

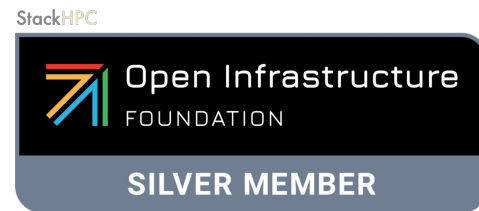
Slurm Appliance

A batch HPC cluster defined as code for
OpenStack clouds

Steve Brasier

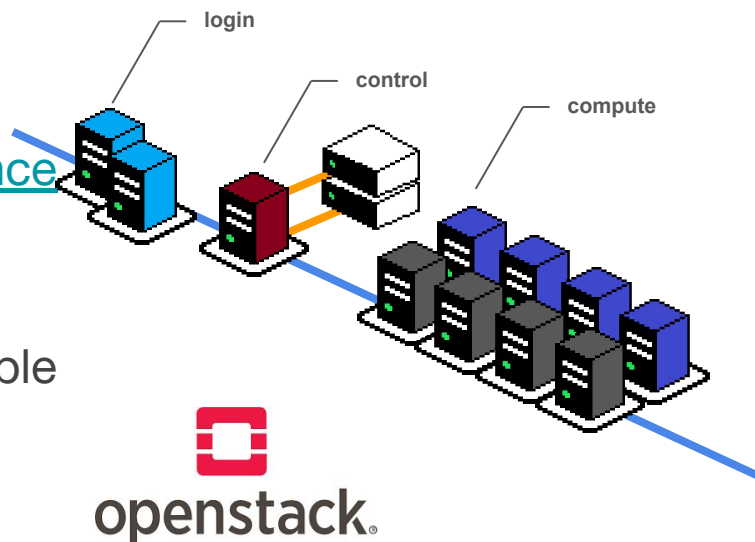
StackHPC Company Overview

- Formed 2016, based in Bristol, UK
 - Based in Bristol with presence in Oxford, Cambridge, France and Poland
 - Currently around 30 people
- Founded on HPC expertise
 - Software Defined Networking
 - Systems Integration
 - OpenStack Development and Operations
- Motivation to transfer this expertise into Cloud to address HPC & HPDA (AI)
- “Open” Modus Operandi
 - Upstream development of OpenStack capability
 - Consultancy/Support to end-user organizations in managing HPC service transition
 - Scientific-WG engagement for the Open Infrastructure Foundation
- Hybrid Cloud Enablement



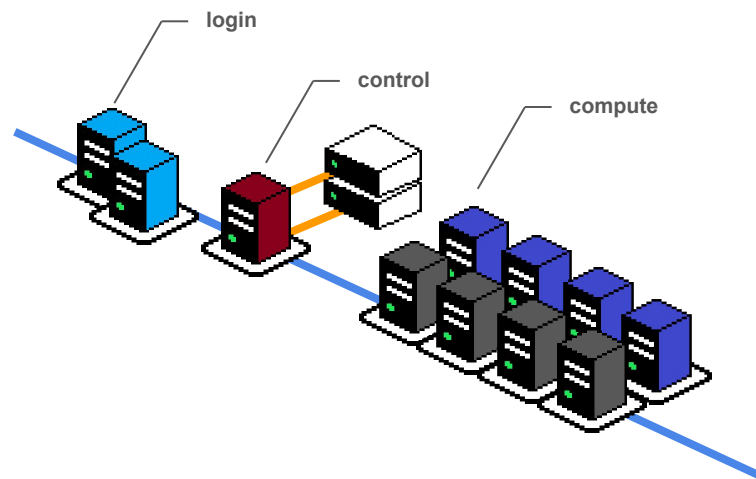
Overview

- Batch HPC environment on OpenStack
- Entirely defined by code
github.com/stackhpc/ansible-slurm-appliance
- Multi-user long-lived site cluster
- Shared, extendable and open-source
- Useable & sensible defaults, but customisable



