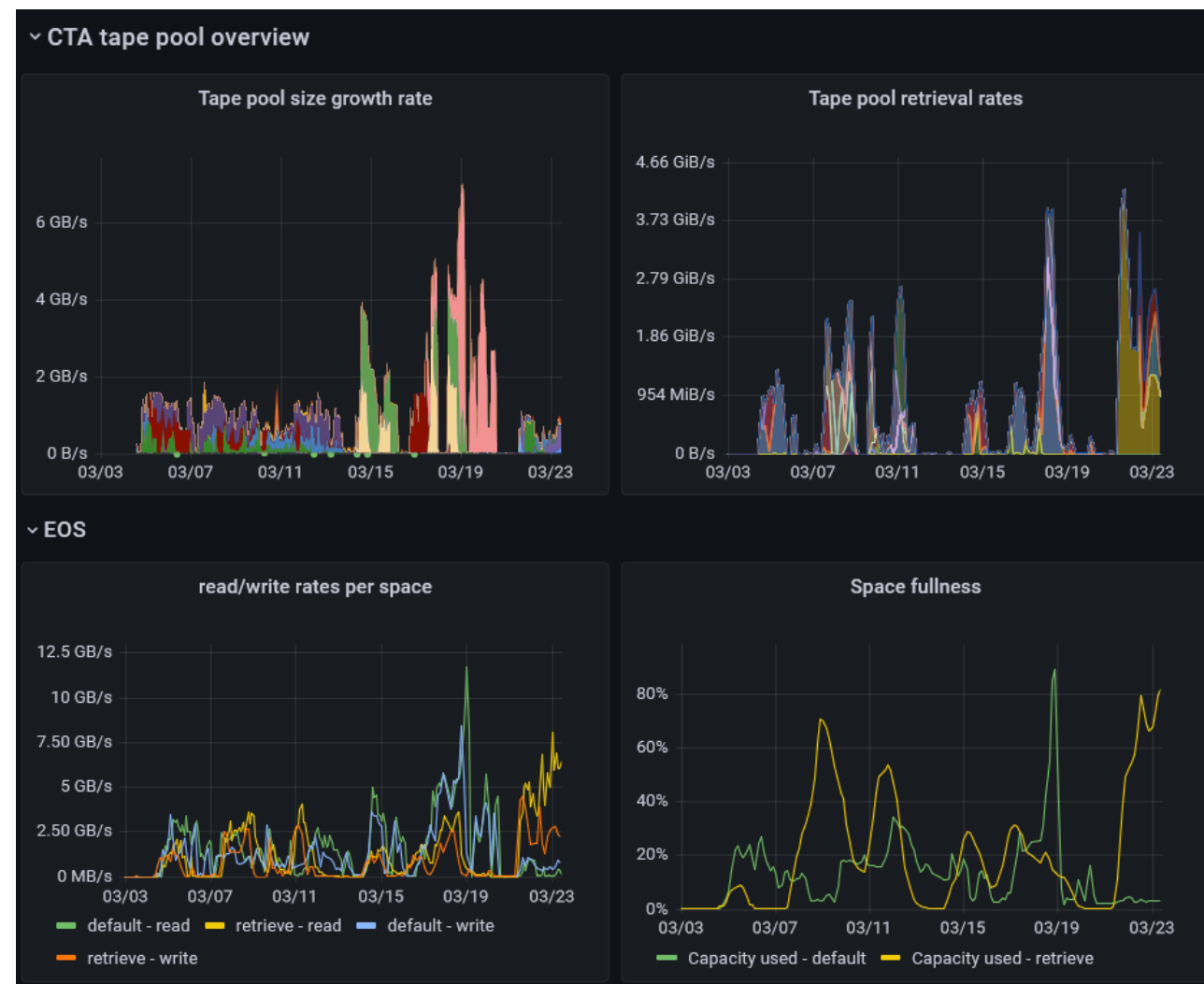# RAL Antares Update

Tom Byrne, George Patargias
23rd March 2022
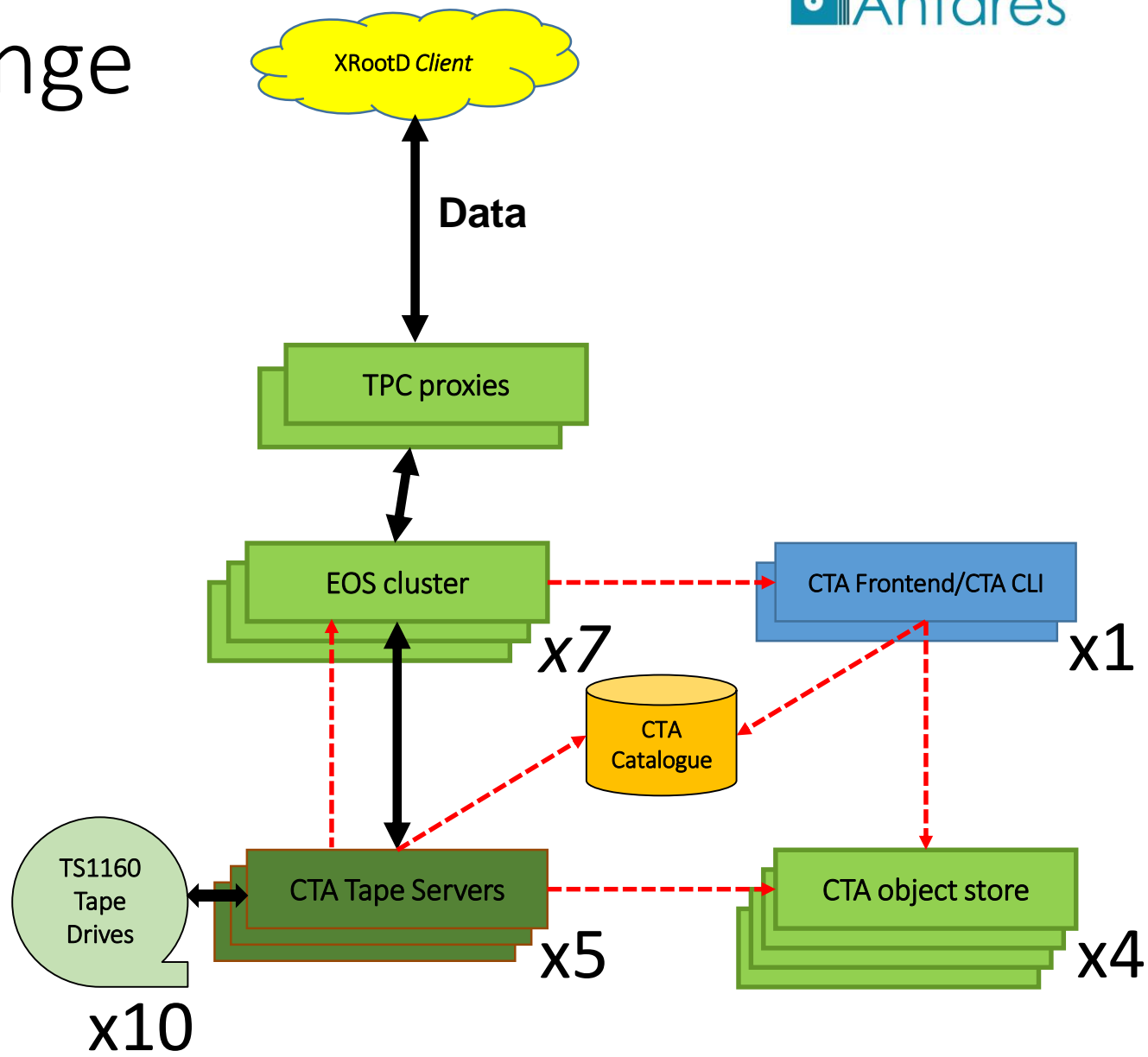
# Recent progress highlights – GridPP47

