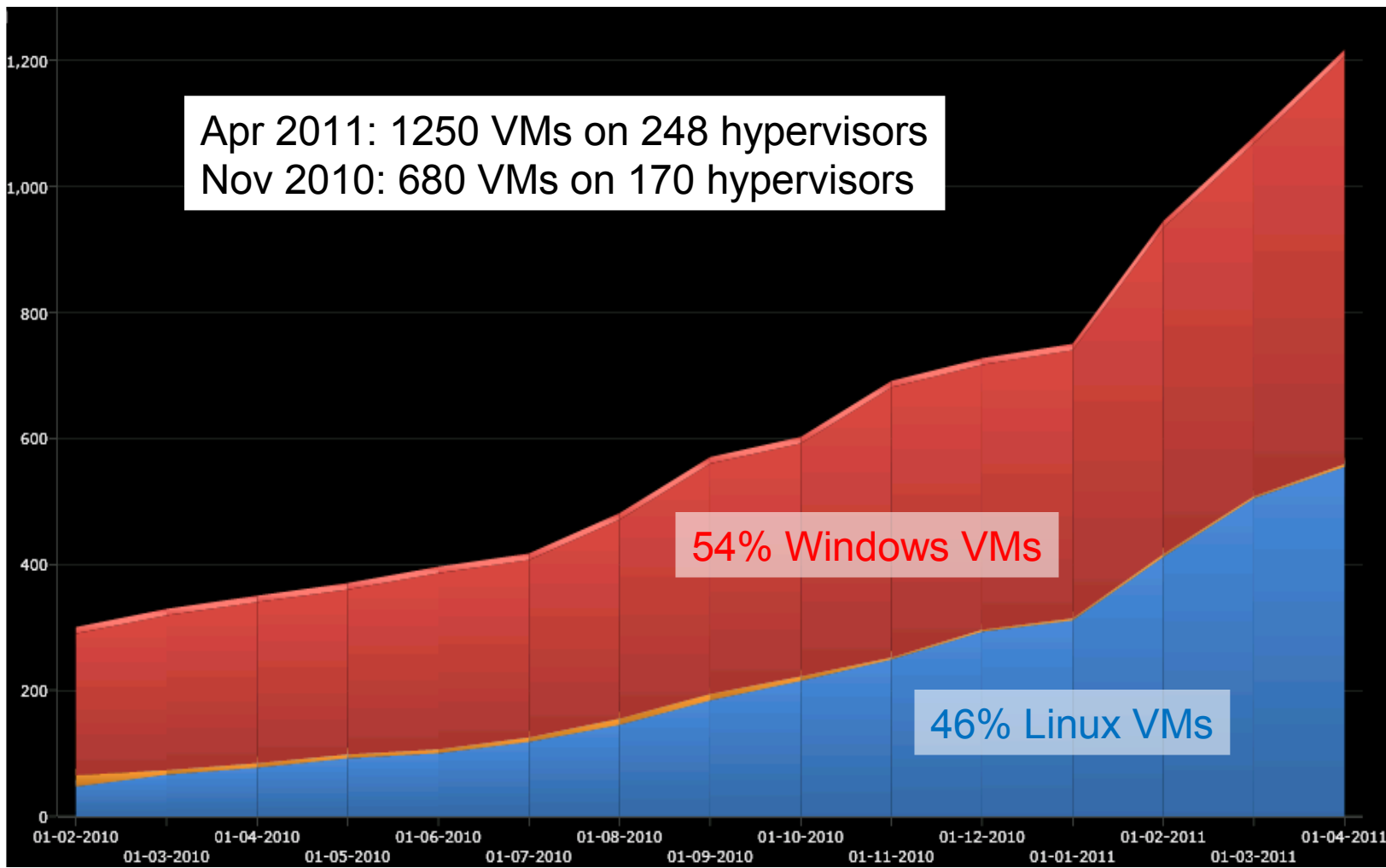# CERN's Internal Cloud Infrastructure - Design and Status

