



# **CEPH**

## **3 ANS APRÈS**

### **(BILAN & PERSPECTIVES)**

**Yann Dupont**

**DSIN Université de Nantes**

[www.univ-nantes.fr](http://www.univ-nantes.fr)

**19 Novembre 2015**



UNIVERSITÉ DE NANTES

# CEPH ?

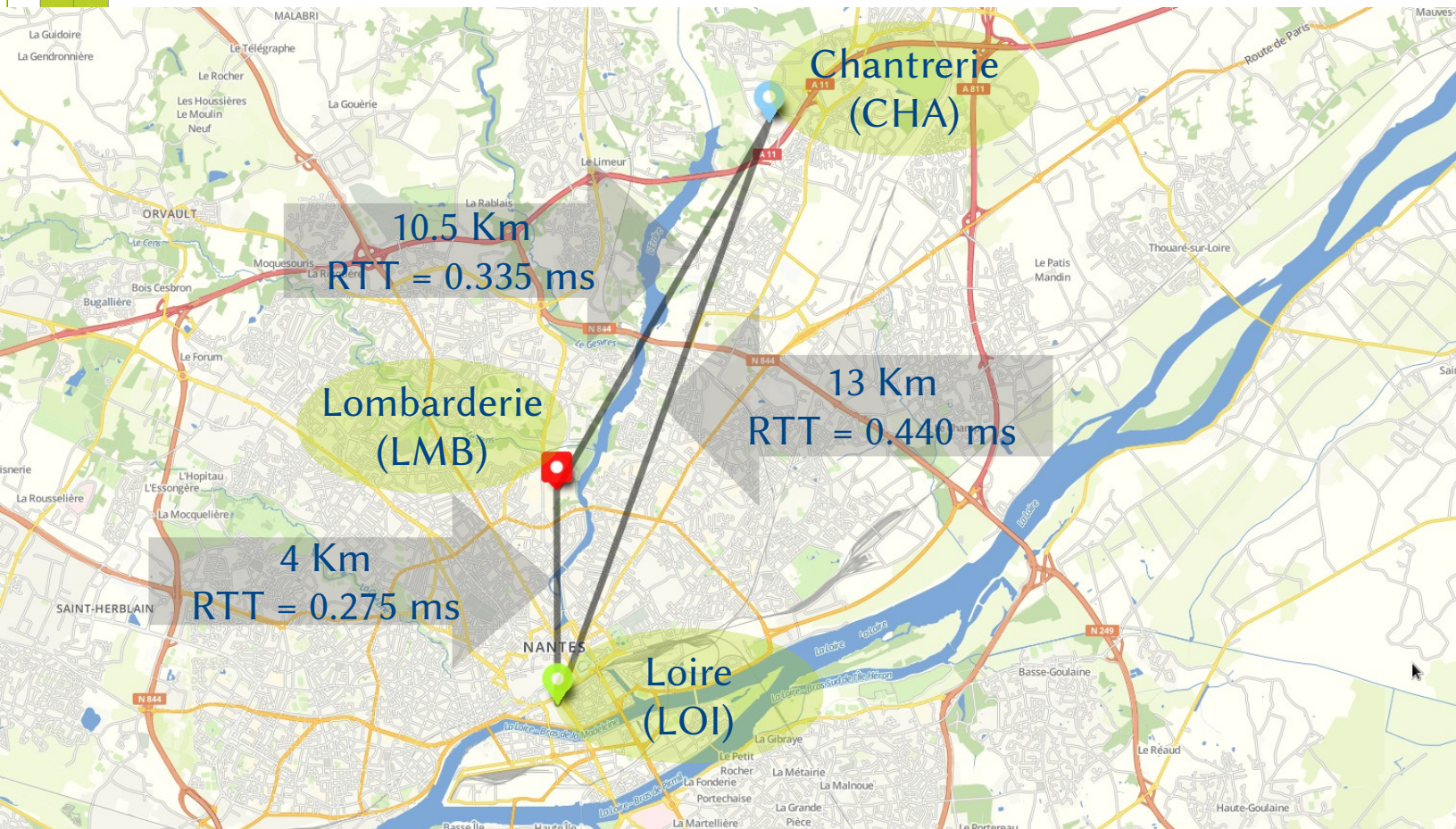


Ceph est UN système complet de stockage distribué  
(**Stockage Objet, Bloc, Système de fichiers**)

Qui fonctionne bien.

Pour une présentation plus complète de Ceph :  
cf Jres 2013, Meetup Openstack Paris, Ceph Days Paris 2014, Bio Ouest 2014...

# TOPOLOGIE RÉSEAU À NANTES



3 points de présence

Fibre noire  
(Allumée  
par la DSIN)

40 Gbit/s  
entre sites

**1 Datacenter**



# MISE EN ŒUVRE OPTIMALE

- 2012 : aucun guide de mise en œuvre...
  - Expériences heureuses puis malheureuses !
  - Erreurs de débutant...
  - Peinture pas encore sèche !
- Partage d'expérience
  - Liste de diffusion
  - Canal IRC
  - Meetups, Ceph Days, Ceph Breizh
  - Jres, OpenStack summit
  - Blogs
- 2015 : Plus simple, assez commun



# HISTORIQUE



UNIVERSITÉ DE NANTES

Tests

Début 2012

47 Préversions (!)

Argonaut (v 0.48)

Bobtail (v 0.56)

Cuttlefish (v 0.61)

Dumpling (v 0.67) LTS

Emperor (v0.72)

Firefly (v0.80) LTS

Giant (v0.87)

Hammer (v0.94) LTS

Infernalis (v9.2.y)

