



**XXV International Symposium on Nuclear Electronics & Computing**

28<sup>th</sup> Sep – 2<sup>nd</sup> Oct 2015, Budva, Becici, Montenegro

# **CERN LHC Run 2 on OpenStack**

**Sebastian Bukowiec**

CERN Cloud Infrastructure



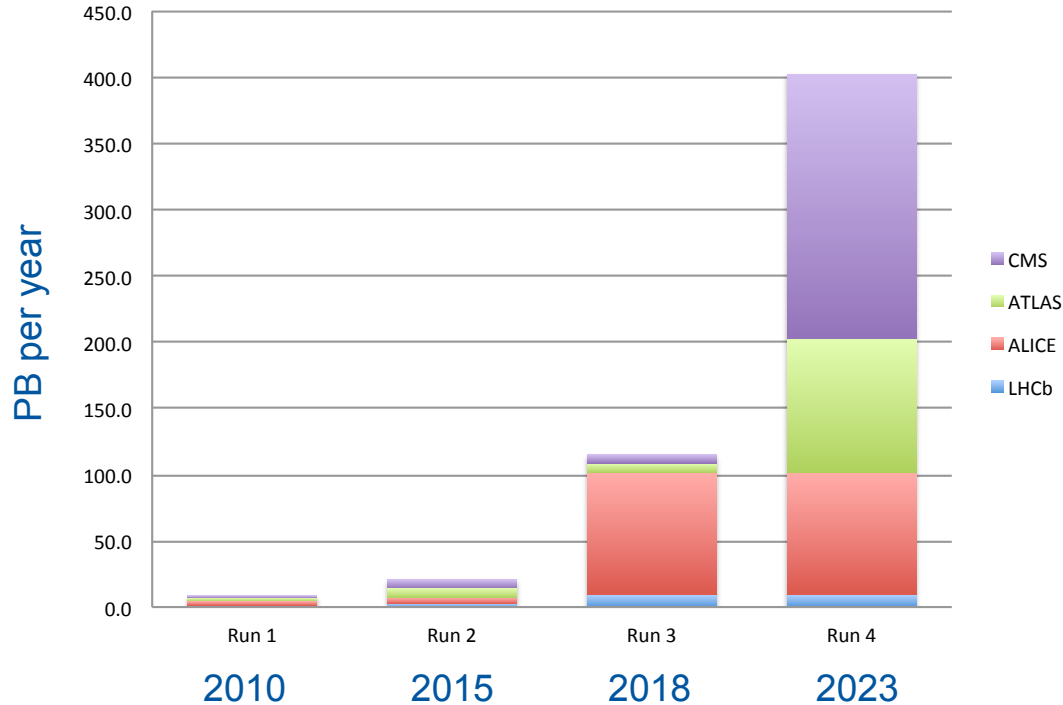
# LHC Run 1 (2010-2013)

- The computing challenge resulted in a great success
- ~65 PB of data produced
- Analyzed in record time proving a theory proposed by Professors Higgs and Englert in 1960s



Source: <http://www.uniovi.es/en/~higgs-englert-y-el-cern-explican-en-la-universidad-de-oviedo-el-camino-hasta-el-descubrimiento-del-boson>

# Planned LHC data growth



- Estimating 400PB/year by 2023
- Compute needs expected to be around 50x current levels if budget available

# Infrastructure growth

- Large increase of computing resources required from 2015
- Constrained budgets
- The same amount of personnel

Enforce to evolve towards more dynamic, efficient and flexible systems.

# CERN's Toolset

Agile Infrastructure project aimed to identify new tools needed to build CERN Cloud Infrastructure and enhance communication within IT department

Configuration Tools

Cloud Software

Monitoring Tools

Storage Solution



# CERN's Toolset

Agile Infrastructure project aimed to identify new tools needed to build CERN Cloud Infrastructure and enhance communication within IT department

