

The CNRS logo, consisting of the letters 'cnrs' in a white, lowercase, sans-serif font inside a dark blue circle.

cnrs

Centre de Calcul

de l'Institut National de Physique Nucléaire
et de Physique des Particules

ANF Ceph 2022 ***Ceph au CC-IN2P3***

Loïc Tortay, tortay@cc.in2p3.fr

- IN2P3
- CC
 - stockage et traitement des données pour les expériences & collaborations auxquelles participent les physiciens de l'IN2P3
 - 3500 utilisateurs, 40 à 50% hors IN2P3 (& CEA)
 - activité d'ouverture
 - 70 Pto (utiles) de disques (& Flash), 125 Pto sur bandes
 - ~1300 machines de calcul, ~500 serveurs de stockage
 - 80 personnes

- 2017 : cluster pour OpenStack @CC (RBD)
- 2019 : cluster hébergé pour l'IFB (OpenStack)
- 2020 : cluster pour LSST (CephFS)
- 2021 : extension cluster LSST (capacité +60%)
- 2021 : cluster pour WoK (RBD+CephFS+S3)
- 2022 : cluster pour remplacement GPFS (CephFS)

- OpenStack : Nautilus (Ceph v14), 600 To, 16 serveurs de disques (disques + SSDs & SSDs uniquement), 258 OSDs, réplication
- *IFB : Pacific (v16), 500 To, 4 serveurs de disques (disques + NVMe), 64 OSDs, réplication*
- WoK : Pacific, 150 To, 3 serveurs de disques, 72 OSDs, réplication
- LSST : Octopus (v15), 7 Po, 32 serveurs de disques (disques + SSDs), 640 OSDs, *erasure coding (8+2)* pour les données & réplication pour les métadonnées
- SPS : Quincy?, 7 Po, 30 serveurs de disques (disques + SSDs), 480 OSDs, CephFS, *erasure coding*
- 3 MONs (physiques) sur tous les clusters

- Dell R730xd, R740xd, R740xd2, R440 : 67 serveurs
- HPE A4200 & DL360 (pas encore en prod) : 33 serveurs
- Dell R740xd2 (cluster LSST), x32 :
 - 2 Xeon Silver 4210, 128 Gio
 - 20 disques 12 To SAS-NL, 4 SSDs 1.6 To SAS
 - 2 x 10G Eth
- HPE A4200, x30, similaire Dell R740xd2, sauf :
 - 16 disques 14 To SAS-NL, 6 SSDs 1 To SAS
 - 2 x 10/25G Eth

- RBD pour OpenStack
- CephFS :
 - très limité pour le cluster OpenStack (Spark, etc.), Manila peu utilisé en production
 - WoK en préproduction
 - utilisation à assez grande échelle pour LSST & SPS (bientôt 2 clusters de ~5 Plo utiles chacun)
- RGW très limité pour OpenStack & WoK

- Migration de 1.1 Pio de données de GPFS & Isilon vers CephFS
- Ajout de 12 serveurs de disques (+1.8 Pio utiles) :
 - 11 jours de redistribution des données (~1.2 Pio)
 - sans interruption des jobs utilisateurs
- Performances accès *mono-thread*
- Bugs mineurs (documentation, MGR)
- Montage statique au lieu de l'automonteur
- Maturité de CephFS (vs GPFS) face à des utilisateurs
- Comportement des *hardlinks* (Anaconda, utilisateurs avancés)

- Stabilité MDS (initialement 2 actifs + 1 *standby*) :
 - blocages détectés mais pas toujours leur source
⇒ redémarrage/basculement autoritaire MDS
 - pics de latence constatés par les utilisateurs
 - *scrub* (MDS pas OSD)
 - depuis début 2022, un seul MDS actif
 - `mds_cache_memory_limit` à 32 Gio
- Réduction du taux de vérification (`scrub_rate`)
- Grafana & monitoring Prometheus
- Liste de diffusion `ceph-users@ceph.io`