Jewel (v10.2.y) LTS

Pré-prod  
Prod

Multi clusters

Archi dédiée

Augmentation  
espace

Tiering/Erasure  
coding

03 / 07 / 2012

01 / 01 / 2013

07 / 05 / 2013

14 / 08 / 2013

09 / 11 / 2013

07 / 05 / 2014

29 / 10 / 2014

07 / 04 / 2015

06 / 11 / 2015

?? / ?? / 2016

Tous les 6 mois

Tous les 3 mois

Tous les 6 mois

Majeure

renumérotation



**Tester CEPH**

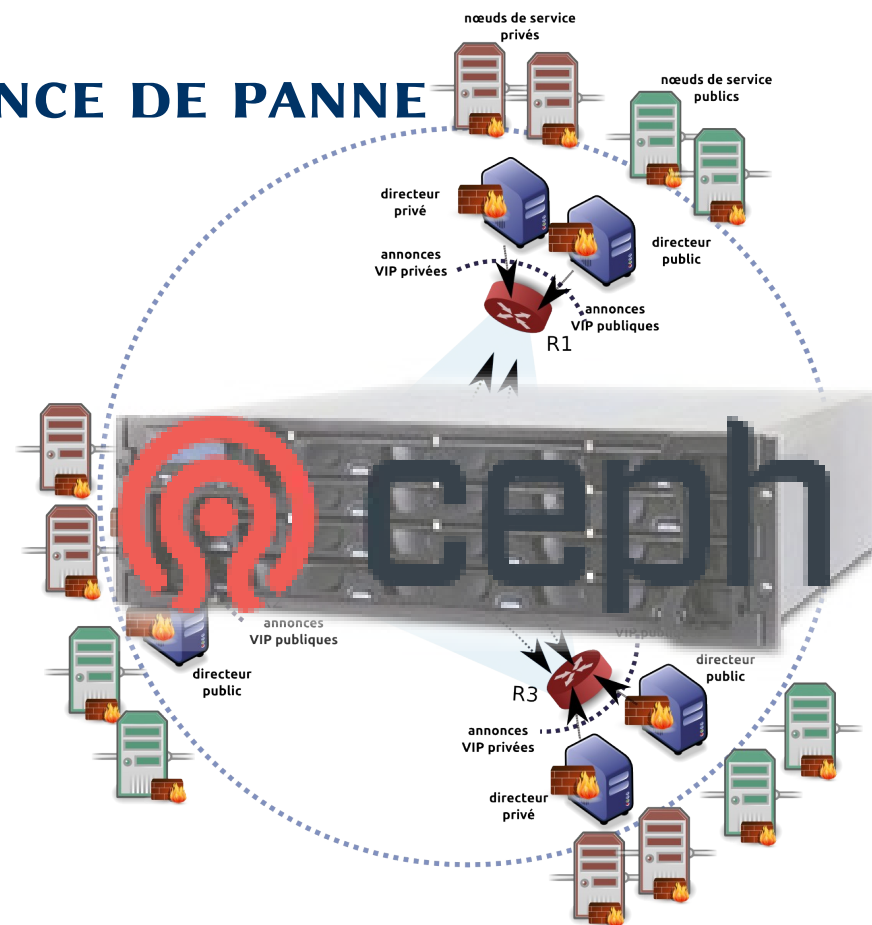
**Valider la plate-forme**

**Démarrer des sauvegardes  
répliquées**

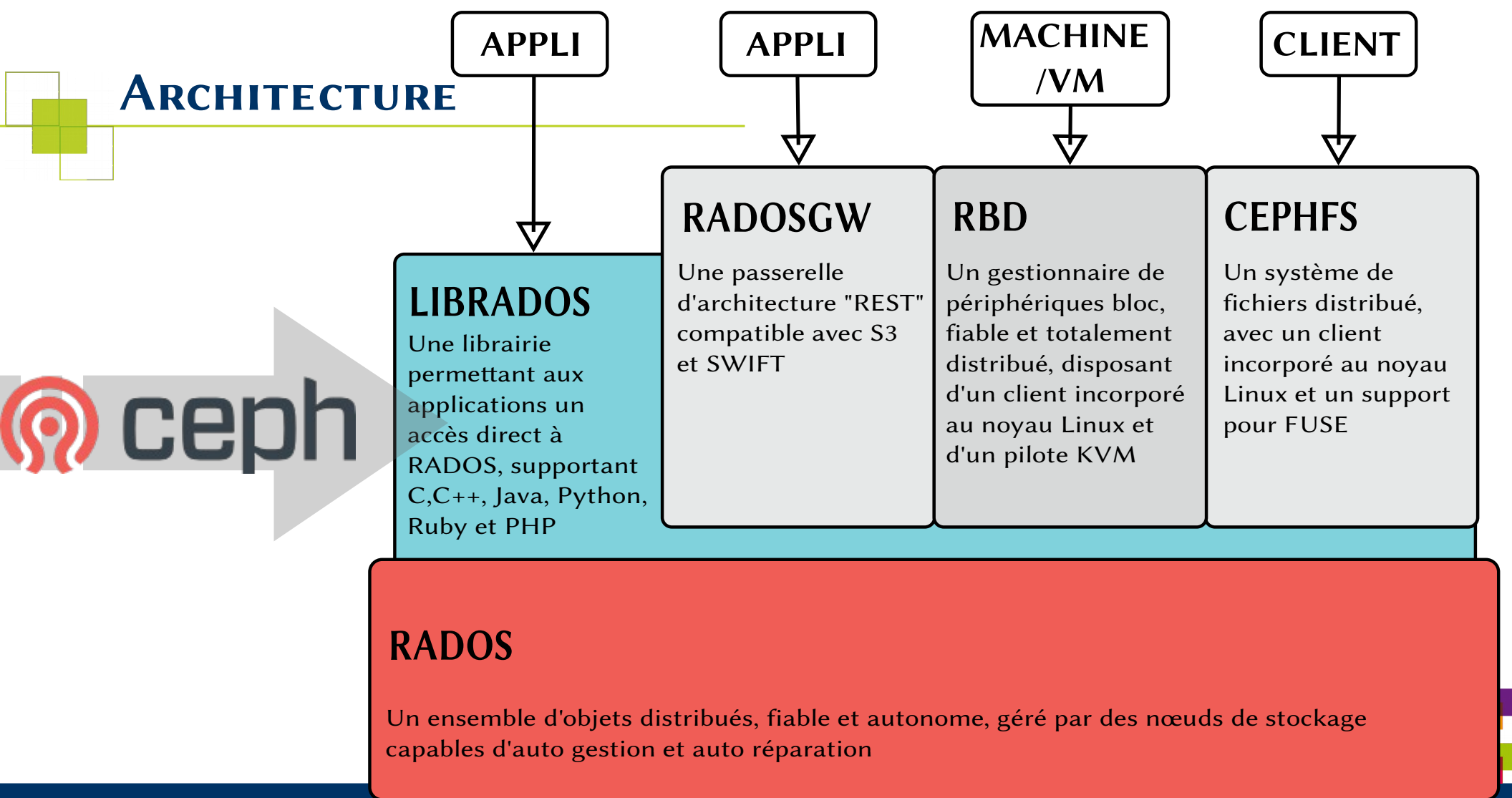


# STOCKAGE DISTRIBUTUÉ À TOLÉRANCE DE PANNE

Après beaucoup d'essais (autres solutions),  
CEPH est validé  
**Début 2012**



# ARCHITECTURE





# ARCHITECTURE DE CEPH

Nombre impair



MON

Nombre impair



MDS

Software  
Defined  
Storage



OSD

Vérifie le bon état du cluster  
Assure la communication initiale  
avec les clients  
Vérifie les droits d'accès

Gère les méta données  
(uniquement pour CephFS)

Stocke les objets  
(généralement des fichiers de  
4 Mo sur filesystem local (XFS))  
Communique avec les clients

# BUGS & DÉCEPTIONS (2012)

## Cluster

Peu d'OSD très volumineux (baies SAN) :

1 unique GROS cluster CEPH

Journaux sur disques : mauvaise idée, Leeeent !

