**Centre de Calcul**
de l'Institut National de Physique Nucléaire
et de Physique des Particules

# Tape Challenge 2021 @ CC-IN2P3
## *By Aresh Vedaee & Pierre Emmanuel Brinette (18.11.2021)*

- Introduction to the tape challenge 2021

- Results and assessments

- CMS vs ATLAS staging performance during A-DT

- Conclusions

- Tape Data Challenge orchestration:
  - Google doc for a better coordination between VOs & sites: https://docs.google.com/document/d/1rUhHdhlSgpU_Doam3Muox9XxnPQJEmqaf5_ZtzWsI5E/
    - TC Objective: *the validation of the **maximum tape bandwidth** needed for reads and writes to tier0 and tier1s tapes. This will imply a **realistic RUN 3 load** from DAQ systems, experiment T0 activity and exports. This test might include the validation of SRM-HTTP activity (TBC with Alessandra Forti).*
    - Timeline of the tape challenge week (October 11~15, 2021) for each VO:
      - Data Taking (DT) and After Data Taking (A-DT): during DT, migration dominated over staging, during A-DT vice versa
      - RUN3 target throughput per each VO, site, activity (migration & staging) and period (DT & A-DT)
      - VO test programme:
        - For ATLAS and CMS: 2 days of DT and 3 days of A-DT
        - For LHCb and ALICE: 2~5 days of DT
  - Site readiness: https://twiki.cern.ch/twiki/bin/view/LCG/TapeTestsPreparation
    - 16 T1 participated but only 4 sites (CERN,CNAF,CC-IN2P3,RAL) support all 4 VOs
  - FTS Dashboard to monitor all but Alice's activity: https://monit-grafana.cern.ch/d/e5o9PjDnz/fts-status-board-tape-challlenge-with-dt-write-and-a-dt-read-plots
  - "Tapetest" channel on Slack
  - Final report:
    - VOs: https://indico.cern.ch/event/1089983/ (spoiler: ATLAS & ALICE happier than CMS & LHCb!)
    - Sites:
      - https://indico.cern.ch/event/1092988/
      - https://indico.cern.ch/event/1094310/

- Accounting & assessing the TC results vs TC expectations is a tricky matter:
  - **Inconveniences on the TC orchestration:**
    - VO readiness & coordination (LHCb did not test a Run3 scenario, CMS A-DT Rucio/FTS issues + RAL not involved in DT, CMS challenge continued beyond TC timeline)
      - Data volume (not large/sustained enough for stage/migration DT & A-DT targets)
    - Activity focus (inclusion or exclusion of production activity, different targets per activity)
    - Timing (start/end dates global vs per site, on request submission vs done)
    - Others: CRIC tuning (some sites did it, others didn't; some VOs rely on it, others don't), monitoring (ALICE was not in the FTS dashboard)
  - **Inconveniences on the site side (i.e. CC-IN2P3):**
    - Tape accounting still diffucult and limited monitoring tools (e.g. throughput views per VO only for staging, drive usage views only for staging but not per VO)
    - No available list of datasets

# Introduction to the tape challenge 2021 (3/3)

- TC stats @CC-IN2P3 based on following guidelines:
    - Rough estimates based on broad assumptions
    - No distinction between TC and production activity
    - No CRIC tuning
    - ATLAS & CMS stats concern migration only for DT and staging only for A-DT
    - LHCb & ALICE stats concern migration only
    - Reference TC time table is the following

| Time (CEST) | ATLAS | | CMS | | LHCb | Alice |
|---|---|---|---|---|---|---|
| Start | DT | A-DT | DT | A-DT | 11.10.21 at 10h00 | 11.10.21 at 10h00 |
| | 11.10.21 at 10h00 | 13.10.21 at 10h00 | 11.10.21 at 10h00 | 12.10.21 at 22h00 | | |
| End | DT | A-DT | DT | A-DT | 13.10.21 at 23h00 | 15.10.21 at 18h00 |
| | 13.10.21 at 10h00 | 15.10.21 at 17h00 | 12.10.21 at 22h00 | 15.10.21 at 10h00 | | |