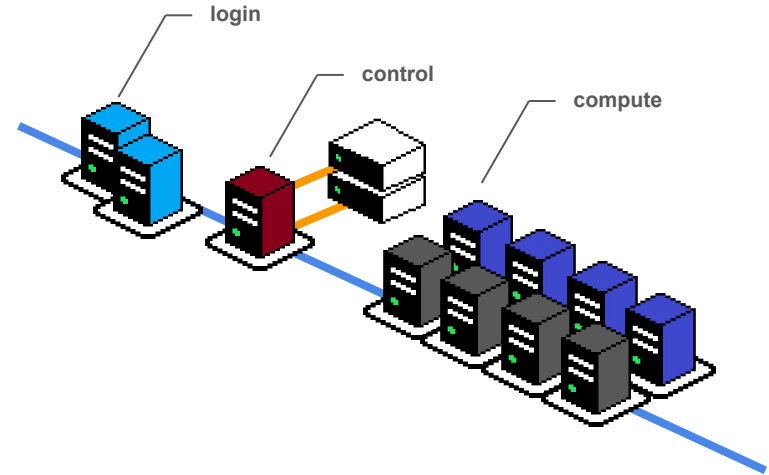
Overview

- Slurm batch scheduler
- Monitoring
- Web-based portal
- Workload support



Overview

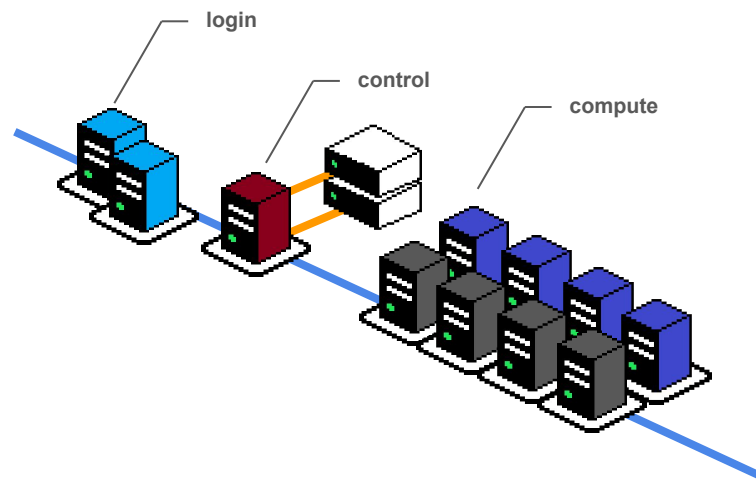
- **Slurm batch scheduler**
- Monitoring
- Web-based portal
- Workload support



- Accounting database
- Login control
- LBNL Node Health Checks
- Topology-aware scheduling

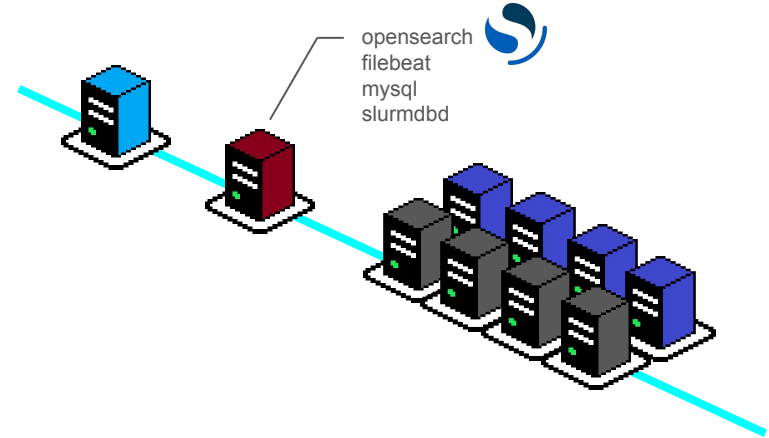
Overview

- Slurm batch scheduler
- **Monitoring**
- Web-based portal
- Workload support



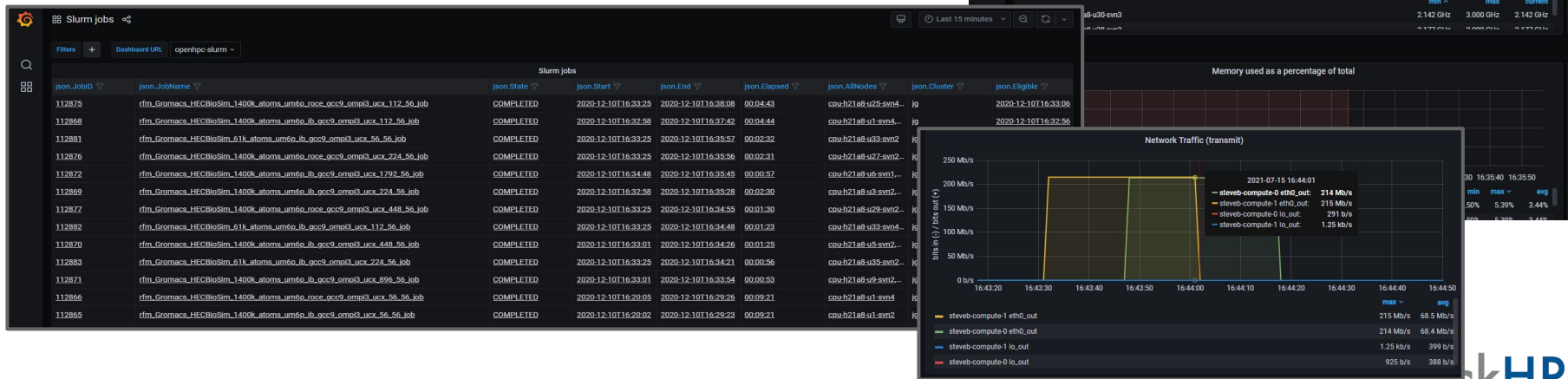
Monitoring

- “Standard” Linux metrics stack of node exporter + Prometheus + Grafana
- SlurmDBD + MySQL + Filebeat + OpenSearch to integrate with Slurm jobs
- Alertmanager



