

# INFRASTRUCTURES AVANCÉES STOCKAGE DISTRIBUÉ AVEC CEPH

## Objectifs

Comprendre la problématique moderne du stockage et la nécessité de Ceph

Comprendre son architecture générale

Comprendre en détails le rôle de CRUSH, des Pools, des PG

Savoir distinguer les 3 services de stockages possibles

Connaître les bonnes pratiques dans l'utilisation de Ceph, et son administration quotidienne

POURQUOI  
CEPH?

# POURQUOI ON PARLE DE CEPH ?

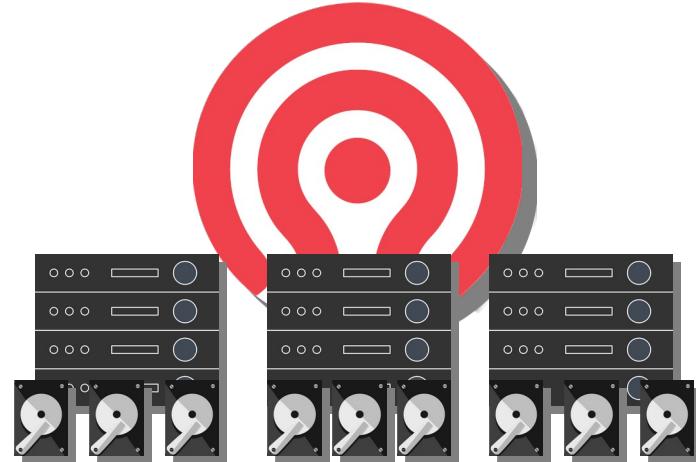
**Ceph** résout un problème très simple : **le stockage doit continuer de tourner même quand un disque, un nœud, ou un switch décède !**

- Et idéalement...
  - **Scalier en ajoutant simplement des machines,**
  - **Sans un point central qui bloque tout** (SPOF!)
  - **En pouvant servir du bloc, du fichier, ou de l'objet** selon le besoin !
  - **En rééquilibrant tout seul quand on ajoute ou retire du stockage !**

# POURQUOI ON PARLE DE CEPH ?

**Ceph** fait tout ça **et le fait BIEN !**

- Il est utilisé **partout en production**
  - OpenStack
  - Kubernetes
  - HPC
  - SI on-prem



VYONS SON  
ARCHITECTURE  
GÉNÉRALE

# CEPH : ARCHITECTURE GÉNÉRALE

Ceph **repose sur une couche logicielle appelée RADOS**

- **3 types de démons principaux intégrés à RADOS**
  - **MON** (ceph-mon) : **gèrent la carte du Cluster !** *Qui est en vie, ou sont les données, quelles règles de placement utiliser ?*
  - **OSD** (ceph-osd) : **Stockent les données**, *en se voyant attribuer du disque, du CPU, de la mémoire de travail...*
  - **MGR** (ceph-mgr) : **fournissent surveillance, métriques, orchestrations complémentaires**