1. Participated in the 2021 LHC data challenge
2. Planned and tested migration from RAL WLCG CASTOR to EOS+CTA
3. Finalised EOS+CTA setup at RAL and rebuilt the production instance at full scale
4. Migrated from CASTOR to EOS+CTA
5. Antares has been in production for 19 days
6. Currently participating in another LHC data challenge



*First 19 days of production*
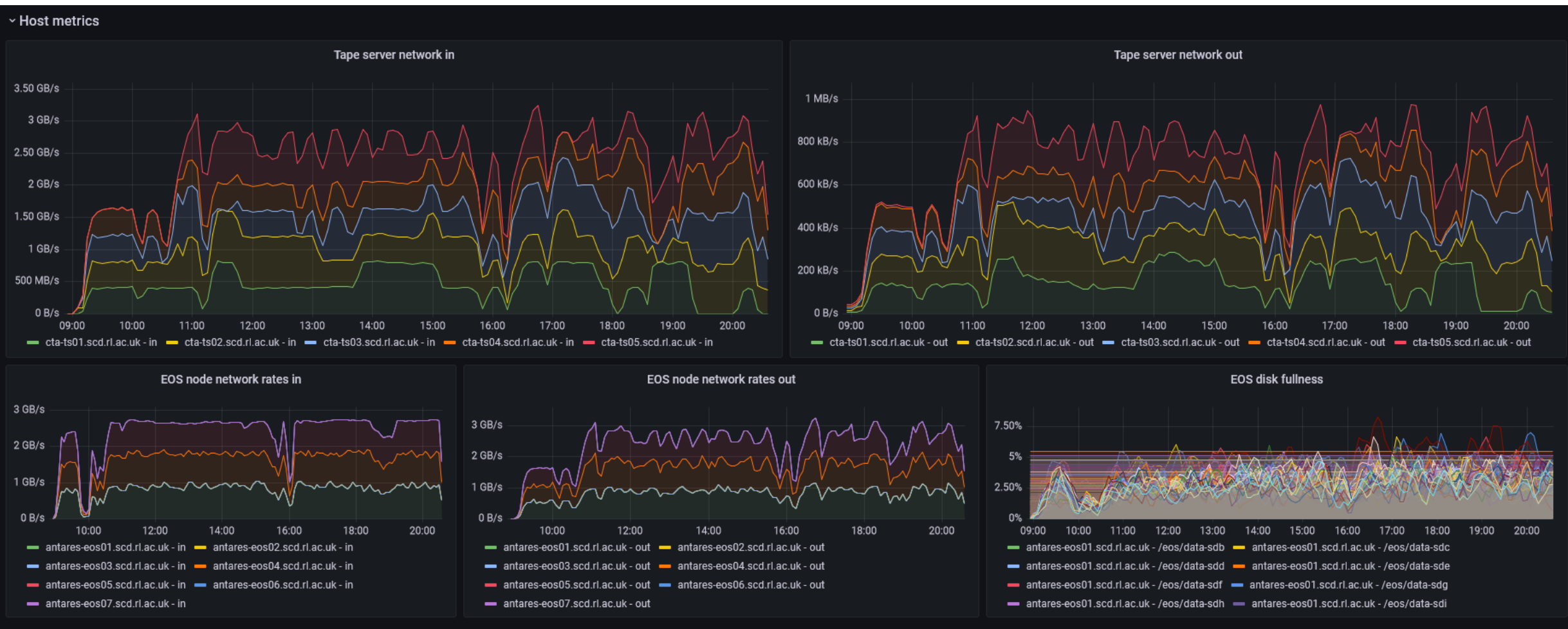
# 2021 LHC data challenge

- The first of several data challenges to ensure custodial data storage systems are ready for Run 3 rates
  - Was a good opportunity to validate CTA performance with these tests

