# Tape operations & experience

*R. Tafirout, TRIUMF*

*with contributions from Tier-1s:*

*BNL (David Yu)*
*CCIN2P3 (Emmanouil Vamvakopoulos)*
*CNAF (Vladimir Sapunenko)*
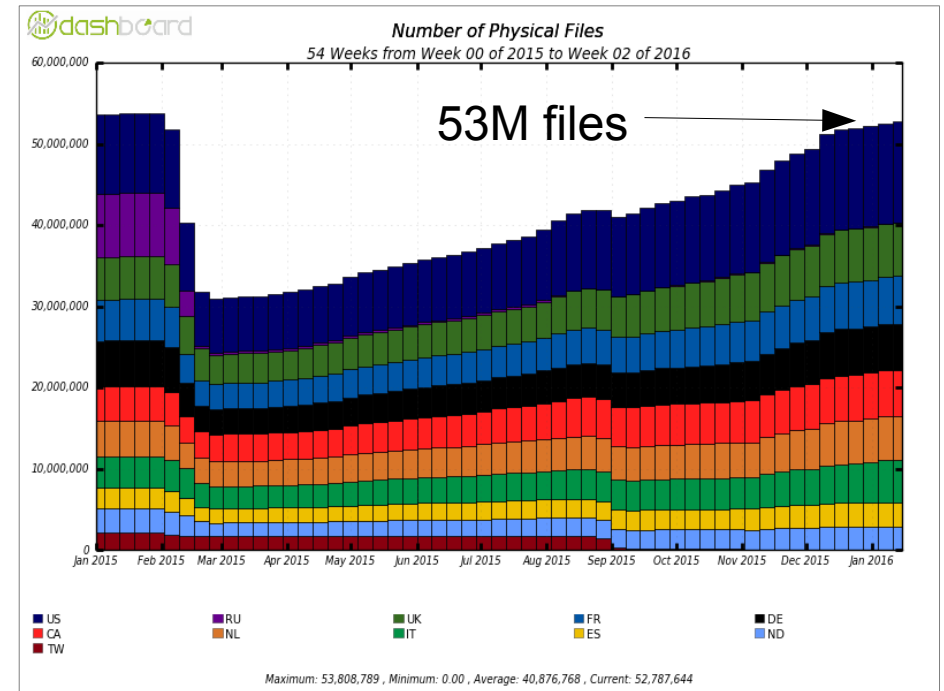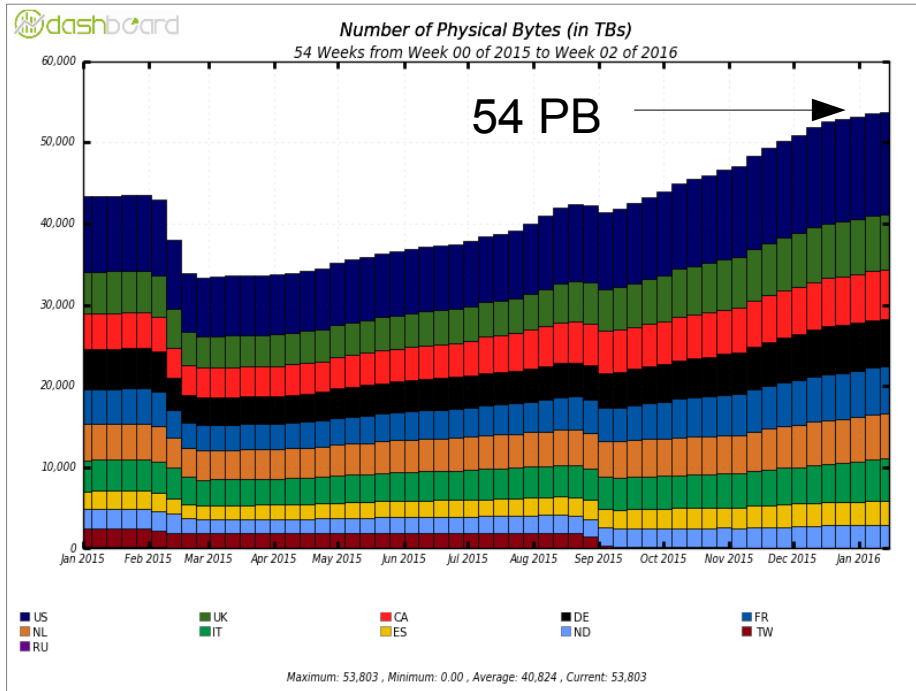*KIT (Xavier Mol)*
*NDGF (Ulf Tigerstedt)*
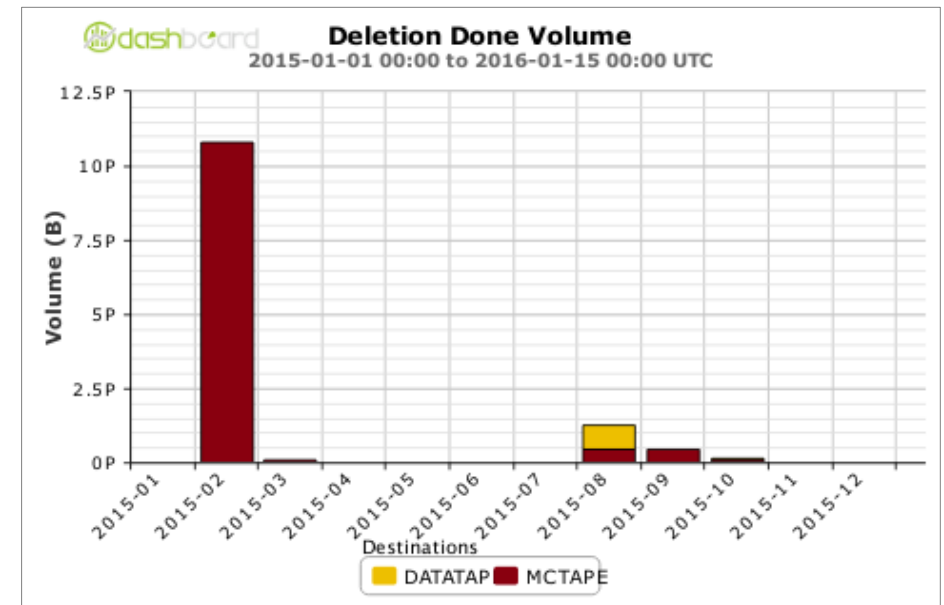*PIC (Josep Flix, Esther Accion, Aresh Vedaee)*
*RAL (Brian Davies)*
*SARA (Onno Zweers)*
*TRIUMF (Simon Liu, Yun-Ha Shin)*