- Bonne résilience à un serveur de disques capricieux
 - `upmap-remapped.py` de Dan van der Ster (CERN) & `pgremapper` (Digital Ocean)
- Export Web de sous-répertoire utilisateur
- Collecte quotidienne métadonnées pour comptabilité & gestion d'espace
- Attributs étendus dynamiques (virtuels) CephFS :
 - `getfattr -n ceph.dir.rbytes $DIR`
 - ```
for attr in rbytes rsubdirs rfiles rctime ; do
 getfattr -n ceph.dir.$attr $DIR
done
```

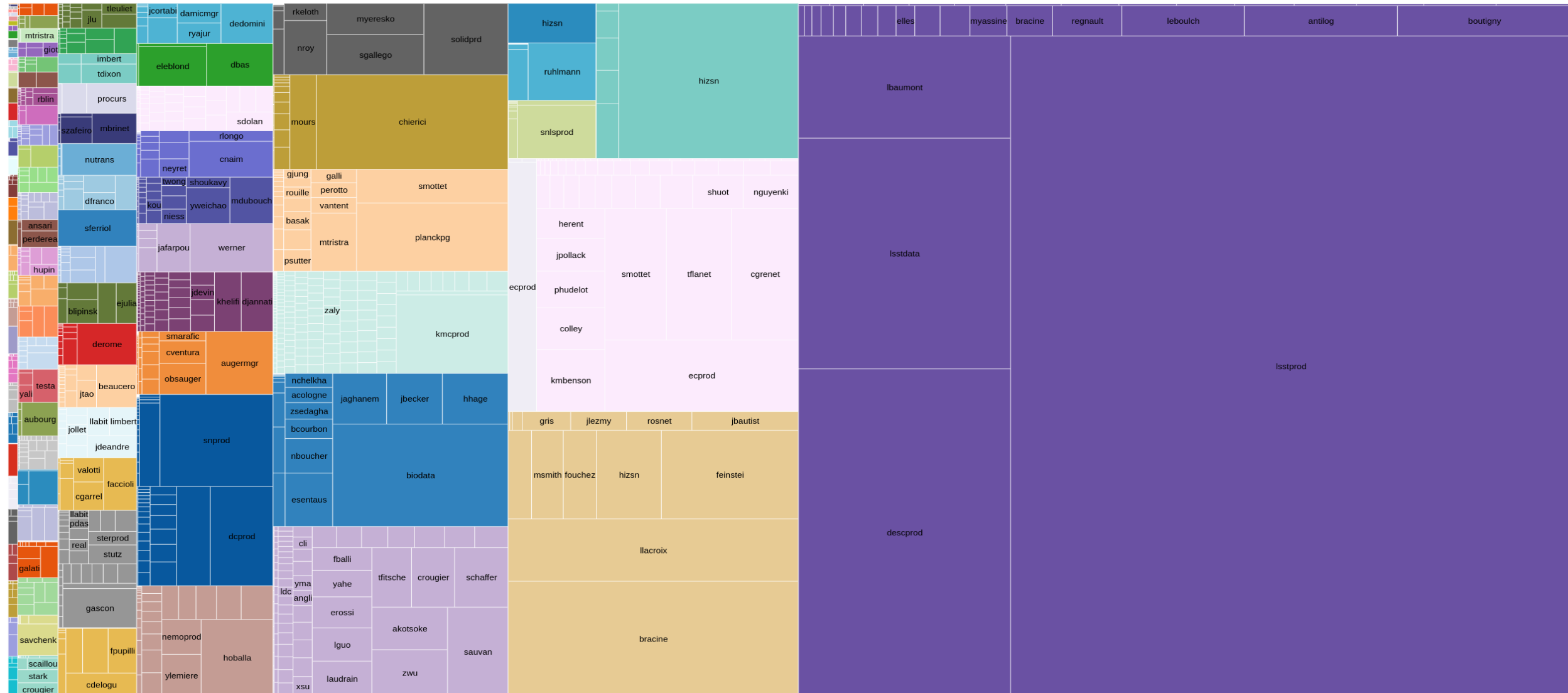
# Ceph @CC : comptabilité/statistiques



Select a metric: Space used

/sps/\* contains 8.841 PiB (+overhead: 899.51 TiB) in 1.693G files (10.1M hard-linked)/140M dirs/96.5M symlinks for 2607 users, max name length: 686, max depth: 52