Configuration Tools

Cloud Software

Monitoring Tools

Storage Solution



# OpenStack

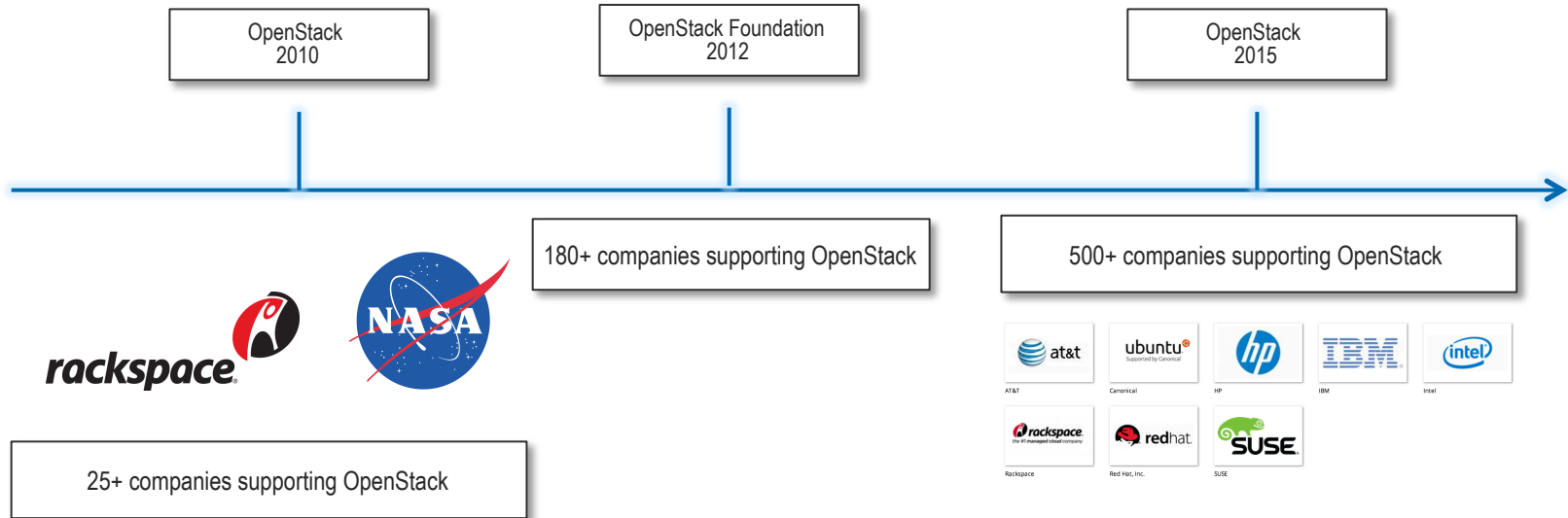


OpenStack software controls large pools of compute, storage, and networking resources throughout a datacenter, managed through a dashboard or via the OpenStack API. OpenStack works with popular enterprise and open source technologies making it ideal for heterogeneous infrastructure.

