# CERN Cloud Infrastructure report

José Castro León
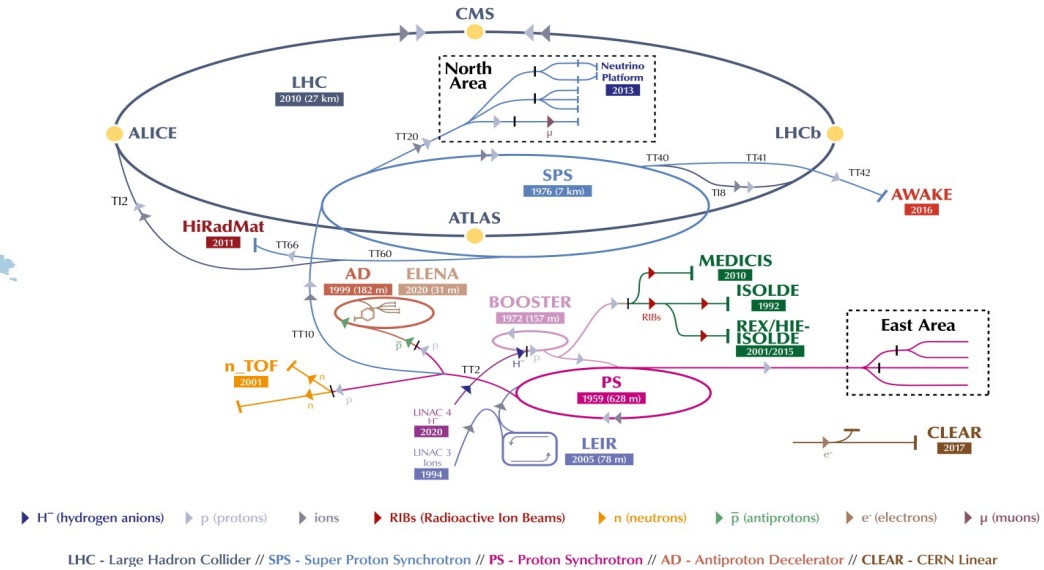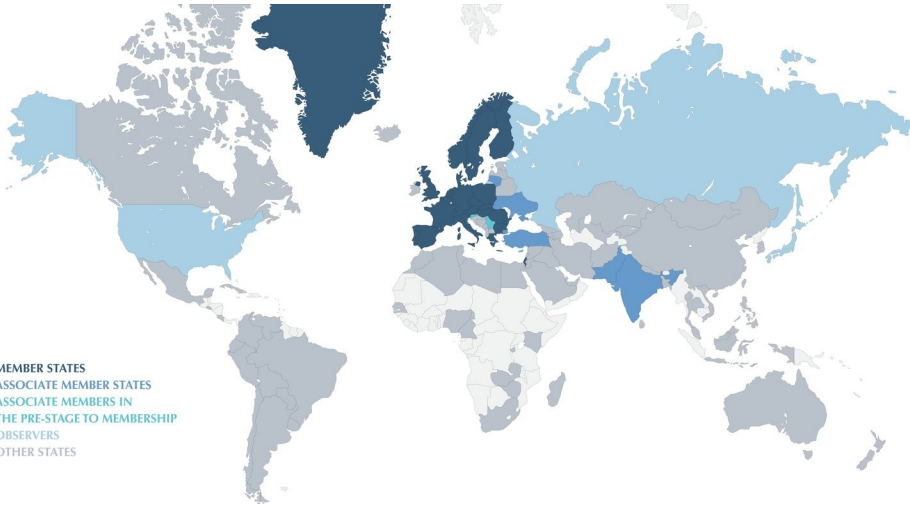CERN Cloud Infrastructure

# Outline

- Introduction

- CERN Cloud service
  - Deployment
  - Monitoring
  - Accounting
  - Identity
  - Probe and Debug

# European Organization for Nuclear Research

- **World largest particle physics laboratory**

- **Founded in 1954**

- **23 member states**

- **Fundamental research in physics**



MEMBER STATES
ASSOCIATE MEMBER STATES
ASSOCIATE MEMBERS IN
THE PRE-STAGE TO MEMBERSHIP
OBSERVERS
OTHER STATES



▶ H⁻ (hydrogen anions)   ▶ p (protons)   ▶ ions   ▶ RIBs (Radioactive Ion Beams)   ▶ n (neutrons)   ▶ p̄ (antiprotons)   ▶ e⁻ (electrons)   ▶ μ (muons)

