# **Virtualisation et approche cloud au CERN**

Manuel Guijarro
On behalf of CERN IT-PES-PS

- Introduction
- Why Virtualisation?. Why Cloud Computing?
- Service Consolidation Service
- Internal Cloud
- Summary and future options

- IT-PES-PS section:

  - CERN Batch service (+4000 WNs)
  - Interactive Login service
  - Grid Computing services
  - Infrastructure Services (DNS Load Balancing, etc)
  - **Service Consolidation Service (aka SCS)** ~500 VMs in ~100 compute nodes
  - **Internal Cloud (aka LXCLOUD)** 432 VMs in 48 compute nodes

- Physical machines: 4300+ in 12000+

- The full list of acronyms includes: Batch Service; BatchVM; **LXCLOUD**; SCAS; LXPLUS; LXADM; GridLFC; GridCE (LCG); GridCE (Cream); GridCE info; GridWMS; GridLB; GridFTS; VOBox; CERNVMFS; ActiveMQ; GridBDII; GridMyProxy; GridMonBox; GridCAProxy; **Virtual Machines (consolidation);** PSSecurityOfficers; GliteVOMS; GliteVOMRS; ROCLayer; WN Publishing; CERN-PROD nagios instance; DNSLoadBalancing; RSysLog

# Low Infrastructure Utilization

Typically one application per server to avoid the risk of vulnerabilities in one application affecting the availability of another application on the same server

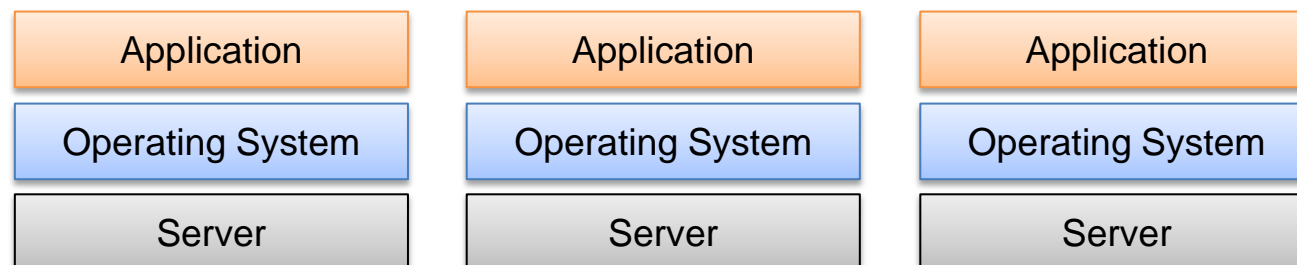# Increasing Physical Infrastructure Costs

Power consumption, cooling and facilities costs that do not vary with utilization levels
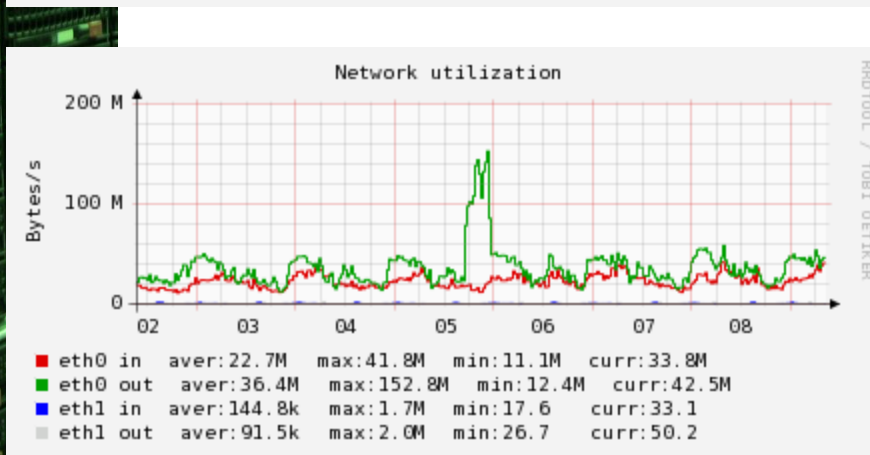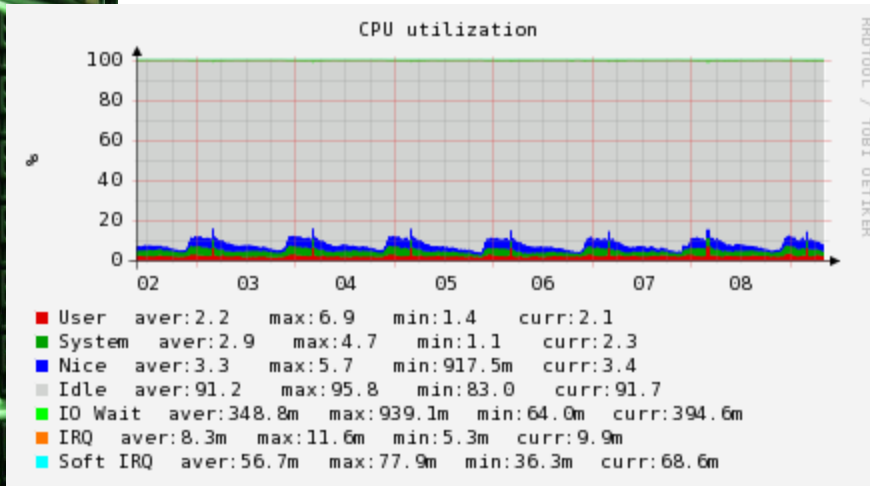
# Increasing IT Management Costs