-- *OpenStack*  
<http://www.openstack.org>

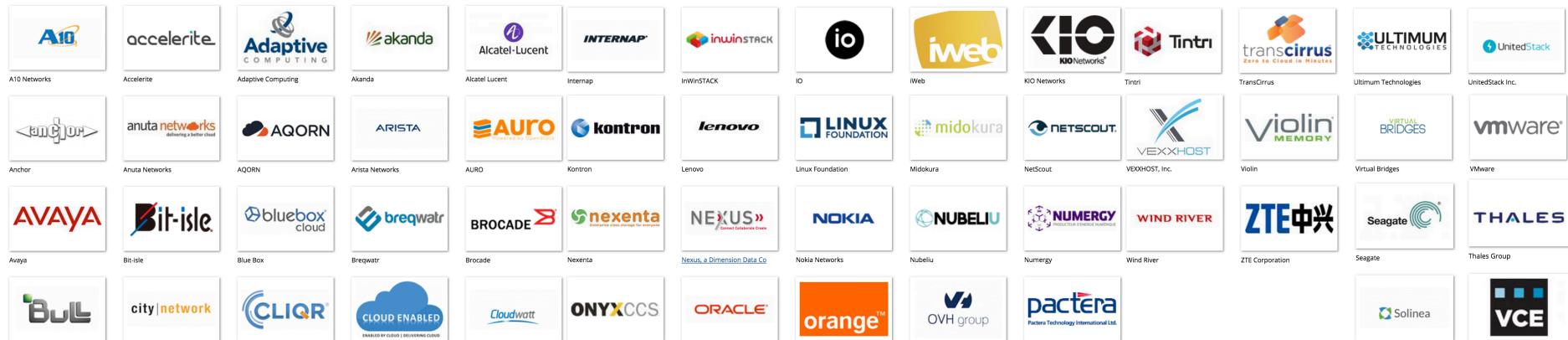


# OpenStack beginning





Gold Members



Corporate Sponsors



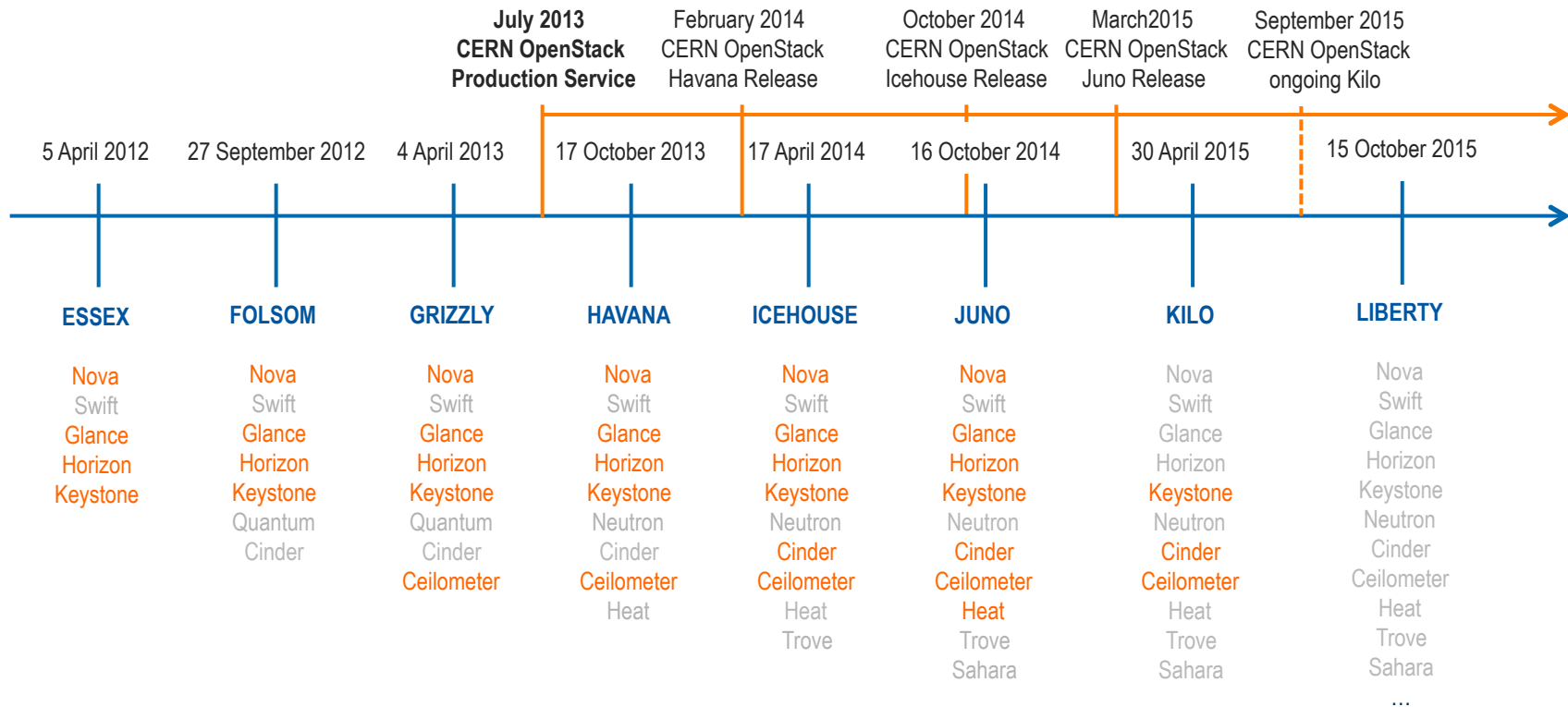
# OpenStack services

## OpenStack® Services



www.rackspace.com 11

# CERN OpenStack project timeline



# Strategy to deploy OpenStack at CERN

- Configuration infrastructure based on Puppet
- Community Puppet modules for OpenStack
- RDO - RPM Packages
- CERN CentOS 7 (CC7), Microsoft Windows 2012 R2 Operating Systems
  - Scientific Linux CERN 6 (SLC6) during phase out

# Strategy to deploy OpenStack at CERN

- Wigner - Budapest Computer Center hardware deployed as OpenStack compute nodes

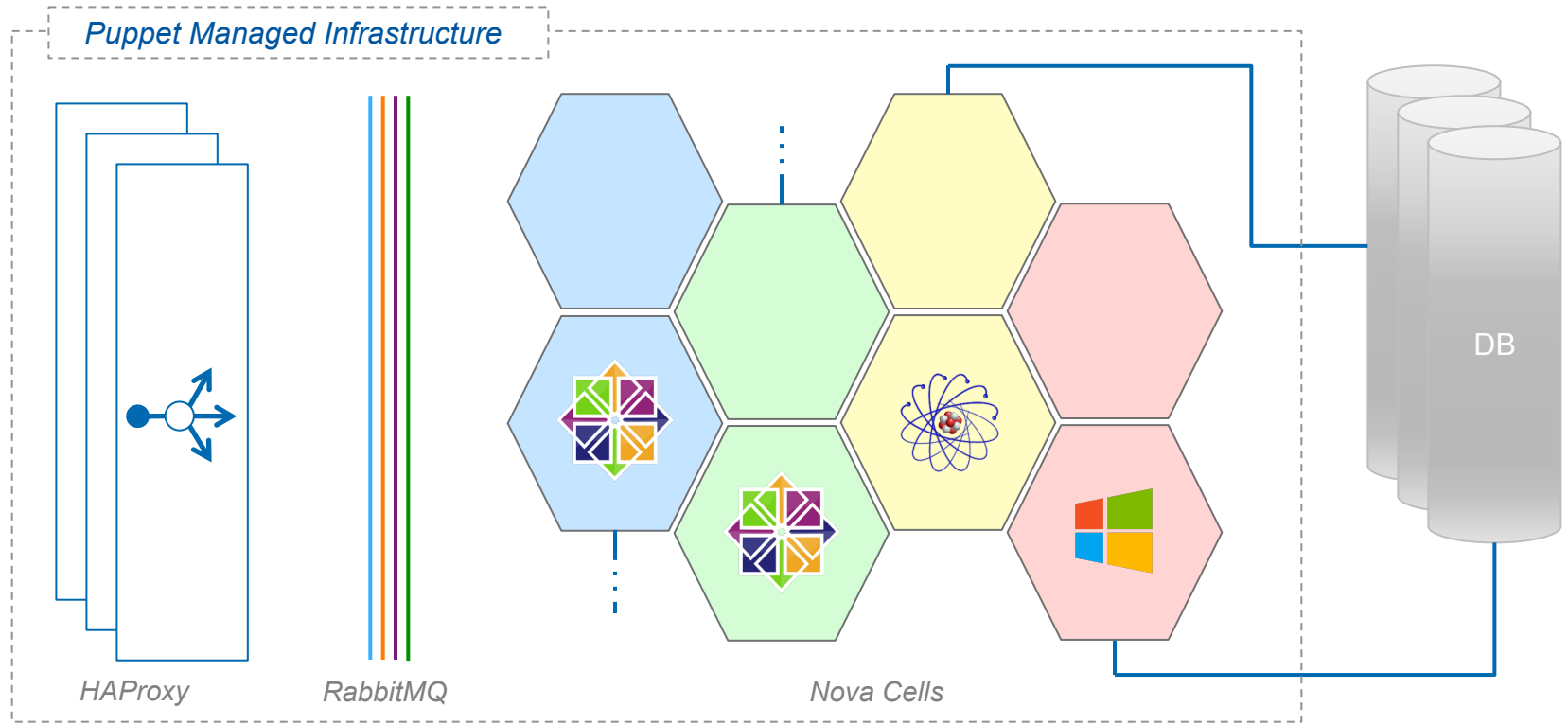


# Size of CERN Cloud Infrastructure

<b>Compute Nodes</b>	5000
<b>Virtual Machines</b>	15 000
<b>Cores</b>	130 000
<b>Tenants</b>	2216

continuously growing ...