## Kernel

Bug mémoire virtuelle, crash fréquent des OSD

## Filesystems

BTRFS ( Lent, se fige, plante...)

XFS (1 Bug sévère et dévastateur) (Vite corrigé)

Cluster presque plein + bugs + effet domino à la reconstruction = « On casse l'incassable »

**Un désastre peut arriver**  
**Choix architecturaux pour atténuer cela**



# NE PAS METTRE TOUS SES ŒUFS DANS LE MÊME PANIER

**Une volumétrie importante**

**Des besoins différents  
Des administrateurs différents**

**Démarrer plusieurs clusters CEPH  
mais de façon virtuelle**

**(Utilisation de virtualisation,  
conteneurs LXC)**



**Plate-formes dédiées**

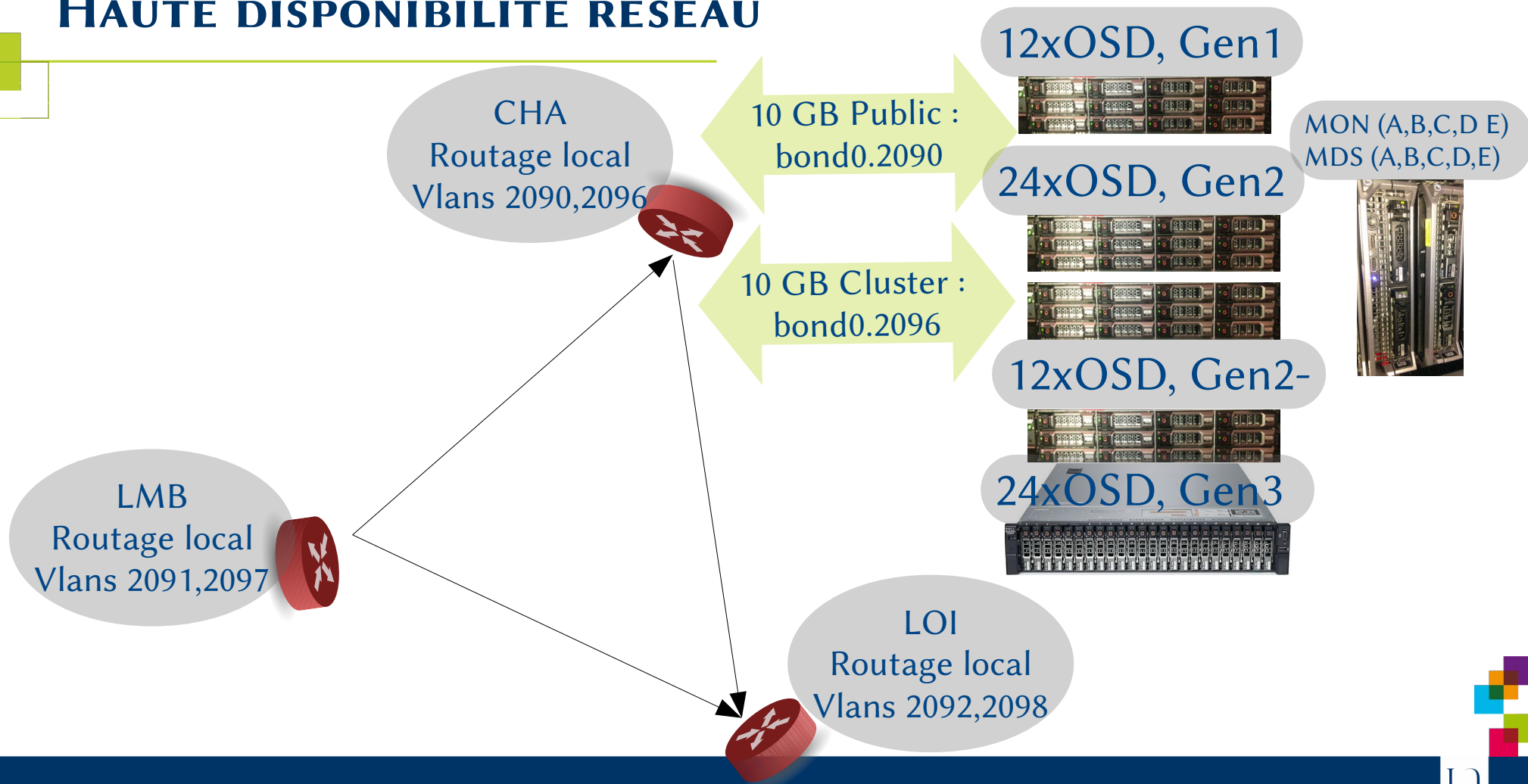
**Plusieurs clusters distincts  
Données placées intelligemment**