Data for Wednesday 5 October 2022 12:01 (UTC+0200)

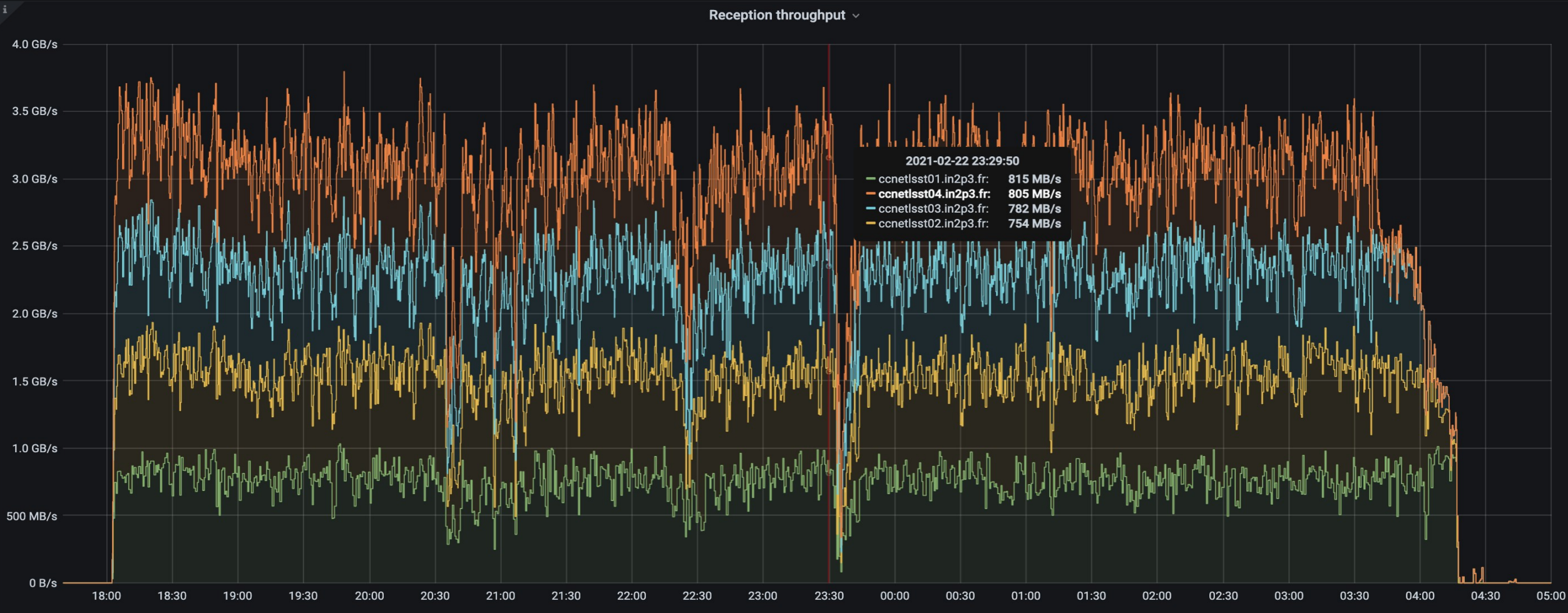


# Ceph @CC : transferts