

The background features a collage of scientific documents. On the left, there's a handwritten diagram with Chinese characters and mathematical symbols. In the center, a blue circular graphic contains mathematical equations like $\nabla \cdot \mathbf{E} = \rho$ and $\nabla \times \mathbf{B} - \frac{1}{c} \frac{\partial \mathbf{E}}{\partial t} = \mathbf{j}$. On the right, there's a snippet of a periodic table showing elements like H, He, Li, Be, B, C, N, O, F, Ne, Na, Mg, Al, Si, P, S, Cl, Ar, K, Ca, Sc, Ti, V, Cr, Mn, Fe, Co, Ni, Cu, Zn, Ga, Ge, As, Se, Br, Kr, Rb, Sr, Y, Zr, Nb, Mo, Tc, Xe, I, Ba, La, Ce, Pr, Nd, Pm, Sm, Eu, Gd, Hg, Tl, Pb, Bi, Po, At, Rn, Fr, Ra, Ac, Th, Pa, U, Np, Pu, Am, Cm, Bk, Cf, Es, Fm, Md, No, Lr, and their atomic weights. There are also some handwritten notes and diagrams scattered throughout.

Optimizing OpenStack Nova for Scientific Workloads

CHEP 2018 - Bulgaria, 2018

Belmiro Moreira

belmiro.moreira@cern.ch

[@belmiromoreira](https://twitter.com/belmiromoreira)

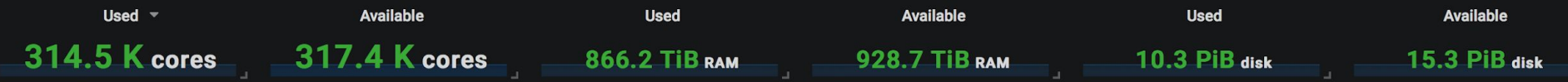
Outline

- CERN Cloud Infrastructure
- Virtualization Overhead challenge
- Resource Utilization challenge
- CERN/SKA collaboration
 - How to Remove Virtualization Overhead
 - How to Increase Resource Utilization

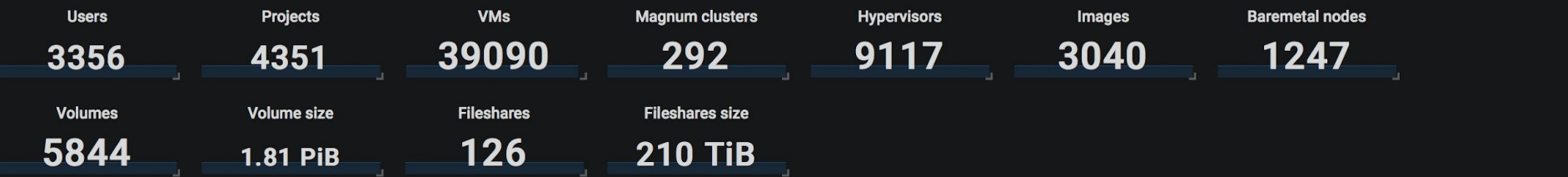
CERN Cloud Infrastructure

- Production service since July 2013
- OpenStack based
 - Offering more than 15 OpenStack projects
 - Keystone, Nova, Neutron, Glance, Cinder, Magnum, Ironic, ...
- Running the last OpenStack release (Queens)
- 2 Data Centres (Geneva and Budapest)
 - One Region
 - 5 Availability Zones
 - Nova Cells (>70 Cells)
- More details about the OpenStack services provided by CERN Cloud
 - (Tue@11:15 - T7, S3) Advanced features of the CERN OpenStack Cloud

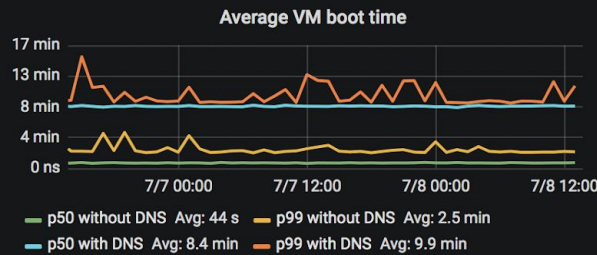
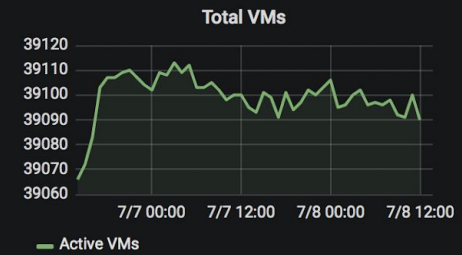
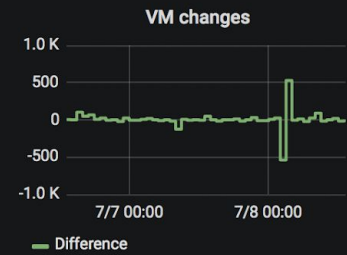
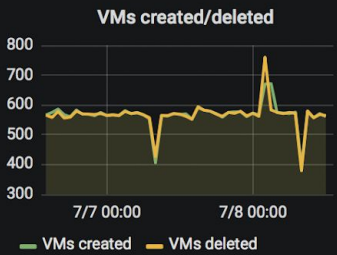