- Deployed the largest EOS+CTA stack we had ever run for this DC

- Lots of learning done – a very valuable experience for a brand new system



**Tier-1 Spectra Logic Tfinity tape library**



XRootD *Client*

**Data**

TPC proxies

EOS cluster   *x7*

CTA Frontend/CTA CLI   x1

CTA Catalogue

TS1160 Tape Drives   x10

CTA Tape Servers   x5

CTA object store   x4

Antares

UKRI
Science and Technology Facilities Council

# 2021 LHC data challenge – Day 1 – ATLAS and LHCb Archiving
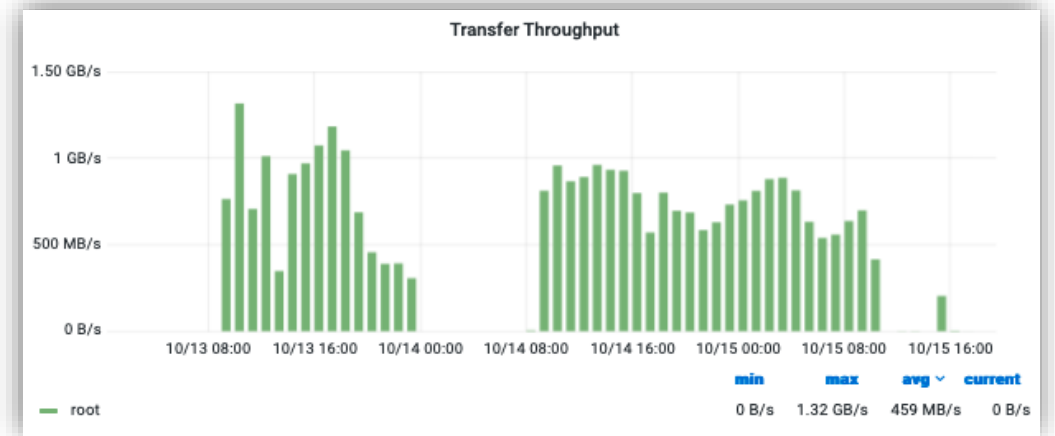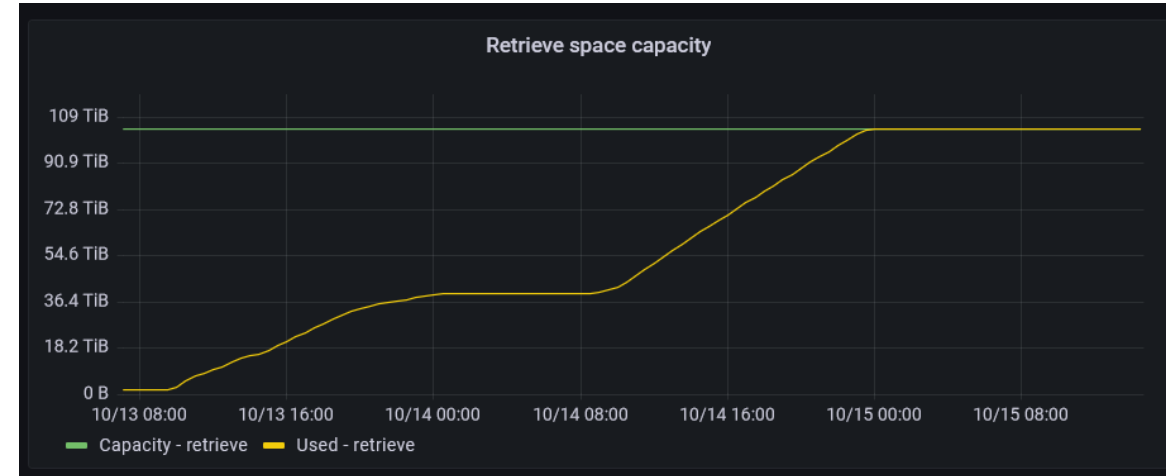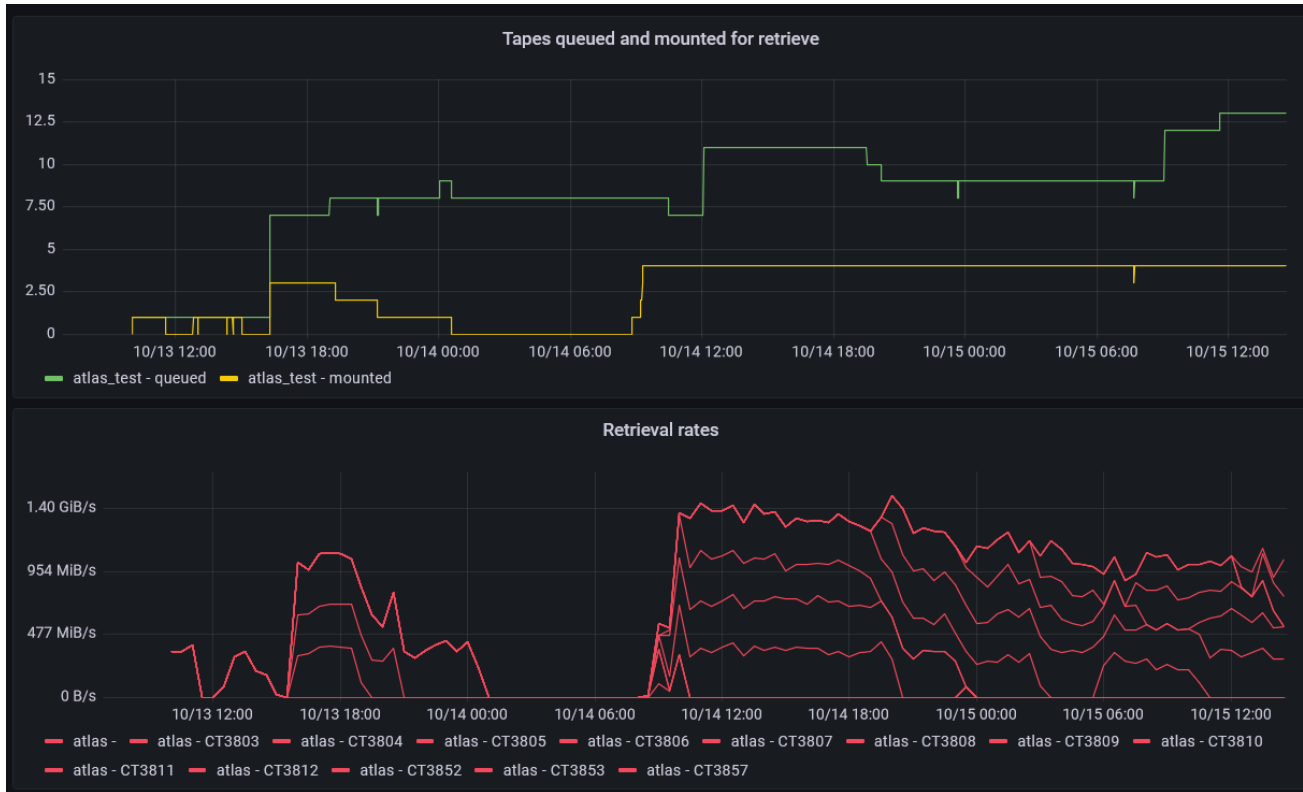
0900 – 2100 11th October