Spend disproportionate time and resources on manual tasks associated with server maintenance, and thus require more personnel to complete these tasks

# Insufficient Failover and Disaster Protection

The threat of security attacks, natural disasters and terrorism has elevated the importance of business continuity

| Application | Application | Application |
| Operating System | Operating System | Operating System |
| Server | Server | Server |

CPU utilization

| | | | | |
|---|---|---|---|---|
| ■ User | aver:2.2 | max:6.9 | min:1.4 | curr:2.1 |
| ■ System | aver:2.9 | max:4.7 | min:1.1 | curr:2.3 |
| ■ Nice | aver:3.3 | max:5.7 | min:917.5m | curr:3.4 |
| ■ Idle | aver:91.2 | max:95.8 | min:83.0 | curr:91.7 |
| ■ IO Wait | aver:348.8m | max:939.1m | min:64.0m | curr:394.6m |
| ■ IRQ | aver:8.3m | max:11.6m | min:5.3m | curr:9.9m |
| ■ Soft IRQ | aver:56.7m | max:77.9m | min:36.3m | curr:68.6m |



Network utilization

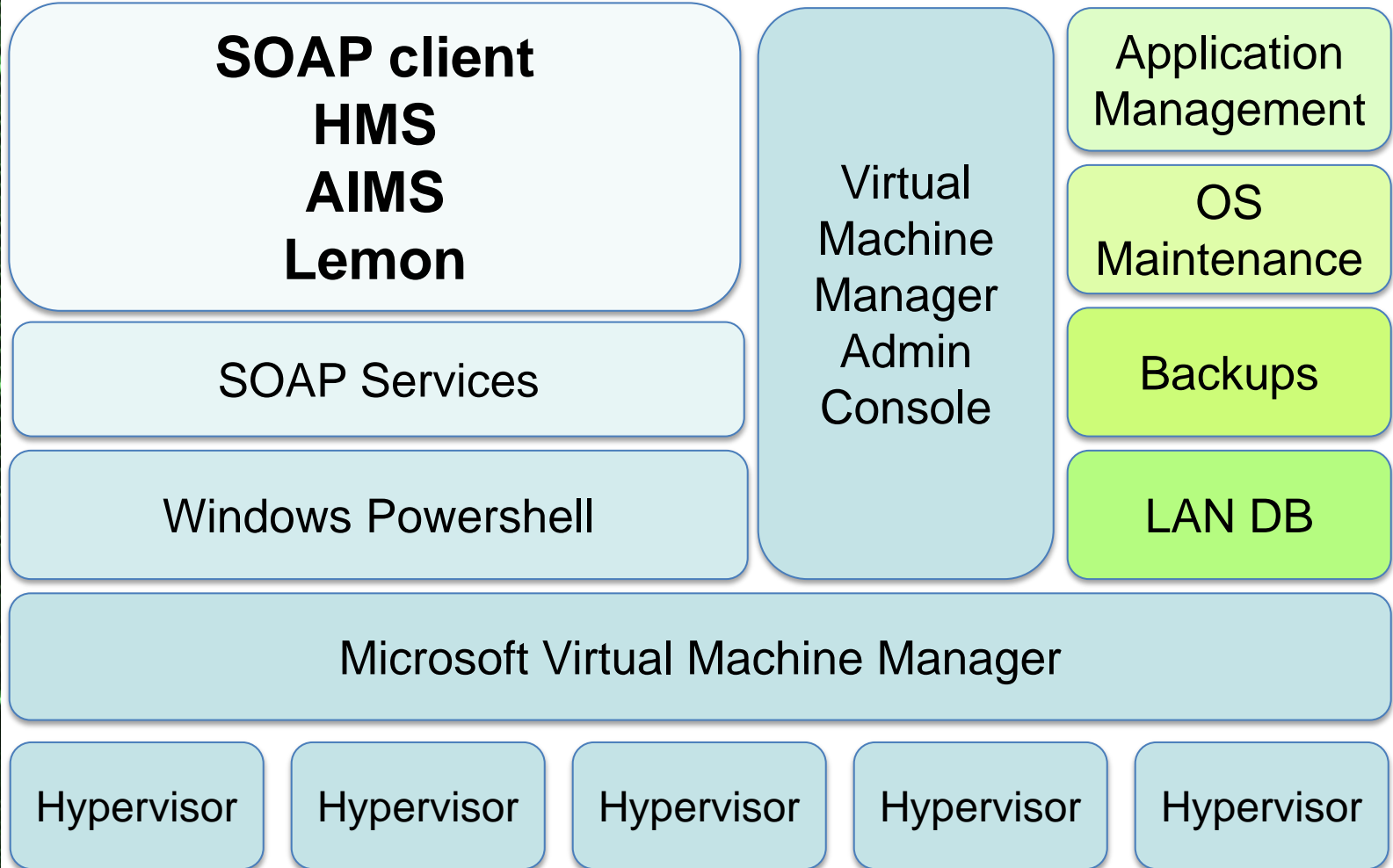| | | | | |
|---|---|---|---|---|
| ■ eth0 in | aver:22.7M | max:41.8M | min:11.1M | curr:33.8M |
| ■ eth0 out | aver:36.4M | max:152.8M | min:12.4M | curr:42.5M |
| ■ eth1 in | aver:144.8k | max:1.7M | min:17.6 | curr:33.1 |
| ■ eth1 out | aver:91.5k | max:2.0M | min:26.7 | curr:50.2 |