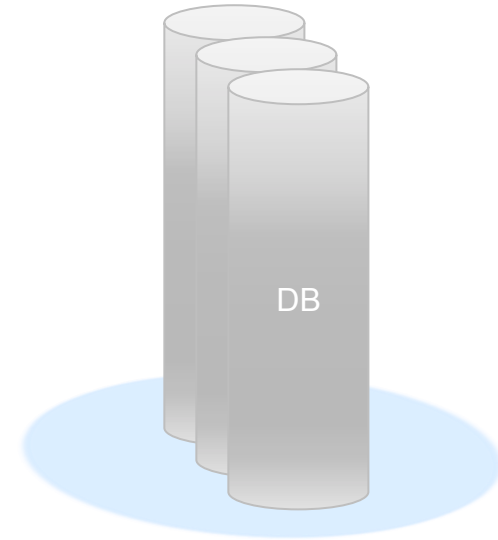
# Infrastructure Overview





# Infrastructure Overview

- **mySQL instance per Cell**
  - provided by CERN DB team (DB on demand service)
  - Running on top of Oracle CRS
  - NetApp storage backend
  - Backups every 6 hours



# Three interfaces

```
$ openstack server list
```

ID	Name	Status	Task State	Power State	Networks
b30d60b3-e4a2-43a1-8889-e022a0c12f20	demo-server-00	ACTIVE	CERN_NETWORK=188.XXX.XX.XXX, 2001:1458:XXX:XX::XXX:XXX		
f79a4a90-9c96-4af8-af94-e89cc35f2fe7	demo-server-01	ACTIVE	CERN_NETWORK=188.XXX.XXX.XXX, 2001:1458:XXX:XX::XXX:XXX		

The screenshot shows the OpenStack dashboard interface. At the top, it says 'CERN Accelerating science' and 'Signed in as: bukowiec'. The main navigation bar includes 'Overview', 'Instances', 'Volumes', 'Images', and 'Access & Security'. The 'Overview' page displays a 'Limit Summary' with five circular progress indicators for Instances (Used 2 of 50), VCPUs (Used 2 of 50), RAM (Used 1.0GB of 100.0GB), Volumes (Used 1 of 10), and Volume Storage (Used 1.0GB of 1000.0GB). Below this is a 'Usage Summary' section with a date range selector (From: 2015-03-01, To: 2015-03-23) and a 'Submit' button. It shows 'Active Instances: 2', 'Active RAM: 1GB', 'This Period's VCPU-Hours: 18.85', and 'This Period's GB-Hours: 0.00'. A 'Usage' table is displayed with columns for Instance Name, VCPUs, Disk, RAM, and Uptime. The table lists two instances: demo-server-01 and demo-server-00, both with 1 VCPU, 0 Disk, 512MB RAM, and 4 months, 3 weeks of uptime. A 'Download CSV Summary' button is located below the table.

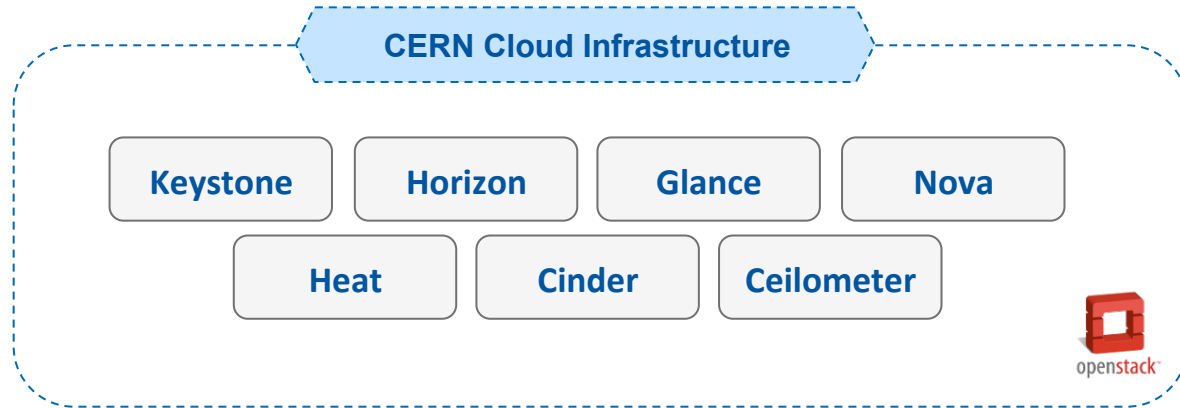
Instance Name	VCPUs	Disk	RAM	Uptime
demo-server-01	1	0	512MB	4 months, 3 weeks
demo-server-00	1	0	512MB	4 months, 3 weeks

```
$ curl -s -X POST http://8.21.28.222:5000/v2.0/tokens \  
-H "Content-Type: application/json" \  
-d '{"auth": {"tenantName": "", "passwordCredentials": \  
{"username": ""$OS_USERNAME"", "password": ""$OS_PASSWORD""}}}' \  
| python -m json.tool
```

```
$ curl -i -H "X-Auth-Token:token" http://8.21.28.222:8774/v2/tenant_id/servers
```

# CERN OpenStack projects

- Modular architecture
- Designed to easily scale out



# Keystone (Identity)

- Main functions:

- Identity: user authentication

Keystone is integrated with CERN AD via LDAP backend

CERN Active Directory infrastructure:

- unified identity across the site
- 44 000 users
- ~200 arrivals/departures per month

- Policies: enforces different user levels: Member, Admin

Keystone

Heat

Horizon

Cinder

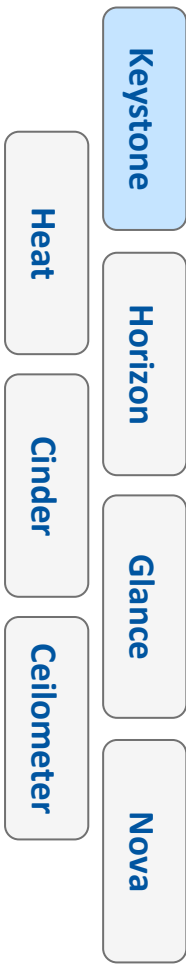
Glance

Ceilingmeter

Nova

# Keystone (Identity)

- Main functions:
  - Token: Provides token to communicate with other OpenStack components
  - Service catalog: service endpoints registered with keystone



# Keystone (Identity)

- CERN user subscribes the "cloud service"
  - Created "Personal Tenant" with limited quota
  - ~ 2000 users
- Shared projects created by request
  - ~300 shared
- Project life cycle integrated with Microsoft Forefront Identity Manager (FIM) to be compliant with standard CERN policies