ATLAS ADC meeting, CERN, 19/5/2011
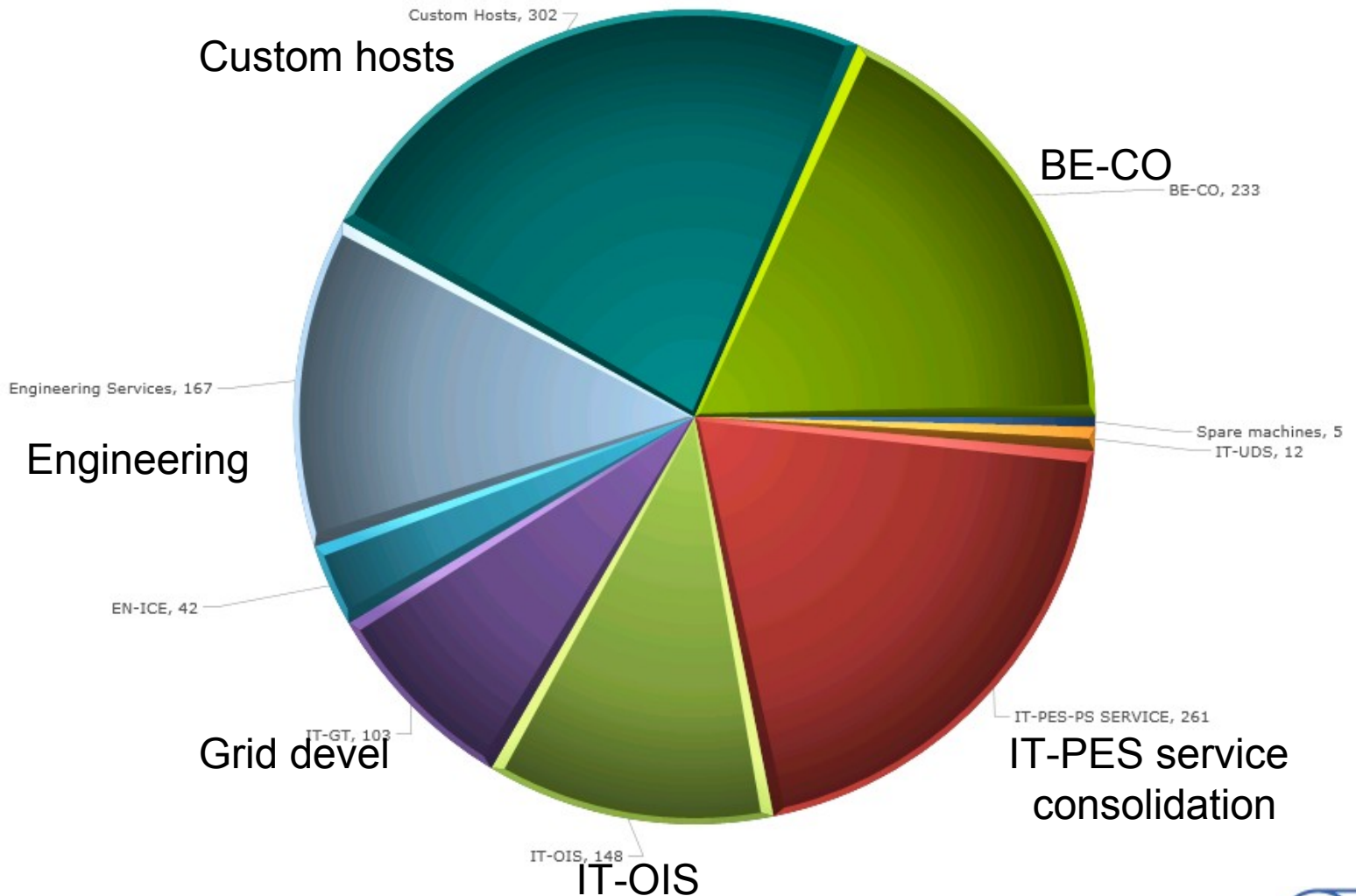
**Outline:**
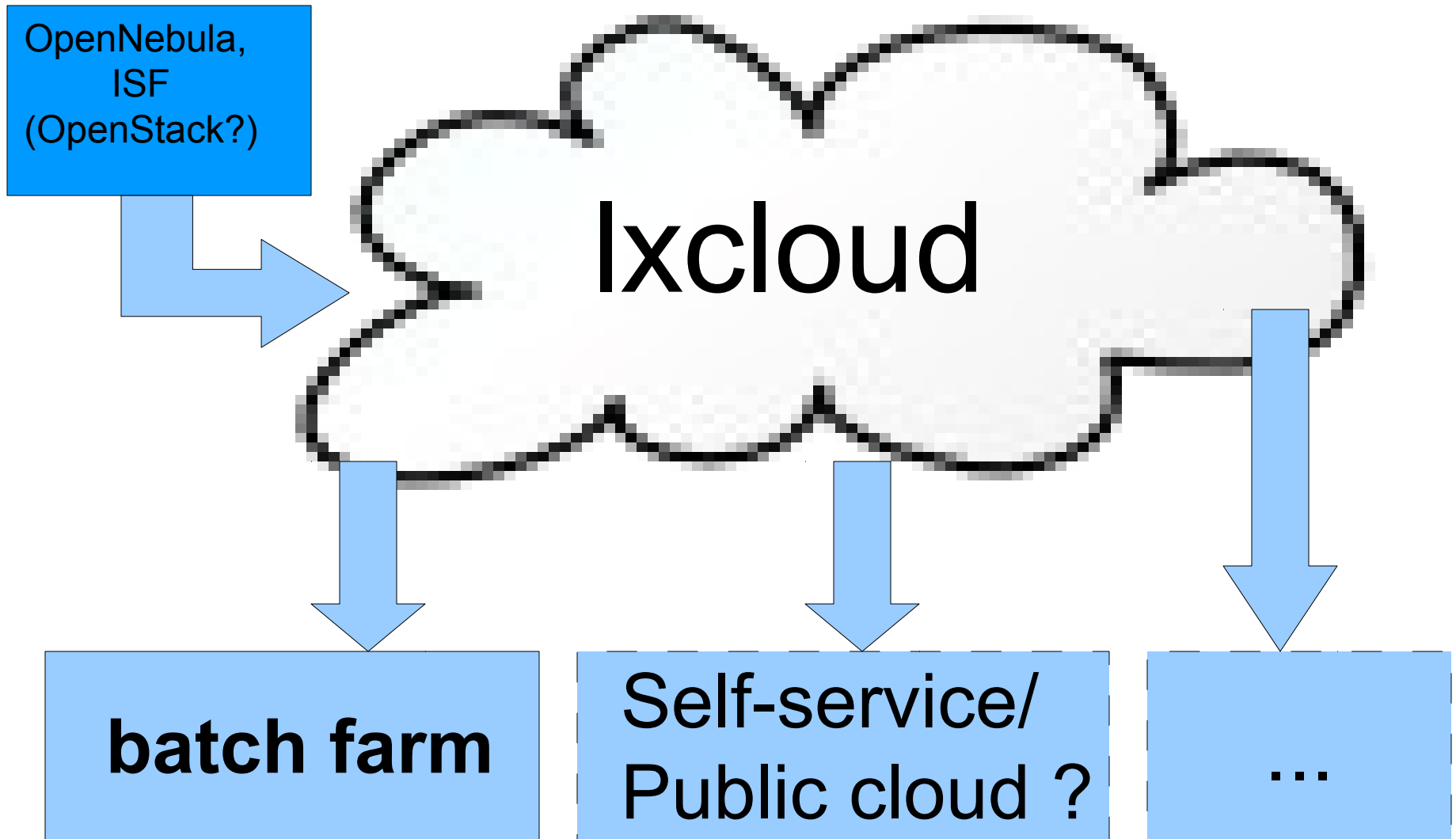- **Part 1: CERN Virtual Infrastructure (CVI) status**
- **Part 2: Internal cloud status and road map**