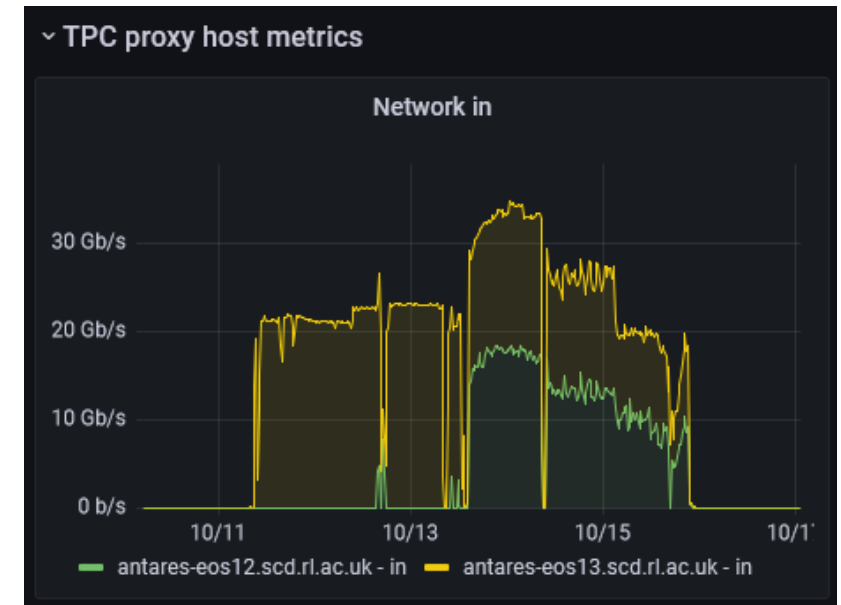# 2021 LHC data challenge – Day 3/4 – Atlas recall

0800 13th – 0800 15th October

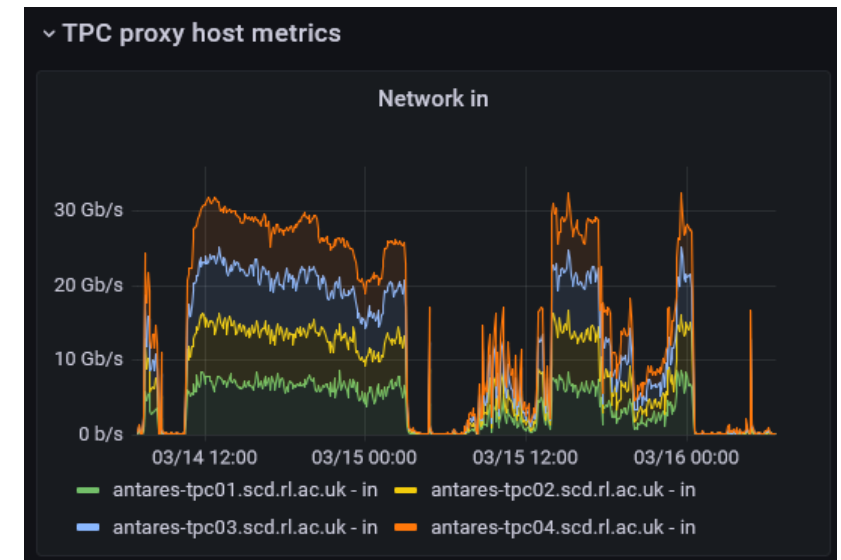**Without specific 'xrootd.site' name specified, FTS will not evict transferred file**





*DDM monitoring of transfers to Echo*

# Data challenge lessons – TPC proxies

- Requirement for these was not well understood going into the first DC
  - CERN CTA do not need these as most TPC transfers are from trusted hosts (CERN EOS)
- Started with one repurposed EOS node as a proxy, added another by the end of the challenge week to cope with demand
- Going into production, we have four dedicated TPC nodes