Keystone

Heat

Horizon

Cinder

Glance

Ceilometer

Nova

# Nova (Compute)

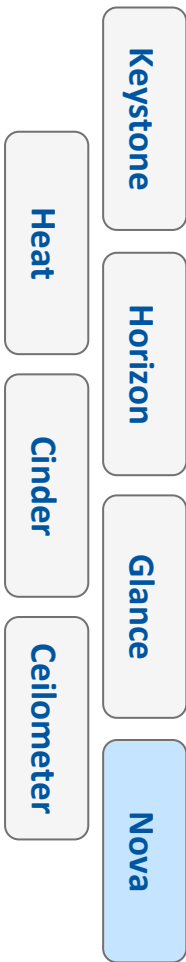
Boots and manages the lifecycle of virtual machine on CERN supported hypervisors:

- KVM (Linux OS VMs)
- Hyper-V (Windows OS VMs)

Hypervisor selection is based on “Image” properties.

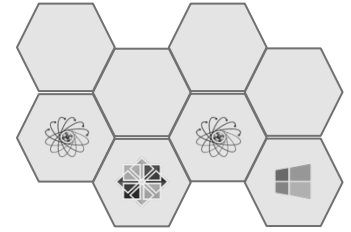
# Nova (Compute)

- Set of compute nodes equipped with hybrid storage for IO intensive workflow
- Various flavors matching ~20 different types of Hardware
- ~5000 compute nodes
  - vast majority QUEMU/KVM (~150 Hyper-V)
  - ~220 HVs on critical power
  - ~800 HVs at Wigner in Hungary
  - ~2000 HVs used by batch - “cattle” use case
  - rest shared by users, services, experiments - “pets” use case





# Nova Cells



Cells are used in order to scale the infrastructure, better fault tolerance and for project distribution.

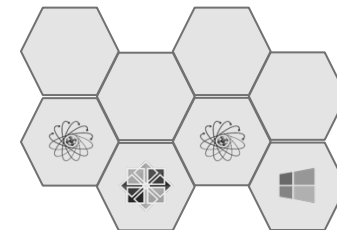
- Allow to scale transparently between different Computer Centers
- Mapped to the availability zones

✓ Any Availability Zone  
cern-geneva-a  
cern-geneva-b  
cern-geneva-c  
cern-wigner-a  
cern-wigner-b

# Nova Cells

## Top level cell

- Runs API service
- Top cell scheduler



API Cell

Controllers



Geneva, Switzerland

## Child cells run:

- Compute nodes
- Nova network
- Scheduler
- Conductor

Compute Cell

Controllers



Compute nodes



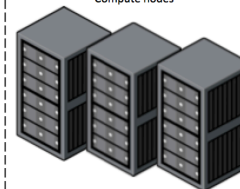
Geneva, Switzerland

Compute Cell

Controllers



Compute nodes



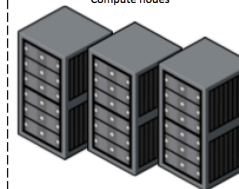
Geneva, Switzerland

Compute Cell

Controllers



Compute nodes

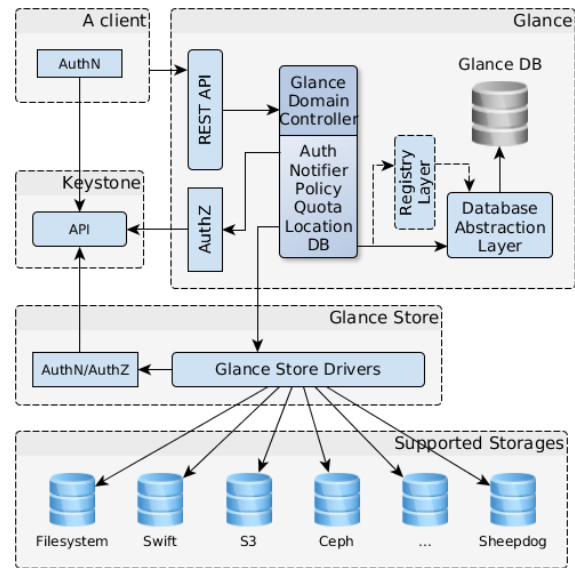


Budapest, Hungary

# Glance (Image service)

Stores and retrieves VMs images  
(virtual machines templates)

- **Glance backend - Ceph**  
(currently ~28 TB of images)
- **Snapshot capabilities**



Source: <http://docs.openstack.org/developer/glance/architecture.html>

Heat

Cinder

Cellometer

Keystone

Horizon

Glance

Nova

# Glance (Image service)

- Maintain a small set of CERN images as default
  - Difficult to offer only the most updated set of images
  - Resize and Live Migration not available if image is deleted from Glance
- Users can upload their own images
  - Users don't pay for storage
  - No quotas per Tenant

Heat

Cinder

Cellometer

Keystone

Horizon

Glance

Nova

# CERN Images

CERN Accelerating science Signed in as: bukowiec Sign out Directory

CERN Cloud Infrastructure Compute ▾ Current Project: Personal bukowiec ▾ Project Settings Submit a ticket Help

