

CERN Cloud Infrastructure status update

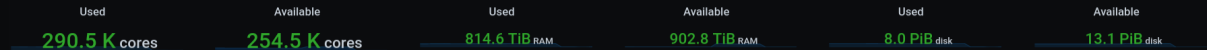
HEPiX Autumn 2020

Thomas Hartland (representing the CERN cloud team)

- What have we done this year?
- What will we be doing next?



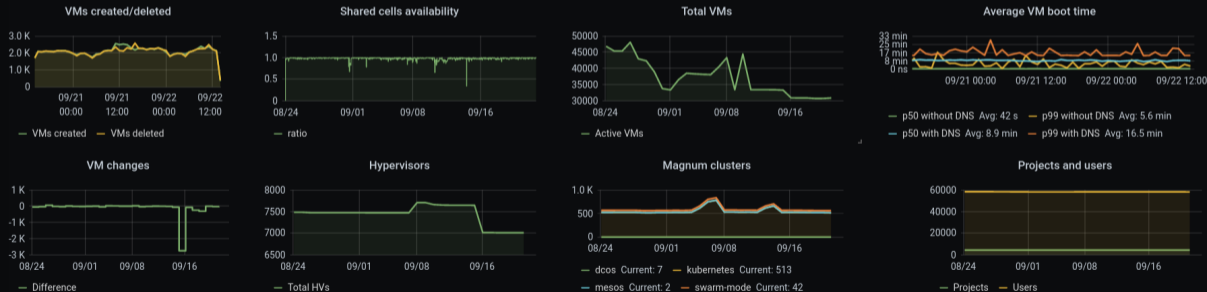
Cloud resources



Openstack services stats



Resource overview by time



Last year's plans

Development areas going forward

- Spot Market / Pre-emptible instances
- Software Defined Networking
 - Introducing LBaaS this month
- Magnum rolling upgrades
- Enrolling all 15K servers in Ironic
 - Containers on Bare Metal
- Evaluation of OpenStack Watcher for resource optimization
- Upgrading to Stein release



Last year's plans

Development areas going forward

- Spot Market / Pre-emptible instances ← **Pre-emptible service running in production**
- Software Defined Networking
 - Introducing LBaaS this month
- Magnum rolling upgrades
- Enrolling all 15K servers in Ironic
 - Containers on Bare Metal
- Evaluation of OpenStack Watcher for resource optimization
- Upgrading to Stein release



Last year's plans

Development areas going forward

- Spot Market / Pre-emptible instances
- Software Defined Networking
 - Introducing LBaaS this month
- Magnum rolling upgrades
- Enrolling all 15K servers in Ironic
 - Containers on Bare Metal
- Evaluation of OpenStack Watcher for resource optimization
- Upgrading to Stein release

In production in 20 projects,
~70 loadbalancers in total



Last year's plans

Development areas going forward

- Spot Market / Pre-emptible instances
- Software Defined Networking
 - Introducing LBaaS this month
- Magnum rolling upgrades
- Enrolling all 15K servers in Ironic
 - Containers on Bare Metal
- Evaluation of OpenStack Watcher for resource optimization
- Upgrading to Stein release

Possible but we encourage users to create new clusters



Last year's plans

Development areas going forward

- Spot Market / Pre-emptible instances
- Software Defined Networking
 - Introducing LBaaS this month
- Magnum rolling upgrades
- Enrolling all 15K servers in Ironic ← **Enrolled 1700 new servers**
 - Containers on Bare Metal
- Evaluation of OpenStack Watcher for resource optimization
- Upgrading to Stein release



Last year's plans

Development areas going forward

- Spot Market / Pre-emptible instances
- Software Defined Networking
 - Introducing LBaaS this month
- Magnum rolling upgrades
- Enrolling all 15K servers in Ironic
 - Containers on Bare Metal ← **Now possible, used by batch to deploy HTCondor with Kubernetes**
- Evaluation of OpenStack Watcher for resource optimization
- Upgrading to Stein release



Last year's plans

Development areas going forward

- Spot Market / Pre-emptible instances
- Software Defined Networking
 - Introducing LBaaS this month
- Magnum rolling upgrades
- Enrolling all 15K servers in Ironic
 - Containers on Bare Metal
- Evaluation of OpenStack Watcher for resource optimization
- Upgrading to Stein release

↑ Some scale issues in our cloud,
improvements contributed



Last year's plans

Development areas going forward

- Spot Market / Pre-emptible instances
- Software Defined Networking
 - Introducing LBaaS this month
- Magnum rolling upgrades
- Enrolling all 15K servers in Ironic
 - Containers on Bare Metal
- Evaluation of OpenStack Watcher for resource optimization
- Upgrading to Stein release

↑ Stein done, Train done, moving to Ussuri



Pre-emptible instances

- New OpenStack service developed by us, “Aardvark”.
- Manages the lifecycle of temporary low priority VMs.
- Integrates with Nova/Placement to free resources when needed.
- Multiple strategies (free on request, free above threshold).