➢ The voatlas cluster – 138 physical machines (in 2009)

➢ Most CPU, network and disk load appears to be from nightly builds

➢ Otherwise very low IO rates

➢ Even a consolidation factor of 2 saves a large amount of resource

➢ In industry a consolidation factor of at least 4 is normally possible for servers with no impact on performance

Reduce hardware needs by providing VOBoxes in the form of Virtual Machines
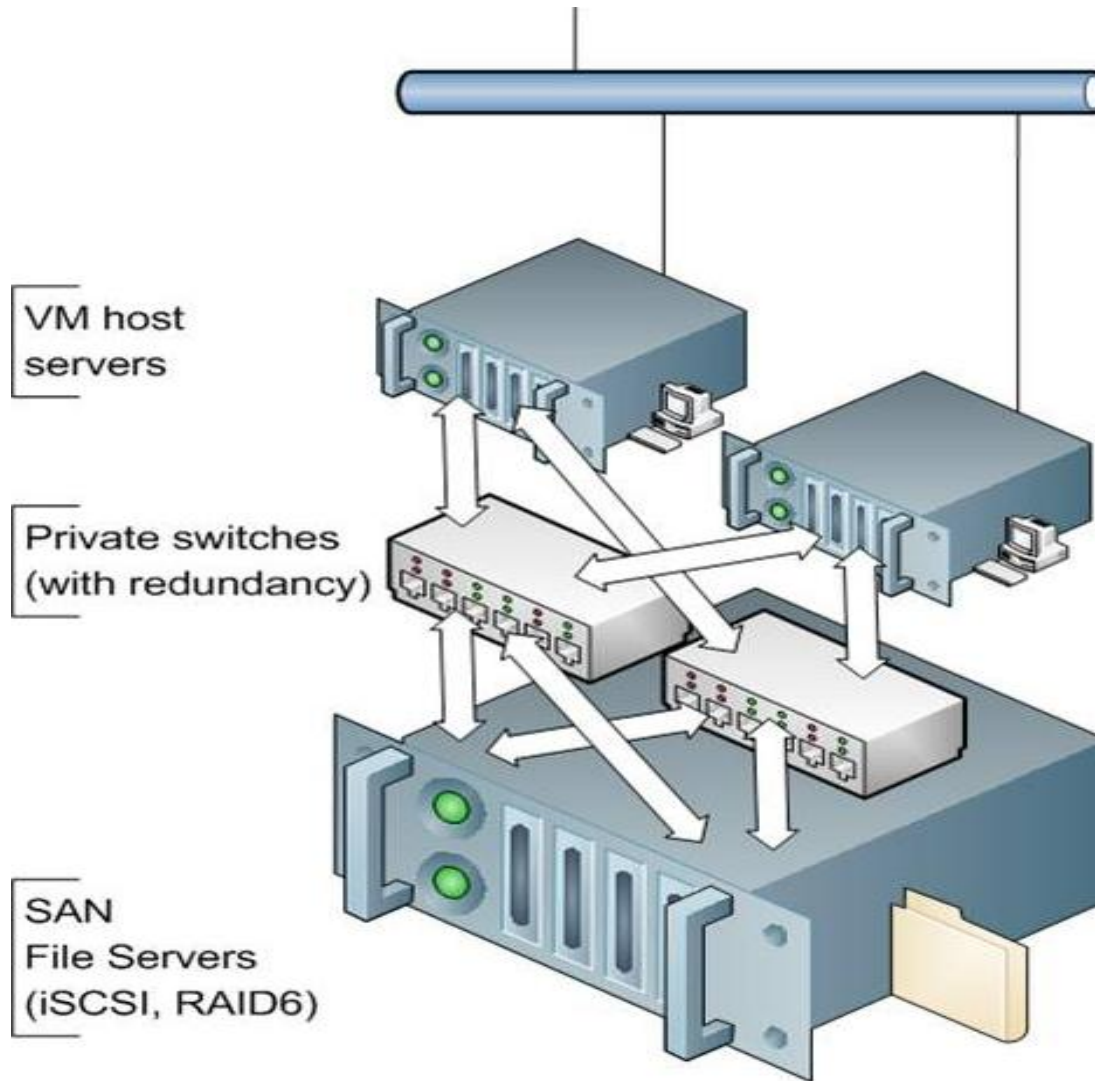
Service Definition: *Integrate the VoBOX service with a virtual machine infrastructure in CERN's computer centre. It enables the provisioning, configuration, administration, monitoring, etc of **Quattor-managed** virtual machines, at least as easily as with physical machines.*

