

The CERN internal Cloud CERNIT

The CERN internal Cloud status report

Sebastien Goasguen, Belmiro Rodrigues Moreira, Ewan Roche, <u>Ulrich Schwickerath</u>, Romain Wartel

HEPIX2010, Cornell university

See also related presentations:

- Related presentations by Tony Cass, Romain Wartel, and myself
- Virtualization at CERN, HEPiX spring meeting 2010
- Batch virtualization at CERN, HEPIX autumn meeting 2009
- Virtualization, HEPIX spring meeting 2009
- Virtualization vision, GDB 9/9/2009 and HEPIX
- Batch virtualization at CERN, EGEE09 conference, Barcelona,



CERN IT Department CH-1211 Genève 23 Switzerland www.cern.ch/it

Outline



CERN

laaS for CERN Principles Components Applications Principles Components **Project status** Summary

CERN IT Department CH-1211 Genève 23 Switzerland www.cern.ch/it





... if we can call it a cloud

CERN

Department

3 CERN

CERN IT Department CH-1211 Genève 23 Switzerland www.cern.ch/it





Requirements: the approach must

be easy to manage

smoothly integrate with the existing infrastructure

be scalable

be flexible for extensions

CERN IT Department CH-1211 Genève 23 Switzerland **www.cern.ch/it**



Principles



Design decisions:

- setup of an internal infrastructure for virtualization
- With the potential to become a public cloud

The initial deployment foresees that

- we deliver laaS internally to IT service managers
- NOT (yet?) to end users
- the first and for now only customer is the batch service

CERN IT Department CH-1211 Genève 23 Switzerland **www.cern.ch/it**





The hypervisor cluster (lxcloud):

Supports XEN and KVM (XEN being dropped now)

Is centrally managed with a unique software setup

Provides the infrastructure required for laaS

Runs the latest OS version to fully support the hardware

CERN IT Department CH-1211 Genève 23 Switzerland **www.cern.ch/it**



The local image cache:

- Is a central place holding all approved images
- Is designed in collaboration with the HEPiX WG

See the presentation of Romain Wartel on VMIC

Internal image distribution

Caching of images on the hypervisors
 Transfer using bit-torrent

CERN IT Department CH-1211 Genève 23 Switzerland www.cern.ch/it



The image provisioning system

- Receives requests
- Selects a hypervisor
- Provisions the machine
- Controls and monitors it during its lifetime

CERN IT Department CH-1211 Genève 23 Switzerland **www.cern.ch/it**

CERN**T** Department





ERN

ISF: InfraStructure sharing facility, Platform computing



ONE: OpenNebula

Modular OpenSource solution
 Now also professional support
 Cloud interfaces

CERN IT Department CH-1211 Genève 23 Switzerland www.cern.ch/it

ISF and ONE relationship

pVMO and ONE do approximately the same thing



Hypervisor cluster (physical resources)

CERN IT Department CH-1211 Genève 23 Switzerland www.cern.ch/it



CER

Department





Batch resources as application

Or: virtualization of batch resources

CERN IT Department CH-1211 Genève 23 Switzerland www.cern.ch/it 13 CERN

The batch application design CERNIT

Requirements

- transparent to the end users
- possible to mix virtual machines and physical worker nodes
 - possible to start small and grow
- scalable
- able to solve the operational issues



Application components

A Golden node:

Is a centrally managed virtual machine

Is a clone of physical machines running the same services

CER

- Does not execute jobs by itself
- Serves as source for virtual machine creation

Adding automation and elasticity(*)

Limit the active life time of the virtual machines (to 24h)

Automatically Monitor the VM population and restart nodes as needed

(*) Batch service managers point of view

CERN IT Department CH-1211 Genève 23 Switzerland **www.cern.ch/it**

The internal cloud infrastructure at CERN - a status report - 16

Department

ERN

Application components

How to manage **intrusive interventions** in a fragmented setup?



Notes:

NEW virtual machines always start with the latest image Image A and Image B can correspond to different OS versions

CERN IT Department CH-1211 Genève 23 Switzerland www.cern.ch/it

ISC-cloud 2010 – Frankfurt / Main



The question is: will it work at the required scale ?

Or: will it ever fly ?

