# Ceph per lo storage Cloud GARR

Alberto Colla

per il Dipartimento CSD GARR

WORK SHOP GARR 2022

NET MAKERS

# Table of Contents

# GARR Storage and Computing Department (CSD)



Role:

- **provider** of resources ("long tail of science")
- resource **aggregator** (**federation**)
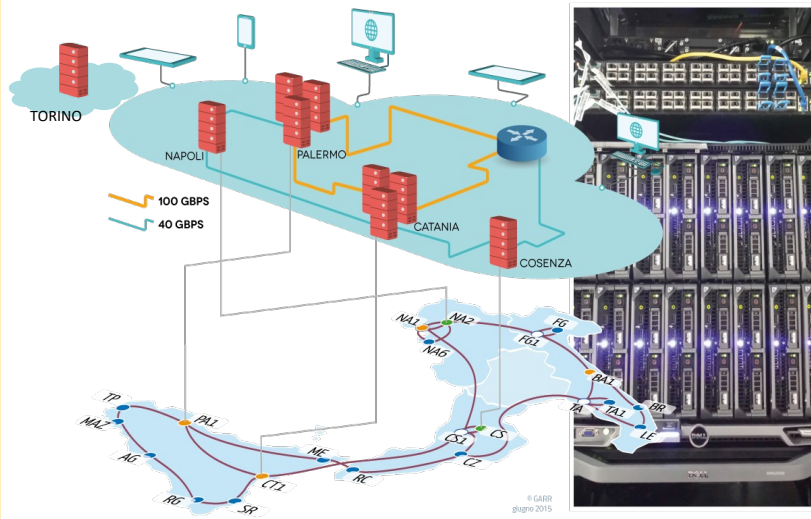- embody a **replicable** model for storage & computing provisioning

Goals:

- **harmonize** (SSO / federation of resources)
- build **secure** and **open** infrastructures
- enhance **user experience**

# GARR Computing and Storage Infrastructure

## Overall

- CPUs: 7600 cores
- RAM: 60 TB
- **Storage: 13.5 PB** (15%) **SSD**
- GPUs: 333 TFLOPS
- Datacentres: 5



## Storage servers

- **DELL PowerEdge R740XD2**
- 48 HT cores, 384 GB RAM
- 3x 1.9TB SSD Mixed Use ➡️ RocksDB
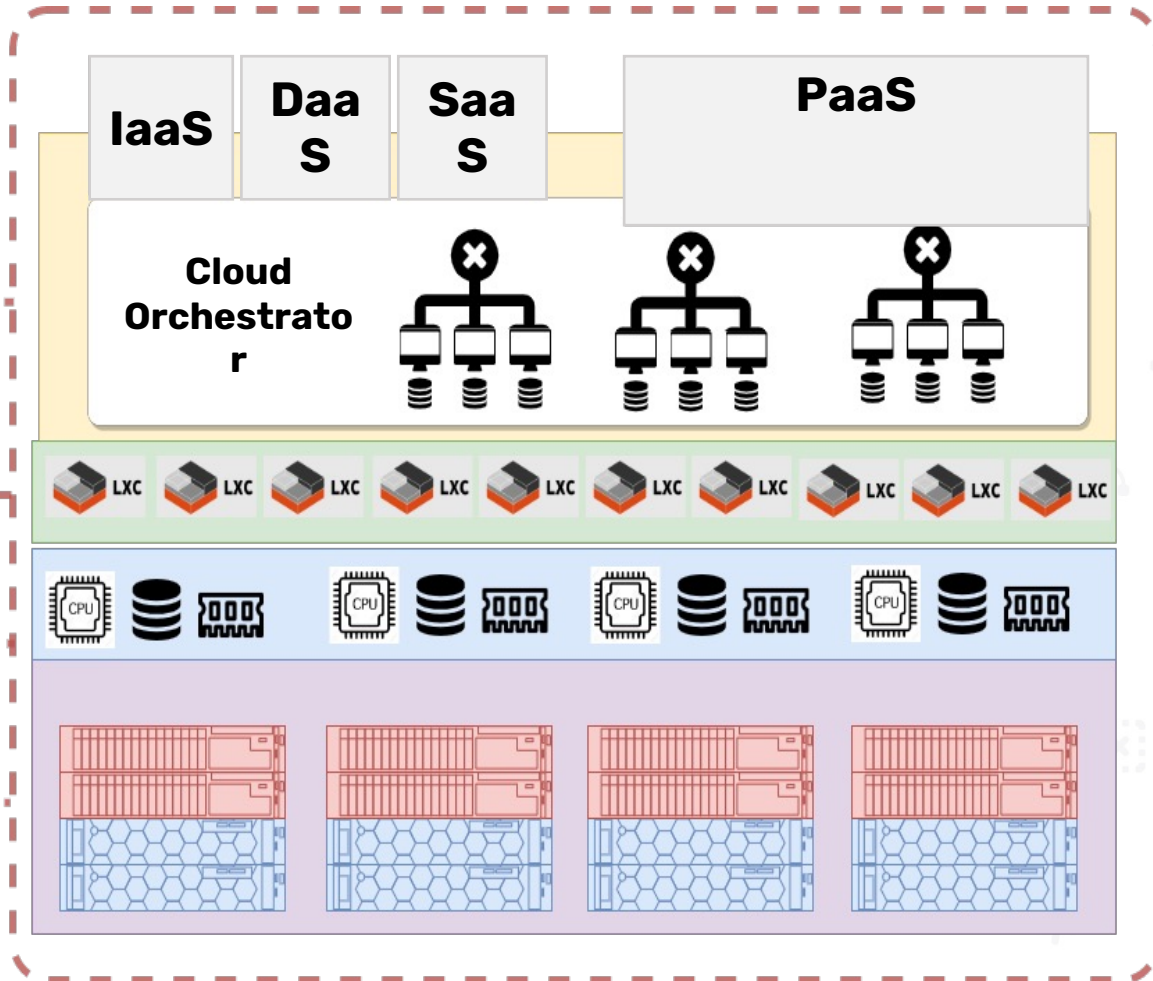- **5x 7.68 TB SSD** Read Intensive
- **14x 18TB HDD** 7Krpm