➢ The base hypervisor technology is provided by IT/OIS and is based upon
- ✓ Hyper-V 2.0
- ✓ System Centre Virtual Machine Manager 2008 R2 (SCVMM)
- ✓ Custom SOAP interface

➢ The highlights include
- ✓ SCVMM is an enterprise class management suite
- ✓ Live migration
- ✓ High availability in the case of hypervisor failure
- ✓ Integration with LanDB via the SOAP interface
- ✓ Full support for Linux guests with up to 4 cores
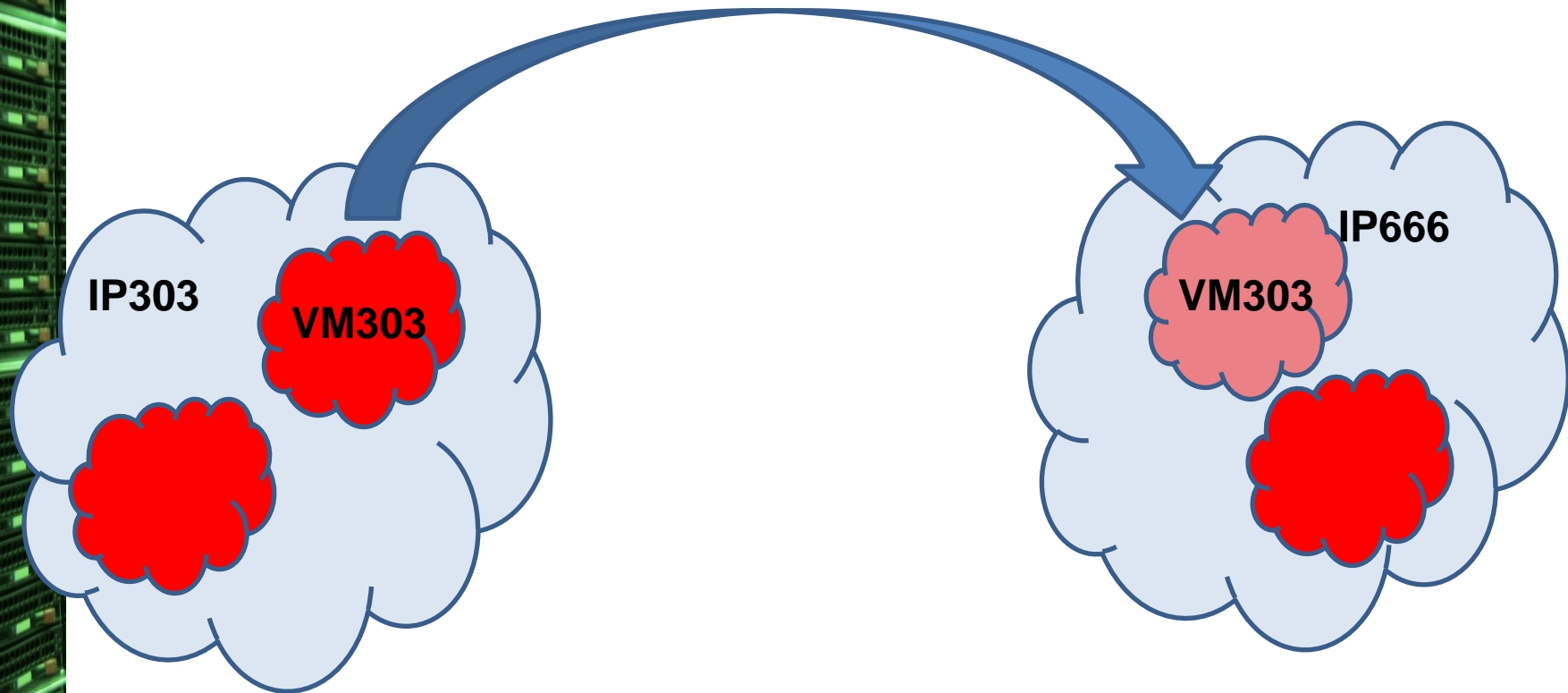- ✓ Paravirtualised I/O drivers integrated into SLC5