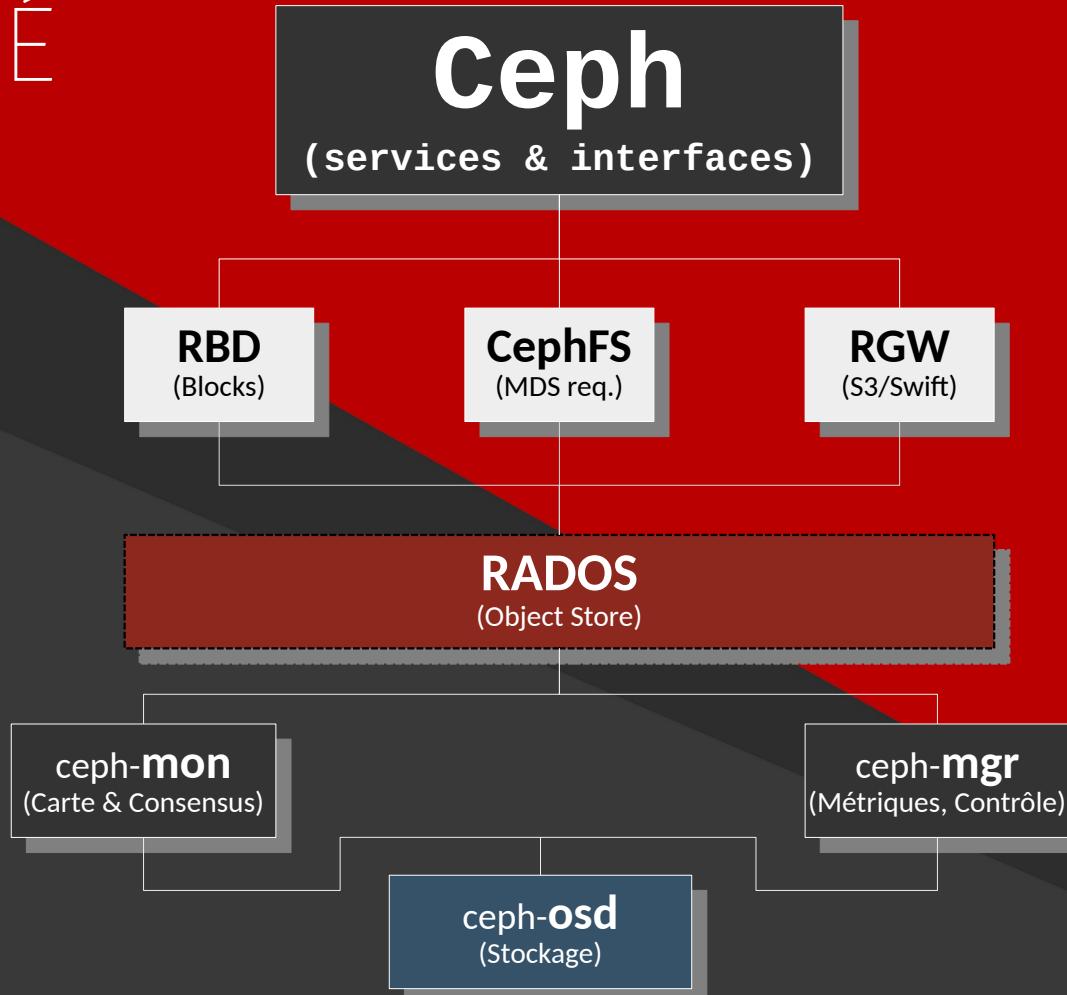
# CEPH : ARCHITECTURE GÉNÉRALE

Il existe des **démons par dessus RADOS**, mais qui n'en font pas partie intégrante

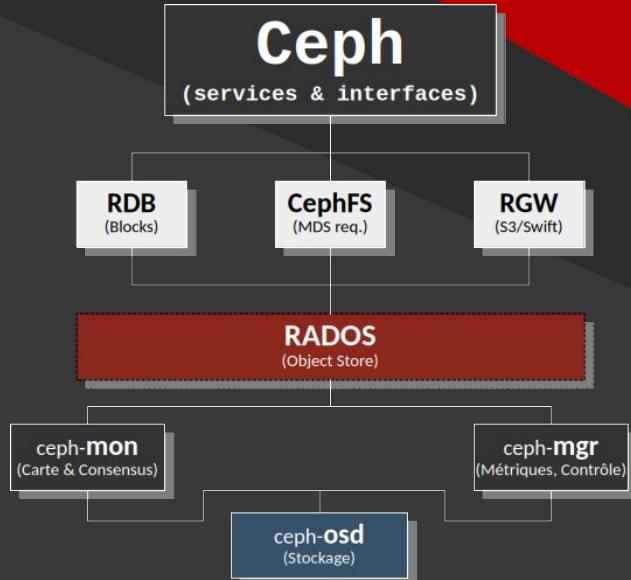
- **MDS** (ceph-mds) : Uniquement **si on stocke en mode système de fichiers**, c'est à dire avec **CephFS** !
- **RGW** (radosgw) : **Passerelle pour stockage objet Amazon S3**
- **RBD** (radosblocks) : **Stockage en mode blocks**



## EN RÉSUMÉ



# EN RÉSUMÉ



**Les clients ne passent pas par un contrôleur central. Ils contactent les MON uniquement pour récupérer des infos, puis ils parlent directement aux OSD.**

**Résultat : pas de SPOF, pas de goulot d'étranglement.**

ET LES DONNÉES SE  
TROUVENT OU,  
CONCRÈTEMENT?