# GARR Cloud Infrastructure: 4-Layers *recipe*
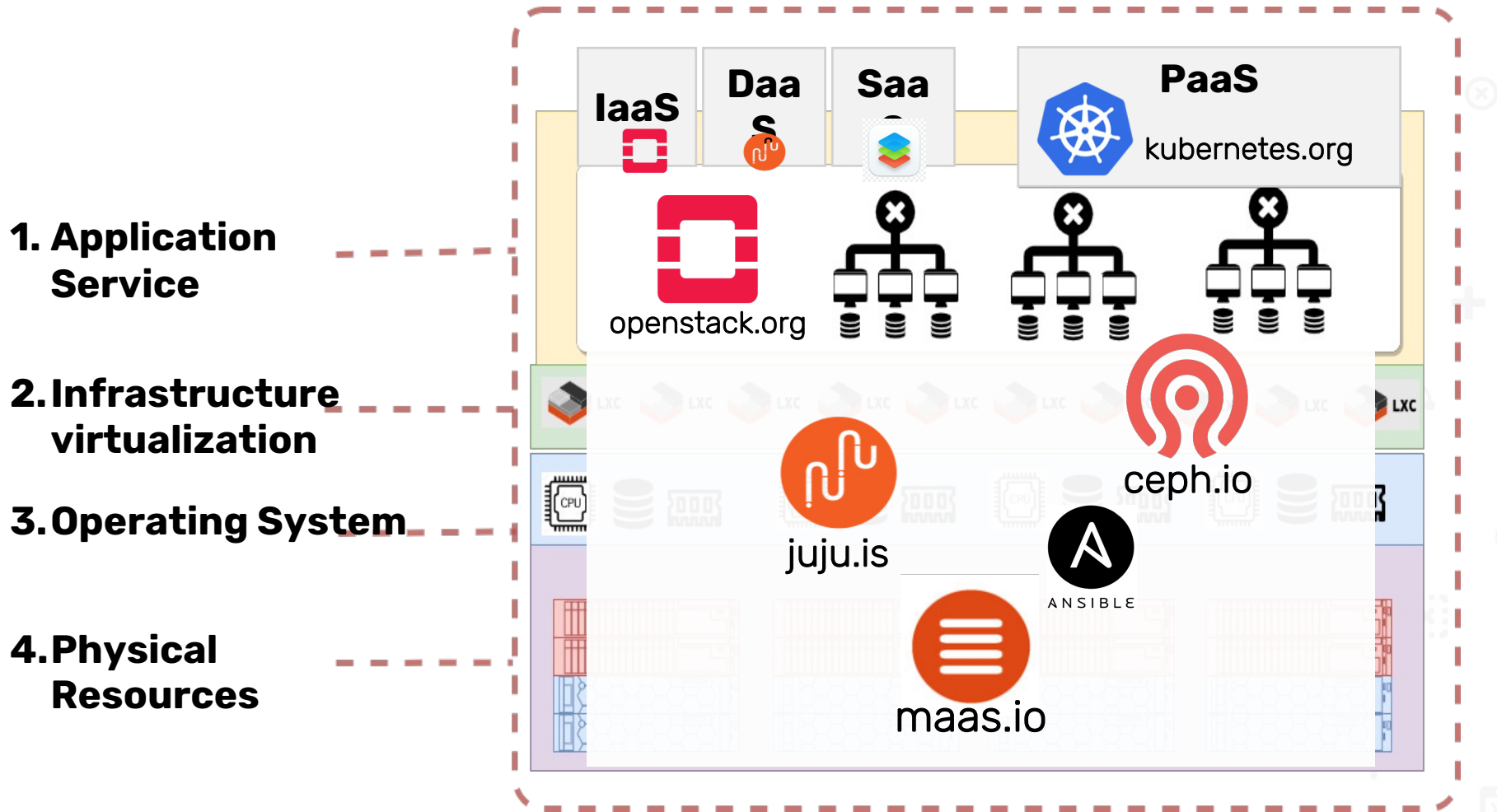


1. **Application Service**

2. **Infrastructure virtualization**

3. **Operating System**

4. **Physical Resources**

# GARR Cloud: the engines

1. **Application Service**

2. **Infrastructure virtualization**

3. **Operating System**

4. **Physical Resources**

IaaS

DaaS

SaaS

PaaS
kubernetes.org

openstack.org

ceph.io

juju.is

ANSIBLE

maas.io

WORK SHOP GARR 2022   NET MAKERS

# OpenStack

**Horizon** (Dashboard)

**Keystone** (Identity Management)

**Nova** (Compute, where VMs are run)

**Glance** (Image Service)

**Cinder** (Block Storage, persistent storage for VMs)

**Swift** (Object Storage, snapshots and not frequently updated data)

~~Neutron (Networking and SDN)~~

# GARR Cloud architecture

**Global services**
- ✓ Identity
- ✓ Images
- ✓ Object Store
- ○ geo-distributed
- ○ DNS HA

**GARR regions**
- ✓ Compute
- ✓ Network
- ✓ Block Store
- ○ Share Identity, Images, Object Store
- ○ Local services

**Federated regions**
- ✓ Compute
- ✓ Network
- ✓ Block Store
- ○ Share Identity
- ○ Local services
- ○ Managed by You

# OpenStack and Ceph

**Horizon** (Dashboard)

**Keystone** (Identity Management)

**Nova** (Compute, where VMs are run)

**Glance** (Image Service)

**Cinder** (Block Storage, persistent storage for VMs)

**Swift** (Object Storage, snapshots and not frequently updated data)
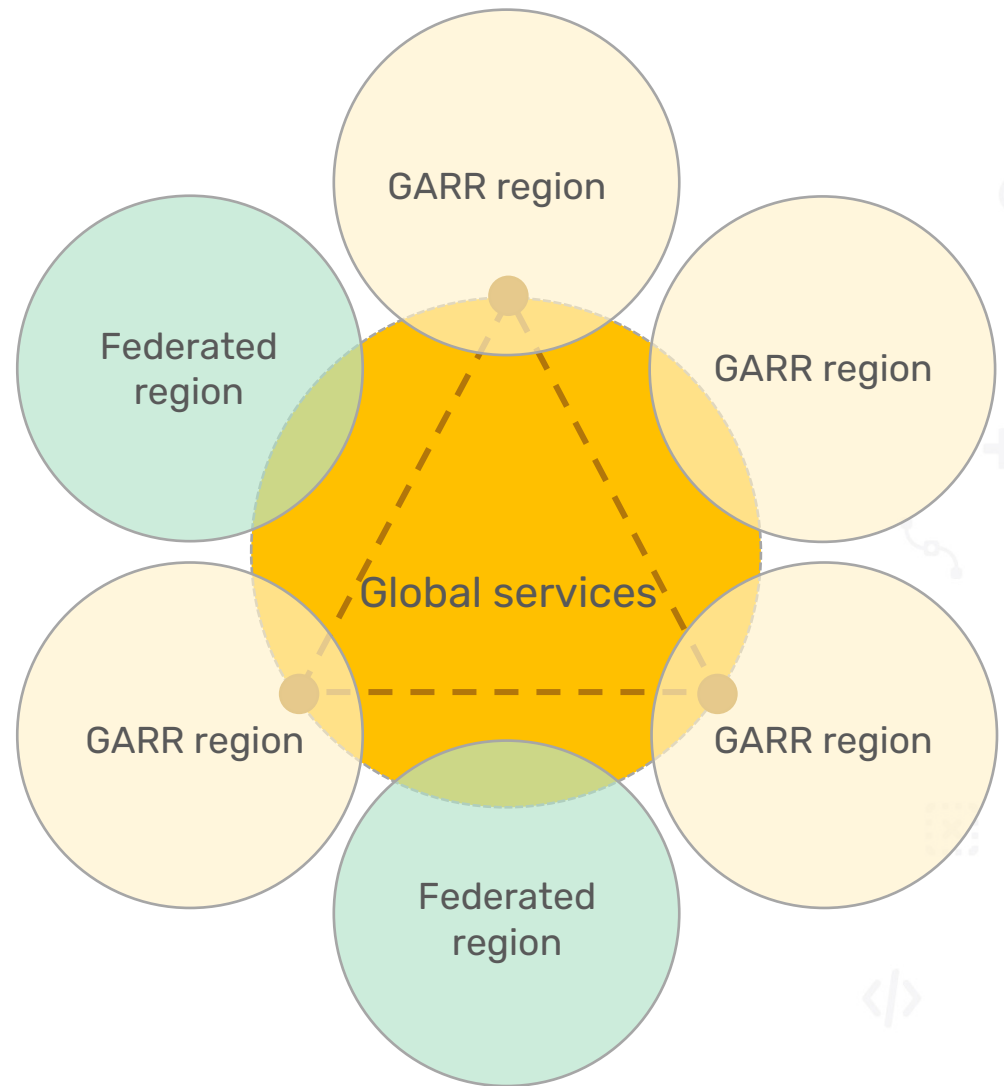