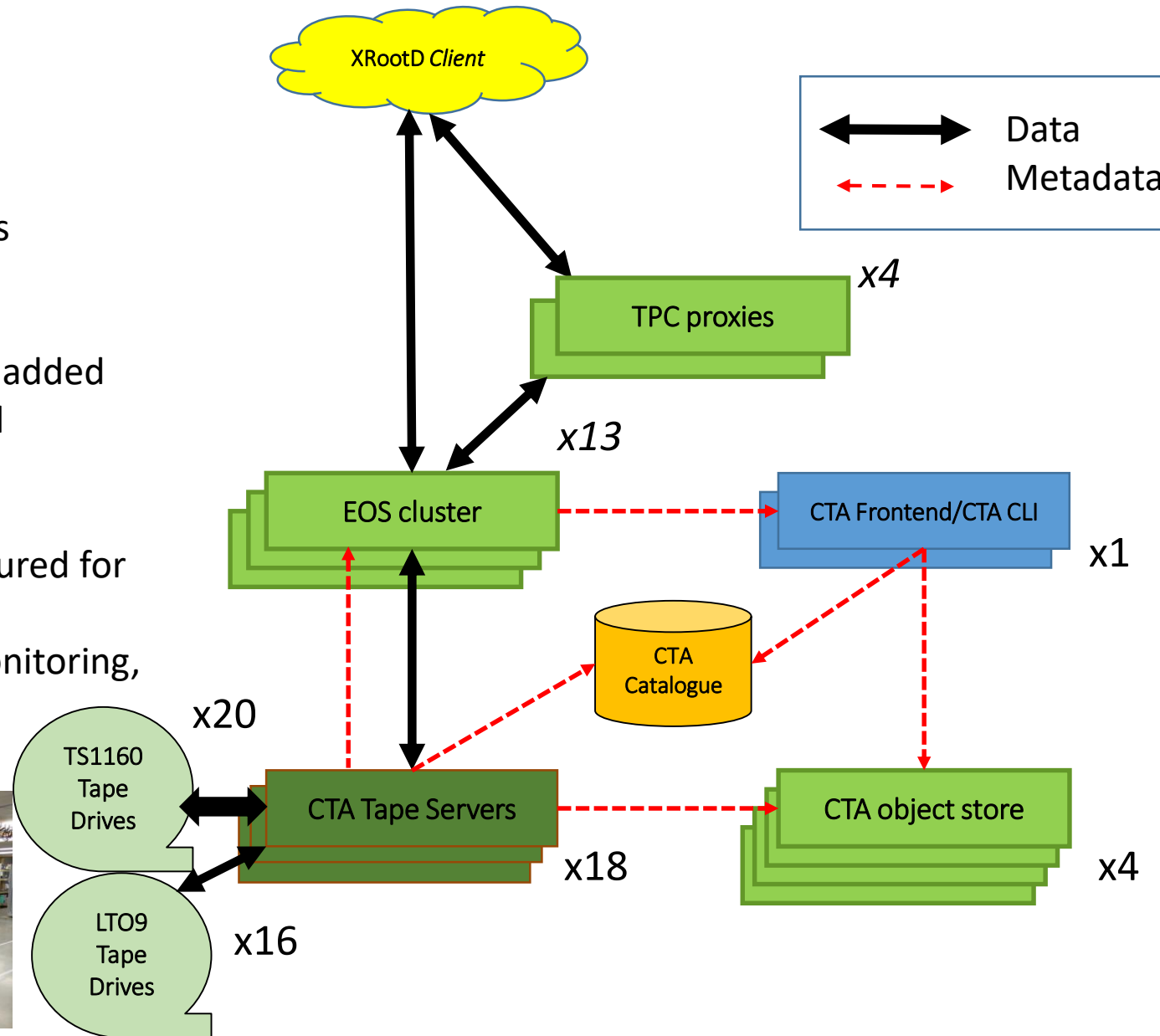*Data challenge 1 setup*



*Production setup*

# Antares production setup

**Hardware changes:**
- Dedicated TPC proxies configured (100Gb/s combined throughput available)
- Full compliment of EOS nodes provisioned
- Another 5 tape servers with TS1160 drives added
- 7 tape servers with LTO9 tape drives added

**Software changes:**
- Alice authentication configured and tested
- WebDav support (incl. TPC support) configured for LHCb needs
- Lots of production readiness changes – monitoring, alerting, DR
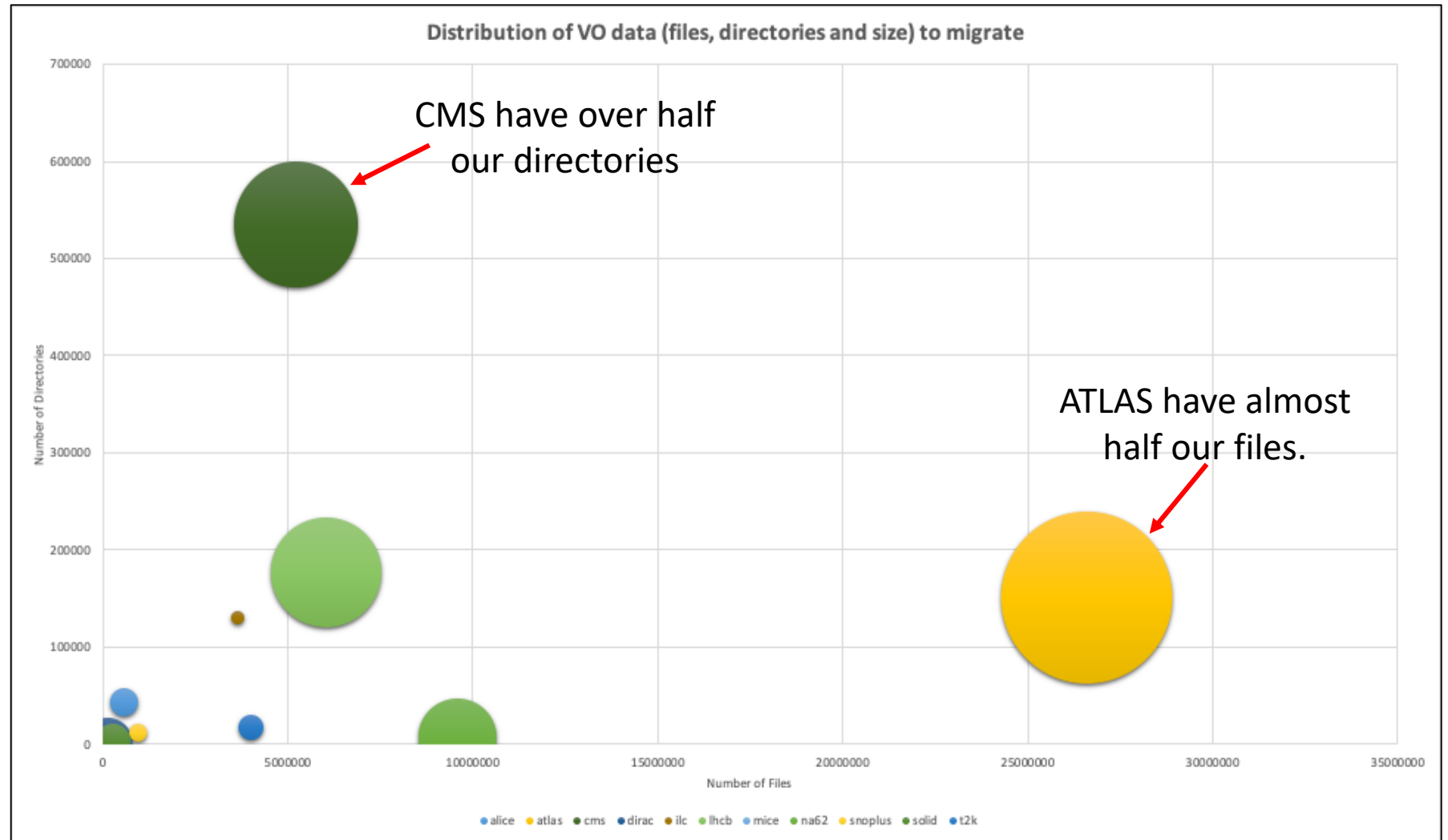
# Tier-1 Castor to Antares Migration