**Virtualisation des clusters**

**Démarrage production**



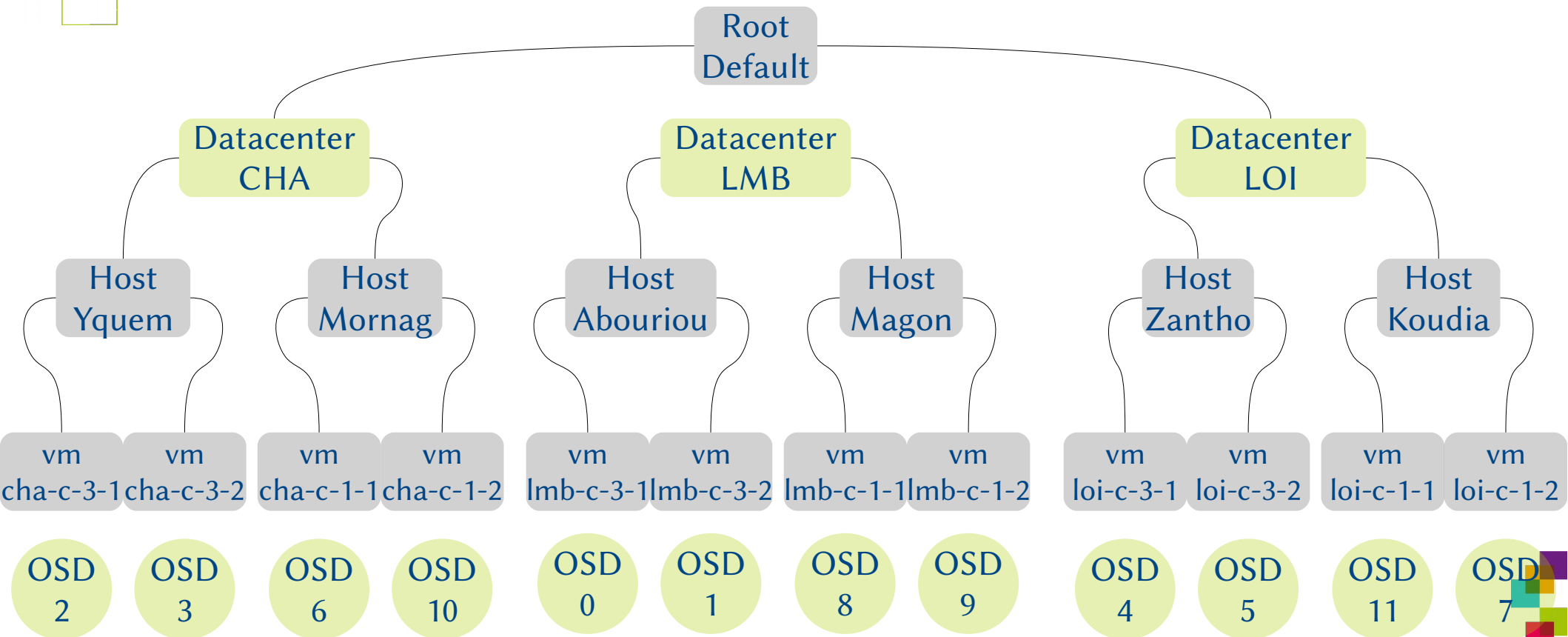
# HAUTE DISPONIBILITÉ RÉSEAU





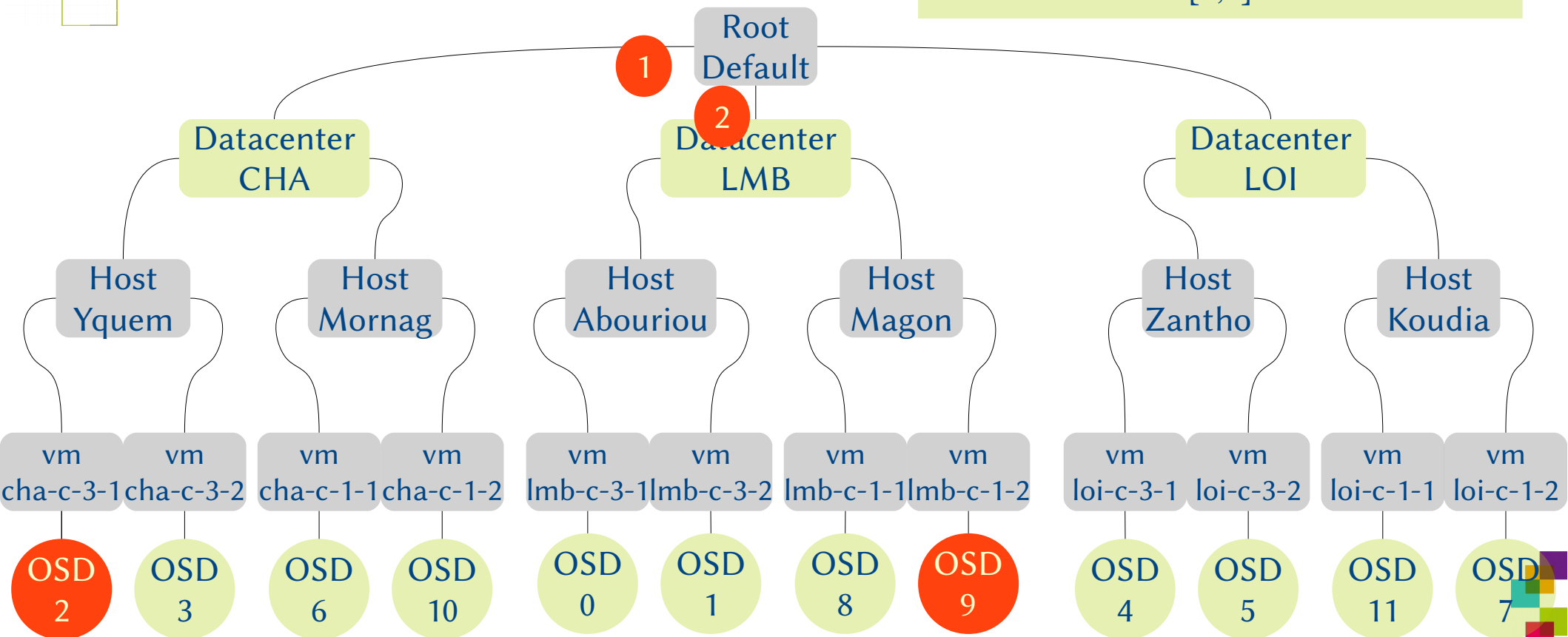
# RÈGLES CRUSH, CLUSTER C

step take default  
step chooseleaf firstn 0 type **datacenter**  
step emit



# RÈGLES CRUSH, CLUSTER C

Pool 3 : mirrors (**2 copies**)  
PG dump 3.1ed  
[2,9]



# DÉPLOIEMENTS (FIN 2013)



1 OSD = 1 Disk = 1 LXC  
plus de RAID Hardware  
2 SSD SLC partagés (Journaux)  
12 OSD / machines  
Capacité Brute : 48 To

Déploiements symétriques  
par plaque.

3 machines identiques.