**SOAP client
HMS
AIMS
Lemon**

SOAP Services

Windows Powershell

Virtual Machine Manager Admin Console

Application Management

OS Maintenance

Backups

LAN DB

Microsoft Virtual Machine Manager

Hypervisor

Hypervisor

Hypervisor

Hypervisor

Hypervisor

VM host servers

Private switches (with redundancy)

SAN File Servers (iSCSI, RAID6)

CERN IT Department

| Name | Sockets/ Cores | Memory | Storage capacity | Disk enclosures | Network | Location |
|---|---|---|---|---|---|---|
| **513critical** lxfssw11[01-08] lxfssw11[09-16] | 2/8 | 48GB | 1x10TB | 2 x baby sumos | GPN | Critical UPS in 513 |
| **General\StorageA** lxfssm38[01-16] | 2/8 | 48GB | 20x14TB | 8 x sumos | OPN | 513 |
| **General\StorageB** lxfssm36[01-16] | 2/8 | 48GB | 20x14TB | 8 x sumos | OPN | 513 |
| **General\StorageC** lxbsm32[01-08] lxbsm32[09-16] | 2/8 | 24GB | 4x14TB | 4 x sumos | OPN | 513 |
| **General\StorageD** cviclr04[01-16] | **2/8** | **96GB** | **6x10TB** | **Equallogic** | **OPN** | **513 *new*** |
| **General\StorageE** cviclr03[01-16] | **2/8** | **96GB** | **6x10TB** | **Equallogic** | **OPN** | **513 *new*** |
| **SafeHost** lxbsm34[01-08] lxbsm34[09-16] | 2/8 | 48GB | 4x15TB | 4 x sumos | OPN | SafeHost |

www.cern.ch/it

**IP303**

**VM303**

**IP666**

**VM303**

Networking is based on secondary (virtual) IP services which are attached to primary services – it is possible to migrate a secondary service between primary ones but downtime is required. **A virtual service cannot span more than one primary service.**

Close to 400 VMs totalling

   1800 GB of memory

   715 virtual CPUs

   23 TB of allocated disk space (out of 140TB committed)

Main clients by number of VMs:

   ATLAS (88)

   PES/PS (57)

   CMS (34)

   Arda (27)

We offer Quattor-managed **SLC5** VMs via the [Snow hardware procurement form](#) with (some) combinations of:

- 1-4 CPUs
- 1-8 GB memory
- 100-2000 GB disk
- 1Gbps paravirtualized network
    - 100Mbps during installation

CPU, disk and network are happily overcommitted
- Typical physical CPU usage on hypervisors < 30%
- Typical physical network usage < 2%
- Real disk usage vs. committed capacity < 20%

Memory is not overcommitted

Our shared-storage clusters are expensive bits of hardware enabling "High Availability" to keep VMs running 24/7

If an hypervisor fails (network, hardware, software failure...) other hypervisors can take over the VMs it was running

Transparent live-migration of VMs between hypervisors to even the load or perform maintenance on hosts

SAN storage with redundancy (network paths, RAID...)

But HA is relative...