- Migration method:
  - ➢ Castor namespace injection to EOS
  - ➢ Castor tape metadata migration to CTA DB

- Migration pre-requisites:
  - ✓ Upgrade to CASTOR to 2.1.19-3
  - ✓ Import CASTOR DB schemas (NS,VMGR,STAGER) snapshot to the CTA DB
  - ✓ Review/modify PL/SQL scripts to be run on the schemas
  - ✓ Further namespace clean up (repack files to the right tape pools)
  - ✓ Set up a migration node to run the migration client tools
  - ✓ Estimate timings to be scheduled in the intervention plan

| Row Labels | Mispl | Count of | Average of FILI | Sum of FILESIZE |
|---|---|---|---|---|
| ⊟ dead | 1 | 1 | 1.36E+09 | 1.36E+09 |
| atlas | 1 | 1 | 1.36E+09 | 1.36E+09 |
| ⊟ dirac | 3 | 42394 | 3.51E+08 | 1.49E+13 |
| ilc | 1 | 39242 | 1.52E+08 | 5.97E+12 |
| lhcb | 1 | 1994 | 4.46E+09 | 8.89E+12 |
| t2k.org | 1 | 1158 | 8.79E+06 | 1.02E+10 |
| ⊟ ilc | 2 | 412076 | 1.12E+08 | 4.62E+13 |
| lhcb | 1 | 1 | 1.07E+09 | 1.07E+09 |
| t2k.org | 1 | 412075 | 1.12E+08 | 4.62E+13 |
| ⊟ lhcb | 2 | 58525 | 1.11E+08 | 6.47E+12 |
| ilc | 1 | 6514 | 1.53E+08 | 9.95E+11 |
| t2k.org | 1 | 52011 | 1.05E+08 | 5.47E+12 |
| ⊟ t2k.org | 1 | 14002 | 8.36E+07 | 1.17E+12 |
| ilc | 1 | 14002 | 8.36E+07 | 1.17E+12 |
| Grand Total | 4 | 526998 | 1.30E+08 | 6.87E+13 |

# Tier-1 CASTOR to Antares migration

**Total:**
- 1,079,217 dirs
- 57,011,928 files
- 70.5PB



Distribution of VO data (files, directories and size) to migrate

CMS have over half
our directories

ATLAS have almost
half our files.