# Overall Tape Activities (ADC)
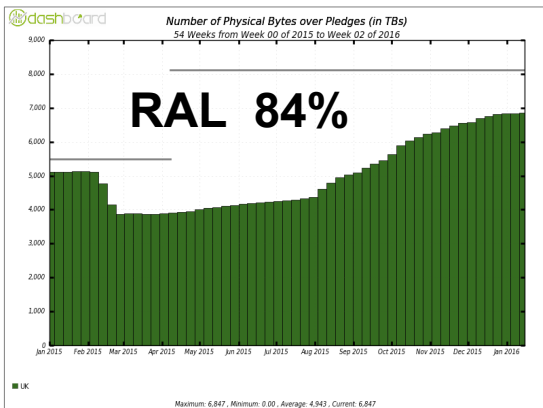


Number of Physical Bytes (in TBs)
54 Weeks from Week 00 of 2015 to Week 02 of 2016

54 PB

Maximum: 53,803 , Minimum: 0.00 , Average: 40,824 , Current: 53,803



Number of Physical Files
54 Weeks from Week 00 of 2015 to Week 02 of 2016

53M files

Maximum: 53,808,789 , Minimum: 0.00 , Average: 40,876,768 , Current: 52,787,644

- Massive deletion campaign in February
- ASGC (TW) decommissioned, no data hosted since ~September
- No tape activity at NRC-KI (RU)
- Includes grouptape (SUSY + TOP) at CNAF, NDGF, SARA, TRIUMF.



Deletion Done Volume
2015-01-01 00:00 to 2016-01-15 00:00 UTC

Destinations: DATATAP, MCTAPE

2

# Tape Usage versus Pledged

# Tape Systems Specs

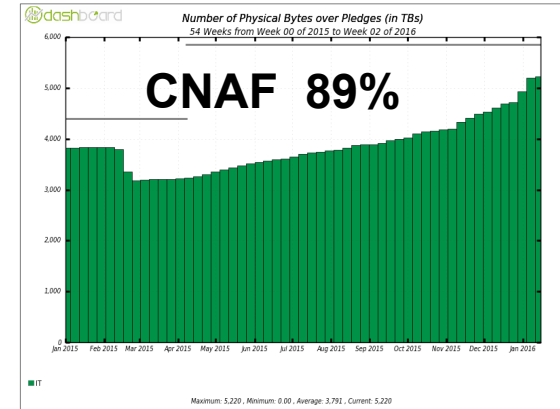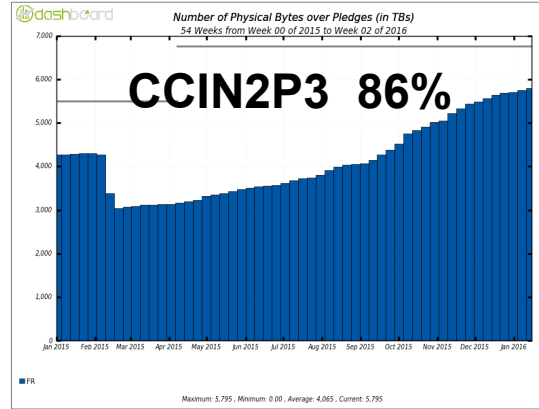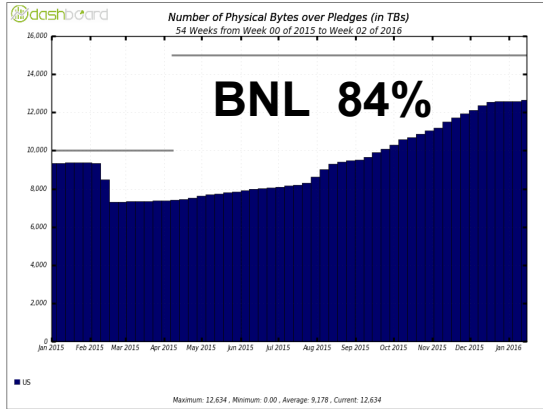| | MoU share (%) | Library | Software | Drives | Disk buffer |
|---|---|---|---|---|---|
| BNL | 23 | 4 x SL8500 | HPSS | 30xLTO4 , 23xLTO6 | 209 TB |
| CCIN2P3 | 10 | 4 x SL8500 | HPSS | T10K: 28xB,22xC,43xD | 120 TB |
| CNAF | 9 | 1 x SL8500 | TSM | T10K: 17xD | ~220 TB (dynamic) |
| KIT | 13 | Multiple / mixed | TSM + ERMM | mix of LTO and T10K | 350 TB |
| NDGF | 5 | 4 sites / 3 countries | TSM | - | ~160 TB |
| PIC | 5 | 1 x SL8500 | Enstore | 14xLTO4, 4xLTO5 , 8xT10KC, 5xT10KD | 60 TB |
| RAL | 13 | 1 x SL8500 | Castor | 14 x T10KC | 360 TB |
| SARA | 9 | 2 x SL8500 | DMF | T10K: 16xB, 8xC, 13xD | 84 TB |
| TRIUMF | 10 | 1 x TS3500 | Tapeguy | 14 x LTO5 , 10 LTO6 | 400 TB |