LHC - Large Hadron Collider // SPS - Super Proton Synchrotron // PS - Proton Synchrotron // AD - Antiproton Decelerator // CLEAR - CERN Linear

**and RUN3 has just started …**

# CERN Cloud Infrastructure

- Infrastructure as a Service

- Production since **July 2013**

- **CentOS 7** based (adding CentOS Stream 8 soon)
  - Based on RDO

- Geneva Computer centre (adding a new DC)

- Highly **scalable** architecture
  - 48 cells on 5 regions

- Currently running **Stein\*** release
  - Some services already in Xena release

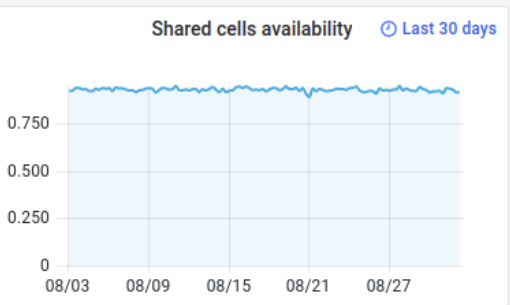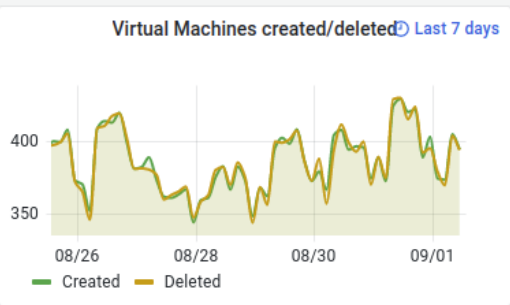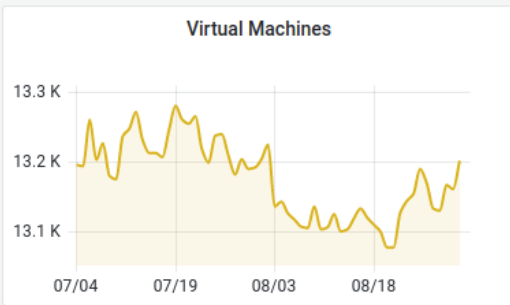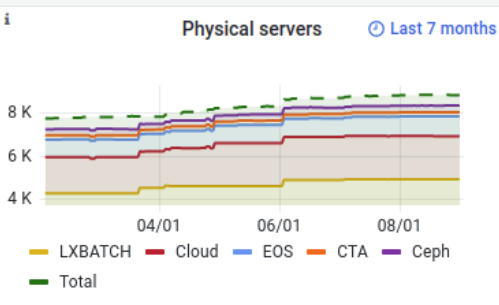## Openstack services statistics

| Users | Projects | Loadbalancers | Images | Volumes | Volumes si... | File Shares | File Shares... | Object Stor... | Object Stor... |
|---|---|---|---|---|---|---|---|---|---|
| 3382 | 4586 | 327 | 4360 | 7329 | 3.78 PB | 5079 | 1.39 PB | 476 | 63.1 TB |

| Servers | | | | Cores | | | RAM | | | Batch | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Physical | Physical in use | Hypervisors | Virtual | Physical | Hypervisors | Virtual | Physical | Hypervisors | Virtual | Servers | Cores | RAM |
| 9112 | 8820 | 2013 | 13674 | 486 K | 58.3 K | 87.9 K | 2.02 PB | 379 TB | 205 TB | 5199 | 281291 | 1.07 PB |

## Time series

# Initial offering



IaaS+

IaaS

User Visible

# CERN Cloud Infrastructure - now

**IaaS+**

| Orchestration | Container Orchestration | Automation | Web |
|---|---|---|---|
| heat | magnum | mistral | horizon |

**User Visible**

**IaaS**

| Network | Compute | | Storage | | | Identity | Key manager |
|---|---|---|---|---|---|---|---|
| neutron | ironic | nova | cinder | manila | glance | keystone | barbican |

**Infra**

| Accounting | Metric aggr | Monitoring | | Automation | Probing | Notifications | Integration |
|---|---|---|---|---|---|---|---|
| reporter | kapacitor | dblogger | collectd | rundeck | rally | rabbitmq | cornerstone |