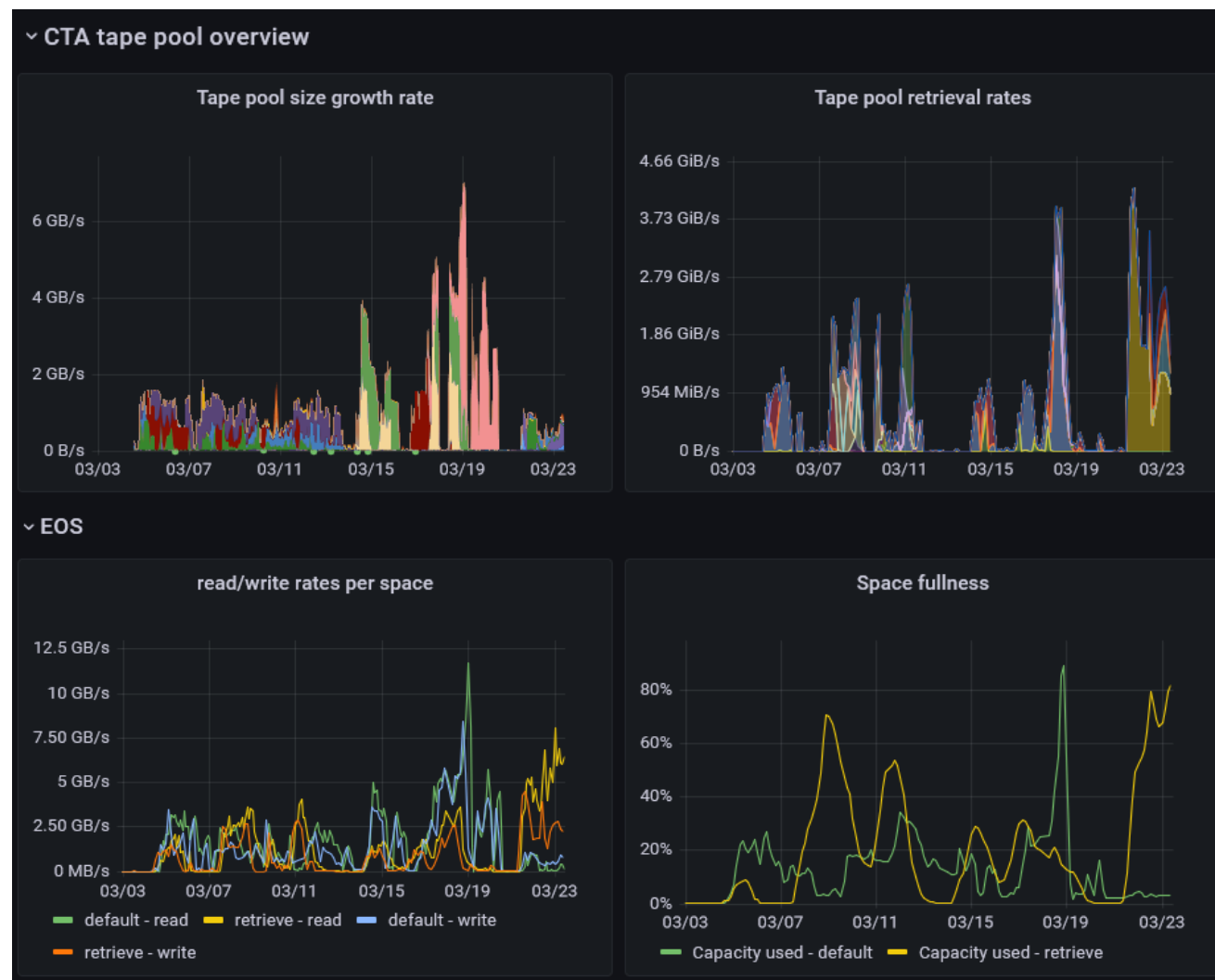Number of Directories

Number of Files

alice  atlas  cms  dirac  ilc  lhcb  mice  na62  snoplus  solid  t2k

# Castor to Antares Migration

- Actual migration – downtime Sunday pm to allow Castor to drain - to midday Thursday

- Backups of Castor taken and file lists created for VOs pre-migration – rollback checkpoint

- Two team members migrating all the VOs one at a time.  File lists in EOS produced for each VO to compare with the Castor file list --> small numbers of anomalies recorded and investigated!

- Had to apply dir extended attributes (ACLs and CTA workflows) on the _whole_ dir structure after migration

- Scripted applying across the whole namespace – ATLAS: 150,000 dirs, CMS: 535,000 dirs, LHCb: 177,000 dirs – and found that 70,000 directories was the maximum namespace size to apply the attributes without hitting the timeout limit
  - Required in a ~24 hours extension of the downtime

- Production traffic for some LHC VOs started on Friday afternoon

| MON Feb 21 | TUE 22 | WED 23 | THU 24 | FRI 25 | SAT 26 | SUN 27 |
|---|---|---|---|---|---|---|
| Antares downtime | | | | | | |
| | | | | | 7:30pm Castor downtime | |
| 28 | Mar 1 | 2 | 3 | 4 | 5 | 6 |
| Antares downtime | | | | | | |
| 7:30pm Castor downtime | | | | | | |
| | | | 12pm Antares production starts | | | |

# Antares in production

- The first 19 days in production have been dominated by the data challenge

- Otherwise, things have been fairly smooth, but busy!
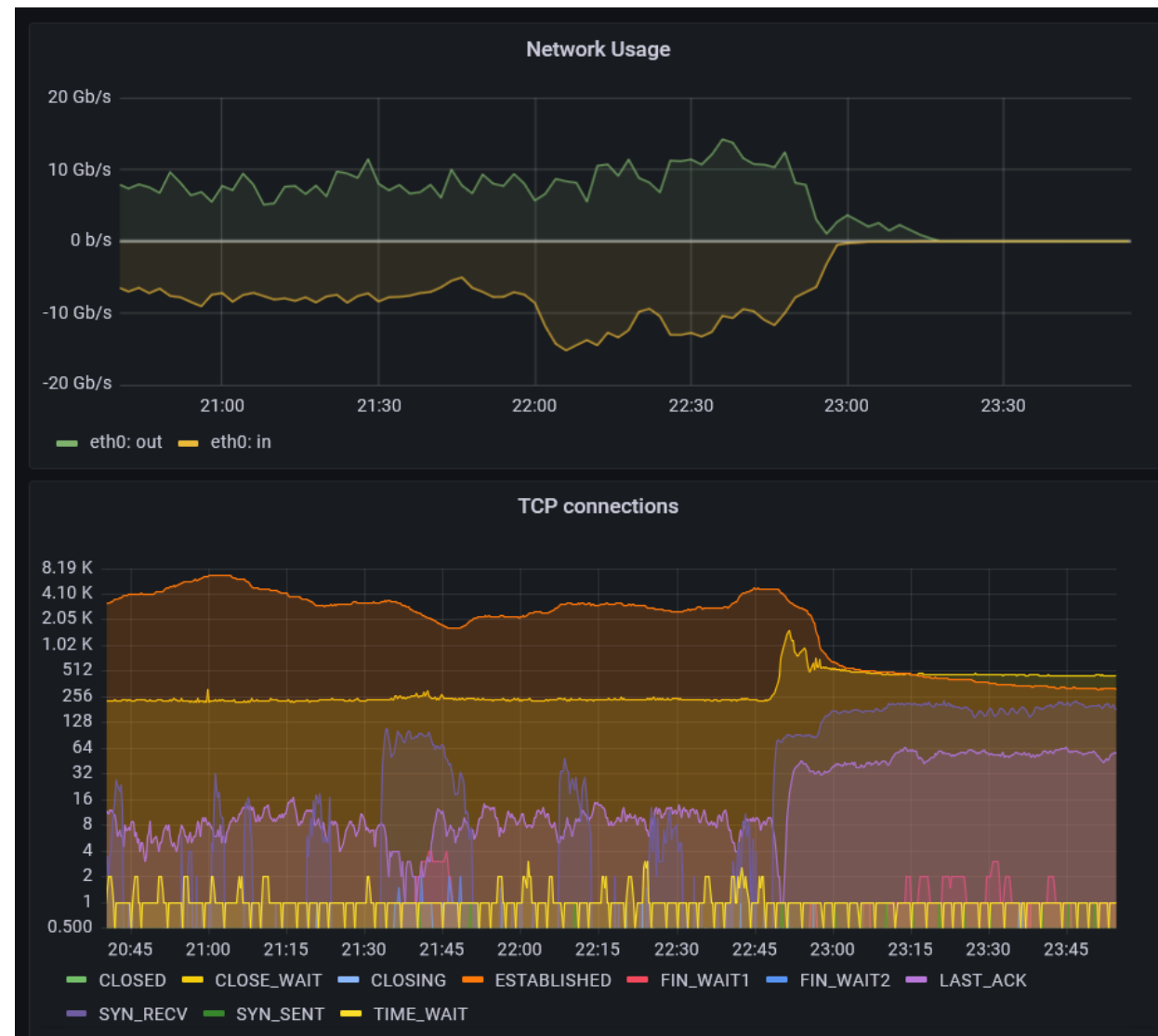
- A few issues to discuss...

