

Une infrastructure mutualisée et recyclée avec Ceph

Bruno Buisson & Philippe Saby
15 novembre 2017



L'Observatoire Midi-Pyrénées

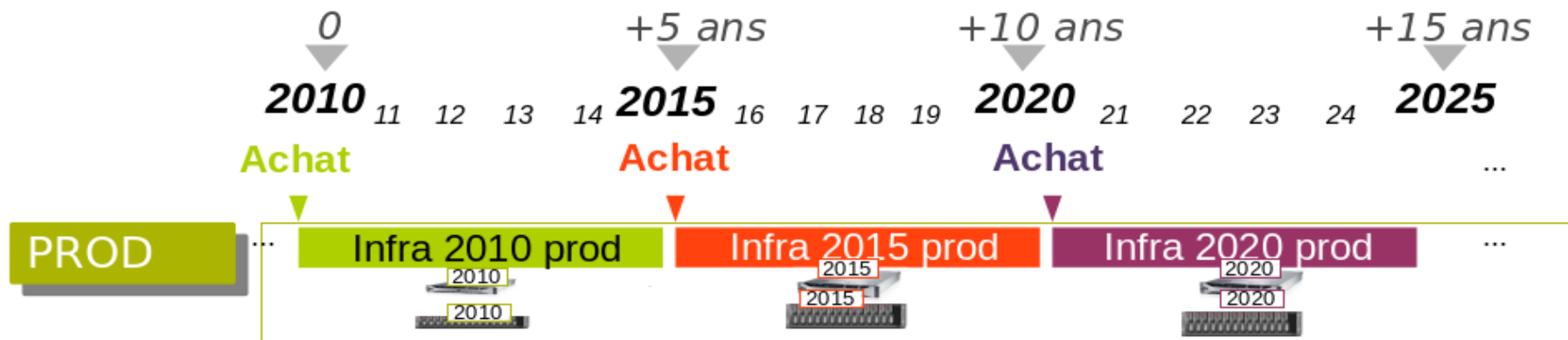
- Observatoire des Sciences de l'Univers (OSU CNRS/INSU, 6 UMR & 1 UMS)
- La donnée scientifique
 - Observations (In situ, satellites, télescopes, etc.) et modèles numériques
 - Volume de ~**2Po** - disques internes et baies de stockage DAS / NAS / SAN
 - **Sécuriser et pérenniser** (stockage secondaire, ~400To)
- **Software Defined Storage** (SDS)
 - Le « silo » de stockage à \pm haute valeur ajoutée devient une infrastructure globale, pilotée par le logiciel
 - Vmware VSAN, Openstack swift, Scality, DataCore, Lustre, GPFS, Ceph, etc.

La valeur ajoutée n'est plus dans le matériel,
mais dans l'infrastructure globale interopérable