Overview Instances Volumes Images Access & Security

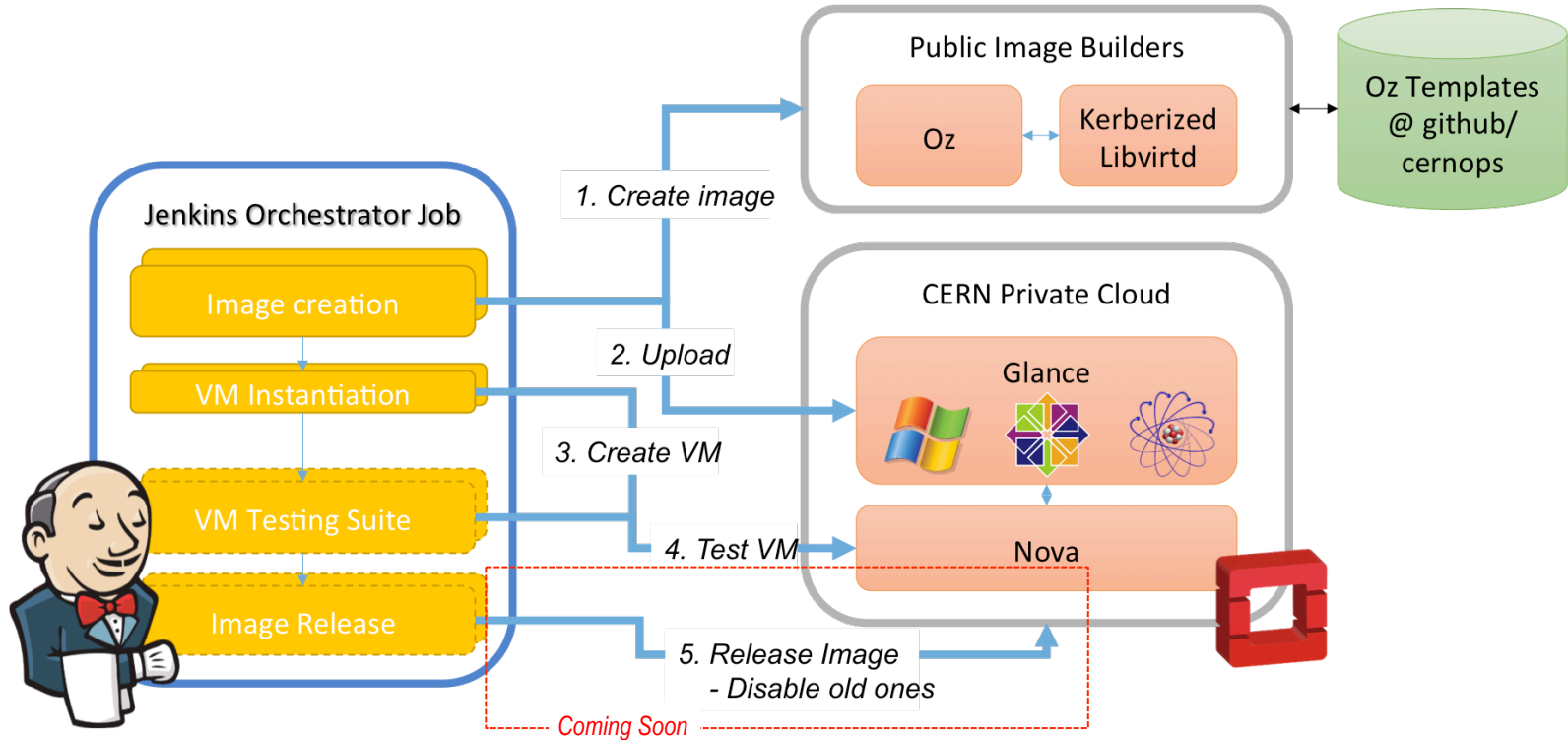
## Images

### Images

Latest All Project (13) Shared with Me (1) Public (56) + Create Image Delete Images

<input type="checkbox"/>	Image Name	Type	Format	Status	Release Date	Actions
<input type="checkbox"/>	<a href="#">CC 7 Base - x86_64</a>	Image	QCOW2	Active	2015-02-10	<span>Launch</span> <span>More ▾</span>
<input type="checkbox"/>	<a href="#">CC 7 Extra - x86_64</a>	Image	QCOW2	Active	2015-02-10	<span>Create Volume Provider</span>
<input type="checkbox"/>	<a href="#">SLC 6 CERN Server - x86_64</a>	Image	QCOW2	Active	2015-02-10	<span>Launch</span> <span>More ▾</span>
<input type="checkbox"/>	<a href="#">SLC 6 Server - x86_64</a>	Image	QCOW2	Active	2015-02-10	<span>Launch</span> <span>More ▾</span>

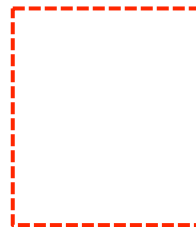
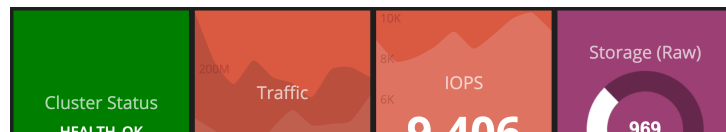
# Image Lifecycle



# Cinder (Block Storage)

Provides persistent disks for virtual machines

- Volumes for Linux and Windows VMs
- Utilizes Ceph and NetApp as storage backends
- Various volume types



Heat

Keystone

Horizon

Cinder

Glance

Ceilometer

Nova

24/09/2015 09:50

# Cinder Volumes

Volume Name \*

Description

Volume Source

Type

Size (GB) \*

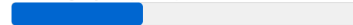
Description:

Volumes are block devices that can be attached to instances.

Type	standard
Usage	default
Platform	Linux
Max IOPS	100
Max Throughput	80 MB/s

Volume Limits

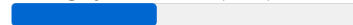
Total Gigabytes (38 GB) 100 GB Available



Number of Volumes (4) 10 Available



Total Gigabytes for standard (42 GB) 100 GB Available



Number of Volumes for standard (4) 10 Available



Volume Name \*

Description

Volume Source

- standard
- cp1
- cp2
- ✓ io1
- cpio1

Size (GB) \*

Description:

Volumes are block devices that can be attached to instances.