- BNL & TRIUMF: ATLAS only

- Disk buffer: includes W + R.

**Drive & Media specs (native):**
LTO-4,5,6 (R/W) = 120 , 140, 160 MB/s
LTO-4,5,6 (capacity) = 0.8, 1.5, 2.5 TB
T10K B,C,D (R/W) = 120, 240, 250 MB/s
T10K B,C,D (capacity)= 1, 5, 8.5 TB

# Tier-1's serving multiple VO's (I)

**Drive access rules for sharing, prioritization & scheduling:**

- **CCIN2P3:**

| Storage class definition for Atlas: | | | |
|---|---|---|---|
| Small Files | → 0-64M | → Titanium 10000B | x 2 drives (for write) |
| Medium | → 64MB – 512M | → Titanium 10000C | x 2 drives (for write) |
| Big Files | → 512MB – 2GB | → Titanium 10000D | x 5 drives (for write) |
| XL Files | → 2GB – 4TB | → Titanium 10000D | x 6 drives (for write) |

- **CNAF:**
  - 8 drives maximum (4 for writing & 4 for recalls), with a global overbooking factor taking into account all VO's.

- **KIT:**
  - at most 10 drives (6 for writing & 4 for recalls), adjusted often for backlogs depending on other VO's activities. Disk reading buffer for ATLAS dedicated (150 TB).
  - tape resources are shared: first-come-first-serve, no way to guarantee a certain number of drives per VO.

# Tier-1's serving multiple VO's (II)

**Drive access rules for sharing, prioritization & scheduling:**

- **NDGF:**
  - no detailed info about drives
  - 4 different sites / 3 countries: each site has one pool for reading and one for writing per VO; each pool size 8-30 TB.

- **PIC:**
  - each VO assigned to a specific tape technology
  - T10KC used for ATLAS with a dedicated disk buffer of 60 TB.
  - each VO can use up to 2 drives per tape family.
  - for reading all drives can be used but system does some balancing.

- **RAL:**
  - 1 T10KC dedicated to ATLAS, the remaining 12 drives shared with LHCb (which has also 1 dedicated drive); no weighting (first come first served)

- **SARA:**
  - Two libraries : T10KC used for ATLAS and up to 5 drives per library.

# Site issues & concerns (I)

- **BNL**: small files is the biggest issue (dragging system performance down), need to watch out for work done near tape library (dust control); unscheduled system maintenance handled carefully to minimize downtime.

- **CCIN2P3:** incident in October (448 TB received in 7 days, ~776 MB/), write buffers filled up; issue fixed by increasing # of drives for migration to tape; system can handle now ~1.2 GB/s)

- **CNAF**: small files slowing down migration and recalls, any file aggregation possible ?. One tape damaged in January 2015, 12 files lost.

- **KIT**: bad file size distribution (5.7 M files registered with average size of 1 GB ± 1.1 GB standard deviation); writing small files to tape not necessarily an issue, but recalls are very slow.

   - issue about tape families: "ATLAS doesn't make use of different tape families" ... "data sets are spread wildly over tape cartridges".

   *(*needs clarification/discussion with ADC...)*

- **NDGF**: have not had any problems with ATLAS. Access through ARC and aCT (more controlled I/O access).

# Site issues & concerns (II)

- **PIC**: no problems observed related to scalability. Tape family via GGUS?
  - usual patterns problems - not observed or affected the site
  - Currently suffering from some tape media integrity for T10KD technology, in contact with Oracle (media taken out of production); files will need a recovery procedure (details to be communicated to ADC asap).

- **RAL:** issues with tape recall policy/algorithm (initially set at 500 files or 32 GB; 32 GB was not working due to timeout; changed to 10 files.
  - Policy to improve drive performance leads to latency: recalls only triggered after 10GB/10files/1hr ; migration only triggered after 100GB/500 files/2hrs (per tape pool).
  - some files are being recalled many times increasing load and churn rate.
  - bringonline request : what lifetime should be aimed for disk buffer ?
  - disk buffer current bottle neck: new hardware being added.
  - being able to control FTS transfers to tape to improve WAN.

- **SARA**: broken pool in Jan. '15 / data loss; lots of SRM timeout issues since dCache upgrade in March,  issue is now understood.

