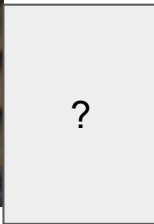# State of Storage

CdG 20 Settembre, 2024
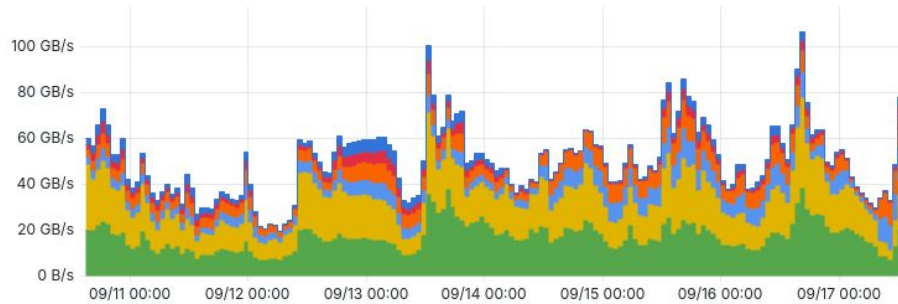
# Business as usual + migration to TP
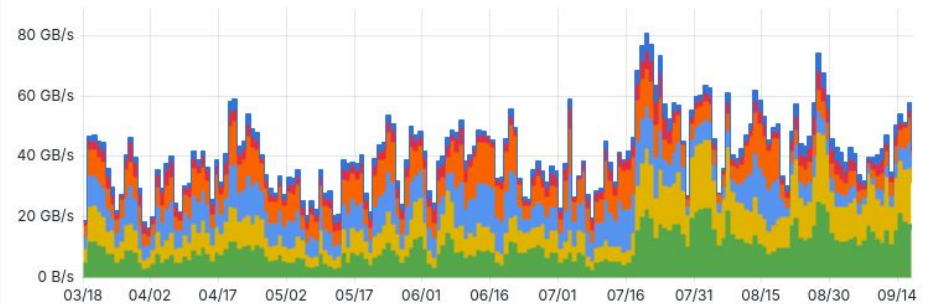


Last month

Last 6 months

# Disk storage in produzione

Installed: 113PB - 33PB (in dismissione)=80.6PB Pledge 2024: 82.08PB, Used: 48.8PB

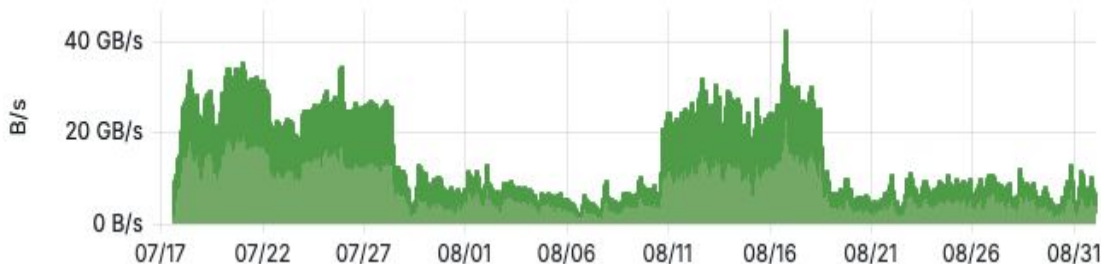| Storage system | Model | Net capacity, TB | Experiment | End of support |
|---|---|---|---|---|
| ~~ddn-10, ddn-11~~ | ~~DDN SFA12k~~ | ~~10120~~ | ~~ALICE, AMS~~ | ~~12/2022~~ |
| os6k8 | Huawei OS6800v3 | 3400 | GR2, Virgo | 07/2024 |
| md-1,md-2,md-3,md-4 | Dell MD3860f | 2308 | DS, Virgo, Archive | 12/2024 |
| md-5, md-6 e md-7 | Dell MD3820f | 50 | metadati, home, SW | 11/2023 e 12/2024 |
| os18k1, ~~os18k2~~ | Huawei OS18000v5 | 960 | LHCb (buffer tape) | 7/2024 |
| os18k3, os18k5, os18k5 | Huawei OS18000v5 | 1200 | ATLAS,ALICE (buffer tape) | 6/2024 |
| ddn-12, ddn-13 | DDN SFA 7990 | 5840 | GR2,GR3 | 2025 |
| ddn-14, ddn-15 | DDN SFA 2000NV | 24 | metadati | 2025 |
| os5k8-1,os5k8-2 | Huawei OS5800v5 | 8999 | ATLAS | 2027 |
| ~~Cluster CEPH~~ | 12xSupermicro SS6029 | ~~3400~~ | ~~ALICE, cloud, etc.~~ | 2027 |
| od1k6-1,2,3,4,5,6 | Huawei OD1600 | 60000 | ALICE,ATLAS,LHCb, CMS | 2031 |