# Results and assessments (1/2)

**MIGRATION**

| VO | # Files | Volume | Avg Migration Rate | Avg file size |
|----|---------|--------|--------------------|--------------| 
| **ALL** | **277,035** | **1068 TB** | **2.31GB/s** | **3.85 GB** |
| ALICE | 111,412 | 233 TB | 622MB/s | 2.097 GB |
| ATLAS | 58,313 | 307 TB | 1.77GB/s | 5.29 GB |
| CMS | 21,947 | 235 TB | 1.8GB/s | 1.074 GB |
| LHCB | 44,798 | 216 TB | 983MB/s | 4.83 GB |
| Others Vos | 40,565 | 77 TB | 205MB/s | 1.91 GB |

**STAGING**

| VO | # Files | Volume | Stage Rate | Avg File Size |
|----|---------|--------|------------|---------------|
| **ALL** | **495,755** | **808.954 TB** | **1.79GB/s** | **1.6GB** |
| ALICE | 867 | 1.37 TB | 10.23 MB/s | 1.5GB |
| ATLAS | 415,704 | 550.028 TB | 1.222 GB/s | 1.3GB |
| CMS | 16,190 | 156.609 TB | 1.061 GB/s | 9.6GB |
| LHCB | 23,810 | 51.439 TB | 1.33 GB/s | 2.16GB |
| Others Vos | 39,184 | 49.508 TB | 112.65 MB/s | 1.2GB |

**MIGRATION vs STAGING
October 11~15, 2021
(CC-IN2P3 view and no distinction
between DT and A-DT)**

**LHC VOs' VIEW**

| Throughput (GB/s) | ATLAS [1] | | CMS [1] | | LHCb [1] | ALICE [2] |
|---|---|---|---|---|---|---|
| | $DT_{(w)}$ | $A\text{-}DT_{(r)}$ | $DT_{(w)}$ | $A\text{-}DT_{(r)}$ | $DT_{(w)}$ | $DT_{(w)}$ |
| Target | 1.4 | 1.2 | 0.9 | 1.5 | 1.26 | 0.4 |
| AVG | 1.52 | 0.89 | 0.57 | 1.22 | 0.78 | 0.54 |
| MAX | 2.22 | 2.63 | 2.35 | 6.73 | 2.42 | |

**Performance gaps**

**discrepancies**

**CC-IN2P3 VIEW**

| Throughput (GB/s) | ATLAS | | CMS | | LHCb | ALICE |
|---|---|---|---|---|---|---|
| | $DT_{(w)}$ | $A\text{-}DT_{(r)}$ | $DT_{(w)}$ | $A\text{-}DT_{(r)}$ | $DT_{(w)}$ | $DT_{(w)}$ |
| Target | 1.4 | 1.2 | 0.9 | 1.5 | 1.26 | 0.4 |
| AVG | 2.45 | 1.02 | 0.65 | 3.2 | 1.0 | 0.59 |
| MAX | - | 2.7 | - | 5.0 | | |

[1] https://monit-grafana.cern.ch/d/e5o9PjDnz/fts-status-board-tape-challlenge-with-dt-write-and-a-dt-read-plots?from=1633903200000&orgId=20&to=1634335199000&var-activity=All&var-bin=1h&var-dst_country=All&var-dst_rse=All&var-dst_site=All&var-dst_tier=All&var-fts_server=All&var-group_by=vo&var-protocol=All&var-src_country=All&var-src_experiment_site=All&var-src_rse=All&var-src_site=All&var-src_tier=All&var-staging=All&var-vo=All

[2] https://indico.cern.ch/event/1089983/contributions/4581916/attachments/2335897/3981407/ALICE%20custodial%20storage%20challenge%20-%20results.pdf