CERN IT Department CH-1211 Genève 23 Switzerland www.cern.ch/it



CERI

Department

Scalability at all levels

- General infrastructure (eg. Idap)
- Networking: number of IPs
- Virtual machine provisioning systems
- Image distribution efficiency
- Application level:
 - Virtual batch nodes
 - \triangleright Batch system scalability \rightarrow dedicated presentation

CERN IT Department CH-1211 Genève 23 Switzerland **www.cern.ch/it**



CER



Virtual machine provisioning system tests



Time

CERN IT Department CH-1211 Genève 23 Switzerland www.cern.ch/it

ட

Nodes seen in LS

OpenNebula:

Up to 16,000 VMs started, all available slots filled

► ISF/pVMO:

Up to ~10,000 VMs started
Limitation is understood and no longer present

CERN IT Department CH-1211 Genève 23 Switzerland www.cern.ch/it

The internal cloud infrastructure at CERN - a status report - 27

CER



Image Distribution



Virtual batch nodes with real userjobs

Currently we have:

- 12 KVM hypervisors ready for production
- hypervisors with private IPs
- And 8 VM slots with public IPs
- 2 KVM hypervisors deployed into production
- 16 KVM VMs running as public batch nodes for ~4 weeks

CERN IT Department CH-1211 Genève 23 Switzerland **www.cern.ch/it**

Virtual machines in public batch



New feature which is:

not seen on physical worker nodes
 caused by network monitoring tools used at CERN
 under investigation

CERN IT Department CH-1211 Genève 23 Switzerland www.cern.ch/it

The internal cloud infrastructure at CERN - a status report - 24

CER

Department

ERN

Changes since last HEPiX

Dropped support for SLC4

- Favor KVM over XEN
 - KVM is the emerging technology
 - XEN support for in RHES/SL6 is questionable
 - CPU usage monitoring works out of the box with KVM



CERN IT Department CH-1211 Genève 23 Switzerland www.cern.ch/it

The internal cloud infrastructure at CERN - a status report - 25

CERN

KVM in SLC5, impressions

Fairly stable

- Some performance issues
 - CPU performance O(10%) low
 - I/O performance even more
- Requires tuning, and eventually a newer version
- KVM processes stay in S mode
- Utilization turns up as system CPU, not user CPU

Impementation status

	Hypervisor cluster	SLC5 virtual batch nodes
Quattor managed	OK	OK, via golden node
Lemon monitored	OK	OK
Auto-registration	OK	OK
Central maintenance	OK	OK
ISF support	OK	OK
ONE support	OK	OK

CERN IT Department CH-1211 Genève 23 Switzerland www.cern.ch/it

The internal cloud infrastructure at CERN - a status report - 27

CERN

Department

CERN

Implementation status

VM kiosk and VM Batch Hypervisor image management application cluster distribution system Initial OK OK OK OK deployment Central **ISF OK, ONE** OK OK OK missing management Monitoring and OK OK OK missing alarming

CERN IT Department CH-1211 Genève 23 Switzerland **www.cern.ch/it**

The internal cloud infrastructure at CERN - a status report - 28



CERN

Plans and remaining technical issues:

VM CPU accounting / CPU factors

- ISF adaptive cluster may come late
- Still some development needed for VM renewal
- Performance considerations

For scaling up:

- Krb5 authentication for root access may not scale
- May need our own lemon instance ?
- LSF doesn't scale up to 15k nodes, redesign is needed

)epartment

Summary



Quite some progress during summer 2010

Only one application so far (Ixbatch) 16 VMs under test with real payload Up to 96 by Christmas Scale up Expected for 2011 Depends on experiences

Any questions

Benchmarking



No tunings applied
 Up to 17% penalty in CPU with current KVM
 Same order or worse for IOzone
 Needs further investigation and tuning

Tunings: kvm kernel module options, CPU affinity

CERN IT Department CH-1211 Genève 23 Switzerland www.cern.ch/it



Department

Benchmarking





CPUs / VM

CERN

Department

Bare metal: 8 Cores, 8 processes: ~96 HS06, so 17% performance penalty

CERN IT Department CH-1211 Genève 23 Switzerland www.cern.ch/it



Benchmarking



VMs/hypervisor

CERN

Department

CERN

Bare metal: ~12 HS06/core, so 12% performance penalty

CERN IT Department CH-1211 Genève 23 Switzerland **www.cern.ch/it**

I/O Benchmarking

CERN

Department

VM lozone read performance KVM **Bare metal** 120000 100000 80000 60000 40000 20000 0

4

CERN IT DE, CH-1211 Genève 23 Switzerland www.cern.ch/it 1

2

3

The internal cloud infrastructure at CERN - a status report - 35

6

7

8

CERN

5

I/O Benchmarking

lozone write performance **KVM** VM 25000 20000 15000 10000 5000 0 2 3 5 6 7 1 4 8

CERN IT Department CH-1211 Genève 23 Switzerland **www.cern.ch/it**

The internal cloud infrastructure at CERN - a status report - 36

VM Bare metal



CERN