# STRATÉGIE 2014

**Augmenter la volumétrie,  
Replication x3**

**Profils d'usages,  
selon cluster, génération de machines**

**Démarrage des tests IaaS**

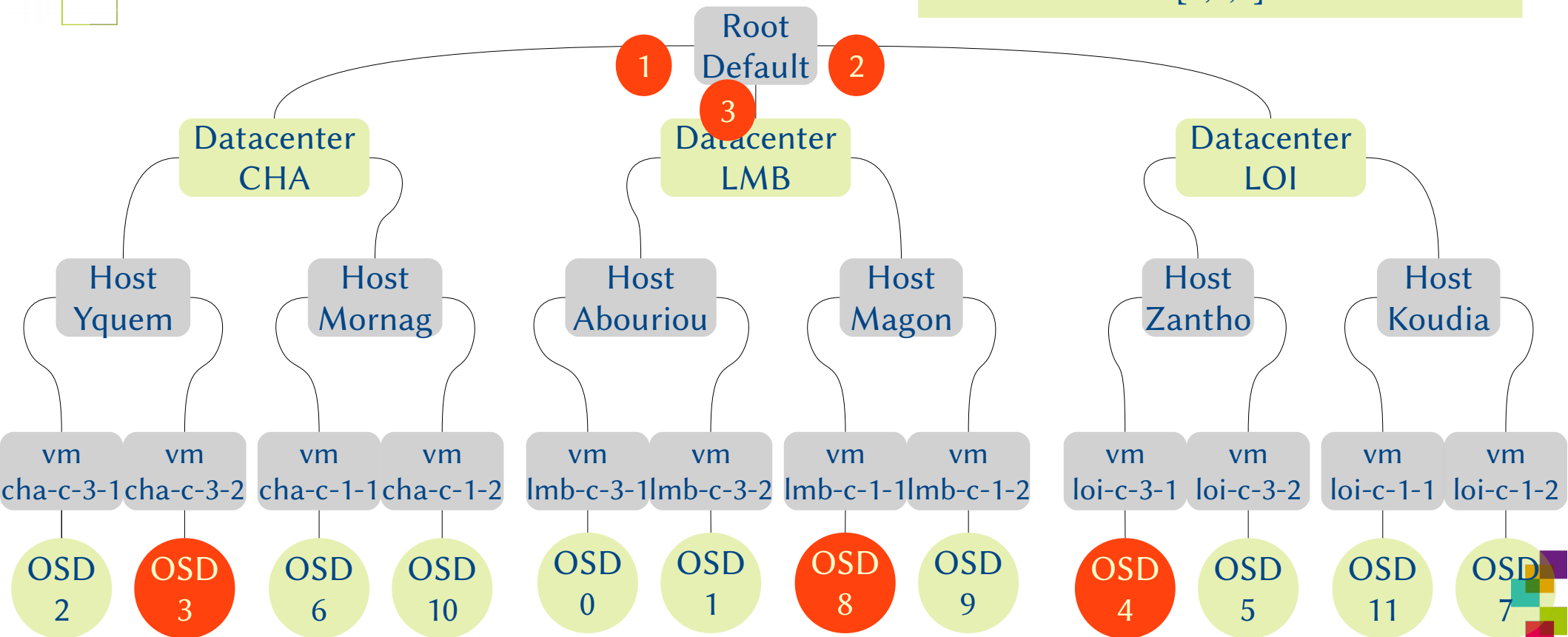
**Commencer à se défaire  
des vieilles baies SAN hors garanties**

**Simplifier les déploiements  
(plus de 200 LXC)**



# RÈGLES CRUSH, CLUSTER C

Pool 6 : OS templates (**3 copies**)  
PG dump 6.3fb  
[3,4,8]

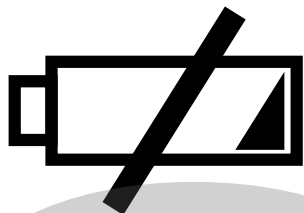




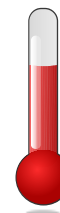
# TOLÉRANCE À LA PANNE : TESTÉ !



(Fibres!)  
CHA



Panne d'onduleur  
et disjoncteur  
LOI



Panne climatiseur  
& arrêt d'urgence  
LMB

**Les problèmes sur machine physique sont rares**  
En général, une salle machine complète est impactée  
3 répliquas permettent d'éviter une reconstruction

# CLUSTERS CEPH ET DESTINATION

Nom	Usage	Depuis	Taille	Version
A	Tiers (labos de recherche)	07/2015	Pourrait être très importante	Stable LTS (Hammer)
B	Incubation	03/2013	Petit	Stable Hammer→Infernalis
C	Production vSAN	05/2014	Moyenne	Stable LTS (Firefly)
D D2	Backups	01/2013 11/2015	Très importante	Stable LTS : Firefly→ D2 : Hammer
E	Expérimental	06/2014	Petit	Versions de développement
F	IaaS	Q4/2015	Moyen à important	Stable LTS (Firefly)
G	Data / Home dirs	Q2/2016	Pourrait être très importante	Stable LTS (Jewel?)

# STRATÉGIE 2015

Augmenter la volumétrie...

Racine de serveurs KVM en prod,  
Serveurs logs, miroirs, ucloud

Se défaire de toutes les baies SAN  
Volumétriques

Augmenter performance  
Ajout de caches (local datacenter)  
baisser les latences  
Augmenter les débits (4x10 Gb)



# CAPACITÉ GLOBALE

Gén	Destination	Nb	Taille	Version	Volume
1	Mixé	3	12 x 3 To 2 SSD MLC	R720xd Juin 2013	108
2	Mixé/Perf	6	12 x 4 To 2 SSD SLC	R720xd, Nov 2013 Juin 2014	288
2--	Stockage de masse, BUDGET	3	12x2 (utilisés), 0SSD	R720xd Juin 2014	72
2-		3	12x4 (utilisés), 0SSD		144
3	Perf lecture	3	24x1 (2,5" sas 10k) pas de SSD.	R720xd Nov 2014	72
4	Stockage de masse	9	16 x 8 To Pas de SSD.	R730xd Nov 2015	1152 (!)
4c	PERF !! Cache, écriture, latence.	2	10x400 Gb SSD write intensive	R630 Nov 2015	0

# ÉVOLUTION (Mi 2015)



LXC D

LXC A

1 OSD = 1 Disk

1 seul LXC par cluster

Pas de SSD. Journal en tête des disques.

24 OSD / machines

Capacité Brute : ~ 24To

«Seulement» 5 LXC max/machine.  
Beaucoup d'axes. Orientation  
performances en lecture  
(Racines IaaS)



# CLUSTER C

GLOBAL:

SIZE	AVAIL	RAW USED	%RAW USED
43483G	11839G	31644G	72.77

14 → 43 To : Replicat à 3

POOLS:

NAME	ID	USED	%USED	MAX	AVAIL	OBJECTS
data	0	0	0		5055G	0
metadata	1	9470	0		5055G	21
rbd	2	0	0		5055G	2
mirrors	3	3038G	6.99		5055G	785006
cloud-perso	4	817G	1.88		3370G	232739
openstack	5	301G	0.69		5055G	79617
os-patrons	6	169G	0.39		3370G	45347
os-dsin-prv	7	2347G	5.40		3370G	605041
os-dsin-pub	8	0	0		3370G	0
logs	9	2460G	5.66		3370G	637182
os-dsin-sig	10	3368M	0		3370G	1189
data-dsin	11	2200G	5.06		3370G	564697
data-tiers	12	180G	0.41		3370G	46167
opennebula	13	26555M	0.06		3370G	6734
esxi-backup	14	8973M	0.02		3370G	2913

Mirroirs Ubuntu, Debian...