## VOs VIEW vs SITE VIEW (per VO & period)

- Assessment based on VOs' view:
  - ALICE target: done
  - LHCb target: not done but not site-related (FTS knobs and EOS gridftp gateways)
  - ATLAS targets: DT done but not A-DT
  - CMS targets: neither DT nor A-DT but A-DT staging throughput better than ATLAS, why? (pattern already noticed during past ATLAS tape stress test)

- Comparison between VO's view vs CC-IN2P3 view:
  - ALICE view vs CC-IN2P3 view stats match
  - LHCb view vs CC-IN2P3 view do not match but not by far (maybe only matter of time window)
  - ATLAS & CMS view vs CC-IN2P3 view: migration/staging stats do not match, and A-DT avg staging throughput from CC-IN2P3 >> avg throughput from FTS dashboard, why? Are we watching the same time series? Indeed as CMS reported, CC-IN2P3 avg throughput is ranked first for staging and migrating (in both cases > 4GB/s). Or maybe FTS dashboard does not count the staging failures (as already noticed also in the past ATLAS Tape Stress Test)?
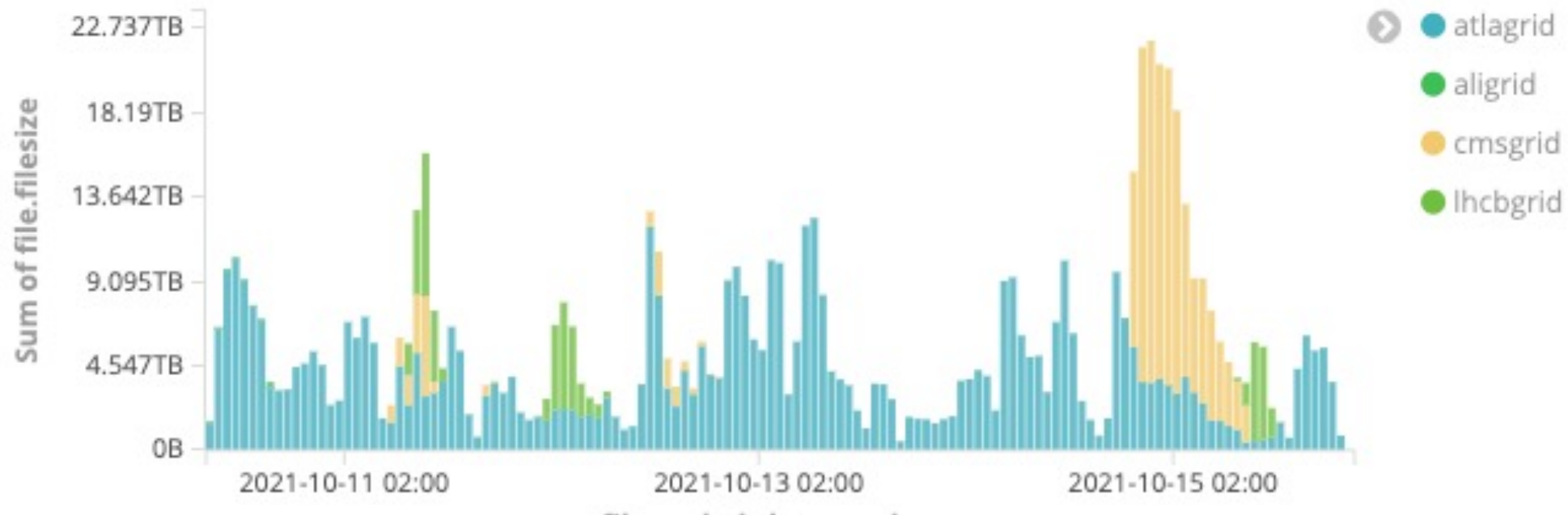
## CC-IN2P3 setup