**Neutron** (Networking and SDN)

# What is Ceph

Ceph is **open source software** designed to provide highly scalable **object-**, **block-** and **file-**based storage under a unified system

# Why Ceph

o Open-source, distributed storage

o Lack of SPOF (single point of failure)

o Runs on commodity hardware

o Aggregates any server with any network and disk setup

o Consistently evolving, new functionalities and several improvements

o Extremely lively and reactive community (ceph-users@ceph.io)

"Which OpenStack Block Storage (Cinder) drivers are you using?"

https://www.openstack.org/analytics

# A little Ceph glossary

- **Object Store Device (OSD)**: the physical disk (plus a slice of CPU/RAM to manage it)

- **Monitors (Mon)**: maintain the map of the cluster state, keeping track of active and failed cluster nodes, cluster configuration, etc.; handle such map to clients

- **Managers (Mgr)**: maintain cluster runtime metrics, enable dashboarding capabilities, provide interface to external monitoring systems
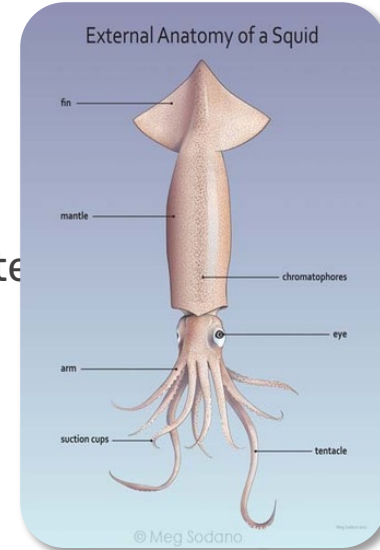
- **Rados Block Device (RBD)**: Ceph's block storage component

- **Rados Gateway (RGW)**: Ceph's object storage APIs (swift and S3)

- **Metadata Servers (MDS)**: store metadata for the Ceph File System

---

- Ceph stores data (objects) within **pools**.

- Pools ensure data redundancy (**Replica** vs. **Erasure-Coding**)

- Within a pool, objects are mapped to **Placement Groups** (**PG**), and placement groups to the OSDs

- OSDs belong to **Device Classes**: default ("hdd", "ssd") or custom ("test", "big")

# A little Ceph glossary

o The **CRUSH** (Controlled Replication Under Scalable Hashing) algorithm determines how to store and retrieve data

   o The **CRUSH map** describes the physical topology of the cluster:
      it is passed to clients who can then interact directly with the cluster

   o The **CRUSH rules** define data placement policy,
      e.g.: what device class to use, how many chunks per rack/host,…

o Two examples of CRUSH rules:

```
rule ssd_rule {
    id 1
    type replicated
    min_size 1
    max_size 3
    step take default class ssd
    step chooseleaf firstn 0 type rack
    step emit
}
```

```
rule default.rgw.buckets.data {
    id 6
    type erasure
    min_size 3
    max_size 10
    step set_chooseleaf_tries 5
    step set_choose_tries 100
    step take default class big
    step choose indep 5 type host
    step chooseleaf indep 2 type osd
    step emit
}
```