OIS

▶ The CERN Virtual Infrastructure  custom virtual machines in the CERN computer centre

> ▶ These VMs have a long-term lifetime of months/years

▶ User kiosk for requesting a VM in less than 30 mins

▶ Based on Microsoft's System Center Virtual Machine Manager (SCVMM)

> ▶ Enterprise class centralized management

> ▶ Rich feature set:

>> ▶ Allows grouping of hypervisors, with delegation of administrative privileges

>> ▶ VM migration, High availability

>> ▶ Checkpointing

>> ▶ PowerShell for administration / scripting

# Grow rate



Apr 2011: 1250 VMs on 248 hypervisors
Nov 2010: 680 VMs on 170 hypervisors

54% Windows VMs

46% Linux VMs

Number of Virtual Machines per Operating System

Custom Hosts, 302

Custom hosts

BE-CO

BE-CO, 233

Engineering Services, 167

Spare machines, 5
IT-UDS, 12

Engineering

EN-ICE, 42

IT-PES-PS SERVICE, 261

IT-GT, 103

Grid devel

IT-PES service consolidation

IT-OIS, 148

IT-OIS

CERN IT Department

OpenNebula,
ISF
(OpenStack?)

lxcloud

**batch farm**

Self-service/
Public cloud ?

...

# Part 2 : internal cloud

## What is it ?
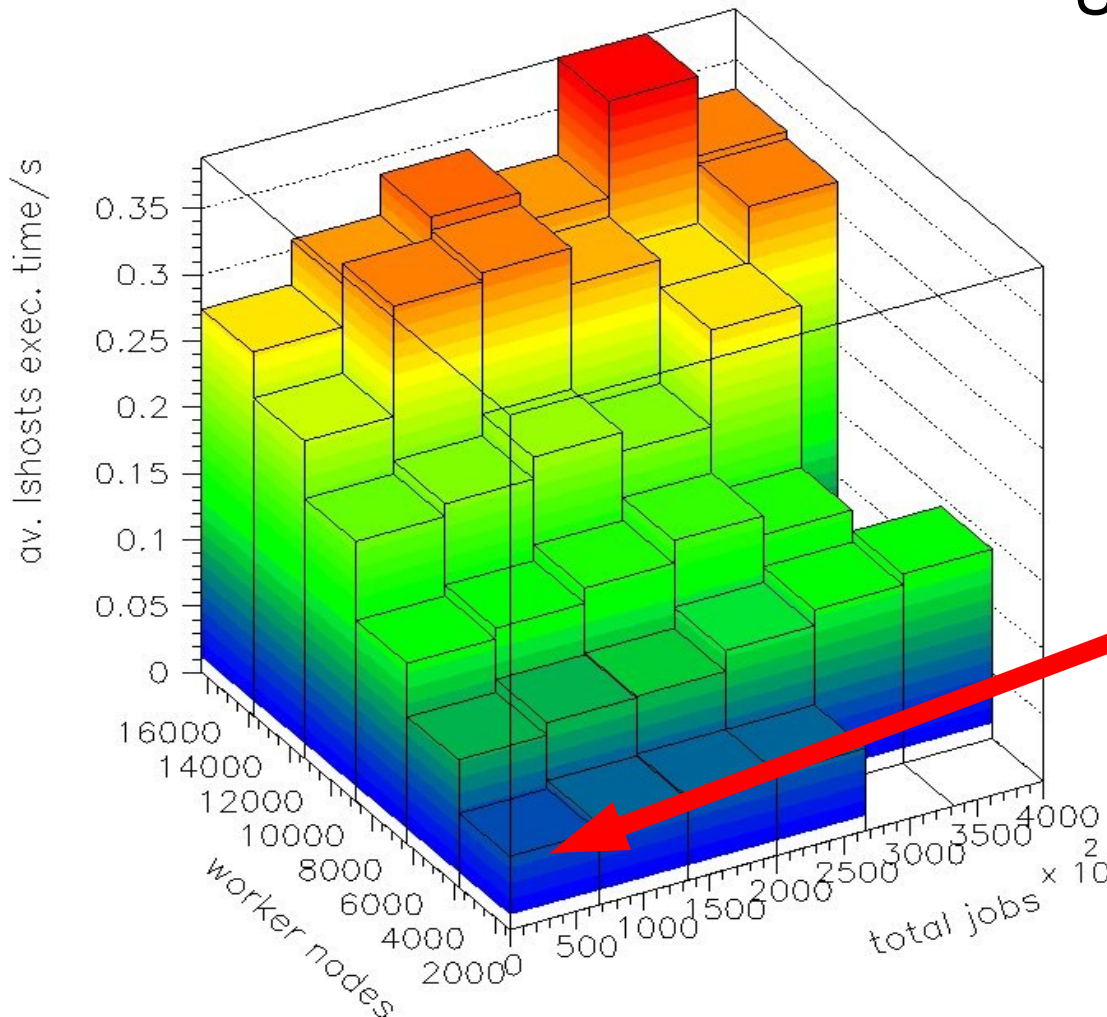
▶ Highly scalable, Linux (KVM) based cloud-like infrastructure

▶ Optimized for efficiency/speed

▶ Main building blocks:

  ▶ Resources with pre-allocated VM slots (lxcloud)

  ▶ Efficient internal image distribution system with bit-torrent

  ▶ Local pre-staged images and LV snapshotting

  ▶ VM management and deployment with OpenNebula/ISF

# History

▶ Idea born during discussions with Platform Computing about batch efficiency in **2008**

▶ First prototypes and feasibility studies (proof of concept) in **2009**, focus on batch

▶ Evolved with time into an internal cloud infrastructure, with lxbatch as one (first) application on top of it

▶ Used for large scale LSF scalability testing in 2010

CERN IT Department

Batch system tests: resource layer

Up to **15,000** nodes
Up to **400,000** jobs

Probed up to more than **3x** of what is officially supported by LSF



Current production system

A single batch system instance works up to 5-10k worker nodes only

# Resource pool details

- **Quattor managed** pool of resources (lxcloud)
- **Hardware**: (cheap) CPU server type, local disks
- **Guest setup**
  - Pre-allocation of VM "slots" in landb
  - Hypervisor "knows" the name if its guests
- **Disk management**
  - Use of LVM snapshots
  - All free disk space in one big LV
  - Pre-stage raw images on LV on the hypervisors
  - Fast instantiation of VMs using LV snapshots