- Tape/drive resources are shared by all VOs (LHC and non-LHC)
- HPSS T10K-D Media migration (repacks) suspended during the TC
- HPSS Staging Configuration (based on TREQS staging scheduler):
  - Jaguar-E/TS1160 (Drive nominal speed 450MB/s)
    - **46 drives available for staging and migration:** staging scheduler (TREQS) requests max 32 drives at each staging pass (so max 14 drives are left for migrations)
  - T10KD (Drive nominal speed 240MB/s)
    - **48 drives only for staging**
  - Pending time for staging requests set up to minimum 4min (but there is no max and the staging file can be served after hours)
  - Migration cycle is every 6h (it applies to files written > 2h ago)
  - File size class setup (more relevant than file family): it determines the number of drives used on migration. For LHC VOs:
    - COS 12 (64MB - 2GB): 5 drives
    - COS 14 (> 2GB): 6 drives

- On migration:
  - All VOs compete for available drives (soft-limited to 14) and with same technology (namely Jaguar-E)
  - File size is also relevant due to HPSS file class drive distribution (the lower the file size the less # of drives)
    - Compared to other VOs, the # files in COS12 for CMS is negligible
- On staging:
  - all VOs compete for 80 drives on staging, BUT with 2 different technologies (so throughput depends on the data distribution)
  - There are less Jaguar-E drives available on staging than for T10K-D
- Additional info:
  - We'll see later that there is more competition for Jaguar-E drives than for T10KD drives
  - Jaguar-E drives underperform on staging wrt migration performance
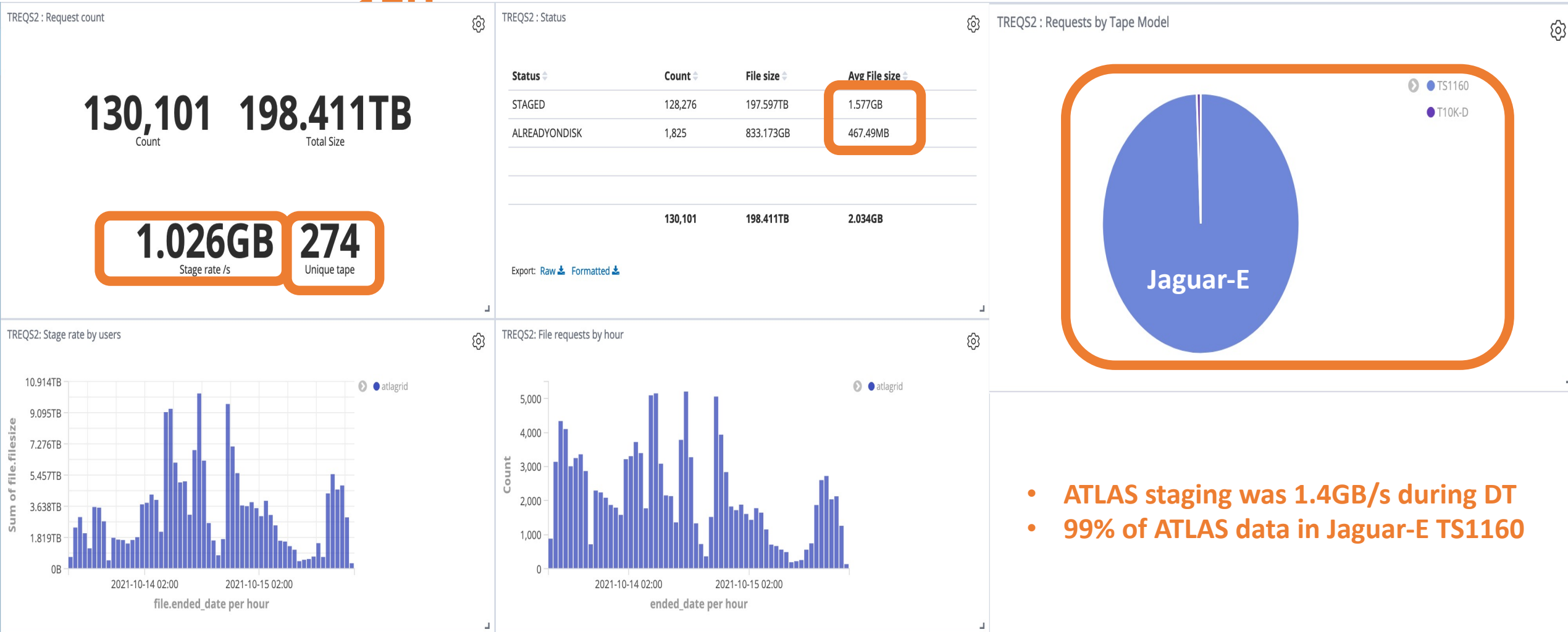
**STAGING activity during the Tape Challenge**