Racines de machines virtuelles  
Dédup (COW)

Vers cluster F

# SNAPSHOTS , LAYERING

```
ceph-mon-lmb-C-1:~#  
rbd info os-dsin-prv/lubuntu-xs2  
rbd image 'lubuntu-xs2':  
size 20000 MB in 5000 objects  
order 22 (4096 kB objects)  
block_name_prefix: rbd_data.3e1292eb141f2  
format: 2  
features: layering  
  
parent: os-patrons/  
lubuntu-2014-10-beta1@20140922  
  
overlap: 20000 MB
```

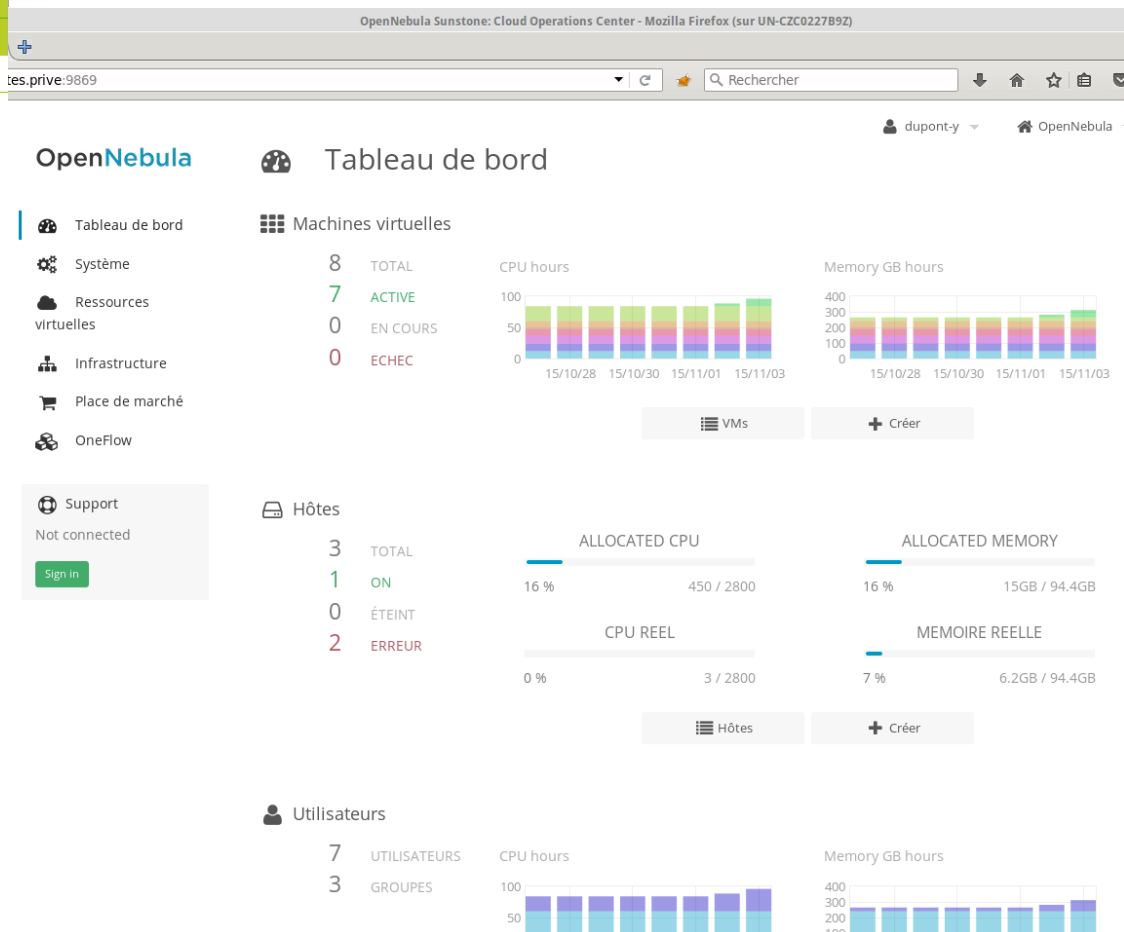
```
ceph-mon-lmb-C-1:~#  
rbd snap ls os-dsin-prv/lubuntu-xs2
```

SNAPID	NAME	SIZE
9	apres-upgrade-wily	20000 MB
6	avant-upgrade-vivid	20000 MB
7	avant-upgrade-vivid-b	20000 MB
8	avant-upgrade-willy	20000 MB

# USAGE DIRECT AVEC KVM (LIBRBD)

```
<disk type='network' device='disk'>
  <driver name='qemu' type='raw' cache='writeback' iothread='8'/>
  <auth username='dsin-prv'>
    <secret type='ceph' uuid='d0a44325-05cb-4bc2-9d95-53af647fe78e'/>
  </auth>
  <source protocol='rbd' name='os-dsin-prv/lubuntu-xs2'>
    <host name='172.20.106.85' port='6789'/>
    <host name='172.20.107.85' port='6789'/>
    <host name='172.20.108.85' port='6789'/>
  </source>
  <target dev='vda' bus='virtio'/>
  <address type='pci' domain='0x0000' bus='0x00' slot='0x05' function='0x0'/>
</disk>
```