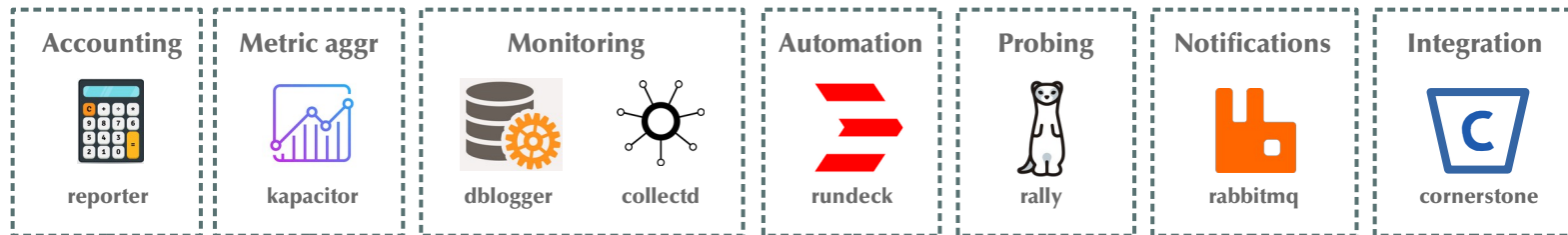# Service deployment

- From shared to "per service" architecture

  - Break dependencies between services

  - Some shared components (rabbitmq, loadbalancers, caches)

  - Freedom to update components under the same API/RPC version

- All deployed in VMs on our own infrastructure: *"eat our own dogfood"*

  - Bootstrap procedure and recovery methods

- Puppet managed running on CentOS 7 (hypervisors) and CentOS Stream 8 (services)

# Service operations

- Deployment upgraded since **July 2013**

- Per-service upgrade model (A/B or in place)

- Compute + Storage availability zones (3 zones each)

- Huge investment on **automation**:

  - Delegate as much as possible administrative tasks (repair team, quota mgmt, end-user)

  - Detect and fix known issues

  - User communication

- Quite some big campaigns:

  - KVM consolidation, Spectre/Meltdown and L1TF, Cold Migration
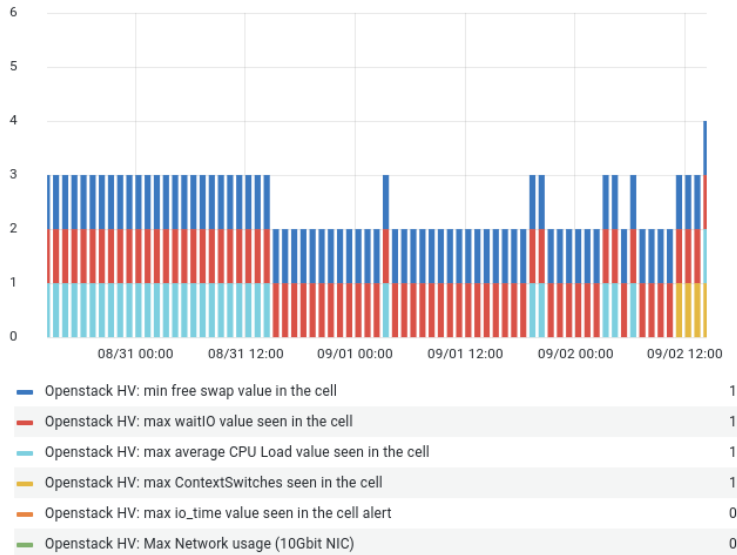
# Cloud Monitoring

- Use same monitoring pipeline as any other IT service

  - Metrics (Collectd => InfluxDB)

  - Logs (Flume => Kafka => ES, HDFS)

