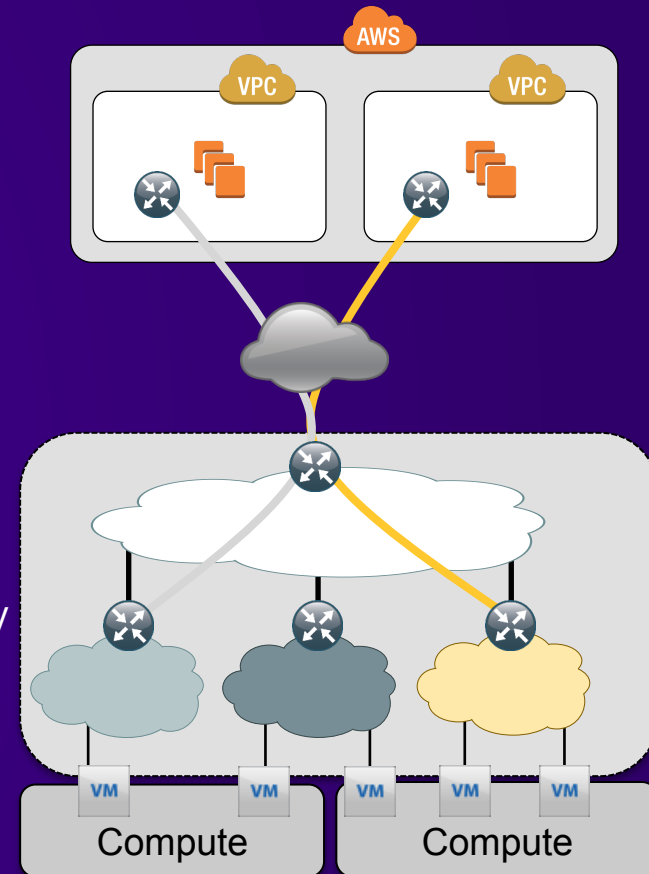
Scheduler Based

- Workers can be distributed across clouds
- Need access to scheduler
- Need access to storage

Network Extension

- VPN
- Direct Connect

Tenant Gateway
Tenant VLANs



Use Cases

#redhat #rhsummit



Use Cases



Application: Seismic Data Analysis

Infrastructure:

- Network: Provider VLANs
- Compute: Rack Servers
- Storage: Object w/local scratch

Scheduling:

- Celery/Mesos

"Burst" to AWS



Application: Genomics

Infrastructure:

- Network: Provider VLANs
- Compute: Blades
- Storage: NFS

Scheduling:

- Grid Engine

"Burst" to AWS

RED HAT
SUMMIT

LEARN. NETWORK.
EXPERIENCE OPEN SOURCE.

#redhat #rhsummit