# Acquisti recenti e futuri

- Gara storage 2022 (14PB netti)
  - Nuova proposta con apparati DDN SFA7990X
  - In attesa per la consegna entro settembre
- Tape Library
  - Installata, collaudo completato
  - Le cassette JF da 50TB sono state inserite nella libreria (7.8PB)
- Gare nastri
  - Nuova gara di acquisto tape JF (96PB)

# Migrazione dati sul nuovo storage a TP

- Con il trasloco a TP abbiamo spostato ~50PB di dati;
  - Pledge complessiva 2024 è 82.08PB;
  - Pledge 2024 per 4 exp LHC (56.2PB) va a coprire quasi tutto il nuovo storage appena installato (60PB);
- Gli spostamenti iniziati il 17/07 con ALICE, ATLAS, LHCb e dal 10/08 con CMS
- Average rate di migrazione ~10-20GB/s per filesystem

- L'attività di produzione ha un impatto significativo sulle migrazioni dei dati.

# Problematiche relativi allo storage AQ 2023-2024

- Huawei OceanStore Micro 1500/1600
  - Sono stati forniti 6 sistemi di 10PB + 32 server
  - Installazione e collaudato andati bene
  - Abbiamo notato i primi problemi (riduzione delle prestazioni) quando lo storage è stato riempito del 95% di capacità.
  - Si è scoperto che lo storage funziona solo in modalità thin provisioning, il che significa che scrive solo su spazi non ancora utilizzati.
  - Per liberare spazio dai dati cancellati nel file system, è necessario eseguire la procedura "reclame space" sullo stesso, seguita dal invocazione del "garbage collector" lato storage.
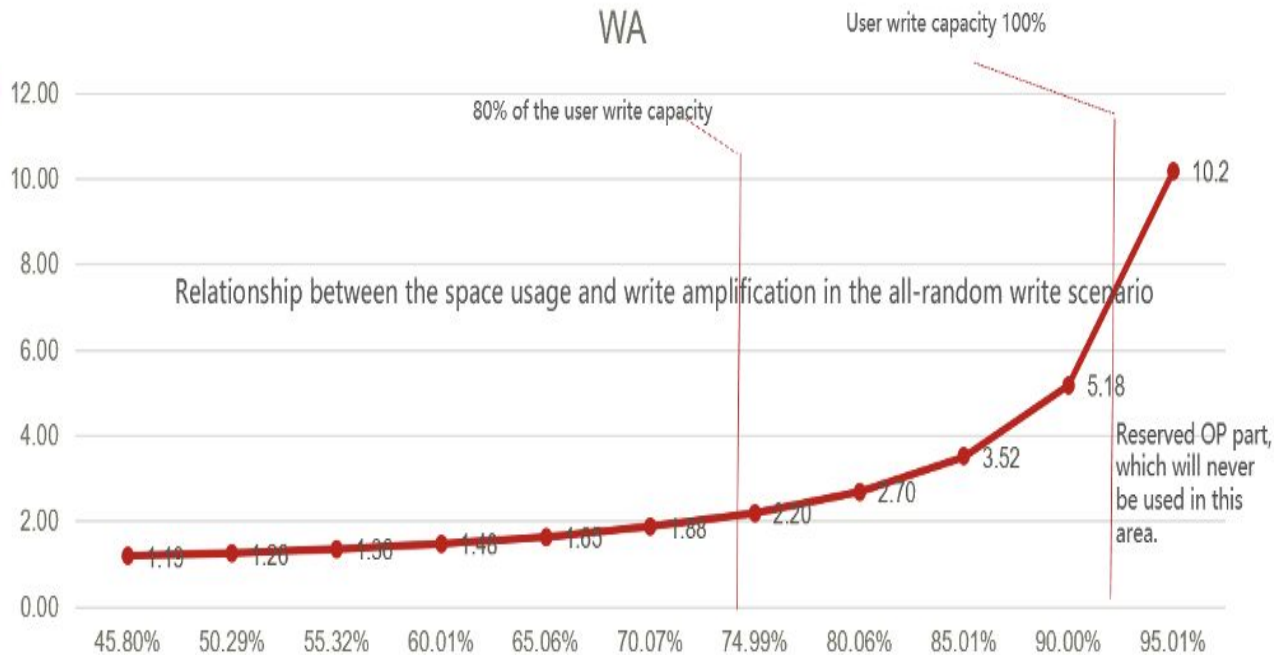
## Always New Write Cost: Garbage Recycling