Transparent patch deployment on hypervisors in Sep-2011(and firmware upgrades in August): OK

"Equallogic incident" in Aug-2011: SAN disconnection with redundant controller not taking over. 72 VMs unavailable for several hours. Firmware subsequently upgraded.

- **The limiting resource is VM memory (in GB)**
- **Current capacity probably sufficient until end 2012**
- **1st-generation hosts will reach end of warranty in January 2013**

- The usual cloud advantages:

  - Maintenance
    - Decoupling from the underlying hardware
    - Easy migration between OS / configurations
  - Dynamic Environment
    - The composition of compute nodes can be rapidly adapted to optimize resources utilization
  - Encapsulation
    - Virtualization allows the separation between the user environment and the hardware

# LxCloud – brief history…

## 2009

- XEN hypervisor;
- ONE;
- Only for Batch resources;

## 2010

- XEN / KVM hypervisor;
- ONE / Platform ISF;
- IaaS concept;
- Scalability Test;
- VM Image Catalog (VMIC);
- Image distribution based on Bittorrent;
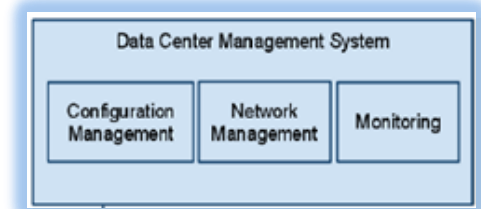- VirtualBatch with 96 VMs;

## 2011

- KVM hypervisor;
- SLC6 migration;
- ONE / OpenStack;
- VMIC improvements;
- EC2 service evaluation;
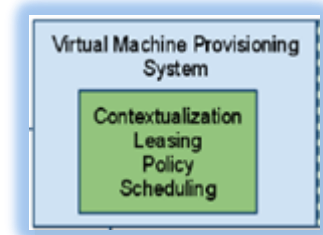- VirtualBatch with 432 VMs;

- LxCloud is running on standard compute nodes used by the batch service
  - 2 x intel xeon L5520 @ 2.27 GHz (8 cores)
    - Nehalem architecture
  - 58 compute nodes
    - 48 nodes in production
    - 10 nodes for tests
  - 24 GB of RAM
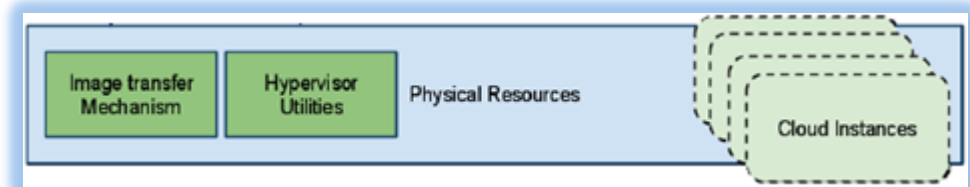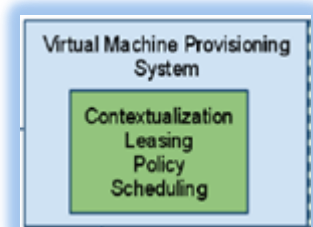  - 3 TB of local disk space

Data Center Management System — Configuration Management | Network Management | Monitoring

- ## Configuration Managment
  - ### LxCloud is integrated with the Fabric Management tools used by CERN
    - Quattor managed pool of resources
    - Alarming with LAS (Operator)
    - Hardware management by sys-admin team
    - "Draining" via sms state management
- ## Network Managment
  - Pre-allocation of VM "slots" in network DB
  - Compute node "knows" the available "name" of its guests
- ## Monitoring
  - All Compute nodes are monitored by Lemon

Virtual Machine Provisioning System

Contextualization
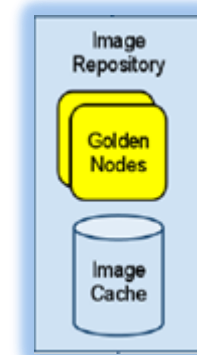Leasing
Policy
Scheduling

- **Virtual Machine Provision System**
  - **OpenNebula**
    - Migration from ONE 2.2 to ONE 3.0
    - Consolidation of different ONE instances in a single server
    - ONE server running on SLC6
  - **Evaluation of OpenStack**

- Compute nodes are running SLC6
  - Updated versions of libvirt and kvm
  - KSM (Kernel Samepage Merging) enabled