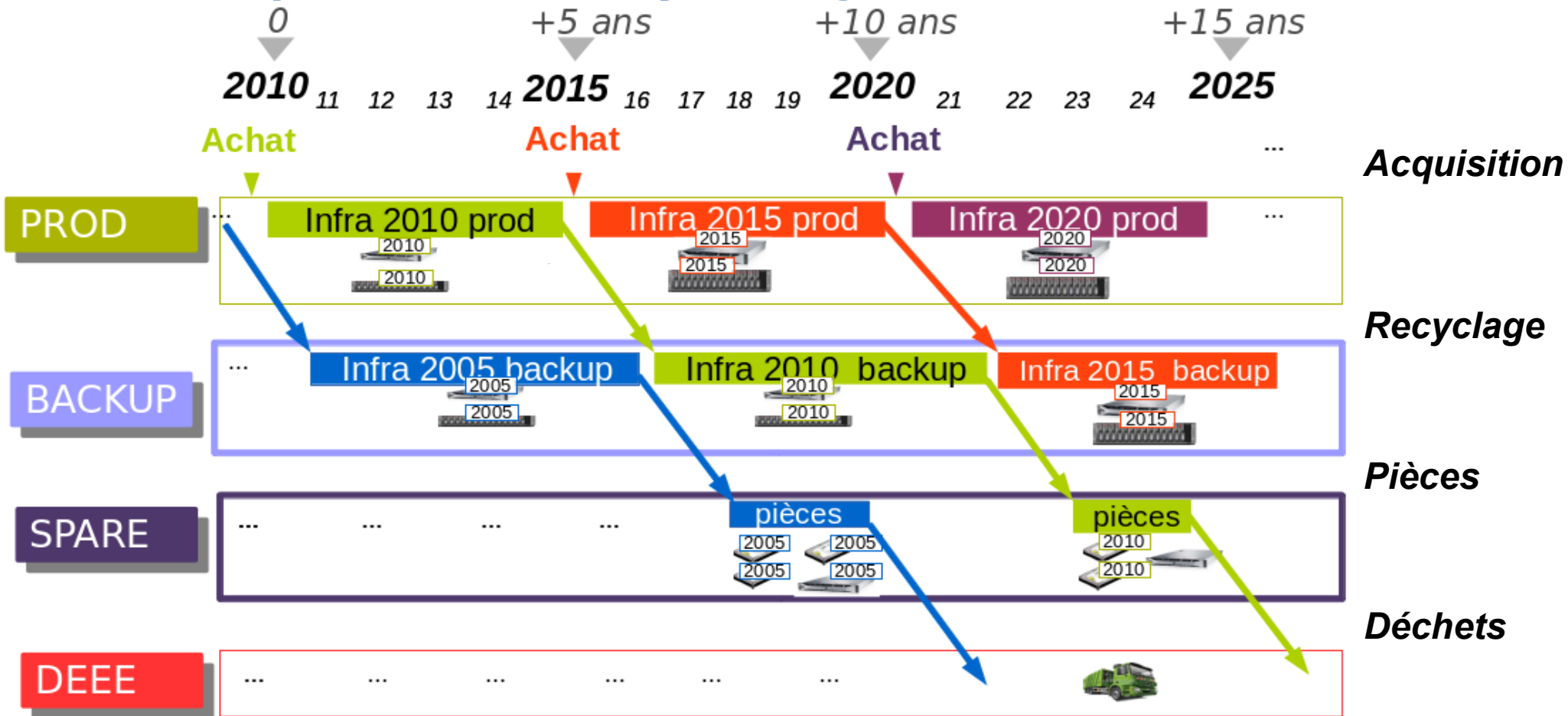
Le projet d'archivage

- › Utiliser du matériel **recyclé** et potentiellement obsolète
 - Serveurs
 - Stockage (disques internes et baies)
- › **Mutualiser** les infrastructures
- › **Répartir** géographiquement les données
 - Trois *datacenters* sur deux sites géographiques éloignés et indépendants
- › **CEPH** pour le « *Software Defined Storage* »
 - Sécurité des données assurée par la résilience du cluster Ceph
 - Indépendance totale vis-à-vis du matériel
 - Volume quasi illimité (stockage objet)

Concept : le recyclage



Concept : le recyclage



Ceph : Proof of Concept

- › Serveurs sous Debian Jessie (actuellement 7 nœuds Dell 2950 → R710)
- › Stockage avec disques (SCSI, SATA, NL-SAS, SAS) internes, baies DAS
- › Répartition géographique sur *trois datacenters*
- › Ceph version stable (*Kraken*)
- › Vlan dédié (flux internes au cluster, administration)
- › Objectifs principaux
 - Valider le déploiement de Ceph sur les matériels
 - Assurer la répartition géographique des données (*crushmap* Ceph)
 - Fournir à chaque unité participante des points d'accès au cluster Ceph, à travers différents types d'interface (LUN, machines passerelles, protocole « cloud », etc.)
 - Éviter la centralisation et faire en sorte que chaque unité puisse être à la fois « consommatrice » et/ou « productrice » de service