Pre-emptible instances

- Used to give additional compute power to batch processing.
- Have to fine-tune placement of pre-emptible batch instances to avoid causing CPU-steal in overcommitted cells.
- Blog post: <https://techblog.web.cern.ch/techblog/post/preemptible-instances/>

CentOS 8

- Migration from CentOS 7 to 8 has started.
- Some control plane services have already moved.
 - For hypervisors we still need to work out a strategy.
- Some staying on C7 for now as their next migration will be to Kubernetes.
- Some Python 2/3 issues that we had to port our own tools/libraries for.



UEFI support

- New CentOS images now support both BIOS and UEFI.
 - Still a single image for VMs and physical machines.
- Had to be done for new hardware with no BIOS support.
 - Plus other benefits.

- Blog post:

https://techblog.web.cern.ch/techblog/post/bios_uefi_cloud_image/

Powercut

- November 2019, CERN wide powercut.
- Data centre ran on UPS power for 9 minutes.
 - Diesel generators started OK for critical power machines.
- One minute away from shutting down all batch machines, the power came back.
- Machines in the LHCb containers (no UPS) lost power.
- These were brought back online an hour later and kept idle.

LHCb containers



Photo: Brice, Maximilien and Ordan, Julien, 13/05/2019

LHCb containers



Photo: Brice, Maximilien and Ordan, Julien, 13/05/2019

LHCb containers

- Hosting returned Wigner hardware (~2000 servers).
- Some issues throughout the year.
 - Lack of redundant power has caused some downtime.
 - Ventilation issues in the summer.
- Overall a good investment for low-SLA batch workloads.
- On the cloud side, these containers are managed as a separate OpenStack region.

Hardware refresh

- Retiring three of the oldest shared cells in the CERN cloud.
- There were ~2500 VMs in these cells.
 - Can't be migrated due to network constraints in these cells.
 - Some have been there for a very long time and are not easy to move.
- VMs recreated using replacement capacity in existing cells.

Hardware refresh

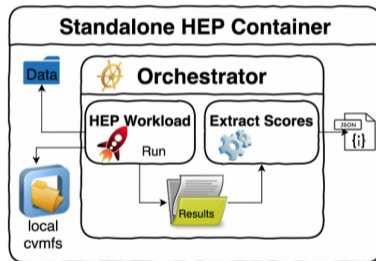
- This is a long and still ongoing process.
 - Have to contact the owner of each VM.
- 75% of the hardware is already decommissioned.
- Aiming to be complete by the end of the year.
- More decommissioning in 2021, should be more transparent.

GPUs

- Previously GPUs were handled by another team, not us.
 - And were distributed by passing access to the physical machine.
- This year we integrated existing GPUs as cloud resources.
- Access is given by allocating specific VM flavors to users' projects.
- Fully supported in Kubernetes.
 - Automatic setup of drivers, monitoring.
- Recently added 64 new NVIDIA T4 GPUs.

Benchmarking

- HEPscore as a replacement for HEP-SPEC06.
- New benchmarking frameworks:
 - Build infrastructure (containers)
 - Running benchmarks in containers
 - Systems to track, compare, share results
- New benchmarking methods:
 - Multi threaded benchmarks
 - GPU benchmarks
 - New architectures?



Ironic

- All new hardware is enrolled into Ironic, but the adoption of existing servers hasn't started yet.
- Integrating benchmarking as part of the Ironic enrollment process.
- Ongoing work for:
 - Hardware inventory tracking
 - Redfish (IPMI replacement)

Kubernetes

- Steady growth of Kubernetes use at CERN.
- New Kubernetes templates: v1.17, v1.18, v1.19.
- Last features to be truly “production ready” this year:
 - SDN gives us external loadbalancers and multi-master clusters.
 - Node groups for splitting clusters across availability zones.
 - We’re confident enough now to move parts of the OpenStack control plane to Kubernetes.

Anomaly detection

- Research into anomaly detection using deep learning methods.
- Two sets of hypervisor metrics for training.
 - Shared cells
 - Batch cells
 - It should be possible to make the data public.
- Integration with grafana monitoring and alerting systems.
- Operators can review/add grafana annotations for anomalies that feed back into the model training.

Plans for the coming year

- Block storage availability zones
- GPUs as managed resources in OpenStack (quotas, vGPUs)
- Integration with CERN identity management
- More hardware replacement
- CentOS 8
- Ironic enrollment of production machines
- More SDN (tenant networks, virtual routing)
- OpenStack upgrades

Questions?