LSST CC ↔ NERSC

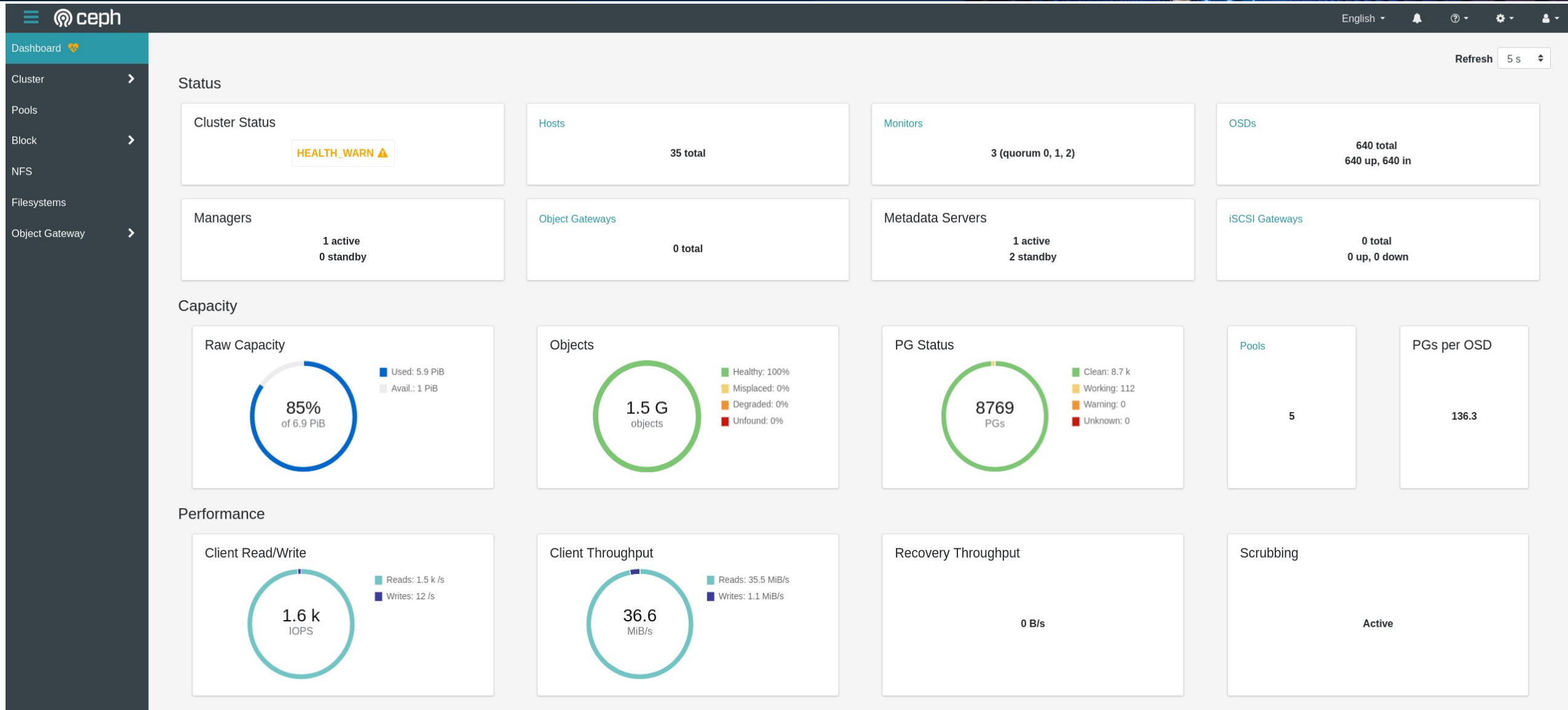


← LSST@CCIN2P3 OVERVIEW ☆ 🔗

📊 📄 ⚙️ < 2021-02-22 