# Ceph deployment - two cases

**Catania, Palermo, Napoli regions**

o Ceph w. **ceph-ansible** on storage nodes OpenStack with **Juju** on compute nodes

o Ceph-OpenStack connection via **ceph-proxy** charm

o Ceph-proxy

    o Is configured with **Mon IP addresses** and **ceph-admin credentials**

    o Is related to the Openstack modules

    o Creates users, pools etc. in Ceph according to the directives of the related OpenStack modules

    o passes config parameters to Ceph clients on the modules

**Torino region:**

o Hyper-converged systems (storage&compute)

    o few servers, but powerful

o Deploy Ceph and Openstack with Juju

o Ceph-mon charm directly connected with OpenStack modules (juju relations)



storage nodes     compute nodes



storage & compute nodes

# From Ceph to OpenStack

| Block store (cinder) | | | | Object store | |
|---|---|---|---|---|---|
| **volume ty** | **defau** | **fa st** | **capaci ty** | **swift / rclone / s3 / ...** | |
| Interfac | rbd | rbd | rbd | rgw | rgw |
| Data Protecti | replica3 | replica3 | data: HDD EC6+4 metadata: SSD repl3 | EC6+4 | EC6+4 |
| Device Class | hdd | ssd | big | big | big |
| Device | | | | ←→ mirro →→ | |

**Ceph Catania**

**Ceph Palermo**

# Benchmarks



SSD-R
EC-R
DEFAULT-R
SSD-W
EC-W
DEFAULT-W

random read+write / 4M chunk

random read+write / 4K chunk

BW (MiB/s)

IOPS

n. processes

# Monitoring

o Currently: Zabbix monitoring configured with custom Ansible scripts

o Ceph recently included native API to Zabbix and Grafana

➡ We will integrate Ceph and Grafana,

to have a single monitoring dashboard for all the cloud systems



cephadm.ct1: Cluster OSD status

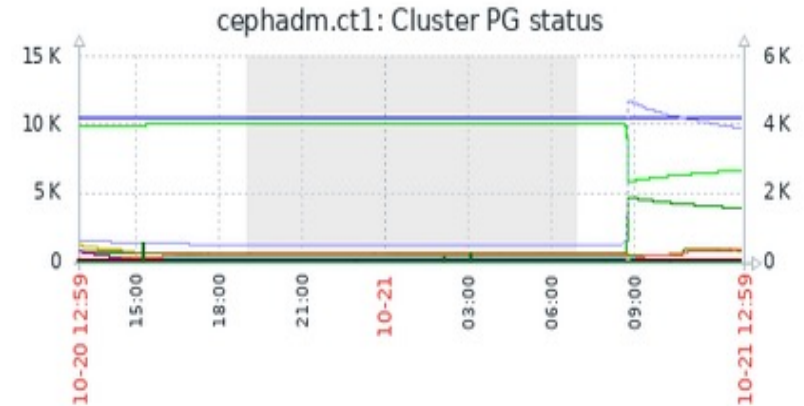| | | last | min | avg | max |
|---|---|---|---|---|---|
| Ceph OSD down | [max] | 0 | 0 | 0.7201 | 26 |
| Ceph OSD up | [max] | 298 | 260 | 293.1017 | 298 |
| Ceph OSD out | [max] | 1 | 1 | 1.3304 | 15 |
| Ceph OSD in | [max] | 297 | 259 | 292.4725 | 297 |

cephadm.ct1: Cluster PG status

| | | last | min | avg | max |
|---|---|---|---|---|---|
| Ceph PG total | [avg] | 10.52 K | 10.52 K | 10.52 K | 10.52 K |
| Ceph PG active | [avg] | 3.92 K | 478 | 1.14 K | 4.7 K |
| Ceph PG active-clean | [avg] | 6.6 K | 5.82 K | 9.39 K | 10.05 K |
| Ceph PG degraded | [avg] | 314.6667 | 176 | 233.5542 | 487 |

# Monitoring - Ceph dashboard

# Ceph management – a little cookbook

o Enable "balancer": can also tune activity window and threshold

```
ceph balancer on ; ceph balancer mode upmap ; # default mode
```