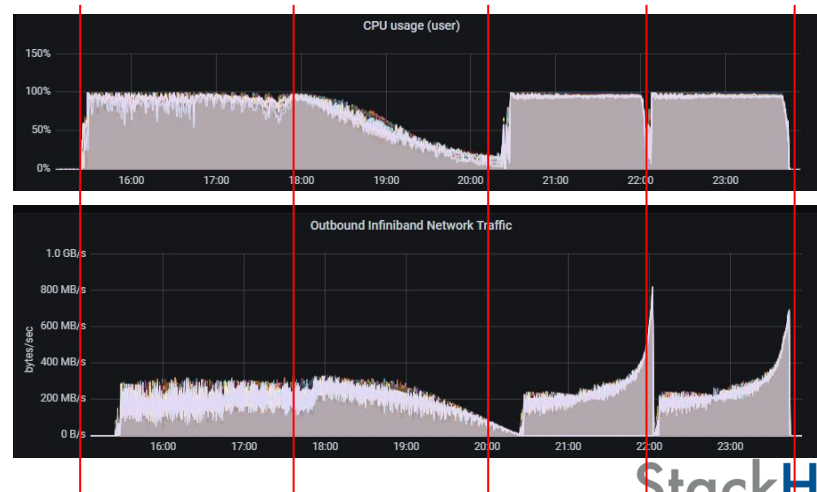
Monitoring

- “Standard” Linux metrics stack of node exporter + Prometheus + Grafana
- SlurmDBD + MySQL + Filebeat + OpenSearch to integrate with Slurm jobs
- Alertmanager



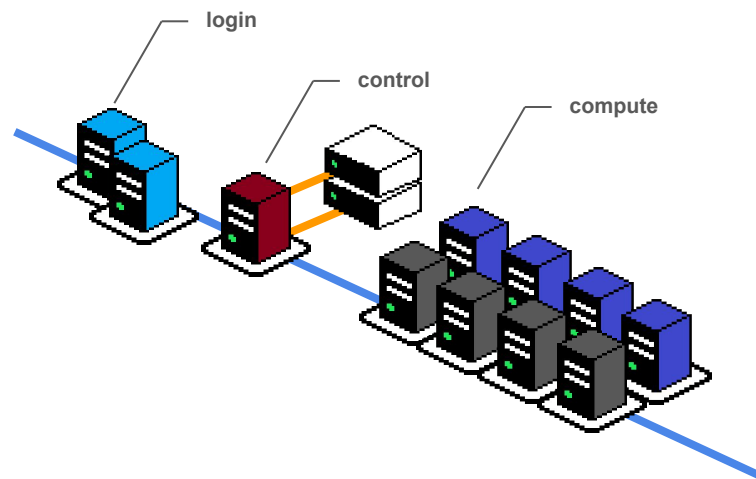
Monitoring

- “Standard” Linux metrics stack of node exporter + Prometheus + Grafana
- SlurmDBD + MySQL + Filebeat + OpenSearch to integrate with Slurm jobs
- Alertmanager
- Plus other prometheus exporters:
 - Open Ondemand
 - Slurm
- LBNL Node Health Checks
 - Set node DOWN if it loses mounts etc ...



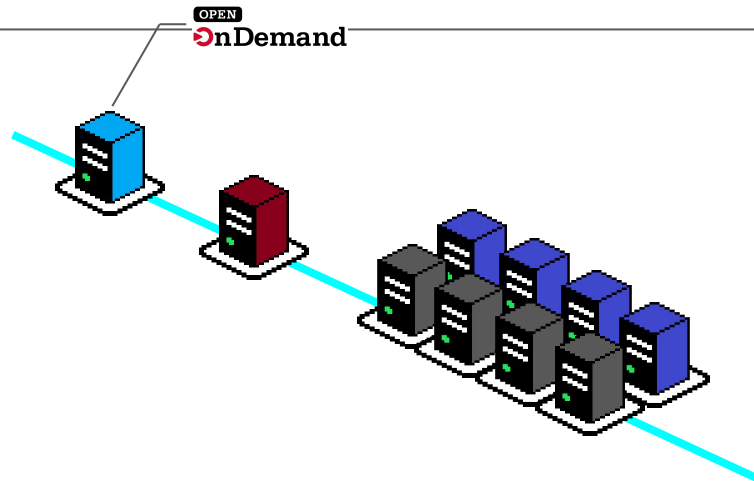
Overview

- Slurm batch scheduler
- Monitoring
- **Web-based portal**
- Workload support



User Portal

- Web-based Open Ondemand user portal
 - Shell
 - File explorer
 - Monitoring (proxied)
 - Job templates & submission etc



The screenshot shows the Open OnDemand user portal interface. The top navigation bar includes "eSearch", "Files", "Jobs", "Clusters", "Interactive Apps", and "Monitoring". The main content area displays "Active Jobs" with a table of job entries:

ID	Name
12	hpl-solo.sh
13	hpl-solo.sh
14	hpl-solo.sh
11	hpl-solo.sh

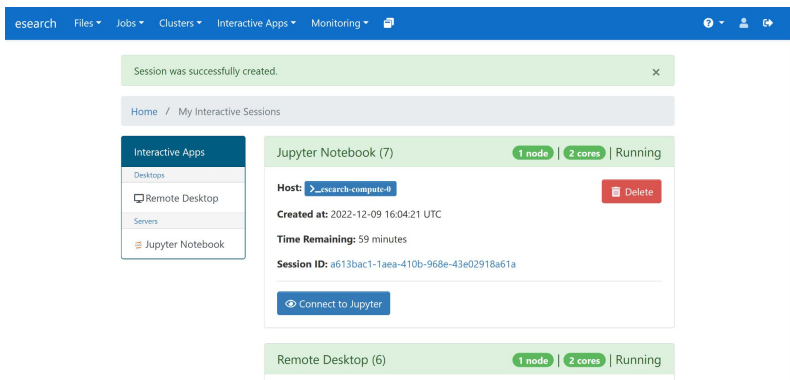
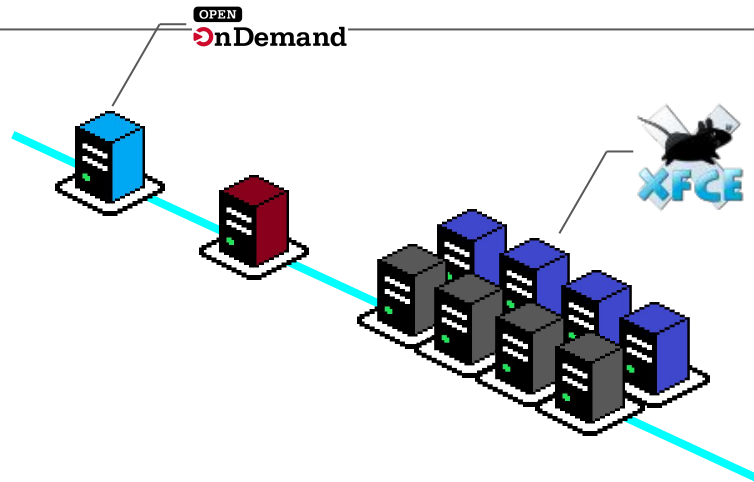
Below the table, there is a "Home Directory" section with a file explorer view showing a terminal window. The terminal output shows the command `sinfo` and its output:

```
Host: esearch-login-0
[testuser@esearch-login-0 ~]$ sinfo
PARTITION AVAIL  TIMELIMIT  NODES  STATE MODELIST
extra      up 60-00:00:00    2   idle esearch-compute-[2-3]
small*    up 60-00:00:00    2   idle esearch-compute-[0-1]
[testuser@esearch-login-0 ~]$
```

The interface also includes a "Message of the Day" section and a "powered by Open OnDemand" logo at the bottom left. The version number "OnDemand version: 2.0.29" is displayed at the bottom right.

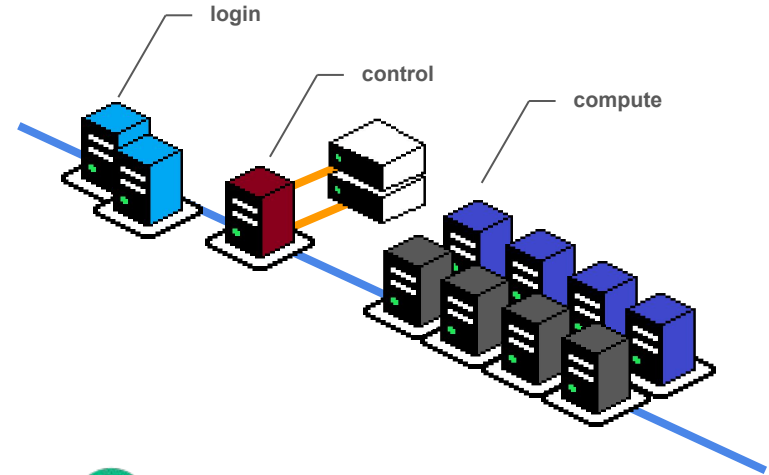
User Portal

- Preconfigured compute node “apps”
 - Remote desktop
 - Jupyter notebook
 - Rstudio
 - Codeserver
 - Matlab*