- **TRIUMF**: no major issues, some minor issues with library; some 10G card issues with HSM pools (cards replaced recently).  Lots of 1 file datasets:
  - Various issues with SUSY group migration (very bursty), FTS damping ?
  - lots of SUSY datasets: 105k vs 50k (datatape + mctape).

# Planned & Tentative System upgrades

- **BNL**: migration to LTO-7 generation, no schedule yet

- **CCIN2P3**: TBD

- **CNAF**: move HSM servers from FC8 to FC16 (February); disk storage system replacement & data migration to new storage (February).

- **KIT**: migration to HPSS as tape management software

- **NDGF**: Not known at the moment

- **PIC**: getting T10KD in production again, finish migration from LTO-4, Enstore upgrade in conjunction with a dCache upgrade.

- **RAL**: move to T10KD media; SRM upgrade & SL6; Castor upgrade to 2.1.15; investigating disk cache with CEPH pool; merging Castor instances into one for WLCG experiments; mainly following CERN's advice.

- **SARA**: <span style="color:red">**major downtime in the fall (moving to another data centre), careful planning with ADC a must / potential data loss at stake**</span>.

- **TRIUMF**: readiness for 2016 pledges / media replacement for 2000 LTO-4 to LTO-6 (migration ongoing); upgrade of 8 LTO-5 to LTO-6 (in February). Tape system software upgrade for various improvements and bug fixes.
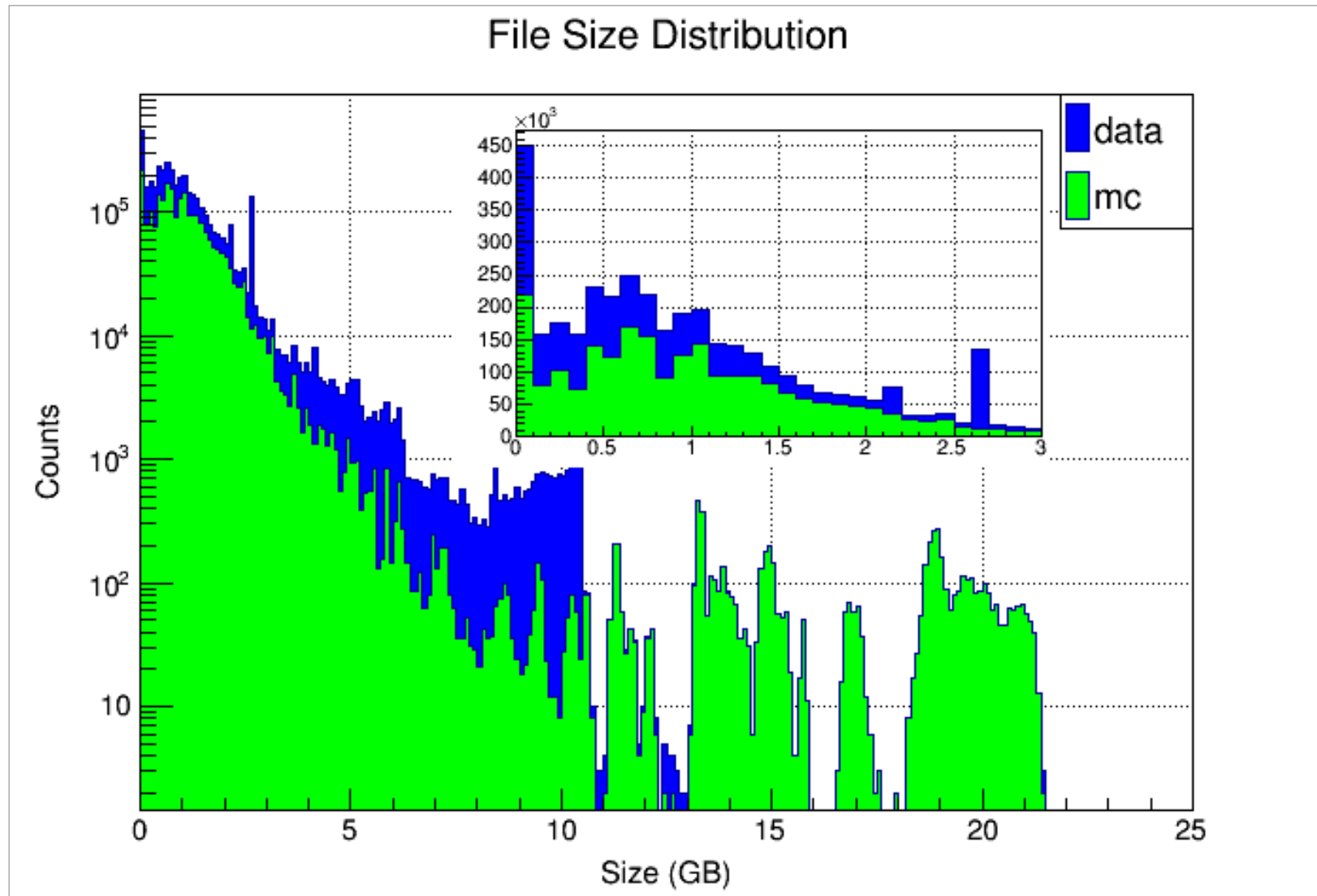
# Other aspects for discussion

- **Various sites have their own monitoring, metrics, etc.**
  - ATLAS tape activity is well monitored and tracked at the sites.
- **Can sites handle more ATLAS activity ?**
  - Based on the information received, it looks like most Tier-1 sites have no major issue beside what was discussed in earlier slides.
  - With both ADC and sites tweaks or tuning, more activity could be handled in principle; drives capacity seems to indicate there is more room.
  - Bringonline / tape recalls strategy
- **How is file deletion handled and tape space reclaimed ?**
  - It is not clear to me which Tier-1 sites have already reclaimed the deleted space; perhaps done automatically (only asked a few), and what strategy will be adopted.
  - Unless space is needed urgently, this is in principle handled automatically when doing media migration / technology refresh.