# Resource pool details

- **Full integration into ELFms infrastructure**
  - Full monitoring with Lemon
  - Alarming with LAS (Operator)
  - Hardware management by sysadmin team
  - Standard tools for installation and management
  - "Draining" via sms state management
- **VM management systems**
  - OpenNebula (version 2.2 with CERN extensions)
  - Platform ISF
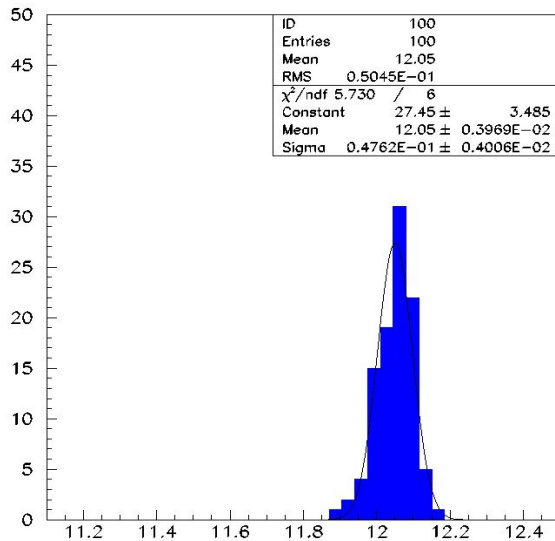  - Evaluation of OpenStack

# Resource pool details

- **Image distribution system**
  - Central image repository of <u>trusted</u> images
  - Distribution using Bit-torrent (rtorrent)
  - Pull model: Hypervisors ask if there are updates
  - Transparent update of images using LV tools
  - Hypervisors advertise existing images
- **Image catalogue (VMIC)**
  - Close collaboration with partners
  - Close collaboration with HEPiX virtualization WG

# Image creation model

- ▶ **"Golden" node(s)**
  - ▶ A <u>fully quattor managed</u> virtual machine
  - ▶ Which runs in lxcloud
  - ▶ PXE installation using a slot on a hypervisor
  - ▶ Usually associated with an existing service
- ▶ **Image creation:**
  - ▶ Halt the golden node
  - ▶ Take a snapshot
  - ▶ De-quattorize and clean up
  - ▶ Move it to the repository
  - ▶ Distribute to the hypervisors

# VM Instance management

▶ Instances are always derived from the **newest available golden node image**

▶ **Customized** at boot time (contextualization phase)

▶ Instances are **no longer known** to Quattor

  ▶ Still possible to manage via the golden node !

  ▶ Still possible to monitor with Lemon

  ▶ Consoles and remote-power-control work
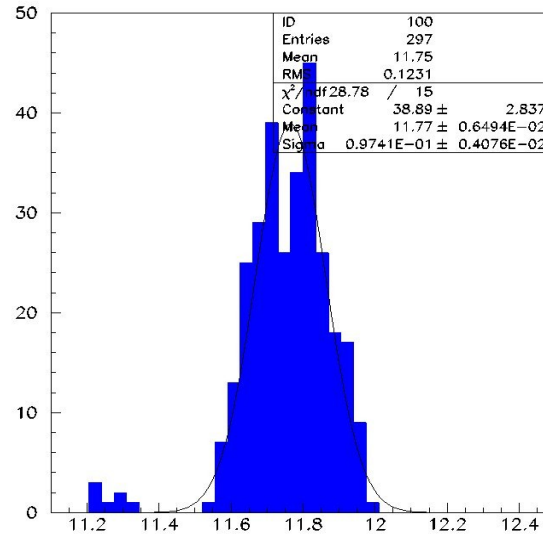    (Remote-power-control for operator only)

CERN IT Department

## HS06 tests



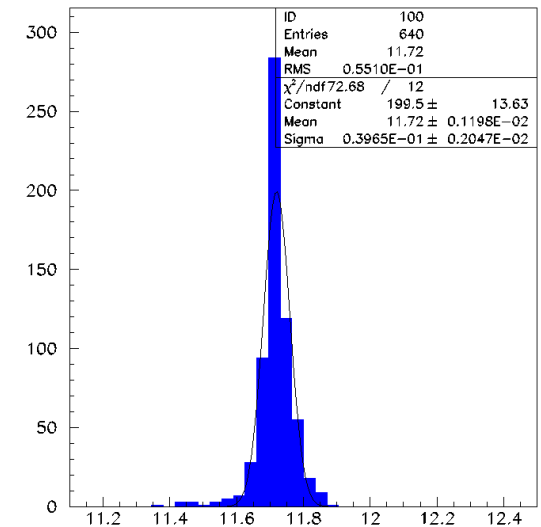**Bare metal:**
▶ SLC5
▶ 2x L5520 Intel Xeon
▶ 2.27GHz