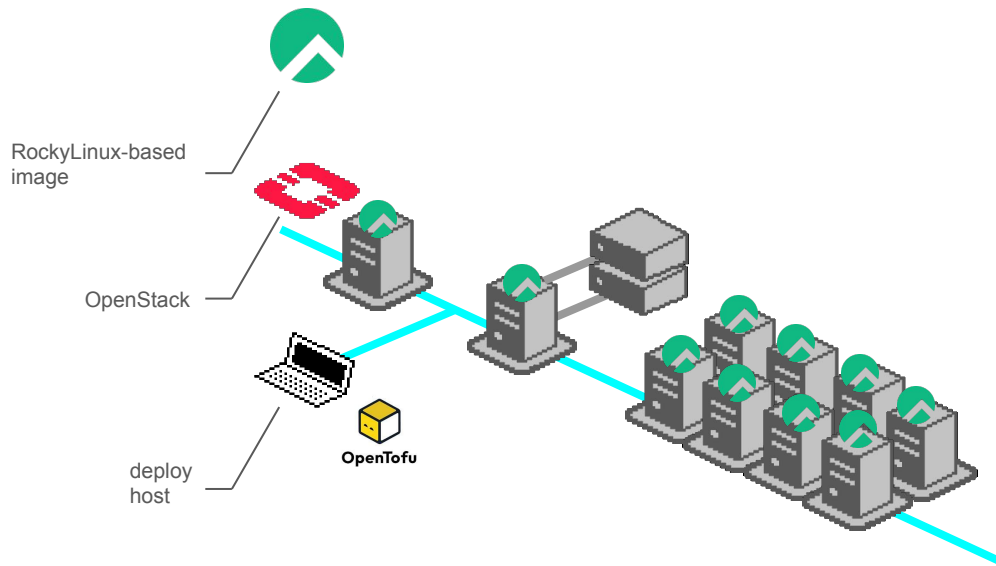


Overview

- Slurm batch scheduler
- Monitoring
- Web-based portal
- **Workload support**



Deployment Process

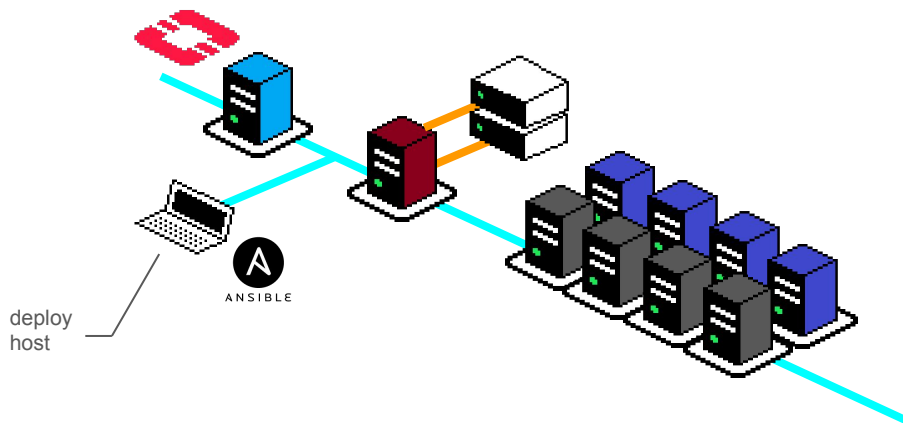


OpenTofu creates infrastructure and Ansible inventory file

```
module "cluster" {  
  source           = "../../site/tofu/"  
  environment_root = var.environment_root  
  cluster_name    = "staging"  
  login = {  
    head = {  
      nodes       = ["login-0"]  
      flavor      = "vm.cpu.himem.eighth"  
      fip_addresses = ["198.51.100.19"]  
      fip_network  = "external"  
    }  
  }  
}  
  
compute = {  
  saphrapid_himem = {  
    nodes       = ["standard-0", "standard-1", "standard-2"]  
    flavor      = "vm.cpu.himem.v2.full"  
  }  
  a100_gpu = {  
    nodes       = ["a100-0"]  
    flavor      = "vm.gpu.4x.a100.full"  
  }  
}
```