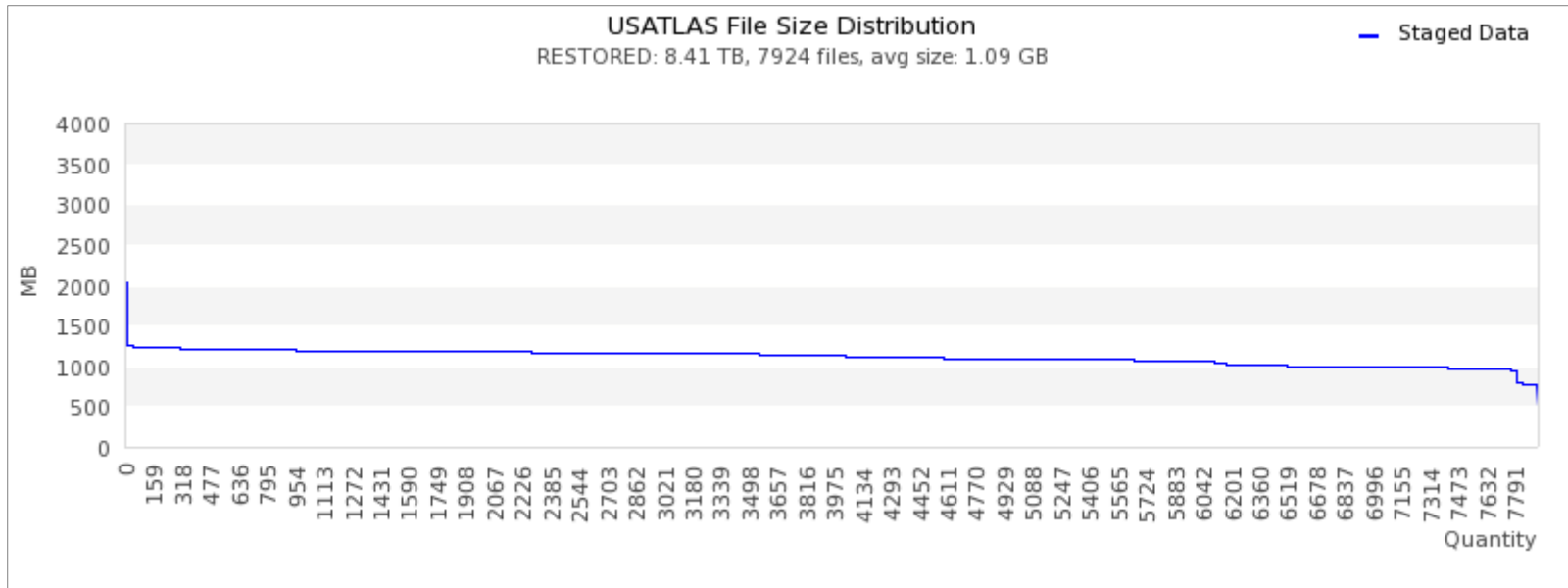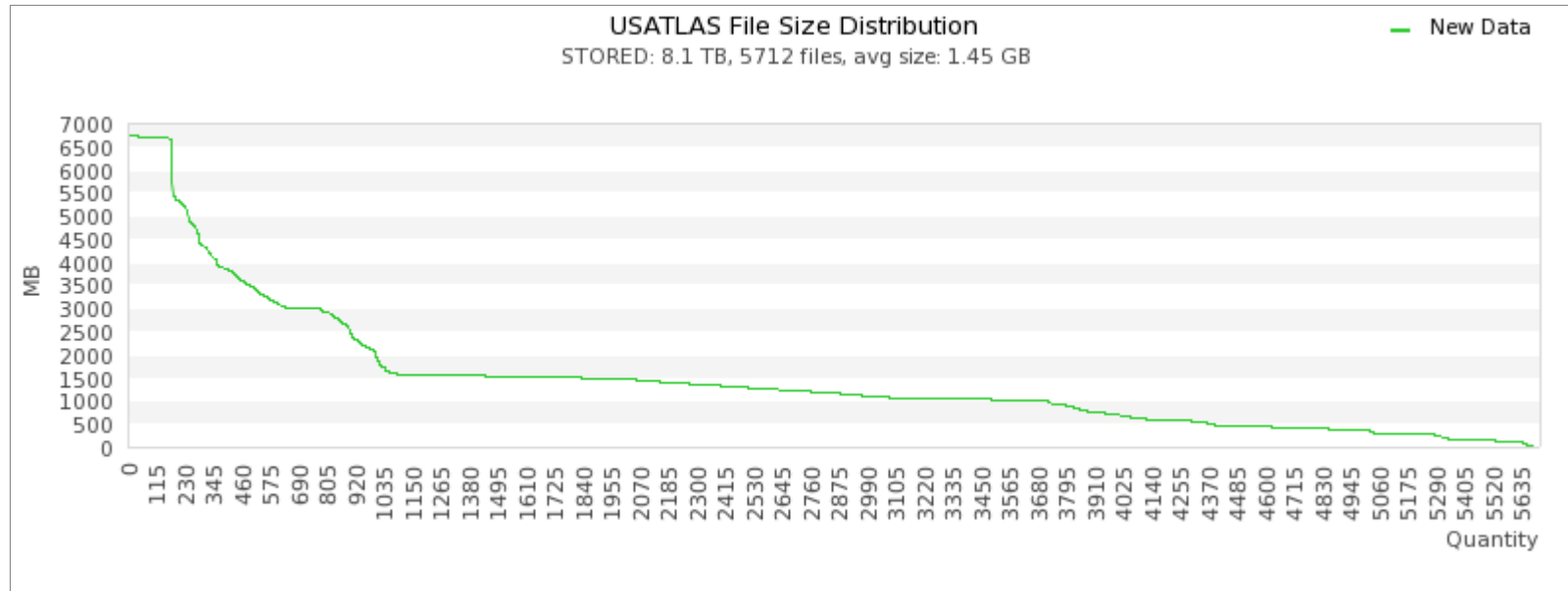# EXTRA MATERIAL
## (monitoring & stats plots)