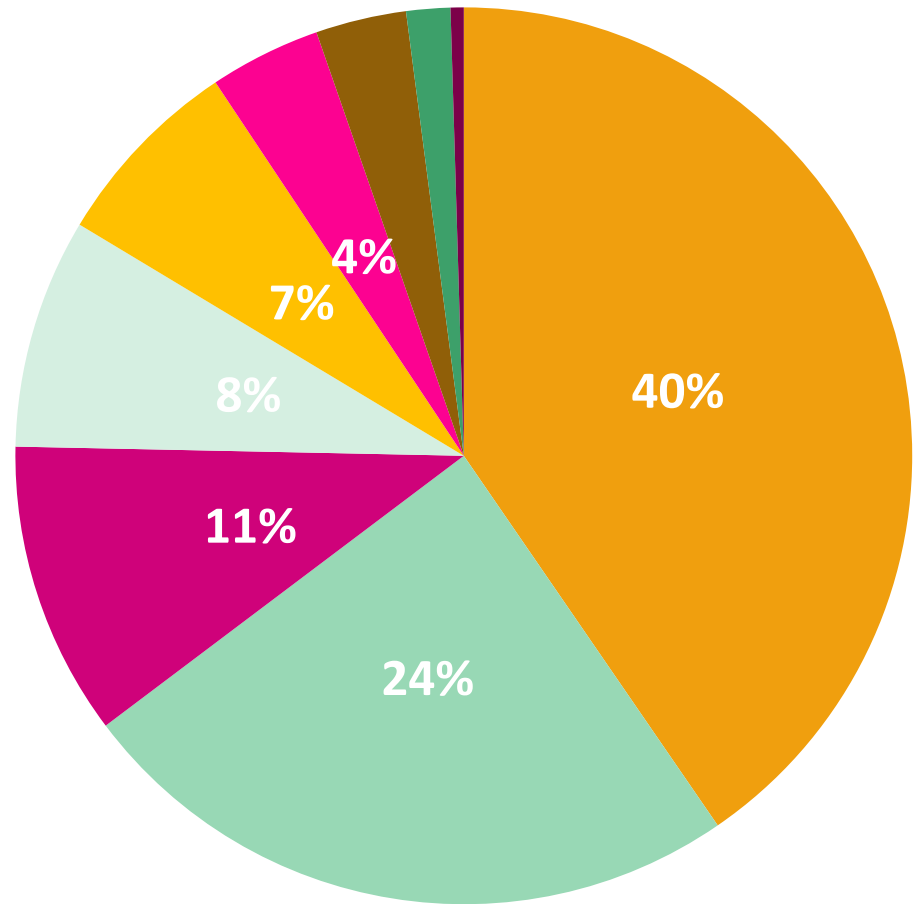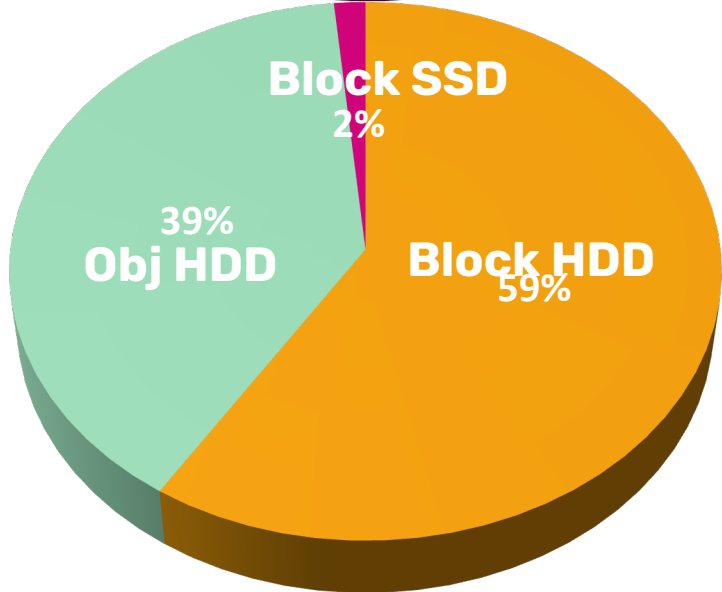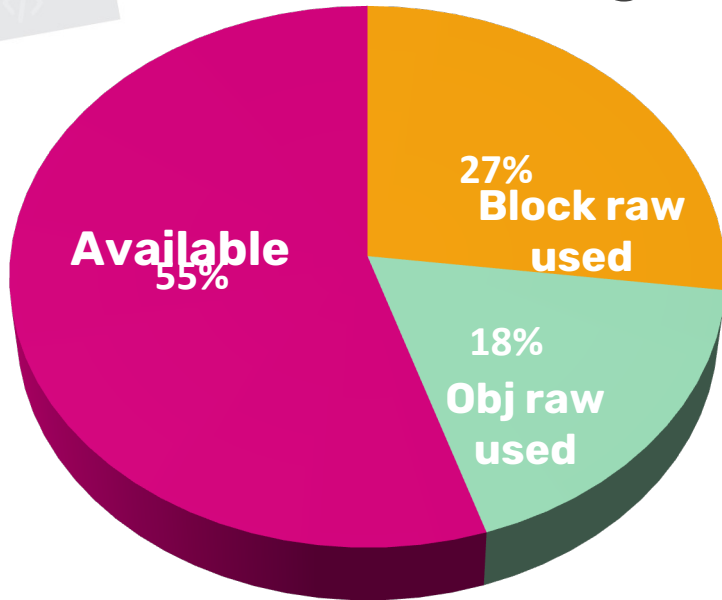o Occasionally, "balancer" fails and one OSD may "drift"