# CEPH : L'EMPLACEMENT DES DONNÉES

## Ceph utilise un algorithme appelé CRUSH

- Il décide où mettre les données en fonction d'une carte décrivant la **topologie** (racks, noeuds, disques, etc...)
- Placement déterministe : Tout le monde arrive à la même conclusion indépendamment.
- **Pas de base SQL cachée**, pas de lookup massive, **rien**



11

# CEPH : L'EMPLACEMENT DES DONNÉES

Grâce à CRUSH :

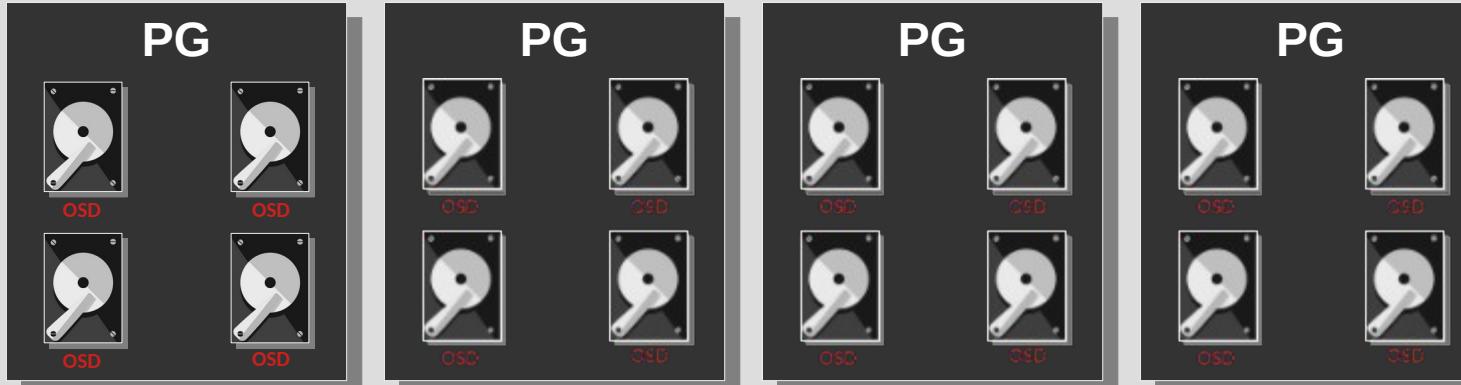
- **Une réPLICATION des données est garantie !**
- **Un disque meurt ? CRUSH recalcule ou recréer les copies**
- **Un disque s'ajoute ? Il redistribue automatiquement, la charge est équilibrée !**



# NOMENCLATURE & ÉLÉMENTS DU STOCKAGE DE CEPH

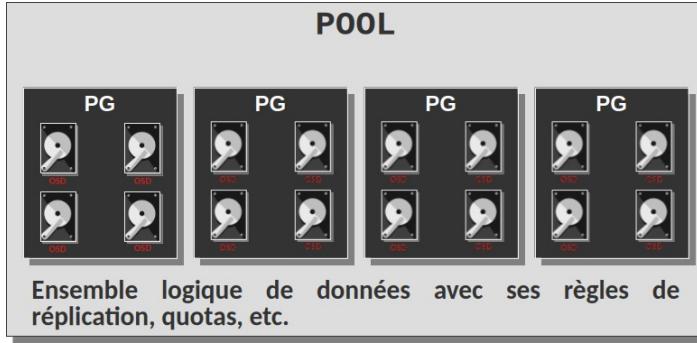
# POOLS & PG DANS CEPH :

## POOL



Ensemble logique de données avec ses règles de réPLICATION, quotas, etc.