# OPENNEBULA (IAAS)



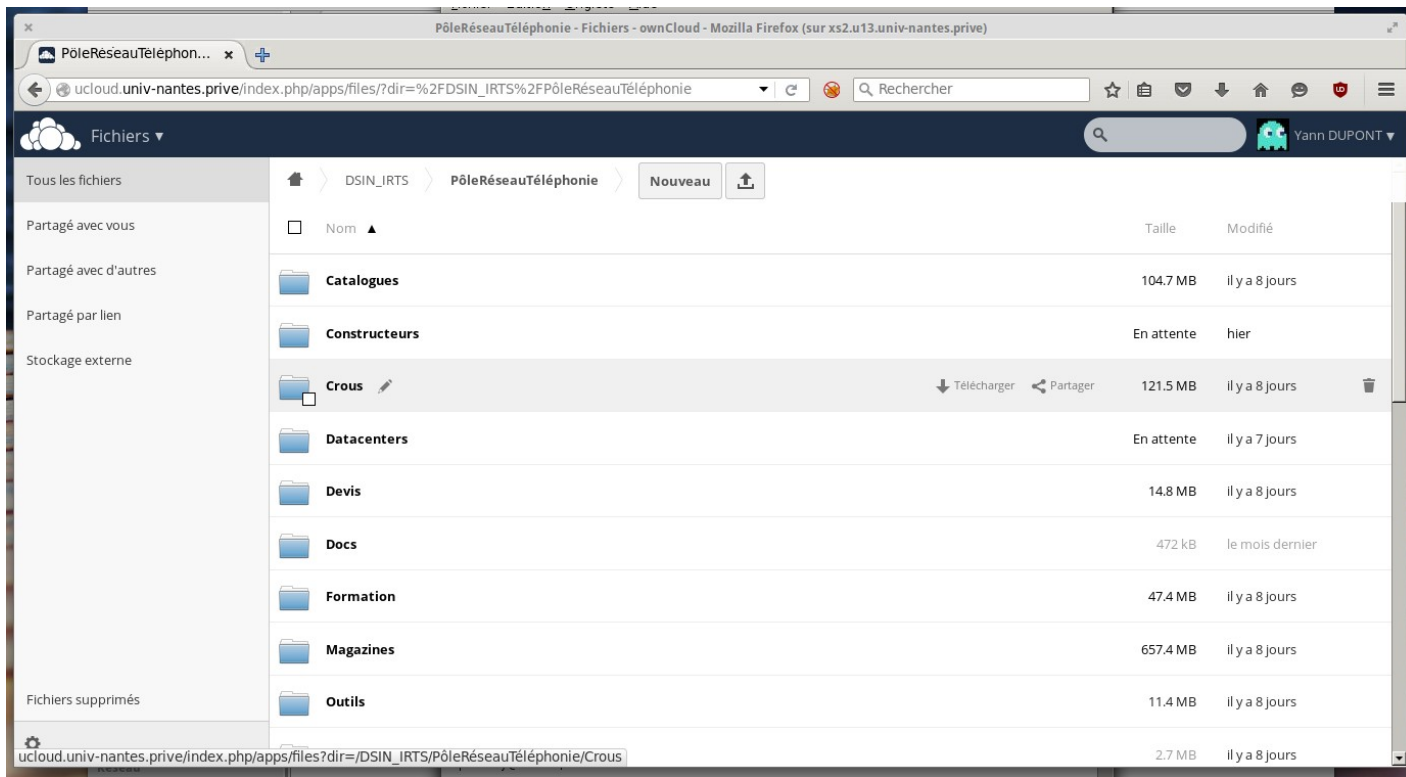
OpenNebula dans un  
Premier temps

Openstack éventuellement  
Comme seconde offre.

Les deux reposent sur  
Les mêmes bases

(Libvirt + KVM + Ceph)

# OWNCLOUD (SYNC & SHARE)



Actuellement, utilisation  
En mode bloc.

À terme,  
multiple frontaux +  
Stockage objet



# CLUSTER D

## GLOBAL :

SIZE	AVAIL	RAW USED	%RAW USED
326T	115T	211T	64.72

326 To, 110 To utiles !!  
(~75 To utilisés)

## POOLS :

NAME	ID	USED	%USED	MAX AVAIL	OBJECTS
data	0	0	0	38561G	0
metadata	1	9470	0	38561G	21
rbd	2	0	0	38561G	0
bacula-D	3	26623G	7.95	25707G	6815746
NFS	4	7480G	2.24	38561G	1961891
bacula-SIG	5	16259G	4.86	25707G	4162837
bacula-NB	6	23495G	7.02	25707G	6027262
MIGRATIONS	7	16	0	25707G	3
backuppc	8	897G	0.27	25707G	229890

Rep = 2

Rep = 3

11/2015 : Cluster D2 = Migration vers Hammer (LTS) + Erasure coding (2+1)  
Choix de ne pas faire de migration à chaud. On ne joue pas avec les données.

