



MINISTÈRE DE L'INTÉRIEUR



I Context

A Cloud offer

B Dev to production strategy

C Focus on storage

II Ceph implementation

A The numbers

B Implementation model

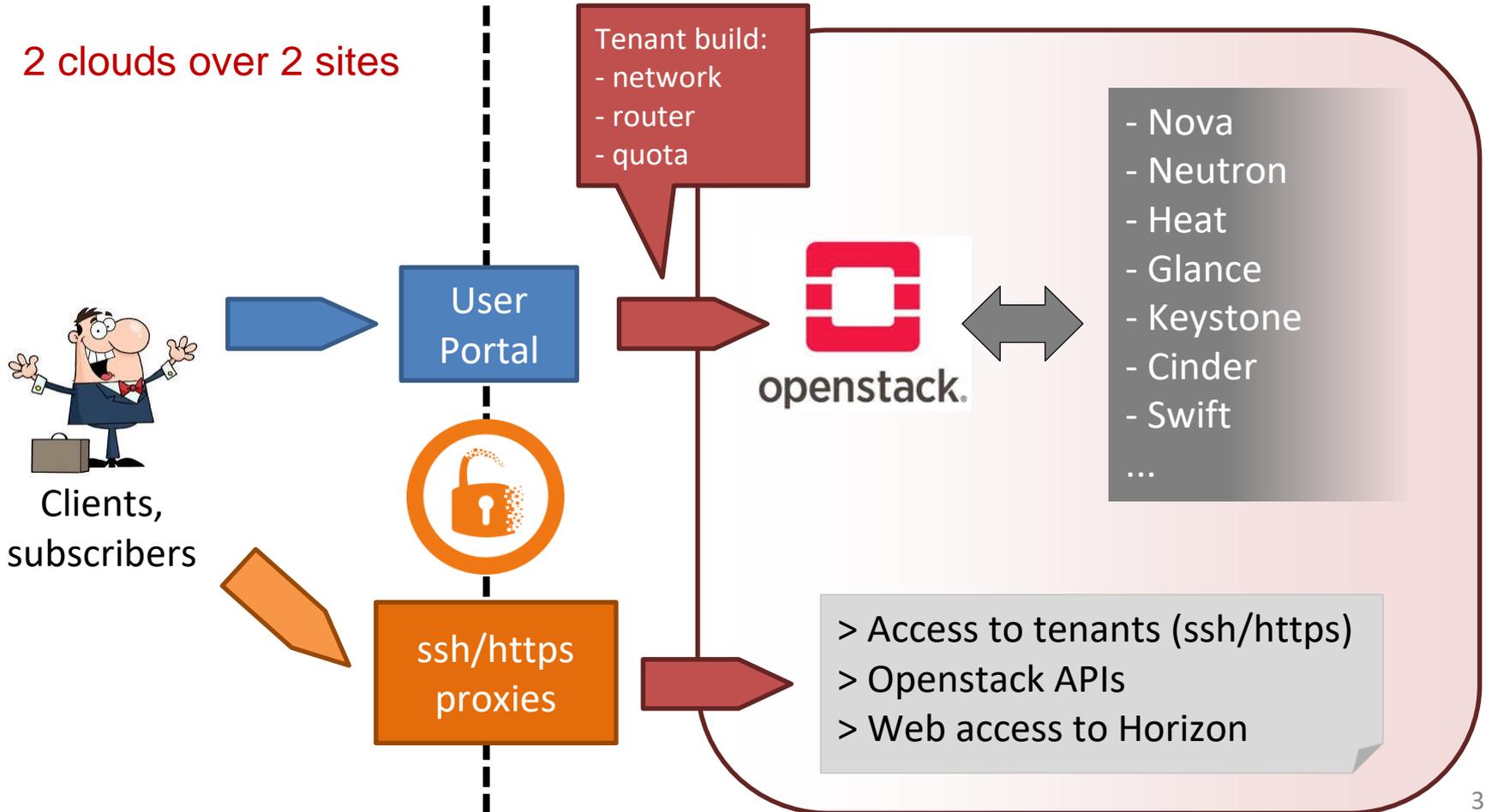
C Storage deployment

D Initial benchmarks with RBD

I. Context

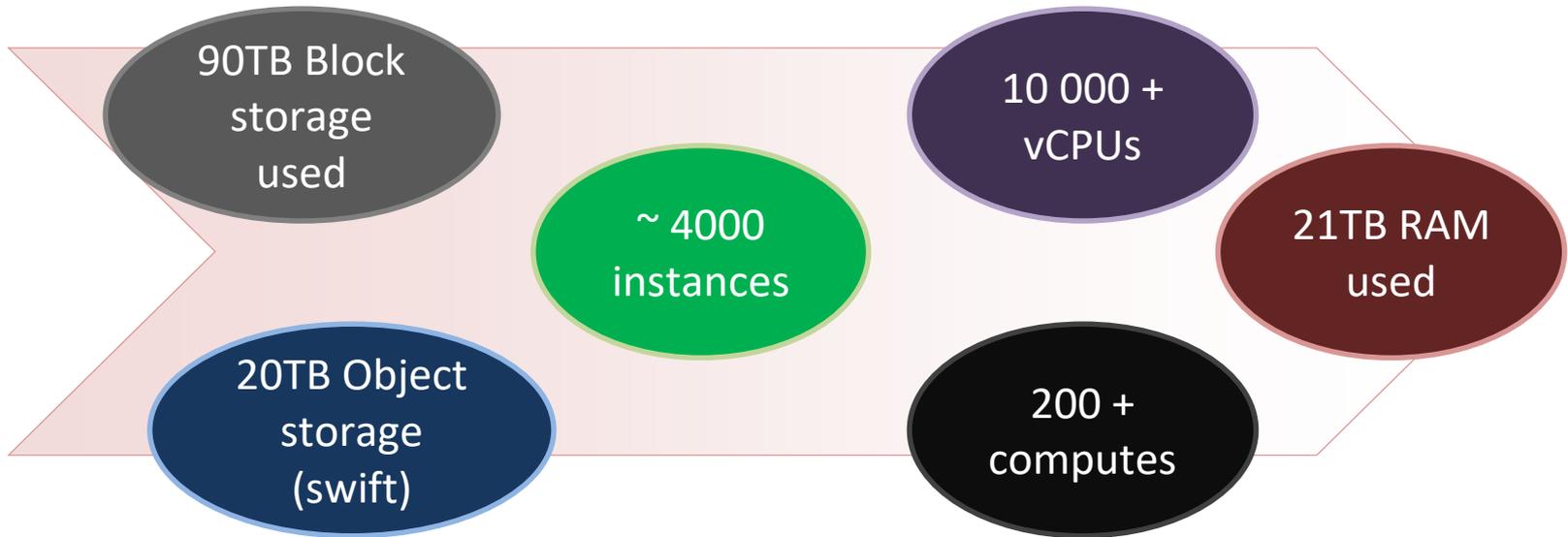