Deployment Process

Ansible configures instances into a cluster



```
- name: Ensure OpenHPC repos
  ansible.builtin.yum_repository : "{{ item }}"
  loop: "{{ openhpc_repos }}"

- name: Install required Slurm packages
  ansible.builtin.dnf :
    name: "{{ openhpc_slurm_pkglist }}"

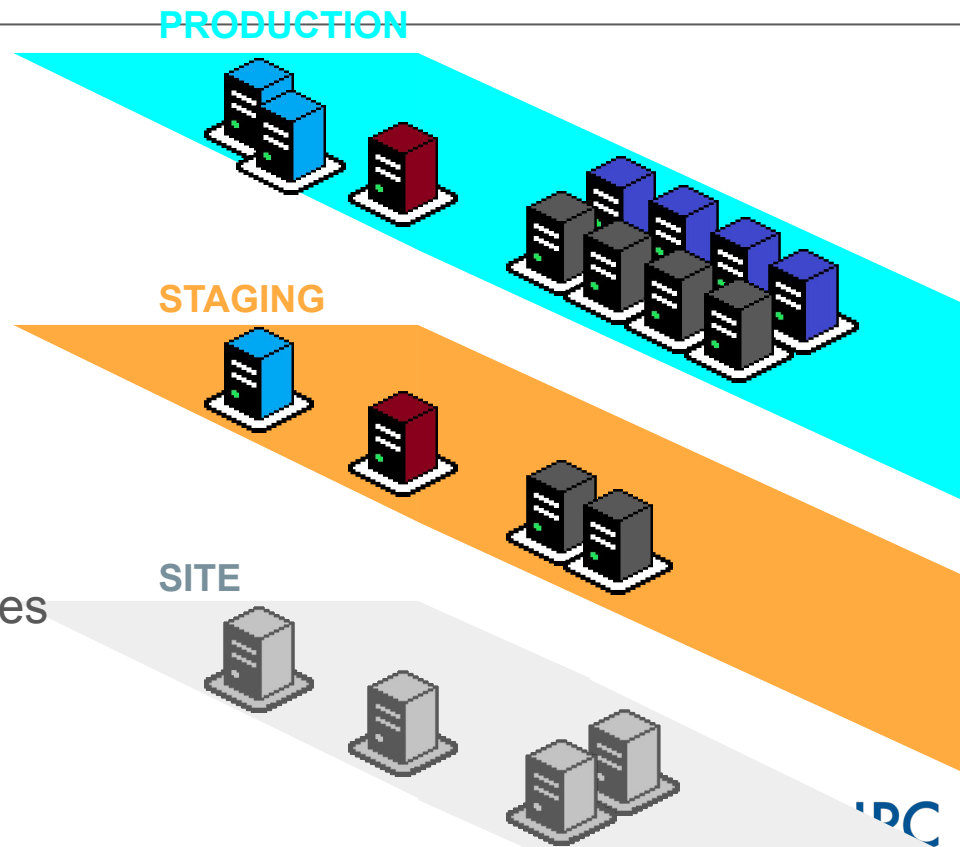
...

- name: Template slurm.conf
  template:
    src: "{{ openhpc_slurm_conf_template }}"
    dest: "{{ openhpc_slurm_conf_path }}"
    owner: root
    group: root
    mode: '0644'
  when: openhpc_enable.control | default(false)
  notify:
    - Restart slurmctld service

...
```

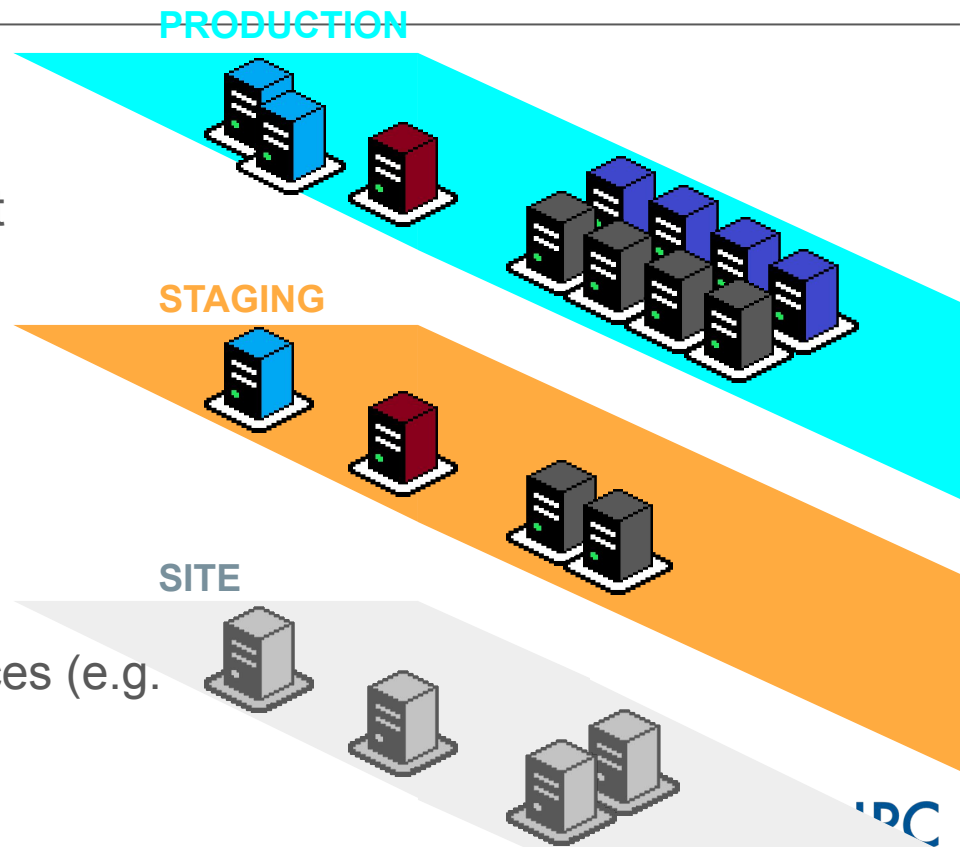
Environments

- Configuration drawn from multiple, layered Ansible inventories
 - *common* - StackHPC defaults
 - *site* - per-site defaults
 - *production, staging, ...*
- Shared aspects only defined once
- Minimal differences
 - Cluster name
 - Floating IPs + external DNS names
 - Numbers of nodes