Garbage collection generates additional write amplification.

The system always selects the CKG with the highest garbage amount for recycling.
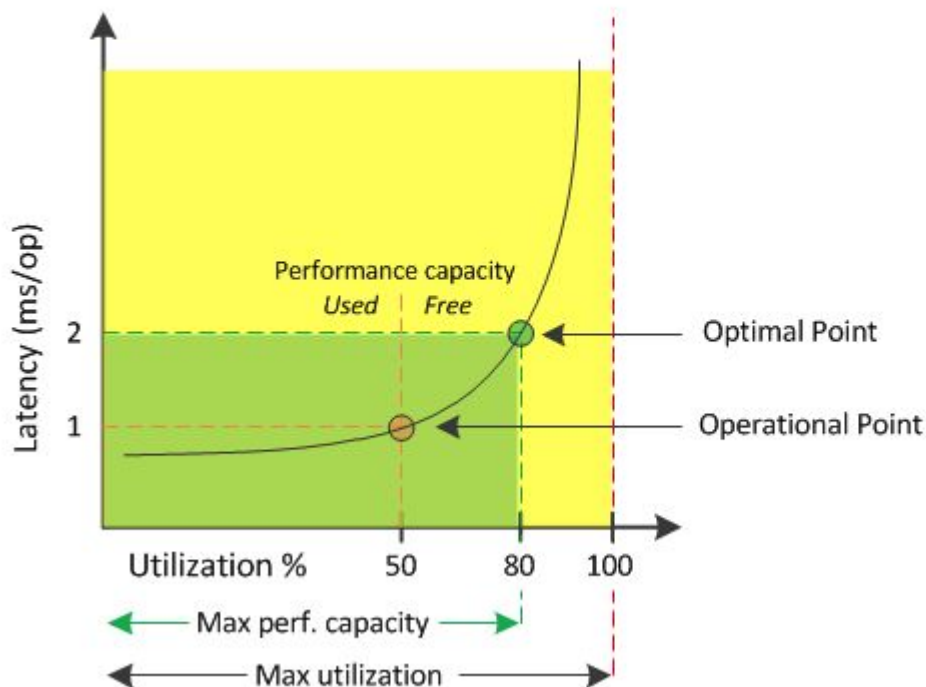


Relationship between the space usage and write amplification in the all-random write scenario

# Problematiche relativi allo storage AQ 2023-2024

- File system come configurato adesso non supporta la funzionalità di thin provisioning.
  - L'abilitazione del thin provisioning
    - Comporta significativa riduzione delle prestazioni, poiché la procedura di "space reclamation" scrive zeri sullo spazio da liberare (che crea traffico e richiede tempo non trascurabile)
    - Richiede avere circa 5% spazio riservato (non accessibile agli utenti) [1]
- Risposta del Supporto Huawei:
  - **Suggerimento:** In generale, l'utilizzo del storage pool dell storage non dovrebbe superare il 95%. Pertanto, consigliamo di ridurre l'utilizzo del storage pool del storage a meno del 95%, ad esempio al 90%, per ridurre i rischi relativi all'affidabilità.
- Conclusione: l'uso di thin provisioning non e' praticabile
  - Che con il nostro "use case" tipico (100TB scritti e cancellati in un giorno) significa di restare fermi per 3-4 ore al giorno.
- Un problema riconosciuto anche dagli altri produttori [2],[3]
    - References:
      1. https://www.ibm.com/docs/en/storage-scale/5.1.8?topic=devices-storage-scale-thin-provisioned
      2. https://www.dell.com/support/kbdoc/it-it/000123351/powerstore-alerts-capacity-utilization?lang=en
      3. https://docs.netapp.com/us-en/active-iq-unified-manager-97/online-help/concept-what-performance-capacity-used-is.html

# NetApp:

The optimal used performance capacity is the point at which a node or an aggregate has optimal utilization and latency (response time), and is being used efficiently. A sample latency versus utilization curve is shown for an aggregate in the following figure.