A Our cloud offer

2 clouds over 2 sites



I. Context

A Our cloud offer



I. Context

B Dev to production strategy

Environments

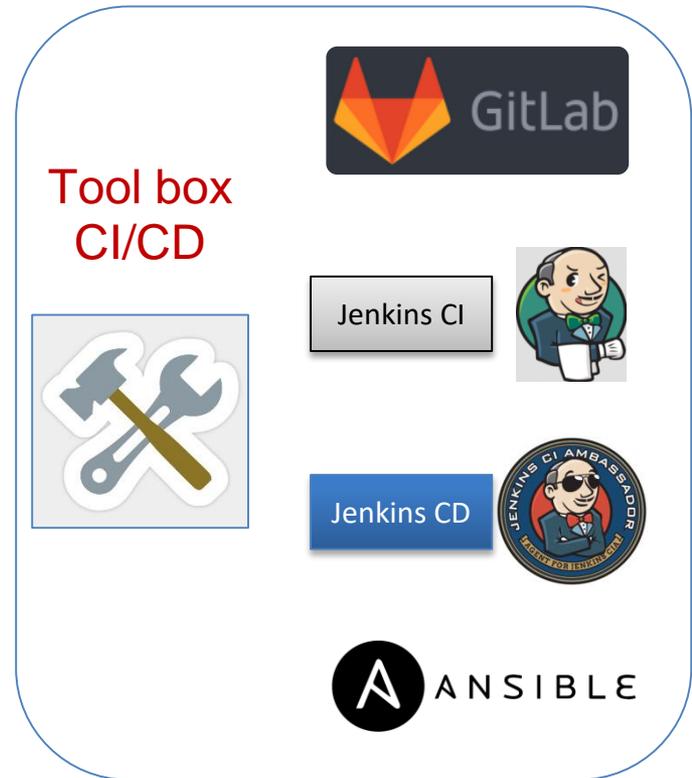
- 2 dev environments for testing
- 2 staging environments (deployed by operations team)
- 2 production environments (deployed by operations team)