Ceph : l'architecture répartie

clusterc4

belinomp

1cnups

d2.belinomp

b2.belinomp

22.1cnups

10.d2.belinomp

12.d2.belinomp

17.b2.belinomp

01.22.1cnups

host5

host1

host6

host7

Mon.2

Mon.0

Mon.1

← Cluster

← Site

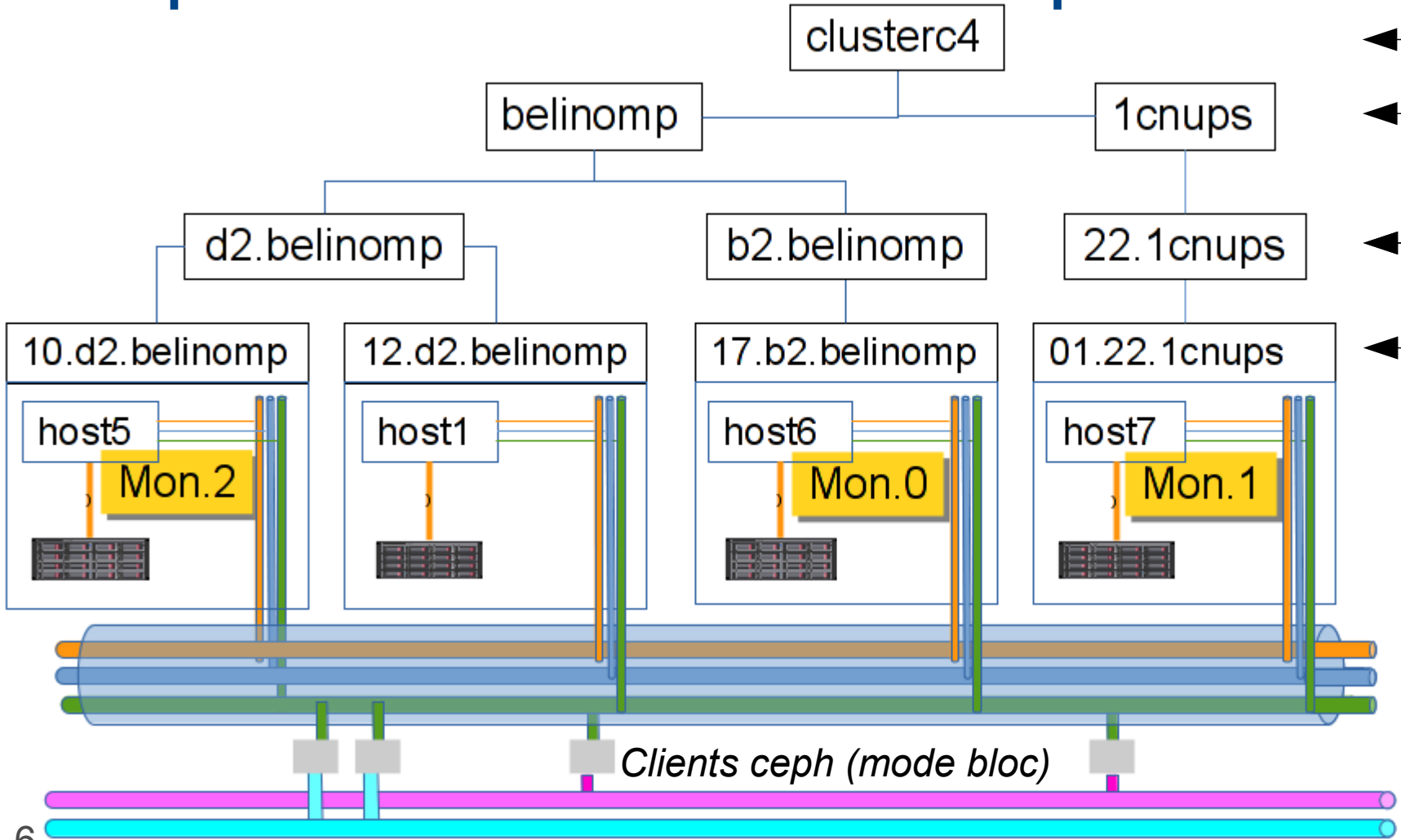
← Datacenter

← Rack

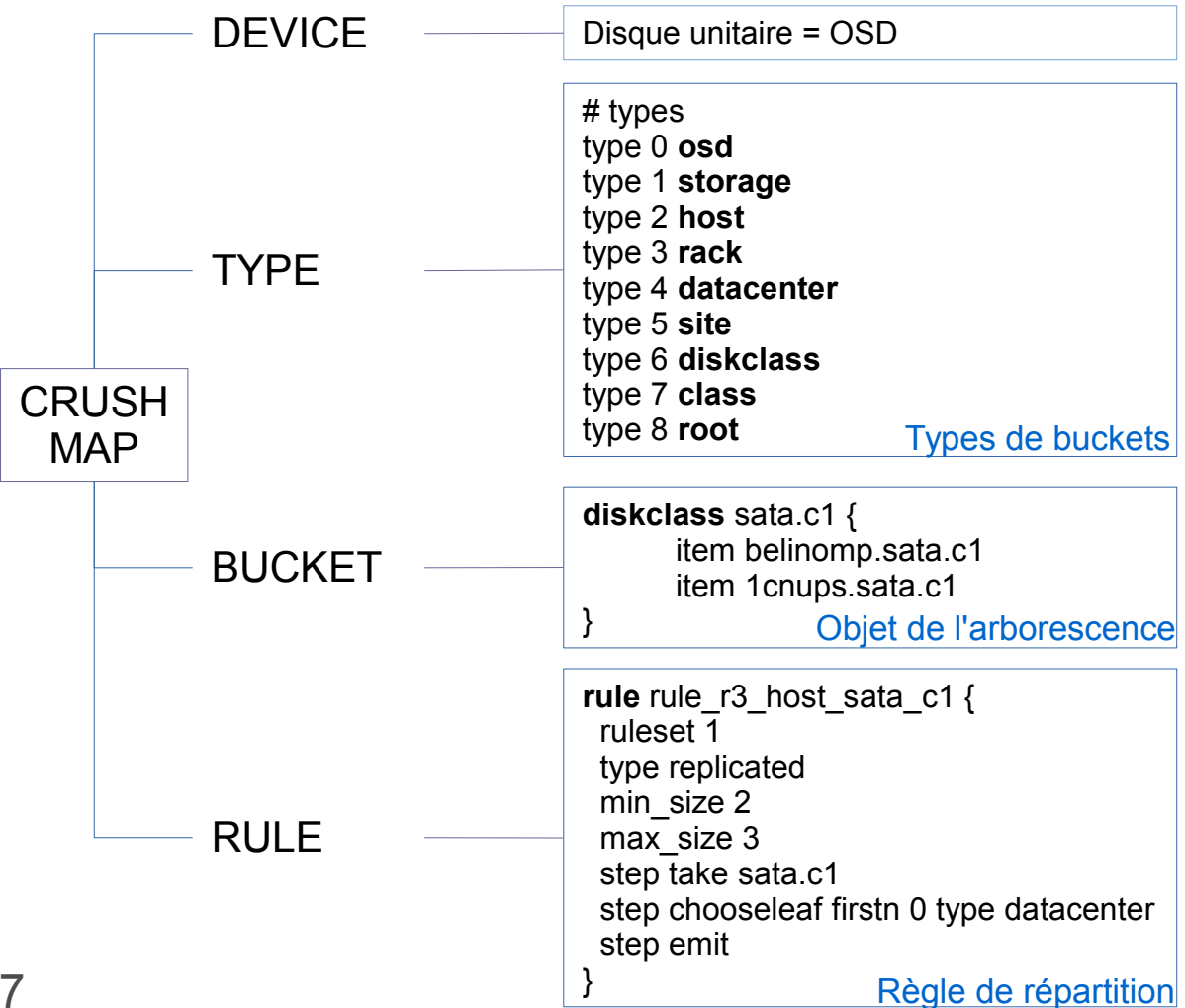
Vlan dédié

Clients ceph (mode bloc)