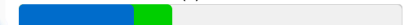
Type	io1
Usage	IO intensive
Hypervisor	Linux
Max IOPS	500
Max Throughput	120 MB/s

Volume Limits

Total Gigabytes (87 GB) 100 GB Available



Number of Volumes (3) 10 Available



Total Gigabytes for io1 (0 GB) 0 GB Available



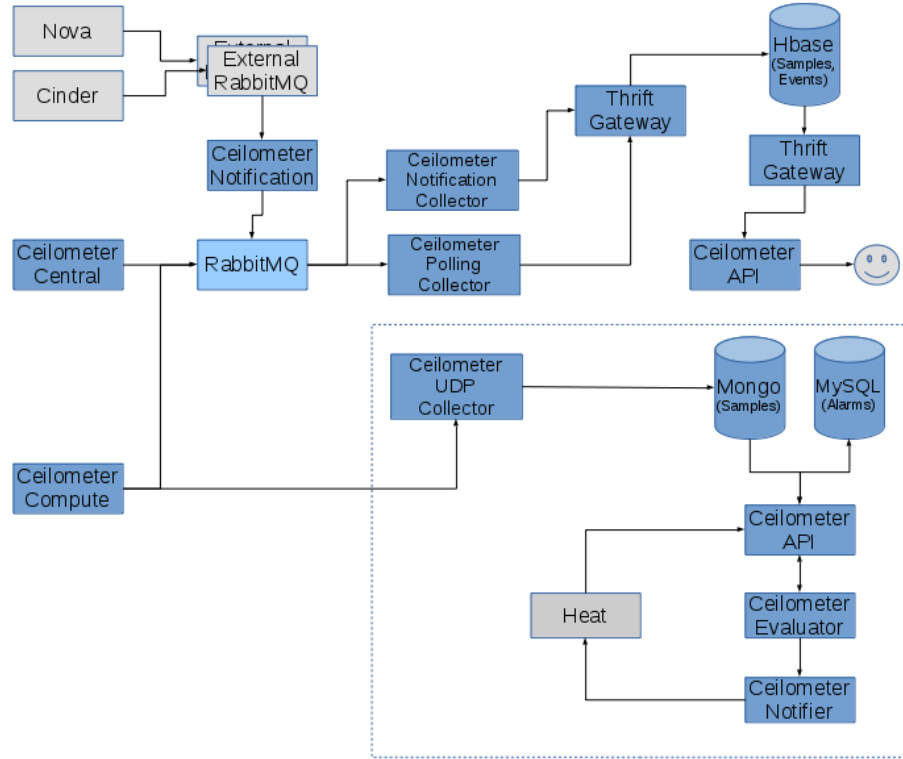
Number of Volumes for io1 (0) 0 Available





# Ceilometer Architecture

- Monitors, collects and stores usage data for all OpenStack Infrastructure
- Provides data used in accounting
- Used for Heat auto scaling scenarios



Heat

Cinder

Ceilometer

Keystone

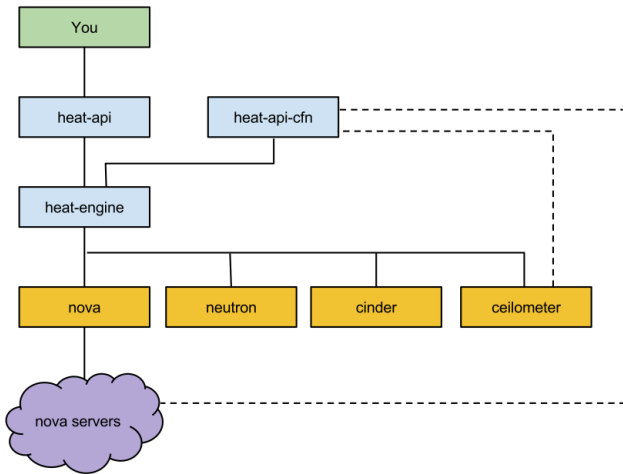
Horizon

Glance

Nova

# Heat (Orchestration)

Heat provides a mechanism for orchestrating OpenStack resources through templates



```
heat_template_version: 2013-05-23

description: Simple template to deploy a single compute instance

parameters:
  key_name:
    type: string
    label: Key Name
    description: Name of key-pair to be used for compute instance
  image_id:
    type: string
    label: Image ID
    description: Image to be used for compute instance
  instance_type:
    type: string
    label: Instance Type
    description: Type of instance (flavor) to be used

resources:
  my_instance:
    type: OS::Nova::Server
    properties:
      key_name: { get_param: key_name }
      image: { get_param: image_id }
      flavor: { get_param: instance_type }
```

# Horizon (Dashboard)

Provides  
a self-service  
web based  
user interface  
to OpenStack  
services  
for end users.

The screenshot displays the Horizon dashboard interface. At the top, it shows the CERN logo and navigation options like 'Signed in as: bukowiec', 'Sign out', and 'Directory'. Below this is a navigation bar with tabs for 'Overview', 'Instances', 'Volumes', 'Images', and 'Access & Security'. The main content area is titled 'Overview' and features a 'Limit Summary' section with five circular gauges representing resource usage: Instances (Used 2 of 50), VCPUs (Used 2 of 50), RAM (Used 1.0GB of 100.0GB), Volumes (Used 1 of 10), and Volume Storage (Used 1.0GB of 1000.0GB). Below the gauges is a 'Usage Summary' section with a date range selector (From: 2015-03-01, To: 2015-03-24) and a 'Submit' button. The usage summary text reads: 'Active Instances: 2 Active RAM: 1GB This Period's VCPU-Hours: 41.39 This Period's GB-Hours: 0.00'. A 'Download CSV Summary' button is also present. At the bottom, there is a table with the following data:

Instance Name	VCPUs	Disk	RAM	Uptime
<a href="#">demo-server-01</a>	1	0	512MB	4 months, 4 weeks
<a href="#">demo-server-00</a>	1	0	512MB	4 months, 4 weeks

Displaying 2 Items

Heat

Cinder

Ceilometer

Keystone

Horizon

Glance

Nova

# Automation



- Friendly and easy interface from where we can organize and launch jobs on our hosts
- Sharing of sensitive tasks to other groups without exposing credentials or procedures
- Operations delegation:
  - **SysAdmins:** Workflows related to hypervisor maintenance (h/w intervention, notify users...)
  - **Cloud-Operations:** Project creation, Health reports, Quota update

# Operations - Rally

- Benchmarking tool for OpenStack
- Performance test
- Cloud verification
- Used for OpenStack Continuous Integration
- Check if services work correctly
  - Rally runs against QA and Production environments regularly
  - We can compare results between the environments.