## *(a small sample from all received material...)*
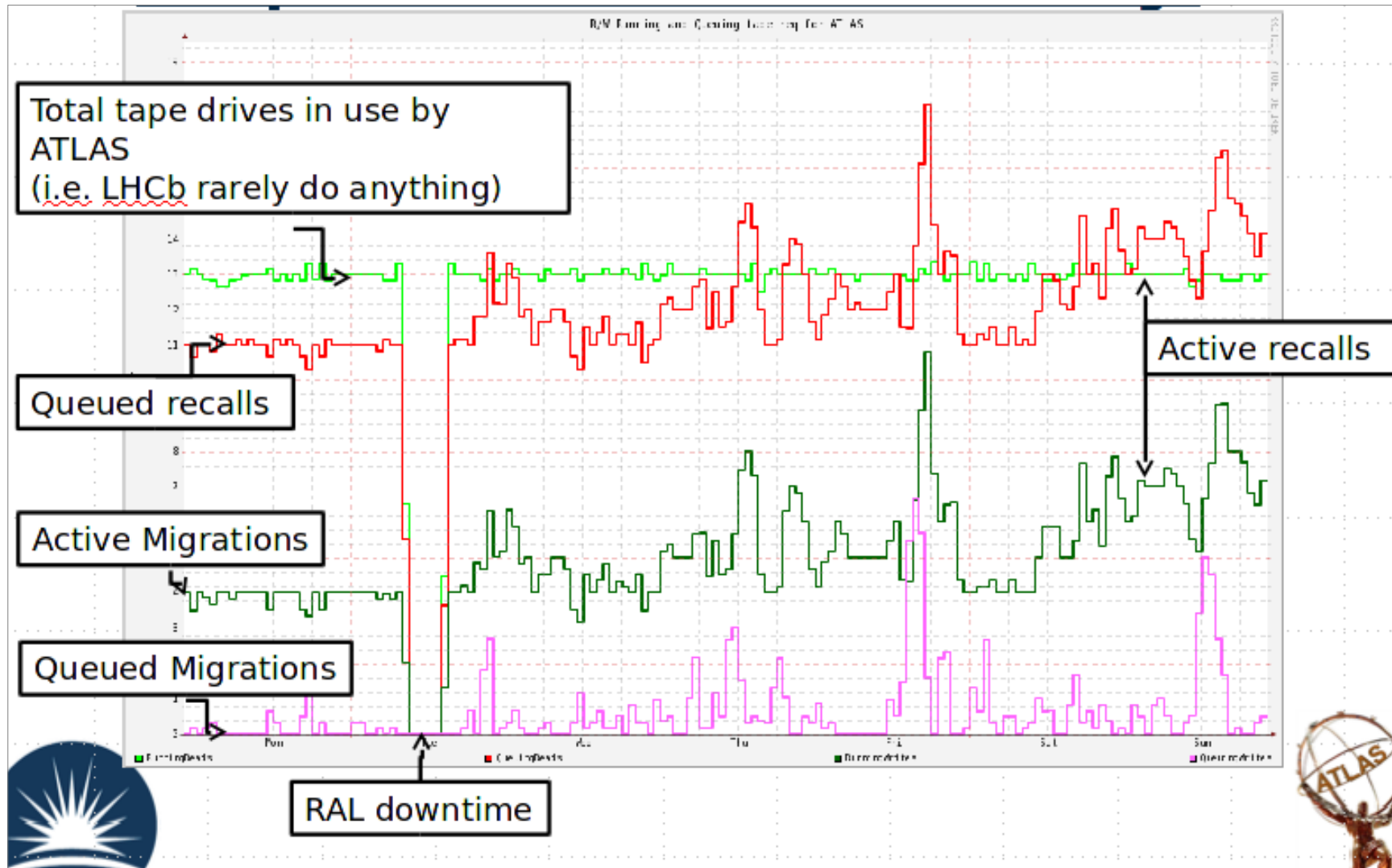
# File size distribution @ TRIUMF