# Come risolviamo questi problemi?

- Per diminuire stress del work flow di LHCb abbiamo migrato buffer tape sullo HW separato
  - Siamo in attesa del feedback dal esperimento
- Consideriamo anche la possibilità di creare un "buffer disco" per i dati "hot"
- Attualmente stiamo discutendo con il team di R&D e Sales di Huawei su come ottimizzare le prestazioni dei loro sistemi.
  - Hanno proposto aumentare capacità del sistema per creare lo spazio "riservato" per restare nella zona con le prestazioni accettabili

# Stato tape

Last 2 months

MSS bytes in/out (per day)



| Name | Mean | Last * | Max | Min | Total |
|---|---|---|---|---|---|
| out traffic (recalls) | 65.8 TB | 109 TB | 241 TB | 6.23 GB | 3.82 PB |
| in traffic (migrations) | 157 TB | 96.1 TB | 297 TB | 31.4 TB | 9.10 PB |

9 PB of new data written to tapes in two month (since last CdG)

# Tapes: Migration from Oracle to IBM library stopped

Repack - Library Space Occupancy



| Name | Last | Difference |
|------|------|------------|
| SL8500 | 36.7 PB | -7.72 PB |
| TS4500_1 | 99.6 PB | 28.7 PB |

Repack - number of free cartridges

380 tapes inserted

Repack - Library Scratch Tape



| Name | Last |
|------|------|
| SL8500 | 4139 |
| TS4500_1 | 56 |

SL8500 has 36.7PB to migrate (2.5PB more than 2 month ago)

# Stato tape

- Liberi 1.1 PB (Scratch tape sulla libreria IBM).
- Usati 136 PB.
- Spazio tape sulla libreria IBM praticamente esaurito
- La nuova libreria IBM non è ancora funzionante a causa di problemi di compatibilità con la versione TSM in produzione.

| Library | Tape drives | Max data rate/drive, MB/s | Max slots | Max tape capacity, TB | Installed cartridges | Used space, PB | Free space, PB |
|---|---|---|---|---|---|---|---|
| SL8500 (Oracle) | 16*T10KD | 250 | 10000 | 8.4 | ~10000 | **36.7** | - |
| TS4500 (IBM) | 19*TS1160 | 400 | 6198 | 20 | 5100+380 | **99.6** | **1.1** |
| TS4500-2(IBM) | 18*TS1170 | 400 | 7844 | 50 | 165 | **0** | **8.2** |

# Current SW in PROD

- GPFS 5.1.2-15, in preparazione migrazione alla 5.1.9-6
  - Supporto RHEL 9 ed architettura ARM
- StoRM BackEnd 1.11.22 (latest)
- StoRM FrontEnd 1.8.15 (latest)
- StoRM WebDAV 1.4.3 (latest)
- StoRM globus gridftp 1.2.4
- XrootD 5.5.4-1
  - LHCb updated to 5.5.5-1
- Ceph 16.2.6 (Pacific)

# Tickets and more

- ALICE
  - Finishing configuration restyling of XrootD GPFS cluster:
    - Finalizing the configuration for the tape cluster (xs-204, xs-304)
      - F. Noferini has a working "tsm" RPM building procedure for EL9
      - Waiting for the migration of servers to EL9 to install and test it
  - Data migrated from CephFS to GPFS
    - Grid transfer via XrootD
    - CEPH XrootD SE dismissed