TREQS2: Stage rate by users  (TB/h)

- **ATLAS staging was continuous during TC with peaks also during DT**
- **CMS staging was more concentrated especially during last day of A-DT**

**ATLAS staging A-DT: Oct 13, 10h – Oct 15,**



TREQS2 : Request count

**130,101** Count   **198.411TB** Total Size

**1.026GB** Stage rate /s   **274** Unique tape

TREQS2 : Status

| Status | Count | File size | Avg File size |
|--------|-------|-----------|---------------|
| STAGED | 128,276 | 197.597TB | 1.577GB |
| ALREADYONDISK | 1,825 | 833.173GB | 467.49MB |
| | 130,101 | 198.411TB | 2.034GB |

Export: Raw  Formatted

TREQS2 : Requests by Tape Model

TS1160
T10K-D

Jaguar-E

TREQS2: Stage rate by users

atlagrid

TREQS2: File requests by hour

atlagrid

- **ATLAS staging was 1.4GB/s during DT**
- **99% of ATLAS data in Jaguar-E TS1160**

## CMS staging A-DT: Oct 12, 22h – Oct 15, 10H

TREQS2 : Request count

**10,623** **135.593TB**
Count      Total Size

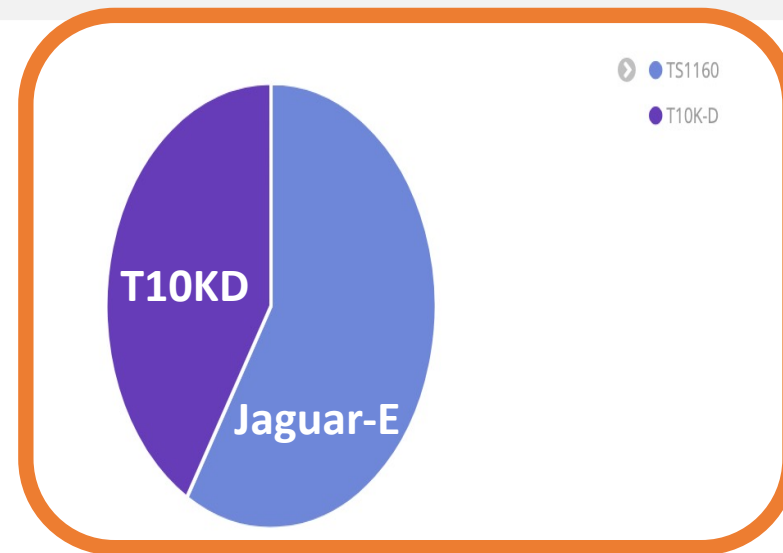**3.214GB** **123**
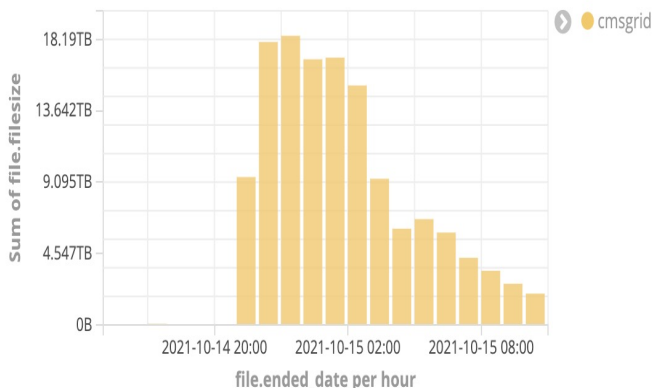Stage rate /s   Unique tape