17:40:00 to 2021-02-23 05:00:00 > 🔍 ↻



# Ceph @CC : Dashboard Ceph



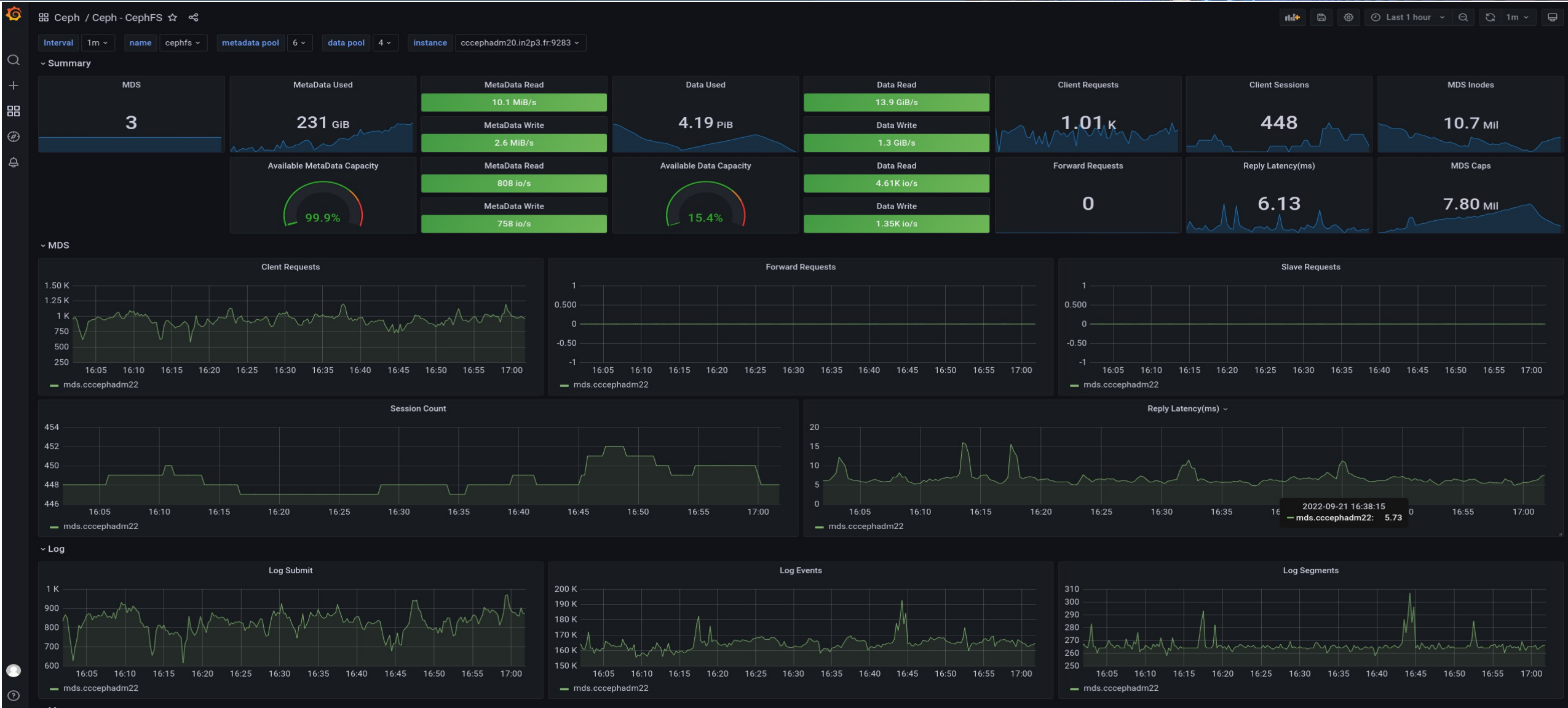