# File size distribution @ BNL

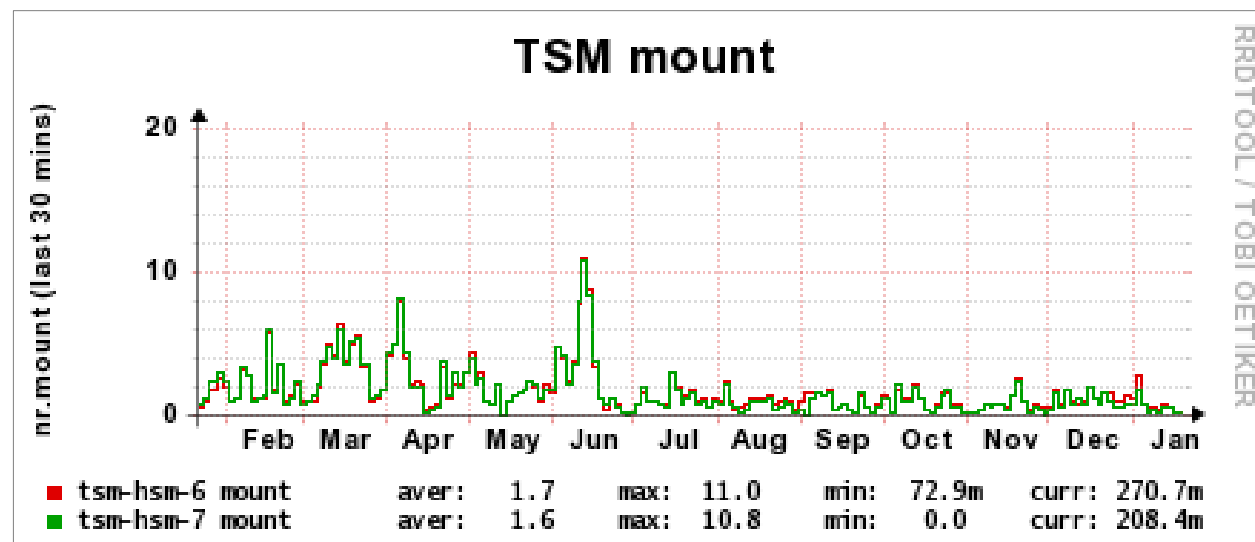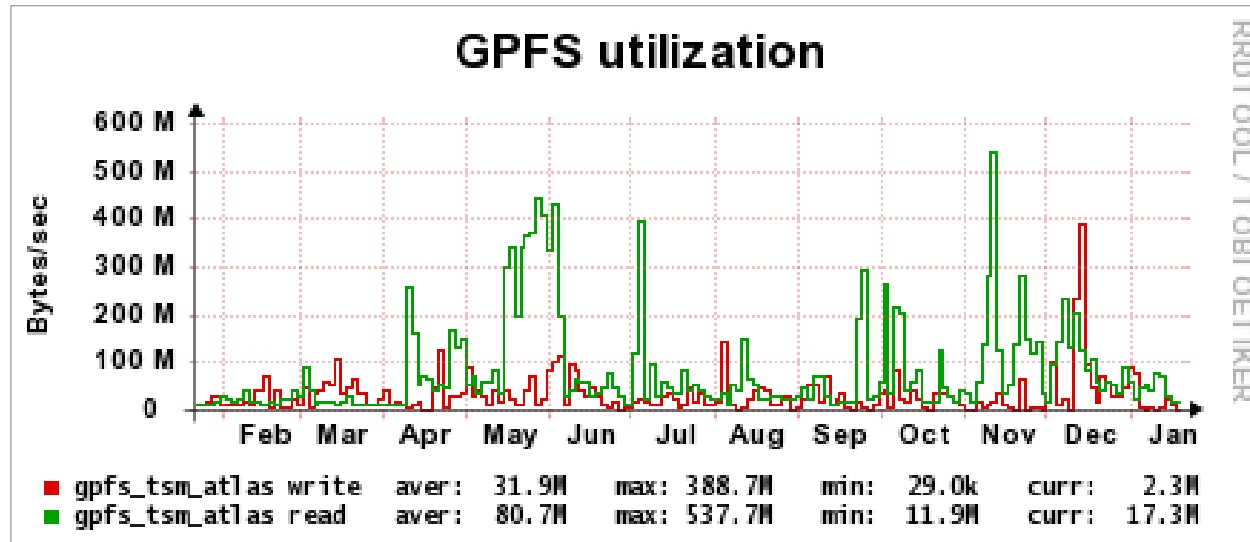- 24 hours activity sampling for stored and staged data:
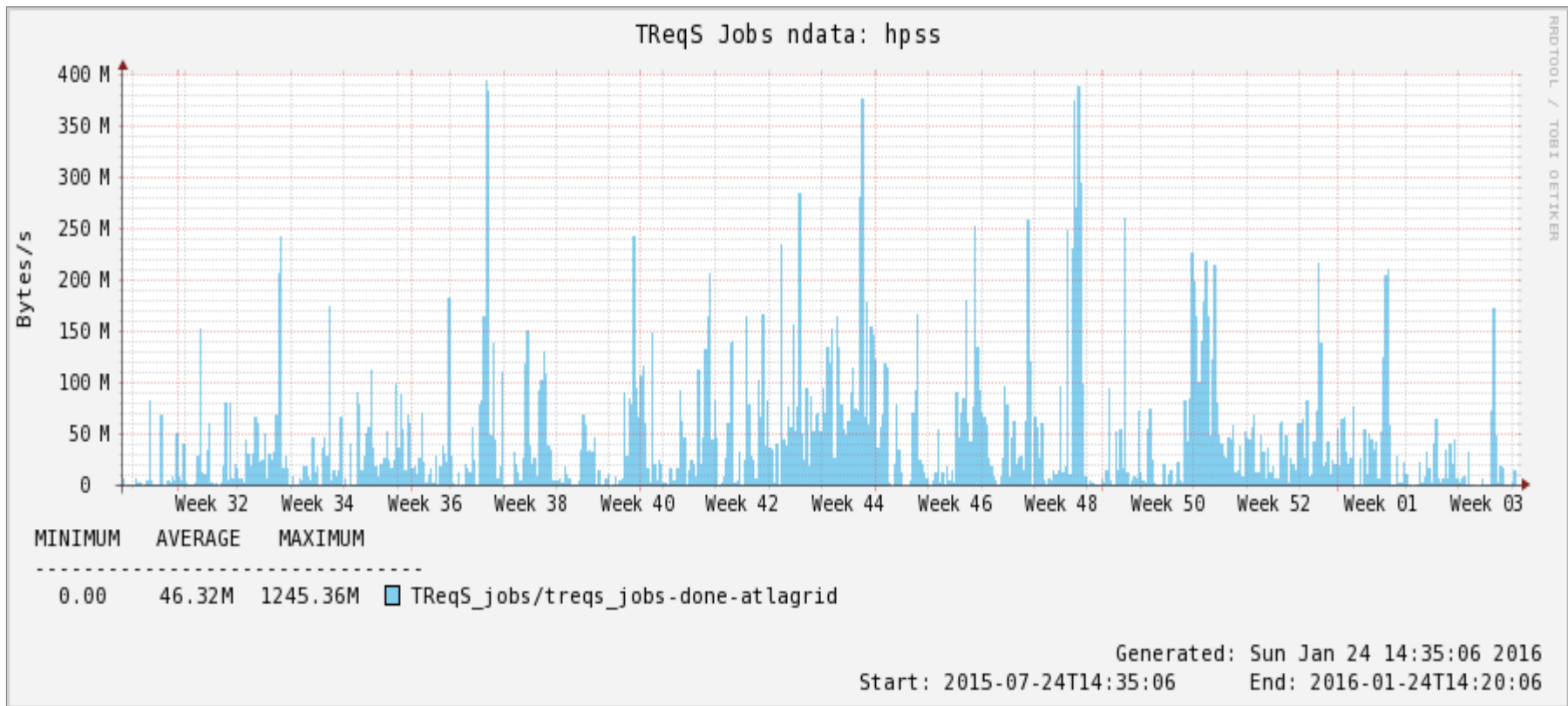
# Tape activity monitoring (RAL)