    o Adjust OSD weight, let it go for ~hour, set it back

```
ceph osd crush reweight osd.XXX 0.9
```

o Reduce impact of rebalancing/recovering operations:

```
ceph tell 'osd.*' injectargs --osd-recovery-sleep=0.1
        --osd-max-backfills=1 --osd-recovery-max-active=3
        --osd-recovery-max-single-start=1
```

    o https://docs.ceph.com/en/latest/rados/configuration/osd-config-ref/

    o ignore message stating restart is required

o Throttle data movement when adding servers

```
ceph balancer off
```

    o add servers, move OSD disks, do whatever…

       (see https://github.com/digitalocean/pgremapper or

       https://github.com/cernceph/ceph-scripts.git)

```
./tools/upmap/upmap-remapped.py | sh
```

    o re-run, to force all "remapped" PGs in their current position

```
ceph balancer on ; # let balancer gradually move PGs
```

    o at any time, can pause data movement by "upmap-remapped.py"

# Object store clients

o rclone
  o swift endpoint
  o cli provides functions equivalent to rsync, cp, mv, ls, ncdu, tree, rm, etc.
  o **rclone mount**: mounts Object store as a disk on many systems
  o (optional) **server side encryption**
  o Works on Windows and Linux
  o (Experimental) GUI available

o s3cmd
  o s3 endpoint
  o client for uploading, retrieving and managing data in s3 bucket

o s3fs
  o s3 endpoint
  o large subset of POSIX including reading/writing files, directories, symlinks, mode, uid/gid, and extended attributes
  o optional server side encryption

# GARR Cloud storage utilization



Left top pie chart:
- 27% Block raw used
- 18% Obj raw used
- Available 55%

Left bottom pie chart:
- Block SSD 2%
- Obj HDD 39%
- Block HDD 59%

Right pie chart:
- 40%
- 24%
- 11%
- 8%
- 7%
- 4%

Legend:
- estensione CED istituti
- ricerca biomedica
- servizi GARR
- ricerca spaziale
- earth observation
- ricerca informatica
- beni culturali
- formazione/didattica
- altro

# Conclusions

o GARR Cloud uses Ceph since 2015 with satisfaction

- o Great direct and indirect support from [ceph-users@ceph.io](mailto:ceph-users@ceph.io)
- o Upgraded seamlessly from Ceph v0.9 (Hammer) to v16 (Octopus)
- o ~10 PB raw storage managed
- o 1.5 PB net / 4.5 PB raw currently used / requested by 250 active users
- o Hot data migration to the new hardware infrastructure done
  - o Not a single bit lost; performance degradation during migrations handled
- o Complete configuration of new hardware by end of November

o Next steps

- o Complete CT⇔PA mirroring of Object Store and Glance (Cloud images) pools
- o Implement Ceph File System provisioning (OpenStack Manila)

Questions
?