Dev Team

- **Developing** new products
- **Maintaining** applications versioning
- **Delivering** new builds to the operations team

Operations Team

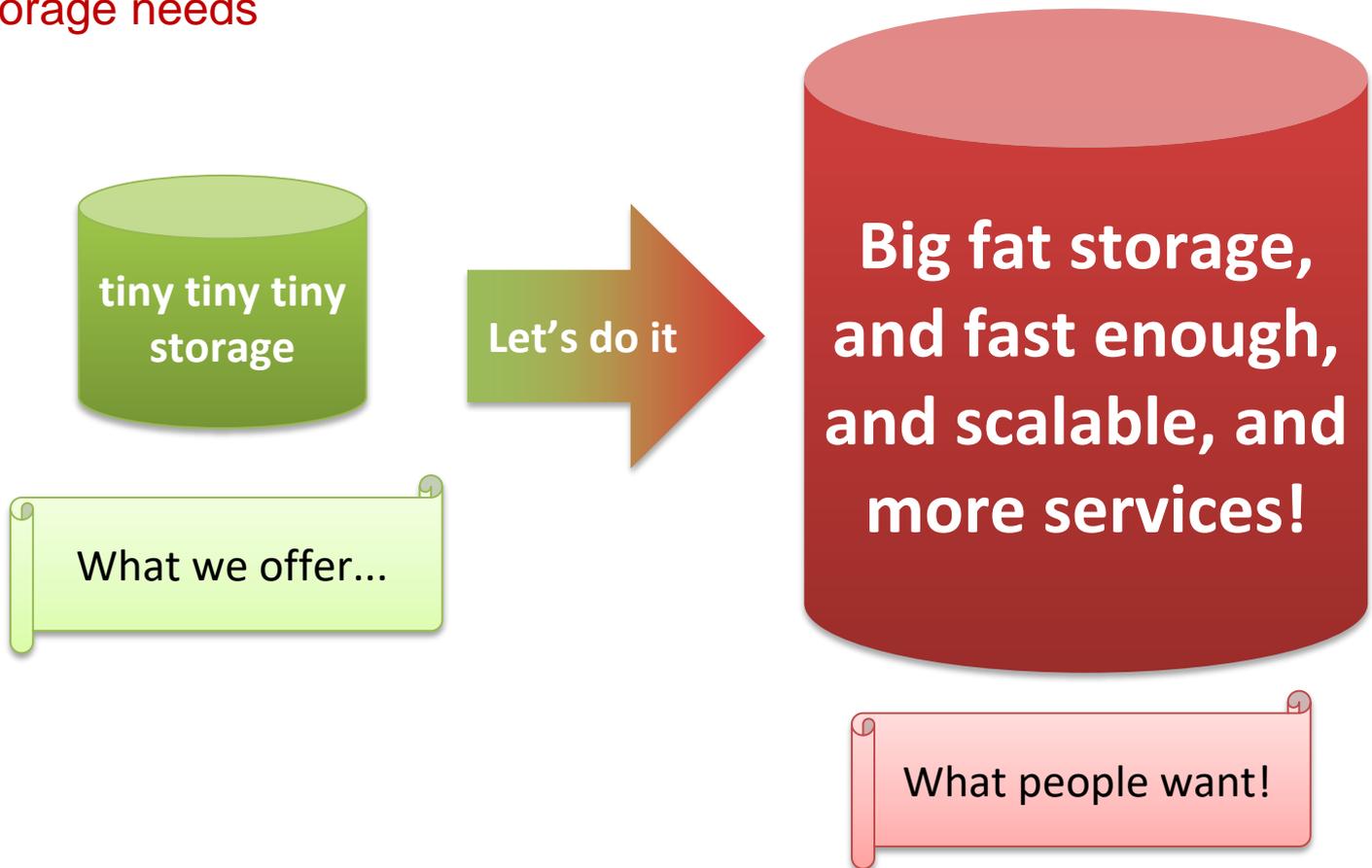
- In charge of staging and **production**
- In charge of services **maintenance**