TREQS2 : Status

| | | | |
|---|---|---|---|
| STAGED | 10,623 | 135.593TB | 13.07GB |
| | | | |
| | **10,623** | **135.593TB** | **13.07GB** |

Export: Raw ⬇ Formatted ⬇

TREQS2 : Requests by Tape Model



● TS1160
● T10K-D

T10KD

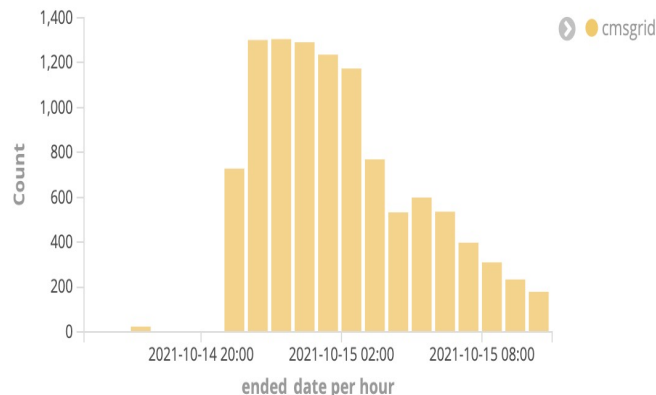Jaguar-E

TREQS2: Stage rate by users