# Tickets and more

- ATLAS
  - GGUS [168159](#) (closed): staging errors due to an expired certificate
  - GGUS [167957](#) (on hold): StoRM WebDAV does not permit the creation of non-existent parent directory even if the scope does it
    - Waiting for the StoRM fix
  - GGUS [167840](#) (closed): low transfer efficiency as a destination (SSL connect errors)
    - Issue caused by GGUS 167725 / 167685
  - GGUS [167725](#) (closed): low staging efficiency and low transfer efficiency
    - Caught more than 100k of recalls up
      - Files with incoherent status on the tape buffer
  - GGUS [167685](#) (closed): on the 23rd of July Atlas wrote 220TB in 20 hours filling the tape buffer up
    - We highly recommend not to exceed the writing rate limit (recalls included) of 1.0GB/s

# Tickets and more

- CMS
  - GGUS [168130](#) (solved): network issues related to a NIC of the core switch at Tecnopolo
    - GocDB DT starting from the 7:55 of Sep 12th up to the 9:00 of Sep 13th (UTC)
  - GGUS [167995](#) (on hold): StoRM WebDAV does not permit the creation of non-existent parent directory even if the scope does it
    - Waiting for the StoRM fix
  - GGUS [167653](#) (closed): SAM tests failed for the issue mentioned in GGUS 167642
  - GGUS [167642](#) (closed): CMS wrote 840TB in 2 days filling the tape buffer up
    - Huge backlog to manage
    - We highly recommend not to exceed the writing rate limit (recalls included) of 1.2GB/s
  - GGUS [167634](#) (closed): CMS WebDAV SSL connection test fails on one server in xfer-cms
    - Issue due to overload of data transfer servers and thread limit reached
    - Servers have been reconfigured to provide a better thread management

# Tickets and more

- LHCb
  - GGUS [167716](#) (in progress): low transfer efficiency with new storage HW installed at TP
    - Performance decreases with the file system occupancy and the pressure of the experiment data flow
      - 6 StoRM WebDAV servers separated from the NSD ones
      - Dedicated HW for tape buffer
      - Following closely the situation via weekly reports to WLCG management board & operations coordination since Aug 30th
  - GGUS [167586](#) (closed): failed data transfers due to overload of the LHCb cluster at CNAF caused by POSIX access and FTS transfers; request to lower submission rate of FTS jobs
  - GGUS [167045](#) (closed): data movement to the new data center facility
    - GocDB DT starting from the 17:00 of Jul 22nd up to the 8:49 of Jul 29th (UTC)

# Tickets and more

- Gsiftp protocol via StoRM backend is still available for two experiments
  - New StoRM release should finally allow to switch GridFTP off (Xenon, CTA-LST)
- AMS
  - HW issues prompted a FS migration to different systems
  - The migration lasted about 7 days during which the jobs submission has been stopped
  - The situation is back to normal since Sep 14th
- Dampe
  - GridFTP "plain" still used
    - TPCs between XrootD server at IHEP and CNAF are working well
    - Rucio+FTS (https) should replace the current gsiftp transfers (WP6-DataCloud)
- DUNE
  - Configuration ongoing to expose data in read mode also via XrootD

# Tickets and more

- FAMU
  - Metadata on tape has been changed in order to expose files via StoRM WebDAV, currently in only write mode with JWT issued by iam-t1-computing, "famu" group
    - davs://xfer-archive.cr.cnaf.infn.it:8443/famu-tape
- Gminus2
  - New StoRM WebDAV storage area pointing to dedicated fileset /storage/gpfs_data/gminus2
    - "fermilab" voms-proxy AuthN/Z
    - davs://xfer-archive.cr.cnaf.infn.it:8443/gminus
- Muone
  - New StoRM WebDAV storage area pointing to dedicated fileset /storage/gpfs_data/muone
    - JWT AuthN/Z - "muone" group of iam-t1-computing
    - davs://xfer-archive.cr.cnaf.infn.it:8443/muone

# Tickets and more

- Newchim
  - Metadata on tape has been change in order to expose files via StoRM WebDAV, currently in only write mode with token issued by iam-t1-computing, "newchim" group
    - davs://xfer-archive.cr.cnaf.infn.it:8443/newchim-tape
- PAuger
  - New StoRM WebDAV storage area pointing to dedicated fileset /storage/gpfs_data/pauger
    - "auger" voms-proxy and JWT AuthN/Z ("pauger" group of iam-t1-computing)
    - davs://xfer-archive.cr.cnaf.infn.it:8443/pauger