I. Context

C Focus on the storage

Storage needs

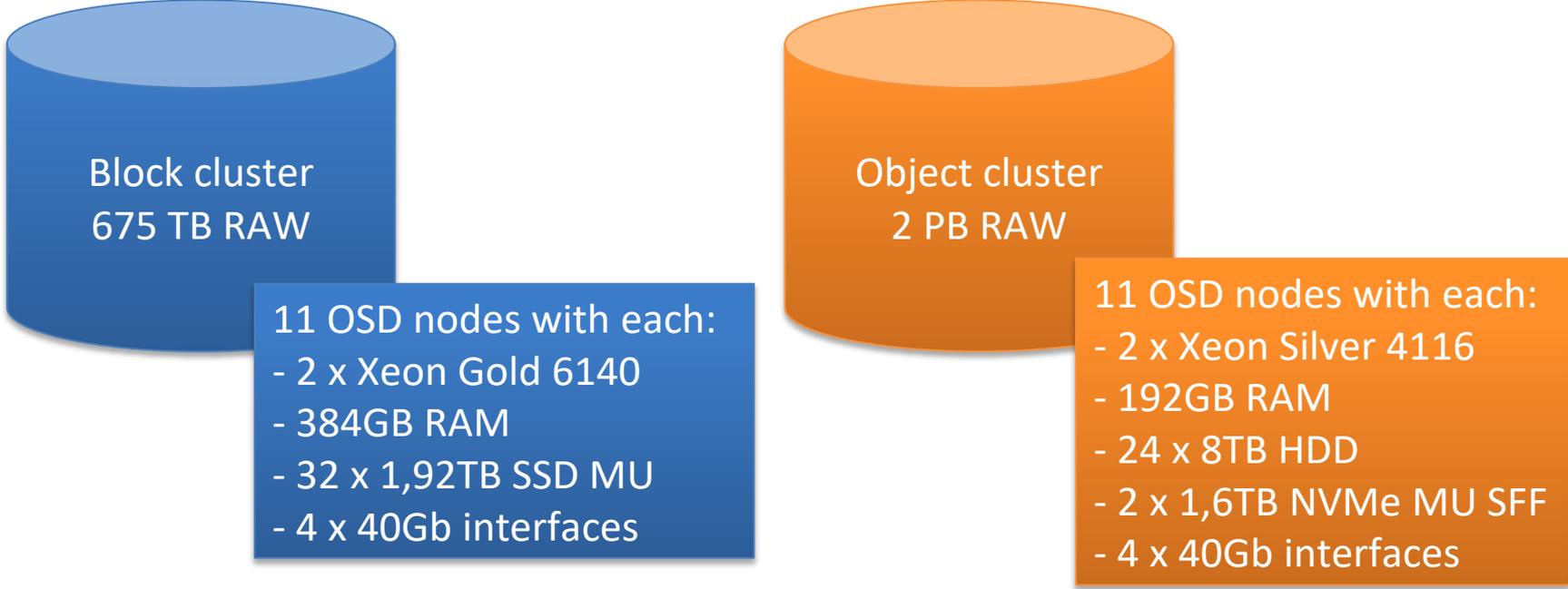


II. Ceph implementation

A The numbers

For each site in production, we are going to deploy:

- 1 Block storage cluster (RBD)
- 1 Object storage cluster
- => smallest clusters for Dev and Staging (4 OSD nodes each)



Block cluster
675 TB RAW

- 11 OSD nodes with each:
- 2 x Xeon Gold 6140