# Issues in production – Low LTO9 tape drive rates

- We found that the write performance of the LTO9 drives appeared to be approximately half of the expected rate

- This was tracked down to very poor network performance to the LTO9 tape servers
  - thought to be the 'unusual' network cables that had been used (Dell rather than Mellanox), but switching cables did not change rates
  - Network (ring buffer size) tuning vastly improved network rates

- Now we have expected network rates to tape servers, testing is ongoing to determine if LTO9 drive performance is similarly improved

# Issues in production – EOS MGM stalls

- Occasional stalls of MGM node have been observed

- XRootD process still running happily, but all incoming connections fail

- Network config/tuning under scrutiny
  - Excessive packet discards seen in some cases

# 2022 LHC data challenge – operational perspective

- The presentation following this will cover how the current data challenge went from a VO perspective:
  - https://indico.cern.ch/event/1128343/contributions/4787155/

- I'd like to present a few thoughts on how Antares handled the challenge from our operational perspective

# 2022 LHC data challenge - archival
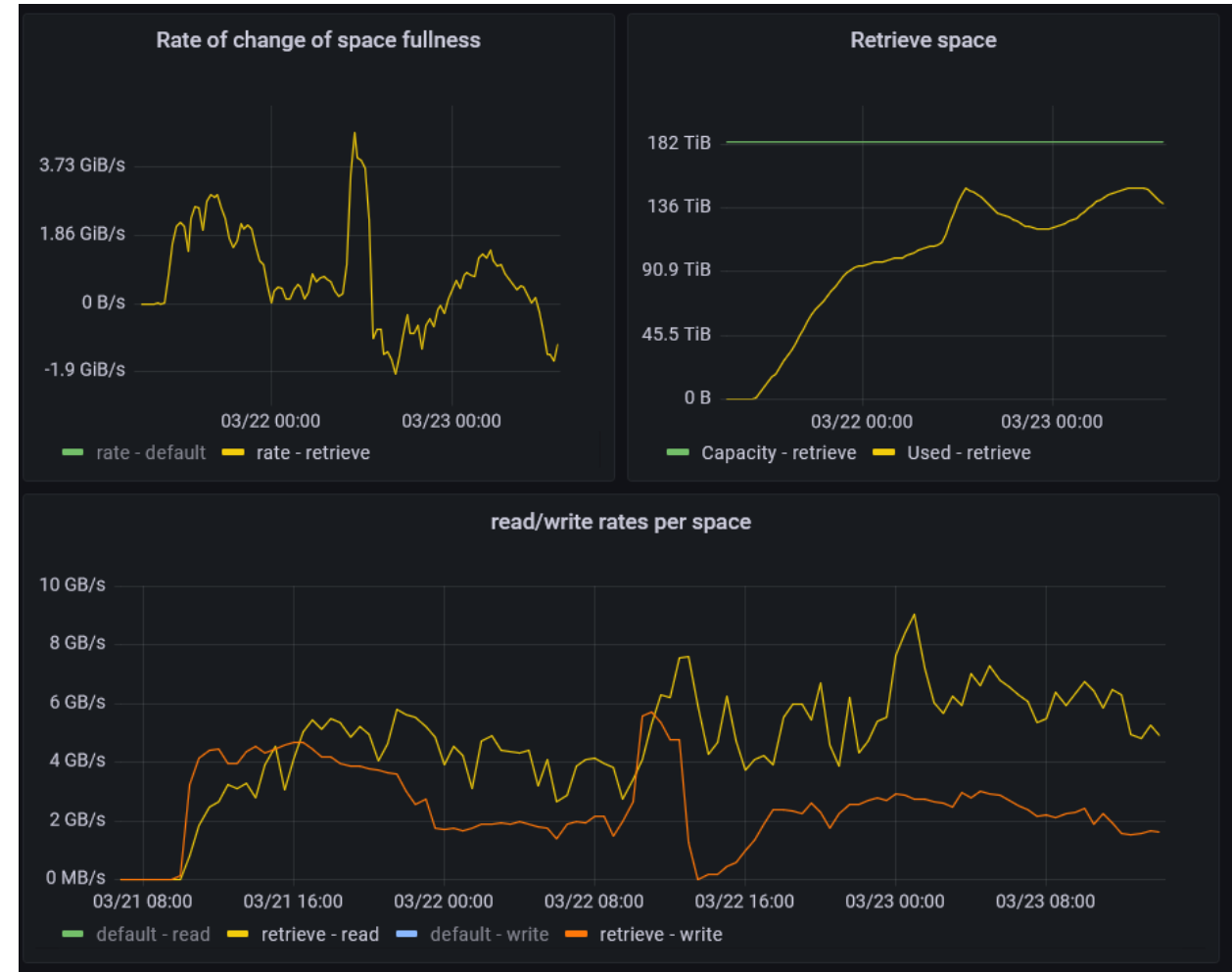
# 2022 LHC data challenge – retrieval

# 2022 LHC data challenge – managing retrieval EOS space

- Archival has been proven to work well with the small, fast buffer model

- Retrieve relies on external systems to monitor status, copy out and clean up files as they come online

  - Any misconfiguration can result in things going wrong very quickly

  - Balancing retrieval pressures between different VOs with different access methods will be an interesting challenge

# Next steps

- Ensure access for non LHC VO's is working as expected (ongoing)

- Upgrade EOS/CTA to the versions deployed at CERN

- Upgrade to EOS5

- Prepare and execute the migration of CASTOR Facilities

- Migrate the CTA Catalogue from Oracle to PostgreSQL

# Thanks

- Questions?