Environments

- Git workflows
- Upgrades via reimaging
 - Confidence on production roll-out
 - But changes are disruptive
- Easy development
- Minimise resources required
 - Fewer nodes
 - Smaller instances
 - VMs instead of baremetal
 - Temporarily move limited resources (e.g. GPUs)



EESSI Integration

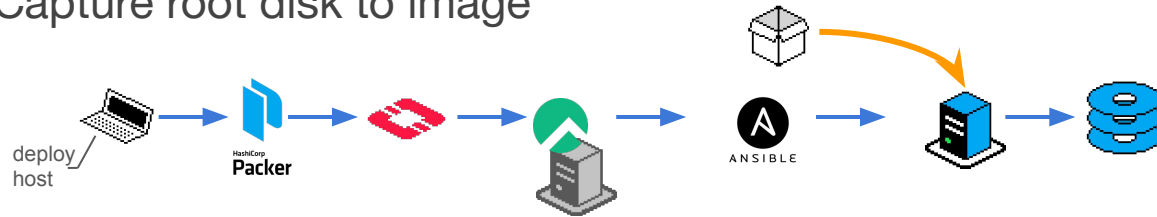
- Since pilot repo days ...
- Docs converted to Ansible

Docs are great - thank you!

- By default:
 - Configure Squid proxy on control node
 - Configure clients to use it
 - Check for NVIDIA GPUs and link host libraries
- Obviously requires outbound HTTP ...
 - Client cluster with only **HTTPS** allowed out
 - Open PR to add CVMFS Stratum 1 server

Image Build

- Using upstream repos isn't repeatable ...
- Hashicorp Packer
 - Launch VM in OpenStack
 - Run (some) Ansible
 - DNF installs
 - Template out files, ...
 - Capture root disk to image
- Adds complexity to Ansible
- *State* preserved via remote storage

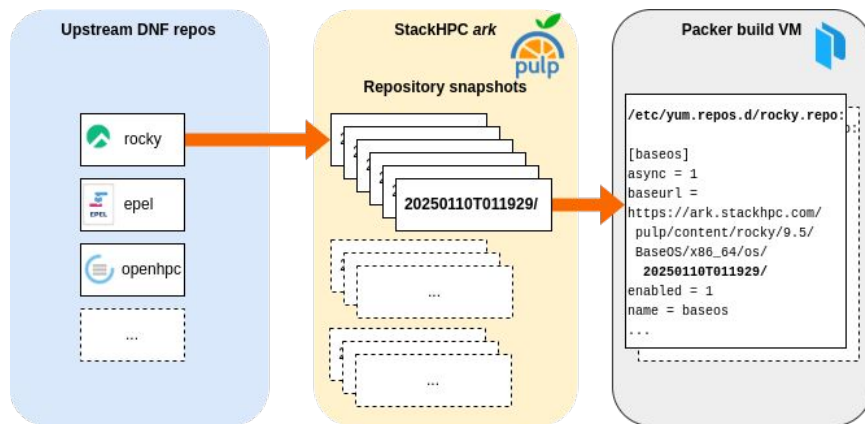


Package Problems ...

- Tests in staging at client site using Lustre: Image missing RDMA verbs header file
 - Latest NVIDIA OFED version clashed with OpenHPC libfabric package required for MPI packages
- New NVIDIA CUDA packages released
 - CI image builds fail - not compatible with the NVIDIA drivers
- RockyLinux 9.5 released
 - Packages from RockyLinux 9.4 disappear from mirrors
 - Our CI tests found issue with latest podman package
- **No version of NVIDIA OFED which worked with RockyLinux 9.5**

Image Build

- Using upstream repos (still) isn't repeatable ...
- Use package mirror with "snapshots" of upstream repositories
- StackHPC "Ark": Pulp server for OpenStack packages