# Ceph @CC : Grafana Ceph



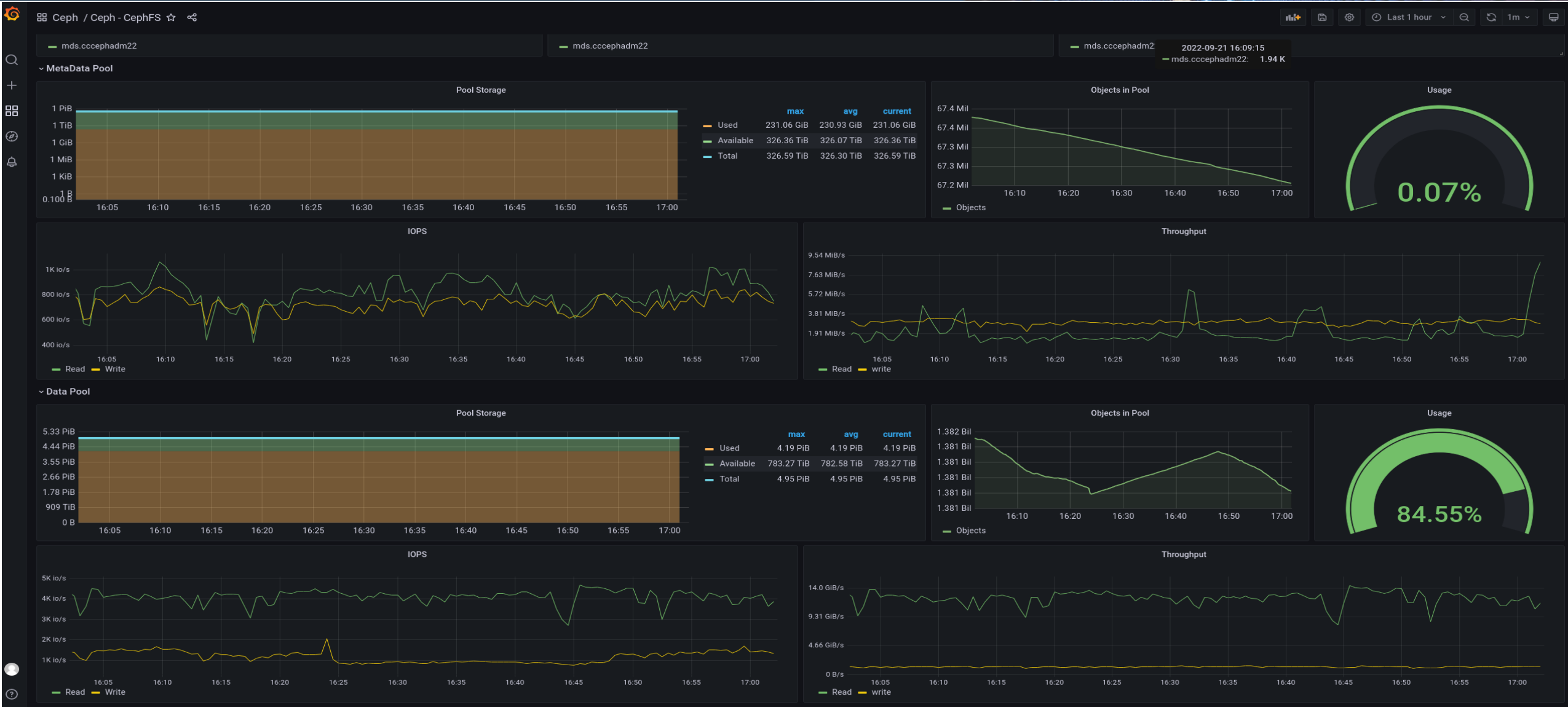
# Ceph @CC : Grafana Ceph activity heatmap