**HS06=12.05/core**

**KVM**
▶ HW as before
▶ SLC5/6 hypervisor
▶ 8 SLC5 guests
▶ No KSM,ept off (SLC5)
▶ Pinned VMs
**HS06=11.4/core (SLC5)**
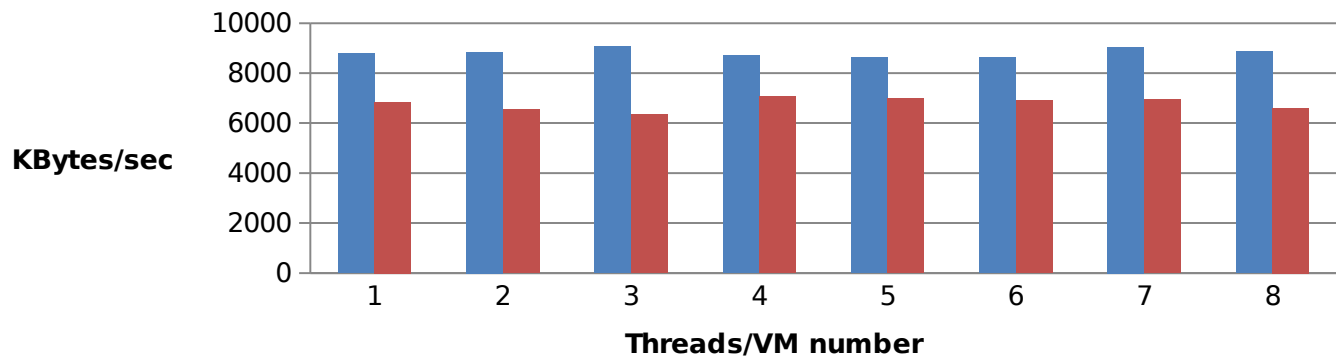**HS06=11.8/core (SLC6)**

**HyperV**
▶ HW as before
▶ 8 SLC5 guests

**HS06=11.7/core**

# I/O Benchmarking

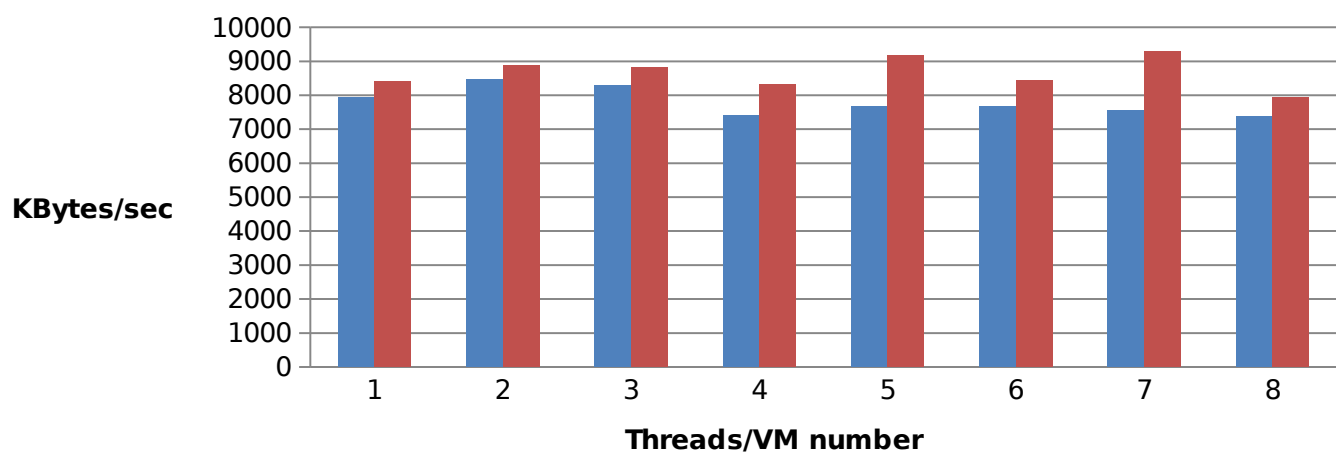-Mce -I -+r -r 256k -s 8g -f /pool/iozone_$i.dat$$ -i0 -i1 -i2