Virtual Machine Provisioning System

Contextualization
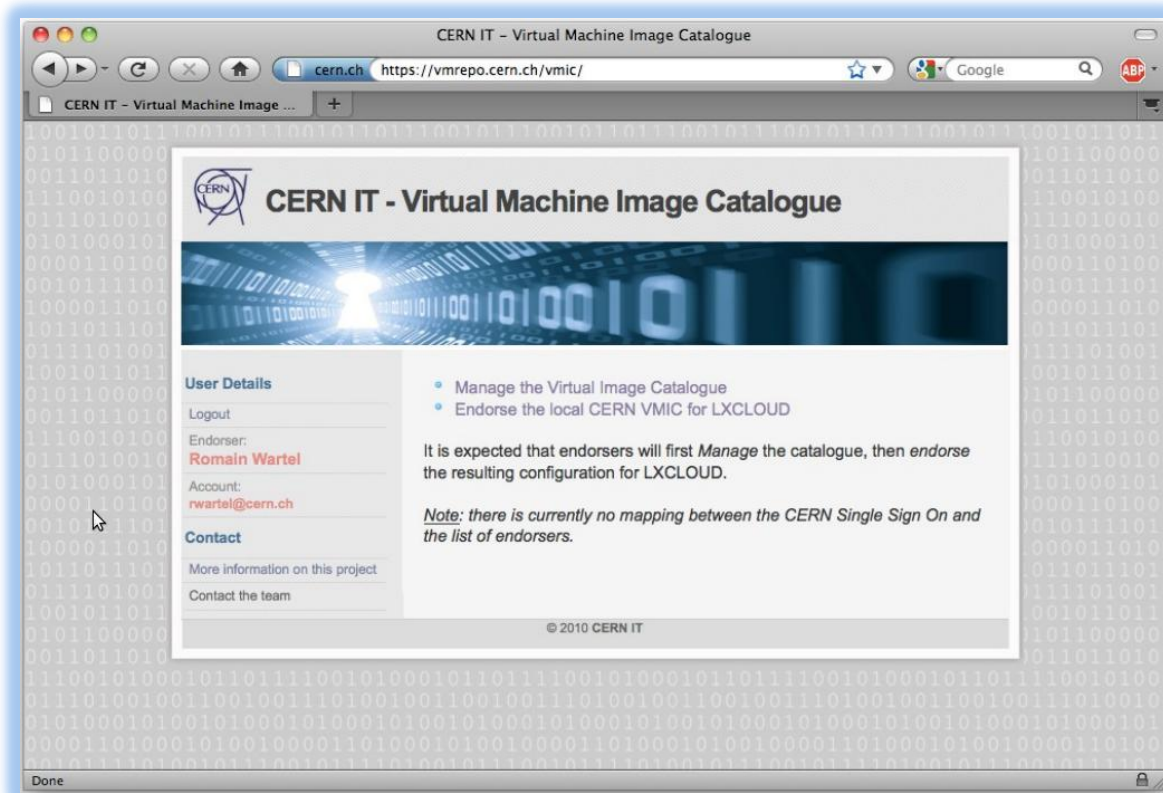Leasing
Policy
Scheduling

- **Contextualization**
  - ISO file with contextualization information is attached into the guest during the boot time
    - Contextualization scripts were defined inside the HEPiX virtualization WG
    - CernVM supports this contextualization model
- **Scalability tests**
  - In 2010 various scalability tests were made to evaluate the provision systems, LSF scheduler, and all LxCloud infrastructure
    - 16000 VMs running in about 500 compute nodes
    - LSF scalability concern

- Image Repository workflow
  - Golden Nodes
    - Images with the desired configuration are installed in a "Golden Node"
      - PXE installation / or not, using a slot on a compute node
    - This image is compressed and transfered to the "Image Cache"
      - If the image is quattor managed it needs to be "de-quattorized" and clean
  - Image Cache
    - It's managed by VMIC – Virtual Machine Image Catalogue

- VMIC – virtual machine image catalogue
  - Using HEPiX virtualization WG specifications
    - Image sharing successful between Clemson University and CERN;
  - See "CloudMan and VMIC projects overview" presentation at HEPiX Fall 2011
- Image distribution management software
  - Responsible for managing images in each compute node
  - Images are pre-staged in the compute nodes;
  - BitTorrent image transfer
    - rtorrent client installed in all compute nodes

- VMIC – virtual machine image catalogue

- BitTorrent protocol for image distribution

- Image flow
  - The compute node query VMIC for the image list
  - The image integrity is verified
  - Check which images are for the node
  - Download the torrent files for the desired images
  - rtorrent downloads the images
  - Image management software in the nodes deploys the new images

- Using RAW format image
  - Images are stored in LVs
  - Using LVM snapshot functionality
- Evaluation of QEMU format image
  - Store images in files
  - Use "COW" functionality of QEMU image format
    - The link between the "main" image and the "snapshots" is hardcoded
      - If the "main" image is renamed the "snapshot" metadata is not updated