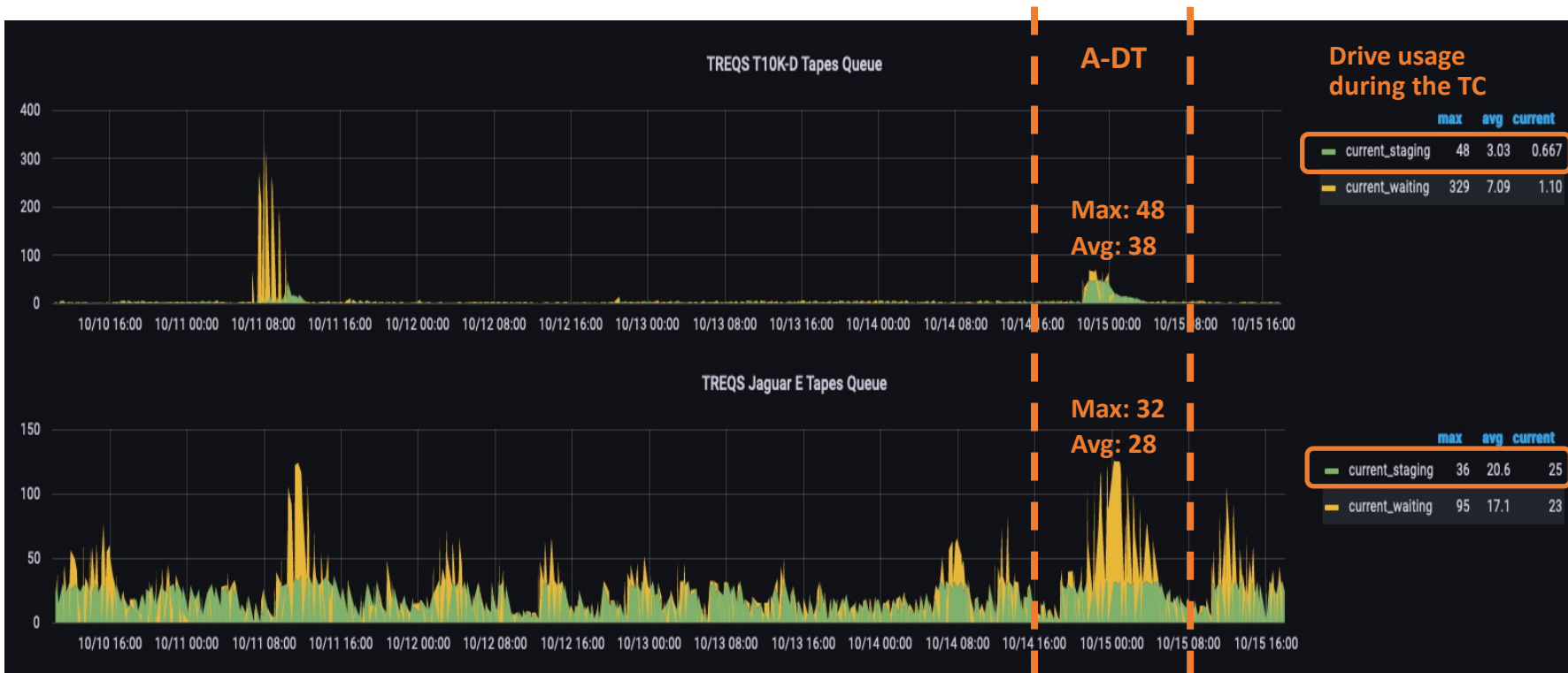


TREQS2: File requests by hour



- # CMS files and unique tapes << ATLAS
- CMS avg file size >> ATLAS avg file size
- 60% of CMS data in Jaguar-E TS1160
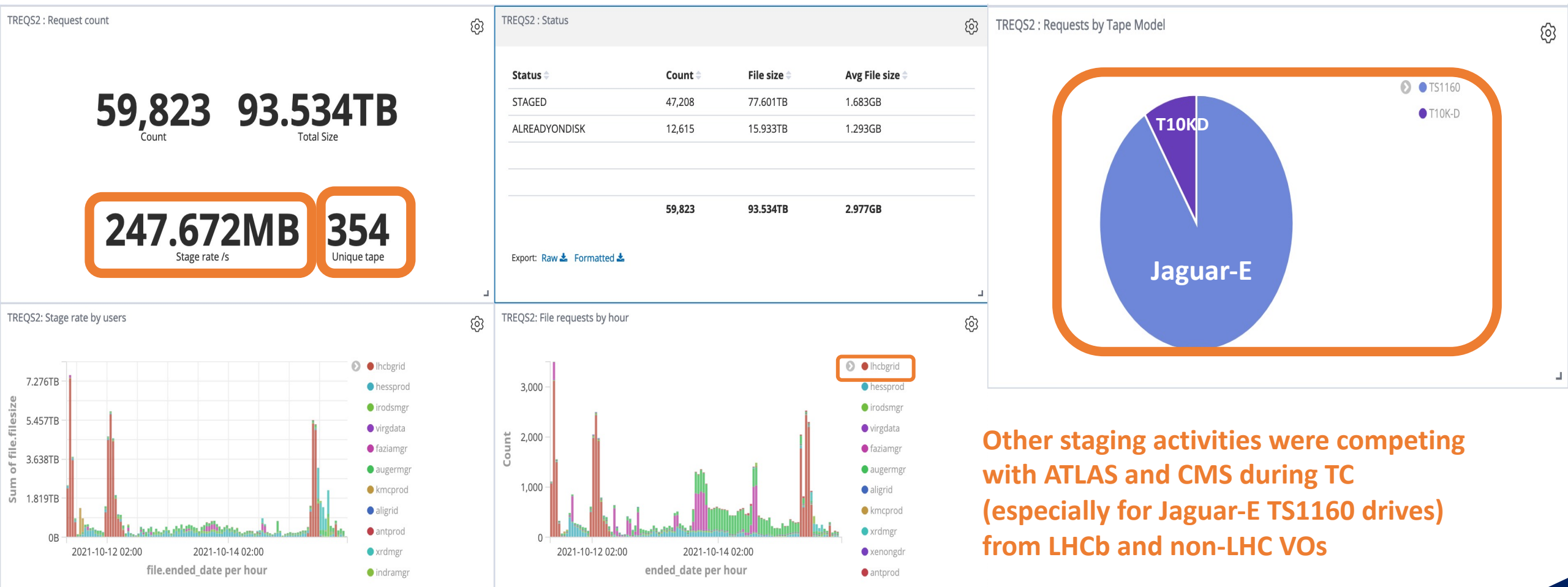
## STAGING activity during the Tape Challenge