 - 384GB RAM
 - 32 x 1,92TB SSD MU
 - 4 x 40Gb interfaces

Object cluster
2 PB RAW

- 11 OSD nodes with each:
- 2 x Xeon Silver 4116
 - 192GB RAM
 - 24 x 8TB HDD
 - 2 x 1,6TB NVMe MU SFF
 - 4 x 40Gb interfaces

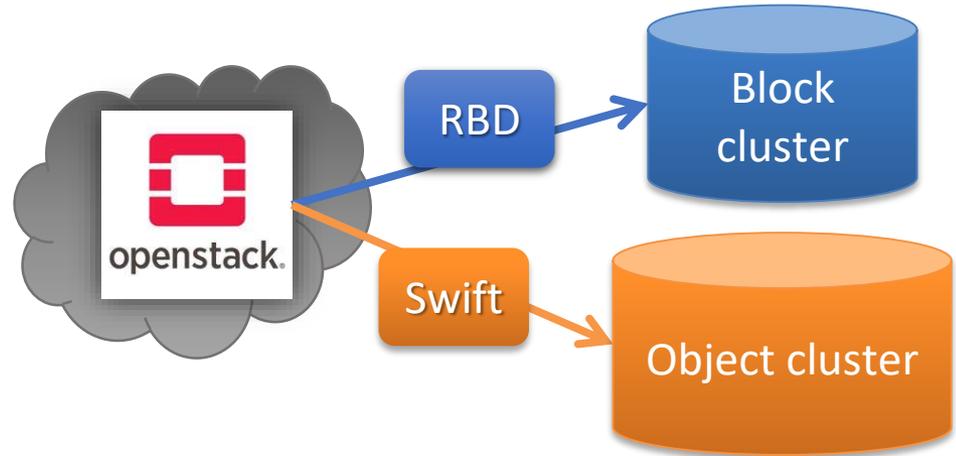
II. Ceph implementation

B Implementation model

Integration with openstack

Each openstack instance is connected to:

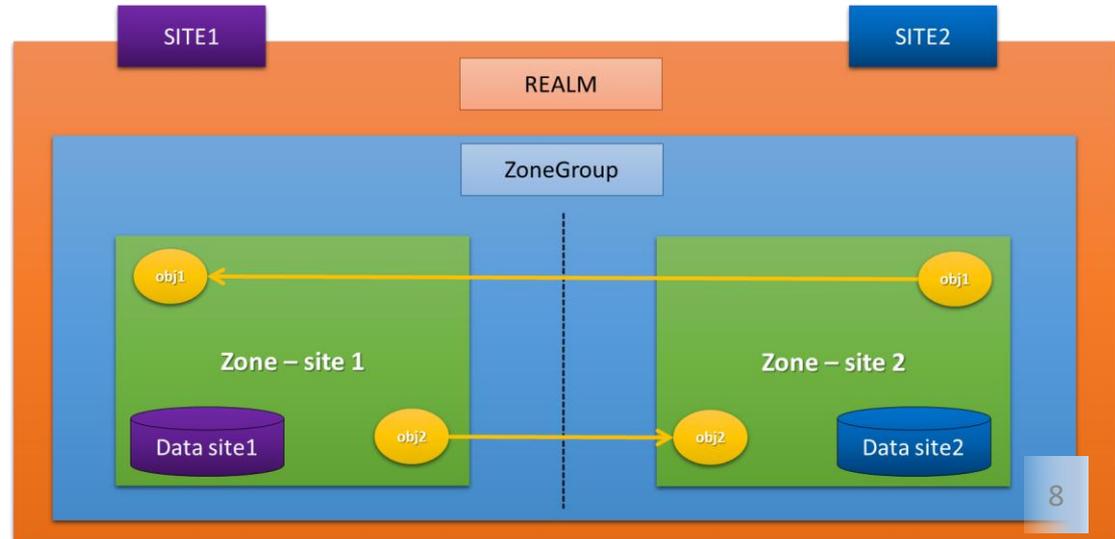
- Ceph block cluster for rbd
- Ceph object storage for swift



Standalone s3 service with replication

Each object cluster is replicated between 2 sites:

- 2 endpoints for clients (1 per site)
- 1 Realm
- 1 ZoneGroup
- 2 zones replicated
- User shared



II. Ceph implementation

C Storage deployment

Each cluster is made of:

- 1 Deployer node (vm)
- **3 Monitors** nodes
- **2 Rados Gateways** for Object clusters
- **4 OSDs** nodes minimum



Ceph version: **Nautilus**

Editor support: **Suse SES6**

OS server: **SLES15-SP1**

Fully automated deployment method from scratch.

Why?



- Control every step of deployment, understand how the product works
- Simplify the deployment for many clusters
- Avoid any human error



II. Ceph implementation

C Storage deployment

Automation!

CD deploys virtual machines :

- cobbler server
- ceph deployer (master)

Deployment steps:

- ceph nodes deployment: cobbler installs every physical node
- ceph deployer configures every node:
 - Ansible
 - ✓ configures sles internal repositories (pulp)
 - ✓ updates servers with the latest staged repositories
 - ✓ installs prerequisite packages (salt-minions...)
 - ✓ updates server configurations
 - ✓ defines salt configuration (nodes roles, custom conf)
 - ✓ configures rados gw and replication between sites
 - Salt / Deepsea
 - ✓ Discovers cluster & configures salt files with defined node roles
 - ✓ Deploys every node with its roles
 - ✓ Deploys a web interface for cluster management and monitoring (you may admin the cluster through the interface)
 - ✓ Removes/Adds nodes in a cluster
 - ✓ ...

II. Ceph implementation

D Initial benchmarks with RBD

Openstack (newton) setup:

- 4 computes dedicated for instances.
- 16 instances (4 per compute)
- 10Gb/s network per compute

Ceph (luminous) setup:

- 4 OSDs nodes (2x10Gb/s pub network and 2x10Gb/s for cluster network)
- 32 OSD, 8 SSD 480GB (SATA) per node
- 256GB ram per node
- 2 Xeon 2650 (14 cores each) per node

Test with FIO:

- Big IOPS: 4MB blocks (for sequential tests writes & reads)

Each compute reached 10Gb/s for writing and reading blocks, so cool!

=> not significantly impact for OSD nodes.

- Small IOPS: 4kB blocks (write & read)

Without RBD cache client: 1600 IOPS / instance (Total of $16 * 1600 = 25600$ IOPS)

With RBD cache client: more than 9500 IOPS / instance ($16 * 9500 = 152000$ IOPS)

=> Some activity noticed for OSD nodes

Experimentation



Questions?