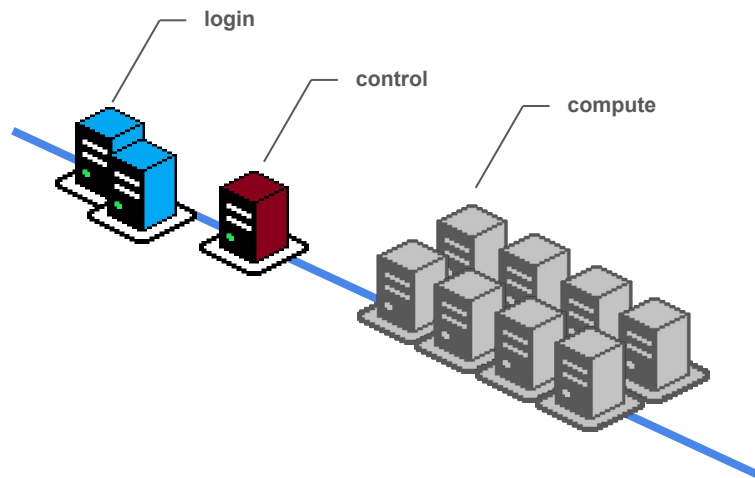


- *Repeatable* builds

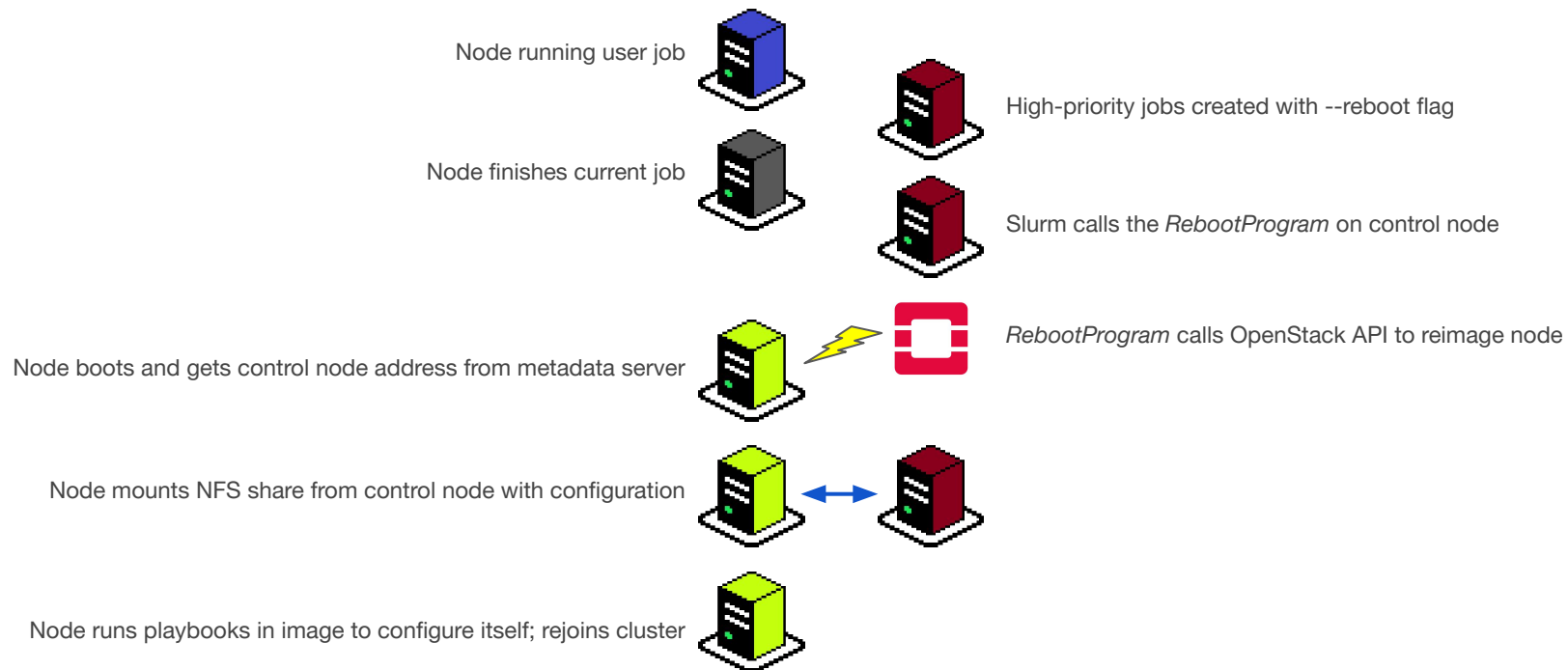
Lower-disruption upgrades

Update via reimage DRAINS compute nodes

- Use the scheduler to reimage compute nodes between user jobs
- Need Ansible to configure nodes?
 - Build code into the image
 - Provide config via metadata + NFS



Lower-disruption upgrades



Features

IAM and access control

- Templated users/groups
- fail2ban & firewalld
- FreeIPA
- LDAP / sssd
- sshd
- opkssh: SSH via OIDC identities

Filesystems

- NFS
- Lustre
- GPFS (Client-specific)
- Arbitrary mount and tmpfs configuration
- CephFS and OpenStack Manila

Slurm integrations

- OpenHPC
- LBNL Node Health Check
- Topology-aware scheduling
- MySQL

Portals

- Open OnDemand

Networking

- /etc/resolv.conf configuration
- Early boot gateway/default route configuration
- Proxy configuration
- /etc/hosts configuration
- Squid

Software

- NVIDIA GPU drivers and CUDA
- NVIDIA DOCA-OFED
- Podman
- EESSI
- Apptainer
- Local Pulp mirror
- Xfce desktop
- Remote apps

System

- journald
- Persistent hostkeys
- Tuned
- CA cert management
- chrony
- mdadm raid root disks
- systemd dropin configuration