**20-30% penalty**



**Write**

KBytes/sec

Threads/VM number

**Read**

KBytes/sec

Threads/VM number

**Analysis by Qiulan Huang** (Chinese academy of science), December 2010
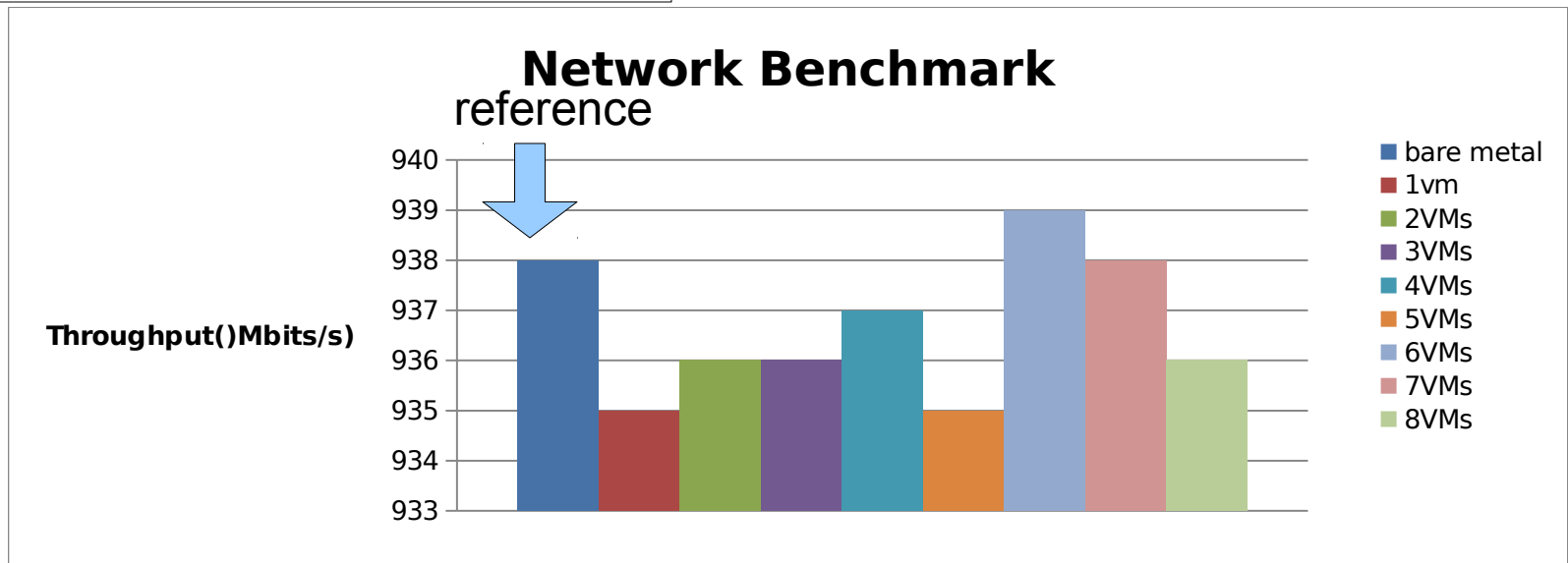
**Notes:**
- Caching off
- SLC5 Hypervisor
- 20% penalty
- write worse
- block device disk on LV, exported to VM

CERN IT Department

**Analysis by
Qiulan Huang**
(Chinese academy of science),
December 2010, CERN

Penalty <=1%

**Network Benchmark**

reference



iperf  with TCP window size of 256k and 60s test time

# Benchmarking conclusions

CPU benchmarking
- ▶ Best results of 2-3% requires tuning
- ▶ EPT=0 has fairly large effect on SLC5, less on SLC6
- ▶ Small effect by using the native CPU (SLC6)

I/O benchmarking house numbers (SLC5)
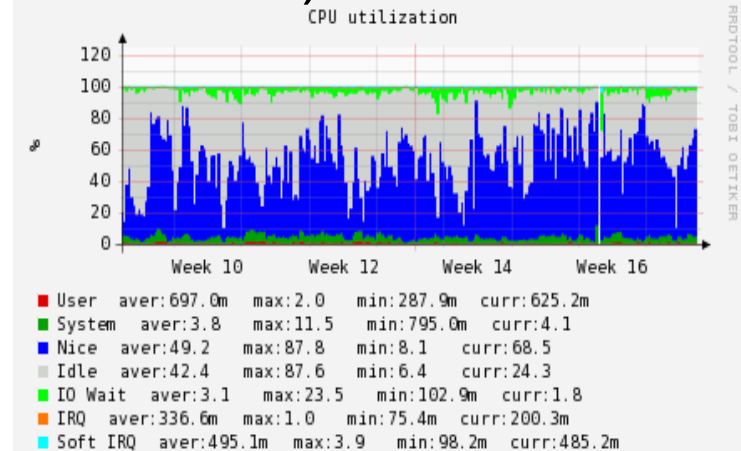- ▶ Read/Write performance penalty 20%-30%