# Ceph @CC : Grafana CephFS



# Ceph @CC : Grafana CephFS





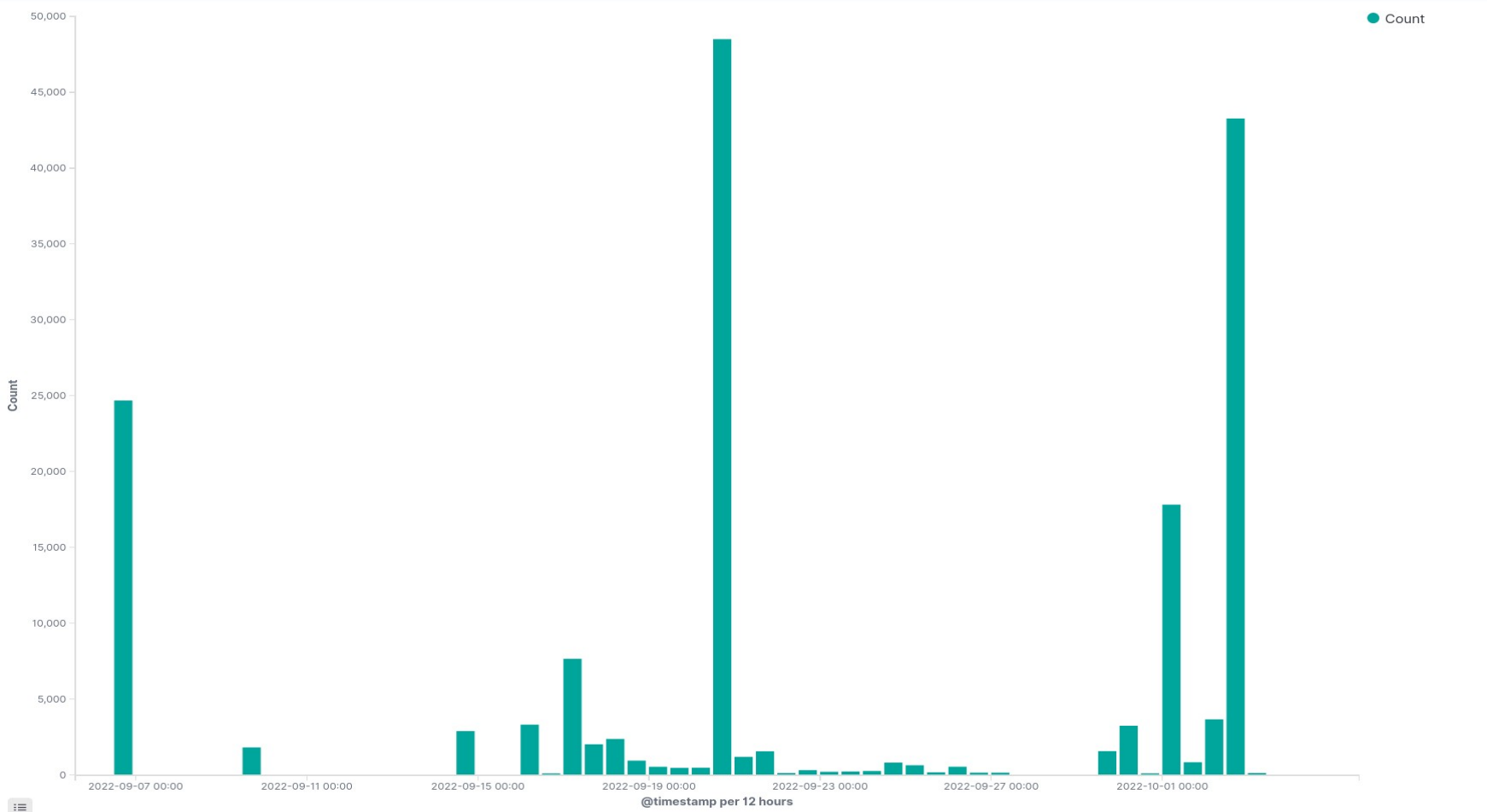
# Ceph @CC : Elastic Search CephFS slow ops



Elastic Visualize / Ceph LSST slow ops

"slow ops" KQL Last 30 days Show dates Refresh

factor.site\_usage: ceph factor.site\_instance: lsst + Add filter



syslog-\*

Data Metrics & axes Panel settings

**Metrics**

- > Y-axis Count + Add

**Buckets**

- > X-axis @timestamp per 12 hours [eye icon] [x icon] + Add

Discard Update

6 octobre 2022

Ceph @CC-IN2P3

# Ceph @CC : Grafana CephFS ?!?