# POOLS & PG DANS CEPH :



Rôle des **PG** :

- Éviter d'avoir un mapping « objet → OSD » gigantesque
- Équilibrer la charge plus proprement
- Simplifier la réPLICATION & la réPARATION

# TOLÉRANCE AUX PANNES DANS CEPH :



Ceph assure la **résilience via** :

- RéPLICATION (x2,x3)
- Ou erasure coding, qui optimise la capacité mais plus lourd en CPU

Quand un OSD tombe, Ceph... :

- Marque la zone inactive
- Recalcule via CRUSH ou reconstituer les copies
- Les OSD restants rééquilibrivent automatiquement

## CHOISIR SON TYPE DE STOCKAGE

# CEPH : CHOISIR SON TYPE DE STOCKAGE

**RBD** : RADOS Block Devices, typiquement adapté pour les cas suivants :

- **Libvirt + KVM**
- **OpenStack Cinder**
- **Kubernetes via CSI RDB**



Fournis des fonctionnalités assez utiles de snapshots, de clones, de thin provisioning...



# CEPH : CHOISIR SON TYPE DE STOCKAGE

**CephFS** : système de fichiers distribué conforme POSIX :

- Nécessite des **MDS** (*MetaData Servers*)
- **Scalable sur les métadonnées**, et meilleur que NFS pour les charges lourdes !

Utilisé en HPC, big data, clusters de calculs, etc...



19

# CEPH : CHOISIR SON TYPE DE STOCKAGE

## **RGW** : Interface objet compatible Amazon S3/Swift

Utilisé pour les applications cloud-native, buckets pour apps web, backups, etc...



# MÉCANISMES INTERNES & VOCABULAIRE TECHNIQUE

# CEPH : MÉCANISMES & VOCABULAIRE

- **Scrub** : vérification de cohérence
- **Backfill** : Reconstruction après panne !
- **Recovery** : rétablissement



22

# BONNES PRATIQUES POUR UN DÉPLOIEMENT CEPH

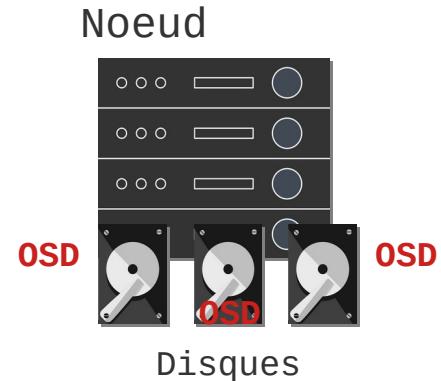
# CEPH : BONNES PRATIQUES DE DÉPLOIEMENT

Ceph est **sensible à 2 choses** :

- **La Latence réseau**
- **L'équilibrage de la topologie**

Donc généralement :

- **3 MON, ou 5, mais jamais pair !**
- **Plusieurs OSD par nœud, et un OSD par disque**
- **Un réseau dédié pour le trafic interne** (cluster network)
- Des disques rapides pour les journaux, comme des SSD ou SSD NVME



# ATTENTION !

**Il faut éviter les machines fourre-tout**

**(OSD + MON + RGW + MDS dans tous les sens sur la même machine)**

OKAI MAIS... CONCRÈTEMENT,  
LES COMMANDES

# CEPH : DÉPLOIEMENT

- La majorité du déploiement se fait par **la commande centrale** :

```
# cephadm
```

- La **sous-commande** permettant de générer le **1<sup>er</sup> nœud** :

```
# cephadm bootstrap --nom-node IP_Noeud1
```

Le **bootstrap** installe :

- **1 MON** sur le nœud en question
- **1 MGR**
- **Un shell Ceph**

# CEPH : COMMANDES ESSENTIELLES

- Vérifier la santé du Cluster Ceph :

```
# ceph status
```

- Vérifier la topologie du Cluster Ceph :

```
# ceph osd tree
```

- Vérifier l'espace :

```
# ceph df
```

# CEPH : DOCUMENTATION

- **La meilleure source d'info pour le déploiement** et la gestion d'un cluster Ceph reste... sa source **officielle** !