- Custom sensors for VM monitoring, service metrics

- Threshold based alarming on individual nodes

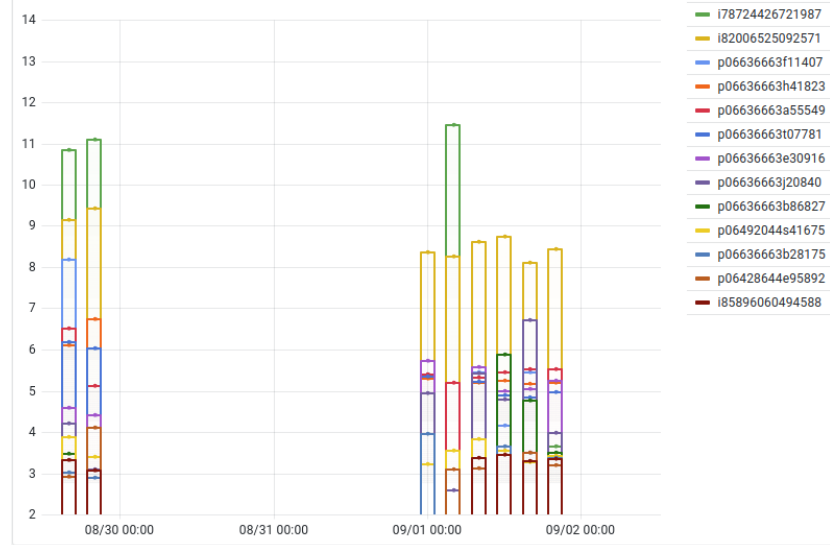- Per-service grafana dashboards

# Find the needle in the haystack

- Threshold based alarming on extreme cases

- Anomaly detection to find misbehaving nodes
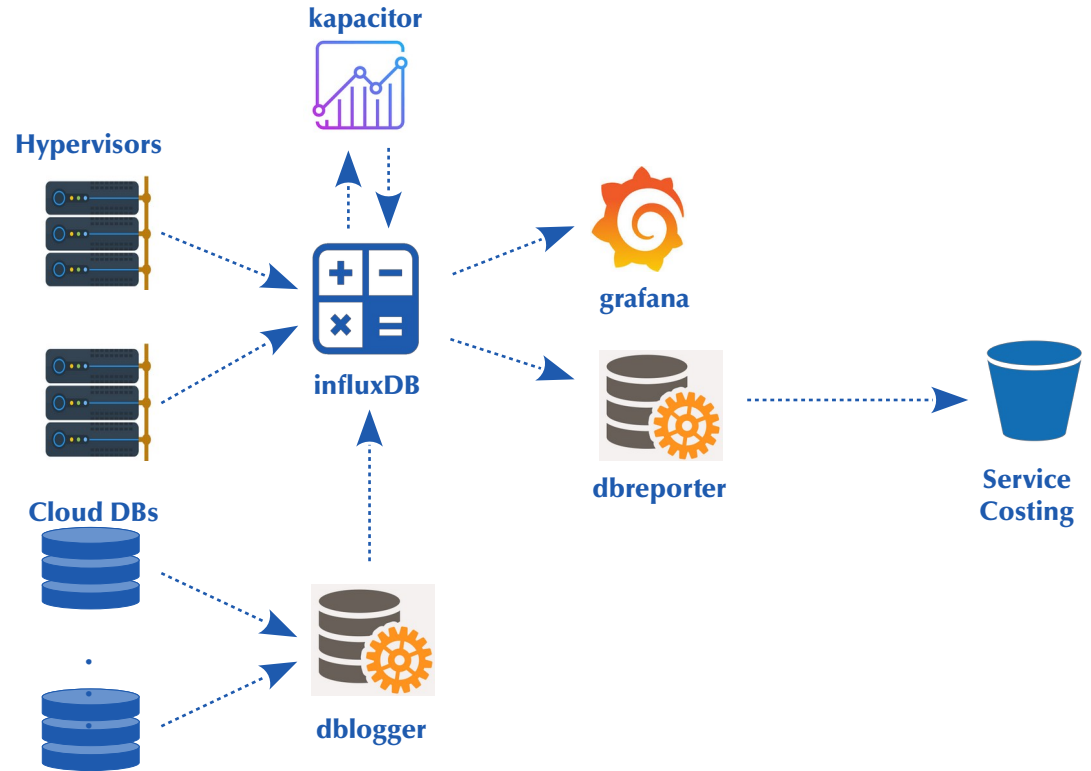
# Cloud Accounting

- All resources grouped by project
  - chargegroup & chargerole

- Producers data
  - VM stats sensor in all HVs
  - Metrics from OpenStack DBs

- Stored in InfluxDB

- Aggregated by Kapacitor

- Exported to Service Costing



kapacitor

Hypervisors

influxDB

grafana

dbreporter

Service
Costing

Cloud DBs

dblogger

# Identity management

- Available to all CERN Users

  - On-demand provision of resources to federation users (based on group membership)

- Types of projects (owned by a CERN primary account)

| | Affiliation Expired | User Disabled | User Deletion |
|---|---|---|---|
| **Shared** | Promote | - | - |
| **Personal** | - | Stop | Delete |

- Provisioning and cleanup in Mistral workflows (inter-dependency handling)

# Resource management for end user

# Security approach

- TLS everywhere (Regular check on TLS security level on endpoints)