Openstack services stats



Resource overview by time



Virtualization Overhead/Resource Utilization

- >80% of the CERN Cloud resources are dedicated for Batch processing
- CPU Performance Optimizations deployed in the CERN Cloud
 - NUMA + CPU pinning
 - 2MB Huge Pages
 - EPT (enabled)
 - KSM (enabled)

Current Virtualization CPU Overhead : **~3%**

- Virtualization Overhead depends on VM size
 - 4 VMs (8 cores each) the Virtualization overhead reduced from 7.8% to 3.3%
 - 2 VMs (16 cores each) the Virtualization overhead reduced from 16% to 4.6%
- Quotas are per Project
 - No Overcommit. Unused quota = Unused resources



HL-LHC



SKA

How to remove the Virtualization Overhead?

- Containers provide a lightweight alternative to VMs
 - Small footprint
 - Increased performance
- OpenStack Magnum
 - Container Orchestration Engine
 - Easy to deploy Kubernetes/Swarm clusters
 - Clusters are per tenant
- Clusters usually deployed using VMs
 - Virtualization Overhead!
 - VMs are used for Security
 - Flexibility VS Performance



MAGNUM
an OpenStack Community Project

Containers on Baremetal

- Containers deployed directly on Baremetal
- OpenStack Ironic
 - Baremetal provisioning using OpenStack Nova APIs
- Magnum to deploy clusters directly in Baremetal
 - CERN and SKA working to add Baremetal support into Magnum
 - Several assumptions required
 - network, location, timeouts, ...
 - Ironic as provision engine
 - Kubernetes and Swarm supported



Containers on Baremetal

- Containers on Baremetal performance
 - Similar to native Baremetal
- SKA testing Containers on Baremetal in ALaSKA prototype
 - Performance and Fast context switching
- CERN evaluating how to run Batch in Containers on Baremetal
 - More RAM available for Jobs
 - More CPU time for Jobs throughput
- Containers on Baremetal can benefit other workloads
 - (Wed@12:00 - T7, S5) Apache Spark usage and deployment models for scientific computing

How to improve Resource Utilization?

- Public Clouds give the illusion of infinite capacity
 - Users pay for resources that they use
- Private Clouds
 - Resource management usually is based in project quotas
 - Prevent resources being exhausted
 - Prevent “over-committing” resources/quota
 - Manage individual projects requirements
 - Reserve resources for operations with higher priority
 - Scientific Clouds
 - Projects have different funding models
 - They expect a predefined number of resources available
 - But not always these resources are used full time



Preemptible/Spot Instances

- Public Clouds
 - Based on different pricing/SLA considering resource availability
 - Reserved instances vs spot-market
- Private Clouds
 - Quotas are hard limits. Leads to a reduction in resource utilization
 - Preemptible instances
 - Projects that exhausted their quota can continue to create instances
 - Opportunistic workloads
 - Low SLA

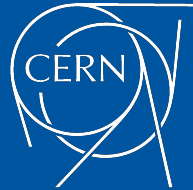
Preemptible Instances - OpenStack Nova

- Preemptible Instances Workflow in OpenStack Nova
 - The creation of a non preemptible VM fails because there aren't available resources
 - Instances that fail with "Nova Valid Host", go to "PENDING" state instead of "ERROR"
 - The Reaper service is notified and it tries to free the requested resources
 - Rebuild the instance
 - Or change instance state to "ERROR"
- CERN and SKA are working with OpenStack Nova team to implement Preemptible Instances
- CERN is deploying a Preemptible Instances prototype
 - Expected to be ready by the end of Q3/2018

Summary

- HL-LHC and SKA will produce an unprecedented amount of data to analyse
- Small compute inefficiencies in large infrastructures translate in a huge number resources underutilized
- Flexibility vs Efficiency
- CERN and SKA working together to build a High Efficient Infrastructure
- How to Remove Virtualization Overhead?
 - Containers on Baremetal
- How to Improve Resource Utilization?
 - Preemptible Instances

belmiro.moreira@cern.ch
@belmiromoreira



www.cern.ch