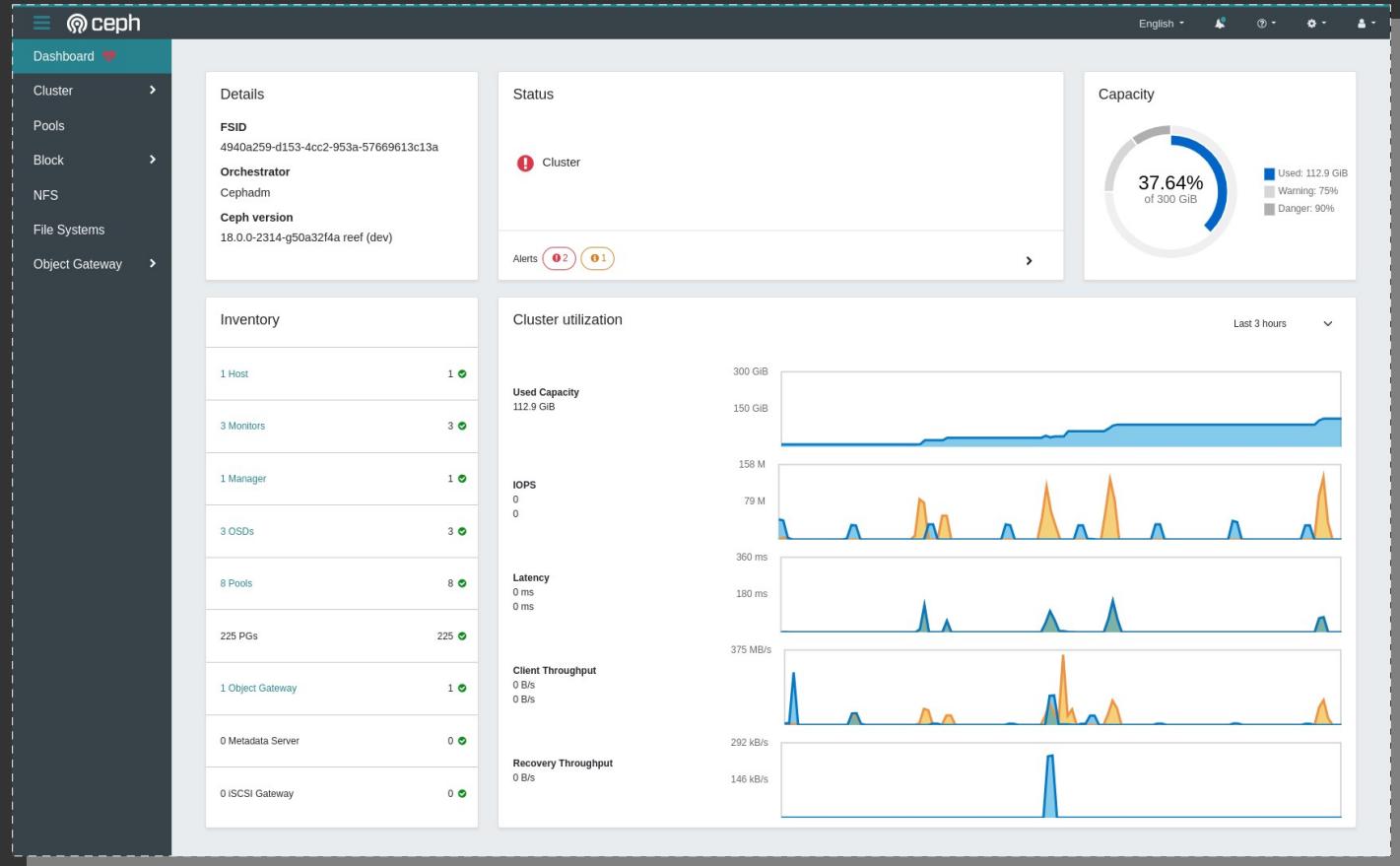


**Docs.ceph.com**

# Ceph Web UI

```
# ceph mgr module enable dashboard
```

Activera la WebUI  
d'un Cluster Ceph  
déjà mis en service



CE COURS ET LES SUPPORTS  
D'EXERCICES QUI Y SONT LIÉS  
SONT LA PROPRIÉTÉ EXCLUSIVE  
DE SON AUTEUR

MERCI DE RESPECTER LE TEMPS ET LE SOIN Y  
AYANT ÉTÉ ACCORDÉS EN NE DIFFUSANT PAS  
SON CONTENU SANS L'AUTORISATION  
EXPLICITE!

Des inexactitudes ? Des outils  
dépréciés ? Des fautes ?  
[contact@scalar-formation.fr](mailto:contact@scalar-formation.fr)

Cédric Surquin  
Admin Système Linux & Réseaux Cisco  
Consultant & Formateur, auteur du  
cours.