- DoS protection on Load Balancers

- 2FA for Administrative operations

- Follow up CVEs on openstack/virtualisation packages with local backports

- Standard vs Audit notifications

- CERN Security team analyses network traffic and controls external firewall

  - Granted additional permissions to stop/lock user VMs

  - Network Isolate any VM/physical node

# Cloud Probing

- Use Rally as automated probe system

- Focus on infrastructure wide issues



Rally: Number of failing tests

| | Mean | Max ⌄ | Min | Last |
|---|---|---|---|---|
| Total | 24.1 | 34 | 5 | 5 |
| nova | 19.8 | 25 | 5 | 5 |
| magnum | 4.28 | 17 | 0 | 0 |
| cinder | 0 | 0 | 0 | |
| glance | 0 | 0 | 0 | |

Passing % in time frame per availability zone

| availability zone | attach-volume | boot-linux | boot-linux-with-dns | ping-linux | ping6-linux | reboot-linux | live-migrate-linux |
|---|---|---|---|---|---|---|---|
| cern-geneva-a | 99% | 100% | 100% | 100% | 63% | 100% | 96% |
| cern-geneva-b | 100% | 100% | 100% | 100% | 58% | 100% | 94% |
| cern-geneva-c | 100% | 100% | 100% | 100% | 50% | 98% | 100% |
| gva-critical | 98% | 100% | 100% | 100% | 41% | 100% | 100% |

Global actions: Passing % in time frame

| deployment | authenticate | boot-from-snapshot-linux | boot-from-volume-linux | create-and-delete-image | list-images |
|---|---|---|---|---|---|
| global | 100% | 100% | 89% | 100% | 100% |

Cinder actions: Passing % in time frame

| deployment | create-and-delete-snapshot | create-and-delete-volume | create-and-extend-volume | list-volumes |
|---|---|---|---|---|
| cinder | 100% | 100% | 100% | 100% |

Manila actions: Passing % in time frame

| deployment | create-share-and-allow-and-deny-ac | create-share-and-delete | create-share-and-extend | create-share-and-shrink |
|---|---|---|---|---|
| manila | 100% | 100% | 100% | 100% |

# Debugging

- Stateless vs Stateful services

    - Focus on reproducing issue

- Use of dedicated "testing" regions

    - Route user requests to extremely verbose setup

    - Connected to other production services

- Probing on Testing Regions

    - Validate minor/major upgrades

    - Introduce feedback with more user scenarios

# Thank you

More info:

https://computing-blog.web.cern.ch/

All our **open source** code is available on:

https://gitlab.cern.ch/cloud-infrastructure
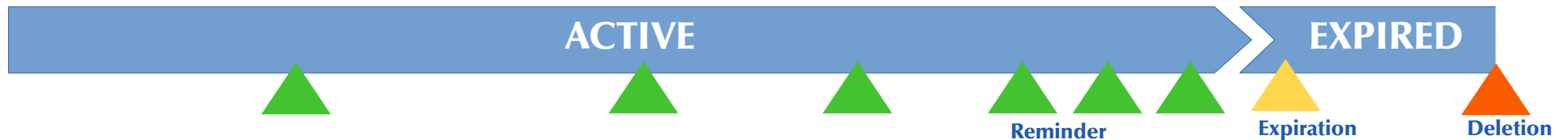
Thank to the work of my team colleagues

# BACKUP SLIDES

# Optimize resource availability - Expiration

- Each VM in a personal project has an expiration date

- Set shortly after creation and evaluated daily

- Configured to 180 days and renewable

- Reminder mails starting 30 days before expiration

- Implemented on a Workbook in Mistral

**ACTIVE** **EXPIRED**

Reminder    Expiration    Deletion

# Task delegation

- Rely on Rundeck for offloading tasks to different teams

  - Repair Team

  - Resource coordinator

  - Cloud operations

- Example: disk replacement

# Why baremetal provisioning?

**ironic**

- VMs not sensible/suitable for all of our use cases

  - Storage nodes, HPC clusters, Batch nodes

- Complete our service offering

  - Physical nodes (in addition to VMs and containers)

  - OpenStack as single pane of glass

- Simplify hardware provisioning workflows

- Consolidate accounting & bookkeeping

  - Machine re-assignments will be easier to track

# HW lifecycle at CERN

ironic

# Ironic Service setup and status



ironic

**Users:**

- Cloud

- Batch

- HPC

- Windows

- DB

- ...