| VO | # uniq tapes (TS1160) | # uniq tapes (10KD) |
|---|---|---|
| **ALL** | **915** | **553** |
| ALICE | 10 | 1 |
| ATLAS | 517 | 9 |
| CMS | 115 | 447 |
| LHCB | 71 | 2 |
| Others Vos | 202 | 94 |

- At a closer look to A-DT, T10KD drive max & avg usage >> Jaguar-E max & avg usage, besides Jaguar-E max usage < 36 (36 is the max available)
  - 48 max used drives during A-DT for everyone's staging activity (including CMS)
  - 32 max used drives during A-DT practically only for CMS staging activity
- The used tape stats show that most of CMS data was in T10KD tapes

18/11/2021

## STAGING COMPETITORS: Oct 11, 10h – Oct 15, 17H



**TREQS2 : Request count**

**59,823** Count    **93.534TB** Total Size

**247.672MB** Stage rate /s    **354** Unique tape

**TREQS2 : Status**

| Status | Count | File size | Avg File size |
|---|---|---|---|
| STAGED | 47,208 | 77.601TB | 1.683GB |
| ALREADYONDISK | 12,615 | 15.933TB | 1.293GB |
| | 59,823 | 93.534TB | 2.977GB |

Export: Raw ⬇ Formatted ⬇

**TREQS2 : Requests by Tape Model**

TS1160 · T10K-D · T10KD · Jaguar-E

**TREQS2: Stage rate by users**

lhcbgrid, hessprod, irodsmgr, virgdata, faziamgr, augermgr, kmcprod, aligrid, antprod, xrdmgr, indramgr

**TREQS2: File requests by hour**

lhcbgrid, hessprod, irodsmgr, virgdata, faziamgr, augermgr, aligrid, kmcprod, xrdmgr, xenongdr, antprod

**Other staging activities were competing with ATLAS and CMS during TC (especially for Jaguar-E TS1160 drives) from LHCb and non-LHC VOs**

- Recall performances on Jaguar E/TS1160 lower than expected (affecting probably more ATLAS than CMS staging throughputs), why?

- We enabled file aggregation on tape for large file (>1GB) in order to improve write performances on TS1160.
  - Aggregation : HPSS feature that aggregates multiple files (up to 50) on a single tape segment.

- During the tape challenge, we noticed that the tape drive position time is greater than expected when reading file within the aggregate.
  - Tape drive positions itself at the beginning of the aggregate segment, then it reads the whole segment until reaching the requested file
  - Fast positioning feature (i.e. Tape Order Recall) is not used when reading files from an aggregate.

- Problem under investigation :
  - Workaround : enable Full Aggregate Recall (after migration to HPSS 8.3 in december)
  - The bug should be definitely fixed in HPSS v9.3 (feature CR 521)

- Accounting and assessing the tape challenge is tricky wrt the TC goal
  - "Shared resources" implies great competition across activities and VOs (LHC and non-LHC) as a suitable TC orchestration could make evident
  - Still some bottlenecks outside T1's perimeter (EOS grdiftp gateways, FTS/Rucio miscommunications)
- CMS staging performance during A-DT better than ATLAS b/c
  - Bigger files (good for both migration and staging)
  - Less competing activities from other VOs
  - Less scattered across tapes
  - Better data distribution across drive sets (more staging drives for CMS and underperforming drives for ATLAS)
  - Déjà vu from the past ATLAS tape stress test

# Merci!

**BACKUP**