Réseaux unités

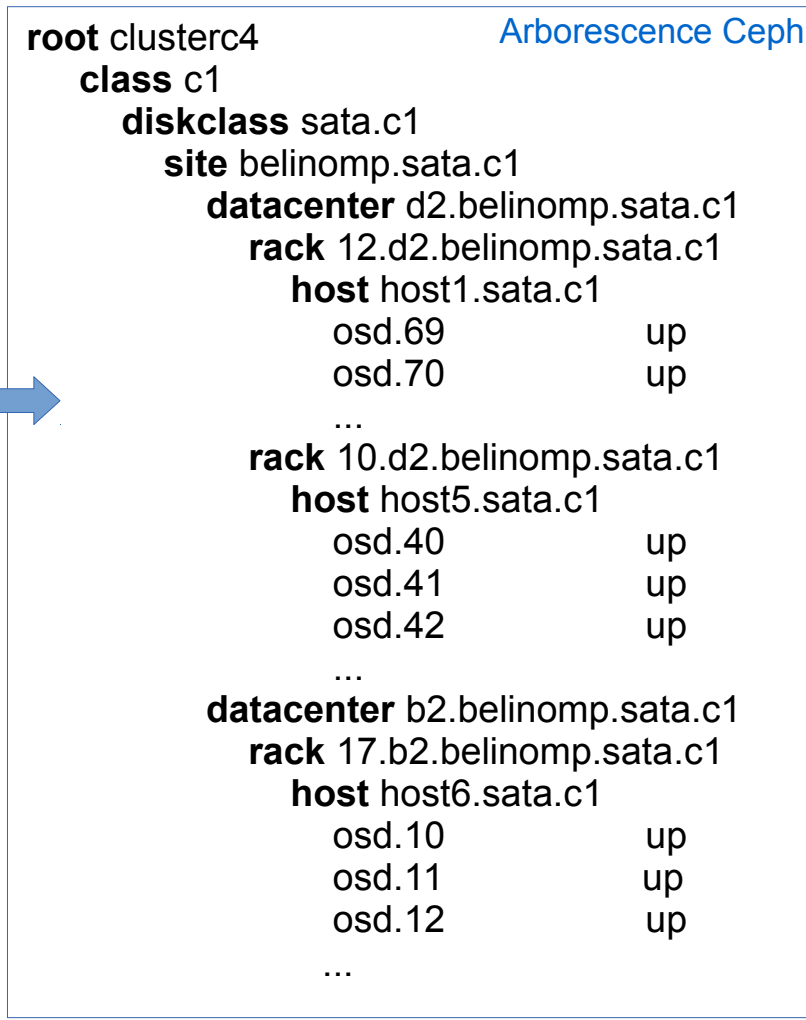
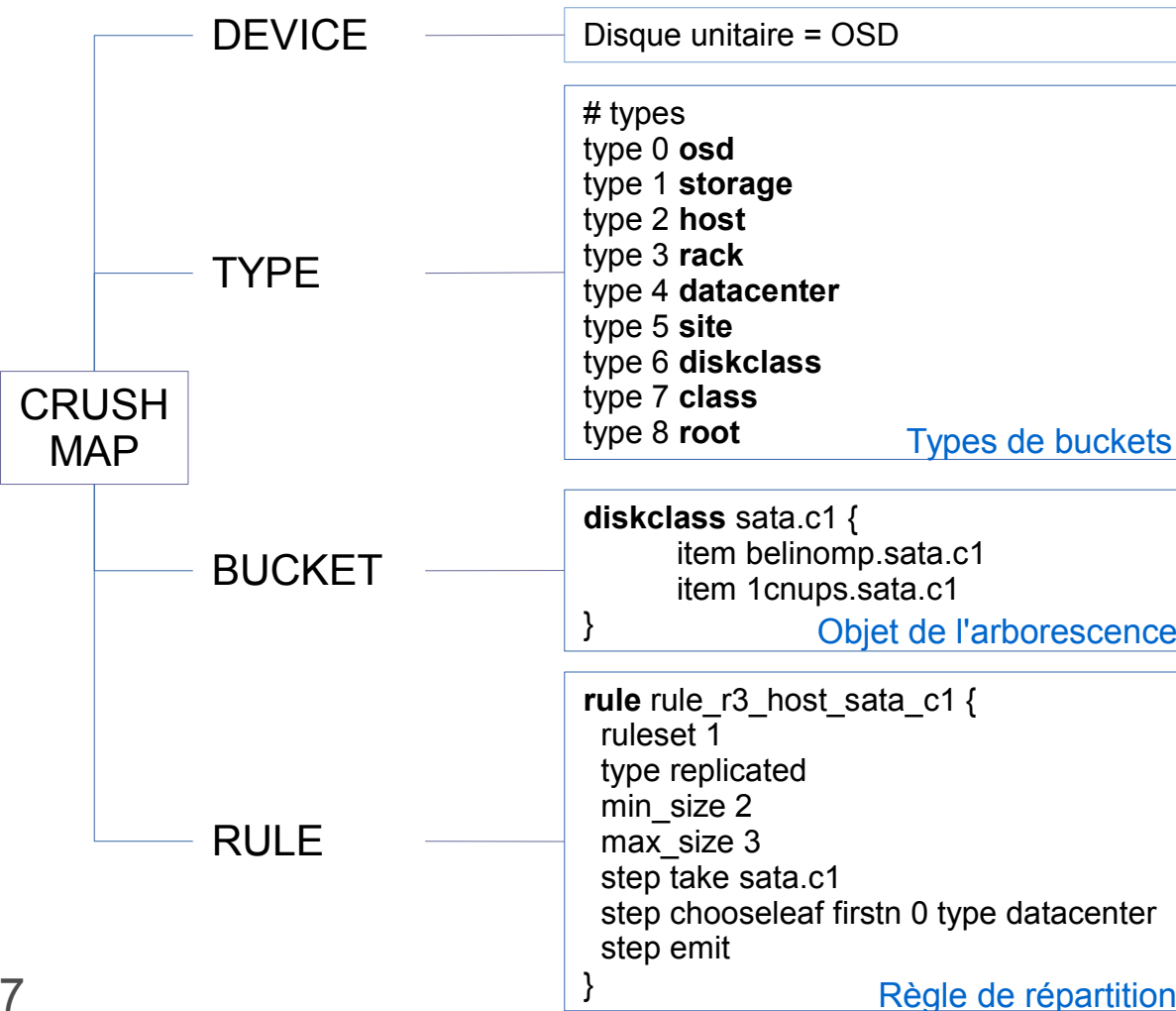