Network benchmarking (iperf, SLC5 GE)
- ▶ Very small performance loss
- ▶ Possibly not significant within statistics

# Virtual batch: production

**96 nodes (12 hypervisors) in full production in public batch**

▶ 6 hypervisors controlled by ISF (version 2.0)

▶ 6 hypervisors controlled by ONE (version 2.2)

▶ Short public and GRID jobs

▶ 1 VM / core and 1 job per VM

▶ Ramping up to 384 now



**Notes:**

▶ updated via Quattor (Golden Node) at boot or external trigger

▶ Image change only required for intrusive updates

▶ 12 identical physical nodes for job throughput comparison

# Virtual batch: observations

**Observation**: More short jobs scheduled to virtual batch nodes

**Job success rate:**

Virtual nodes   : 88%
Physical nodes: 82%

**Delivered wall clock time:**
(Ratio of time and total wall clock time seen by jobs)
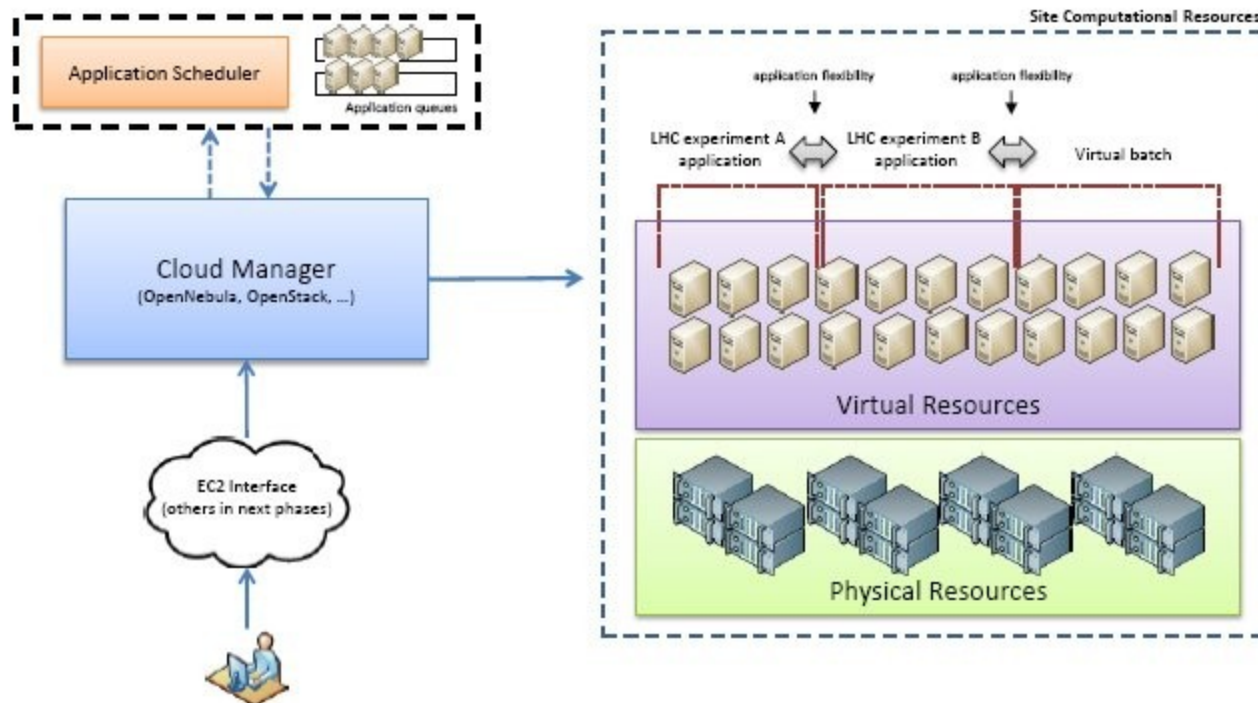
Virtual nodes   : 76%
Physical nodes: 81%

**Possible improvements:**
- Life time of VMs can be increased now
- Several improvements in boot sequence
- Machine renewal

Test instance setup:

▶ ONE server on SLC6 with EC2 interface enabled
▶ Access for restricted users only, on request
▶ Only trusted images (user can't upload their own image)

*The internal cloud infrastructure at CERN - a status report - 20*

# Conclusions

- CERN-IT strategy is going for virtualization everywhere

- Production experience so far looks very promising

- Evaluation of cloud computing options has started and is ongoing

- Input from the experiments is appreciated