Heat

Keystone

Cinder

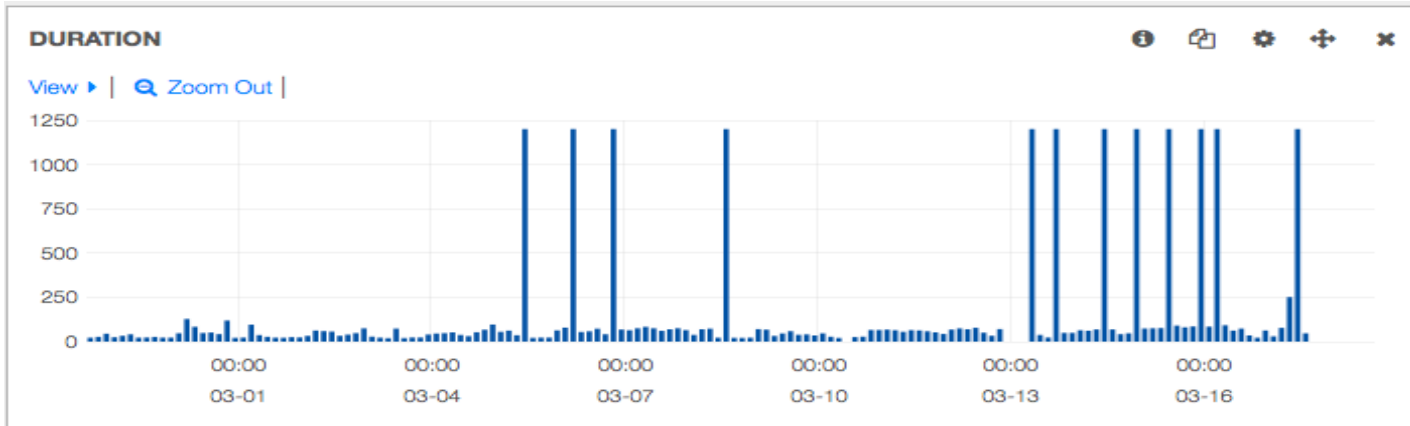
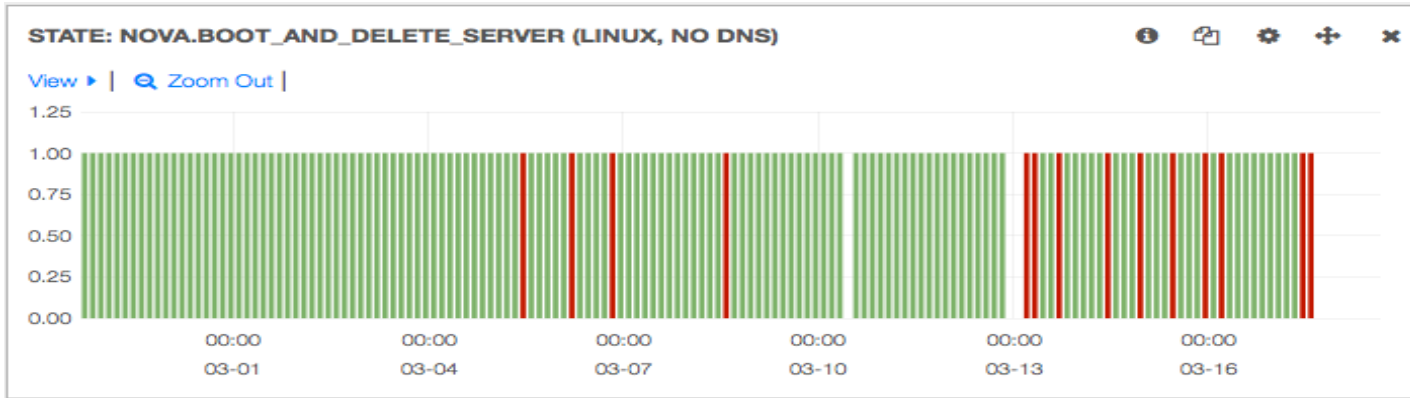
Horizon

Ceilometer

Glance

Rally

Nova



- Heat
- Keystone
- Cinder
- Horizon
- Ceilometer
- Glance
- Rally
- Nova

# OpenStack Federations

Evaluating cloud identity federation available in OpenStack ecosystem that allows for on premise bursting into remote clouds with use of local identities (i.e. domain accounts).

A way to Hybrid Cloud architectures - virtualized infrastructures layered across independent private and public clouds.

# Challenges

- Rapidly growing infrastructure
- Hardware lifecycle
- Neutron deployment
- Federations
- Keeping pace, rolling new updates every six months
- Integration of LXC containers into infrastructure
  - internal compute use cases





# OpenStack at CERN – Summary

- OpenStack scales and meets our needs for LHC Run 2
- Geneva and Budapest Computer Centers run on OpenStack
- Cells deployment allows us to scale out easily adding new groups of nodes
- Ceph backend for images and volumes, NetApp backend for Windows volumes
- Working with upstream for new features and bug fixes

# OpenStack at CERN

- CERN code publicly available on GitHub:  
<https://github.com/cernops>
- Hints and tips from the CERN OpenStack cloud team  
<http://openstack-in-production.blogspot.com>



[www.cern.ch](http://www.cern.ch)