## Cluster health at 10:16:14



OK

## history

2015-11-18 10:14:33 : OK

## Monitors status

b

a

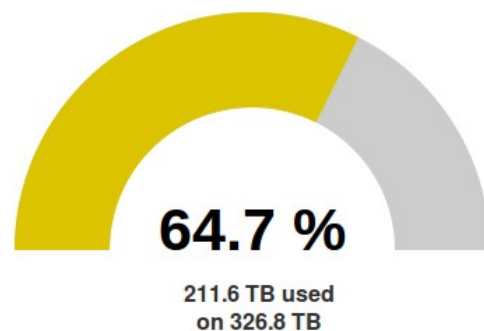
d

## 10432 Placement groups

- active+clean+scrubbing+deep
- active+clean
- active+clean+scrubbing



Avail. capacity :  
115.3 TB



## 117 OSD

	UP	DOWN
IN	117	0
OUT	0	0
Full	near	full

## pools

clean	unclean
9	0

## MDSs

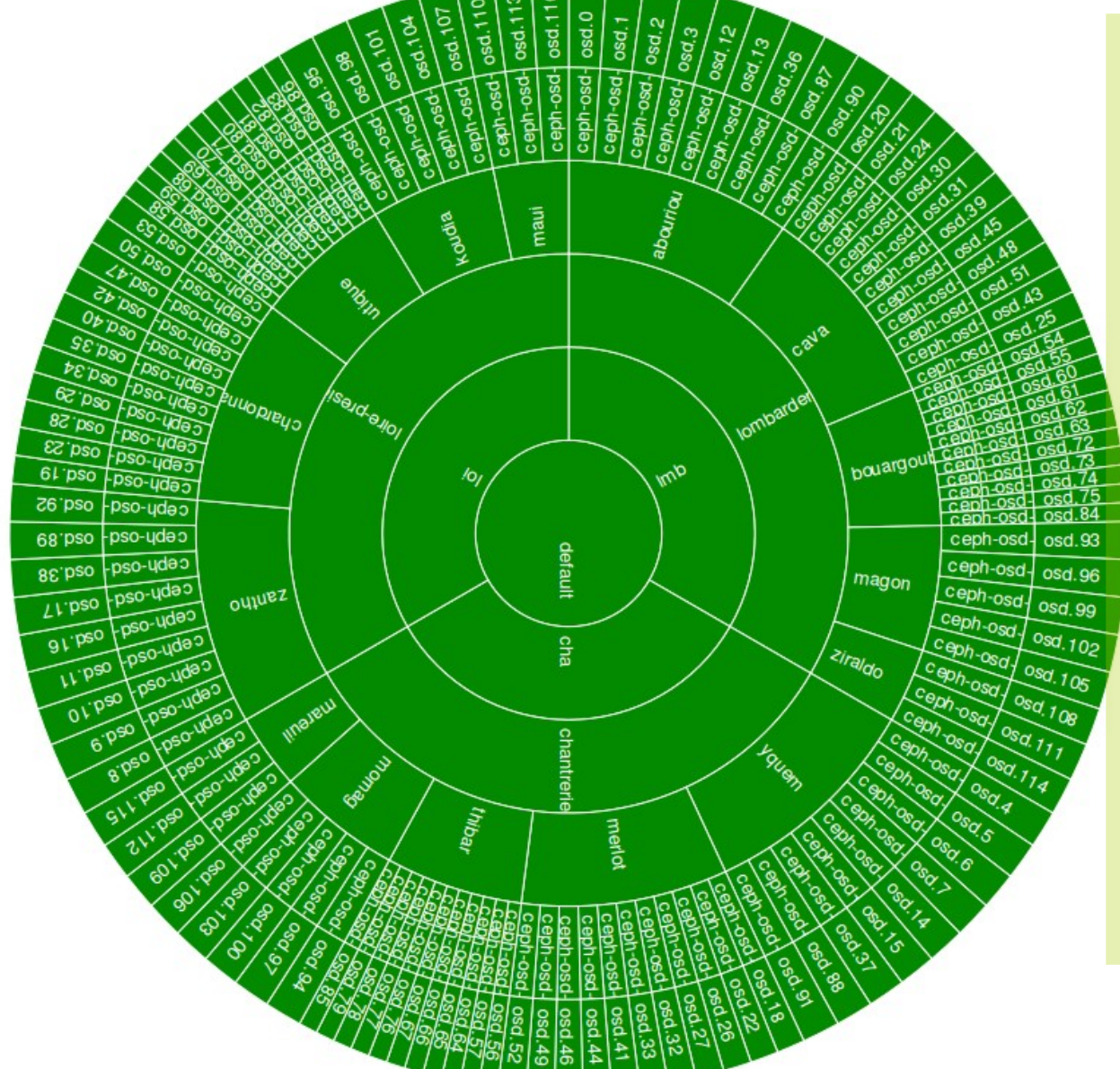
up:in	1	1
up:standby		
max		1

Cluster B  
(vue inkscope)

Répliqua x3

~75 To  
utiles

9 pools



## Cluster B (vue inkscope)

117 OSD  
répartis hiérarchiquement

3 lieux distincts  
1 salle par lieu  
5 serveurs physiques par salle

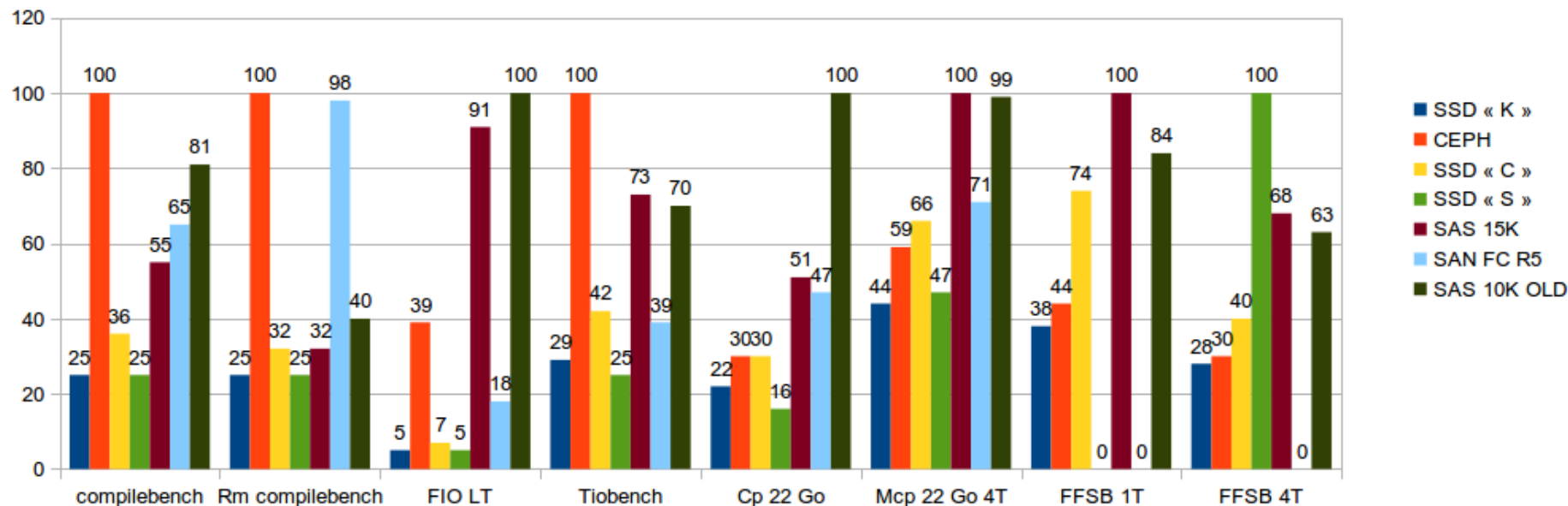
2 à 9 OSD par serveurs volumétriques  
(performance non cruciale)

## ACCÈS POSSIBLES

Type	Bloc	Fichier	Objet
Linux	Krbd (kernel) Kvm+librbd (userspace/kernel)	Cephfs (kernel)	API Swift API S3 API rados
Autres	Rbd-fuse (*BSD ?) Kvm+librbd (*BSD ?) Portages Windows (? ? ?).	Cephfs-fuse	
Via un Tiers	ISCSI : tgt (module ceph, multipath)	NFS : Linux Knfs Ganesha pNFS (HA)  Samba : Linux + krbd/librbd	

# PERFORMANCES (DUMPLING)

Rapidité relative  
(100 % = maximum du temps)



Attention aux benchmarks !!! (Tuning, multithread, RAM, etc...)  
À refaire avec Erasure coding, cache/Tiering SSD et Hammer/Infernalis.

# STRATÉGIE 2016, PERSPECTIVES

Retour sur investissement :

Erasure coding → espace utile x2

Tiering → Utiliser au mieux les machines

Nouveaux usages → performance, volumétrie

**Support direct de l'laaS, cloud → SSD**

Nouveau cluster : Logs & traitements, sécu

Rationaliser/Automatiser les déploiements  
(Ansible, Saltstack, Puppet, Foreman)

Maintenance clusters, évolutions en douceur.  
Clonage du cluster ou évolution à chaud





# MEILLEURES PRATIQUES

- Pas de RAID !
- XFS comme FS de base
- Un Linux de confiance
- Tous les SSD ne se valent pas
- Plusieurs salles machines
- Bien tailler les machines
- Plusieurs clusters : production/sauvegarde/tests
- Les machines doivent être considérées comme des boîtes noires
- Le réseau est déterminant
- Suivre le développement, se documenter, se faire épauler





## CONCLUSIONS

Très stable : plus de 3 années de production et marche remarquablement bien.  
Des possibilités d'évolution infinies, simples à mettre en place. Le Po est facilement atteignable.

**Mais c'est encore un travail d'expert, en particulier au niveau Debug.**

**Beaucoup de nouveautés excitantes** arrivent avec Infernalis :  
(stabilisation de cephfs, amélioration des performances, nouvelles fonctionnalités....)



**Prochaine étape : Jewel**



*MERCI !*

---



# Questions ?

Crédits : Opencliparts / Openstreetmap / Ceph.com