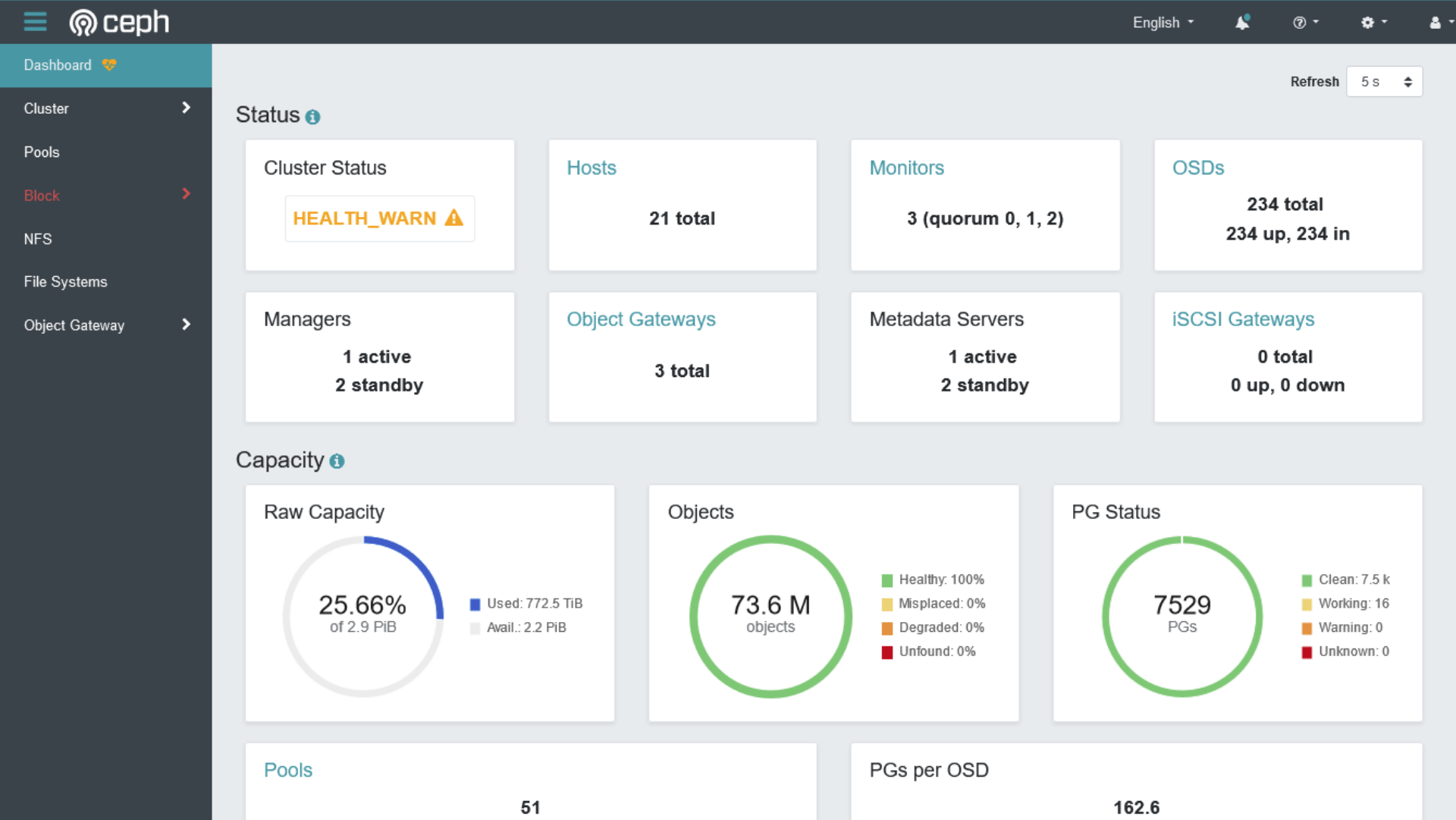photo © Gabriele Colla

# Backup

# Monitoring (Pacific: PA1)

# cloud evolution timeline

## 2015/6

**Tender** GARR X progress.
Preliminary study/tests
**First prod setup**

## 2017/8

Production ready
**Region Pa -> Ct**

## 2019/0

**+Region Na**
Keystone **HA geo**
**Fed GARRCloud** model: INGV
sale sismiche

## 2021

**+Region To**
1st Federated DC: Polito
Federation model
adoption: HPC4AI

## 2022

**Federation spreads**:
EMSO DC ongoing
...
i.e. Uninuvola PG soon

10

2x storage
12

14

HW upgrade: storage Ct/Pa

HW upgrade + GPU
32 A30 + 12 A100

16

O~S juno

O~S mitaka

O~S queens

O~S Stein

O~S Wallaby

*8500 core & 11 PB*
8x blade 16 poweredge m620
8x MD3860/3060e

16

*+2.5 PB + ssd*
+12 power edge R740 XD2

*11500 core & 16PB*
11x 8 poweredge R650/R750

20