Ceph : la *crushmap* (1/2)



CRUSH :
Controlled Replication Under Scalable Hashing

Ceph : la *crushmap* (2/2)



Ceph : intégration / suppression

Intégration

Installer selon choix géographique

matériel

Installer le système (preseed + scripts)

système

Déployer ceph

ceph-deploy install --release ...

ceph

Déployer les osd

ceph-deploy disk prepare ...

ceph-deploy disk activate ...

osd

« Renommer » le nouvel hôte

***ceph osd crush rename-bucket ***

hostx hostx.sas.c1

crushmap

« Positionner » le nouvel hôte

***ceph osd crush move hostx.sas.c1 ***

rack=12.d2.belinomp.sas.c1

Ceph : intégration / suppression

Intégration

Installer selon choix géographique

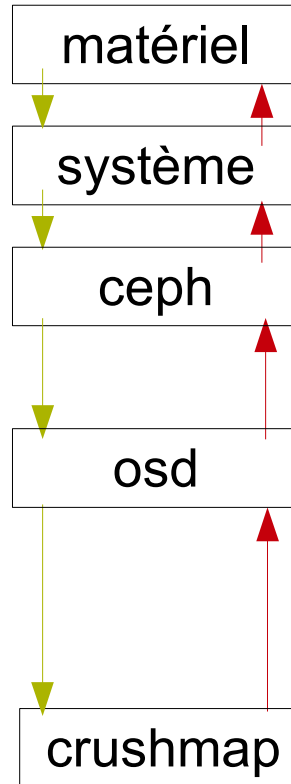
Installer le système (preseed + scripts)

Déployer ceph
ceph-deploy install --release ...

Déployer les osd
ceph-deploy disk prepare ...
ceph-deploy disk activate ...

« Renommer » le nouvel hôte
***ceph osd crush rename-bucket ***
hostx hostx.sas.c1

« Positionner » le nouvel hôte
***ceph osd crush move hostx.sas.c1 ***
rack=12.d2.belinomp.sas.c1



Suppression

Retirer matériel et mise en pièces

Arrêter le système

Supprimer les paquets ceph

Supprimer les osd
ceph osd rm {x ... n}

Arrêter les démons osd
systemctl stop ceph-osd@{x .. n}

Supprimer les osd puis l'hôte
ceph osd crush remove hostx.sas.c1
ceph osd crush remove osd.{x ... n}

Rendre les osd indisponibles
ceph osd out osd.{x ... n}

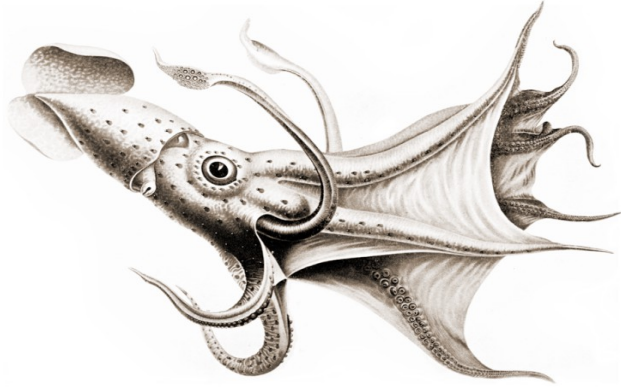
Conclusion et perspectives

- › Le stockage objet sera l'une des solutions de demain
 - L'offre existe aujourd'hui (Scality, Datacore, Ceph, etc.)
- › Ceph s'est avéré être un bon outil
 - Stockage distribué, sécurisé et pérenne
 - *Crushmap* permettant d'appliquer la « bonne » politique de stockage (répartition géographique en particulier)
 - Modèle adapté aux structures distribuées où chaque laboratoire participant peut être « producteur » et/ou « consommateur »
- › Le modèle « matériel recyclé » est :
 - Peu coûteux financièrement
 - Adapté au cas d'usage de l'archivage

L'objectif initial est atteint :
créer un projet fédérateur autour du stockage pour l'archivage

- › Questions ouvertes à plus long terme :
 - Le vieillissement du matériel
 - Le coût humain consacré à l'administration de la solution

Merci de votre attention



Une infrastructure
mutualisée et recyclée
avec Ceph



Ceph : le stockage

Object, Pool, Placement group, OSD