- Services running on top LxCloud

  - virtualBatch
    - Production service since 2010
    - Gradually increasing capacity
      - 96 VMs (2010)
      - 432 VMs (2011);
  - EC2
    - First tests in 2011
    - Restricted set of users
    - Restricted set of images

CERN
IT
Department

- virtualBatch characteristics

  - Instances are always derived from the newest available golden image
    - VM TTL is set to 48 hours
  - Customized at boot time (contextualization)
  - Instances are manageable by Quattor
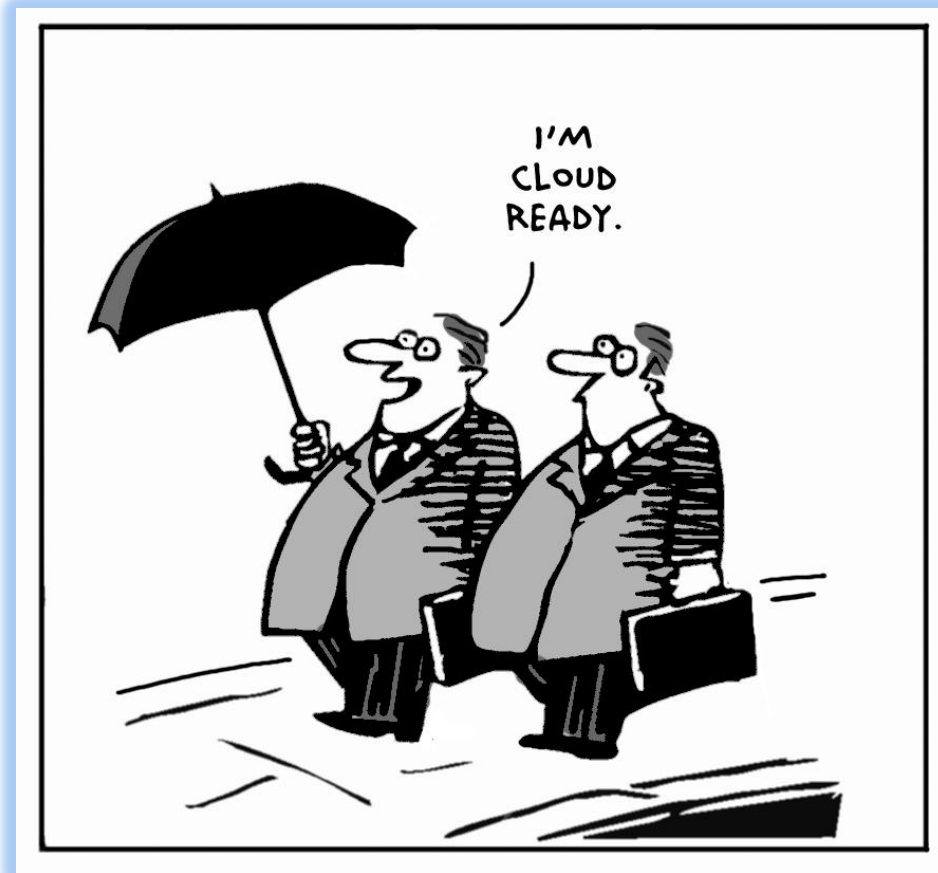
- Running batch jobs:

  - 48 compute nodes (384 cores)
    - 432 VMs (432 job slots)
    - 9 VMs per node
      - CPU pinning
  - 2.6 GB per VM
    - Memory overcommitted
    - Large swap but rarely used

CERN**IT**
Department

- Access for restricted users - only on request
- Predefinied set of images
  - Users can't upload their own images
- Using EC2 driver for OpenNebula
  - econe driver
  - Amazon EC2 API is not totally suported by ONE

- SCS stable production service growing at a good pace
- LxCloud proven flexible and scalable
  - Seamless integration into the existing fabric management tools
  - Running one production service
    - virtualBatch
      - better flexibility and maintenance operations when compared with "lxbatch"
    - New services in test
- Open infrastructure to new tools – OpenStack (a single orchestrator for both services?)

www.cloudtweaks.com – David Fletcher

www.cloudtweaks.com – David Fletcher

R. Wartel, T. Cass, B. Moreira, E. Roche, U. Schwickerath and S. Goasguen; "Image Distribution in Large Scale Cloud Providers"; 2nd IEEE Cloud Computing Conference; Indianapolis; 2010.

HEPiX meeting Spring and Fall 2009, 2010, 2011