# Tape activity monitoring (CNAF)

- Read & Write throughput & tape mounts

# Tape activity monitoring (CCIN2P3)

**Stage-in throughput (last 6 months)**



**Max at 1.2 GB/sec**
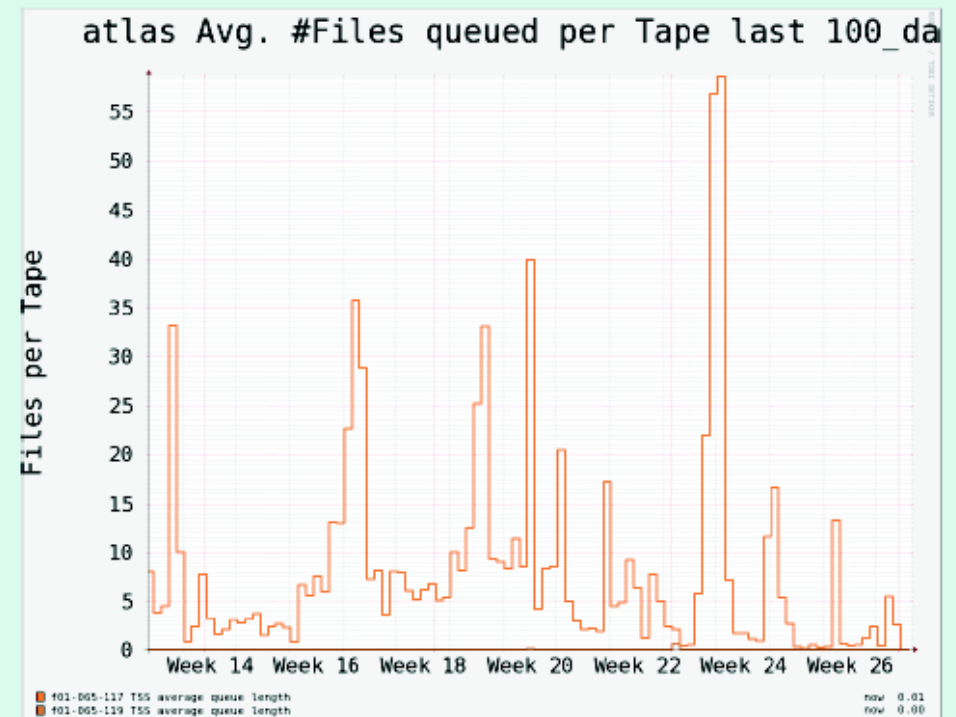**Average media  mount time   ~ 1-2 minutes**

# Tape activity monitoring (KIT)

- Average number of files recalled per tape mount: tape families, grouping issues.
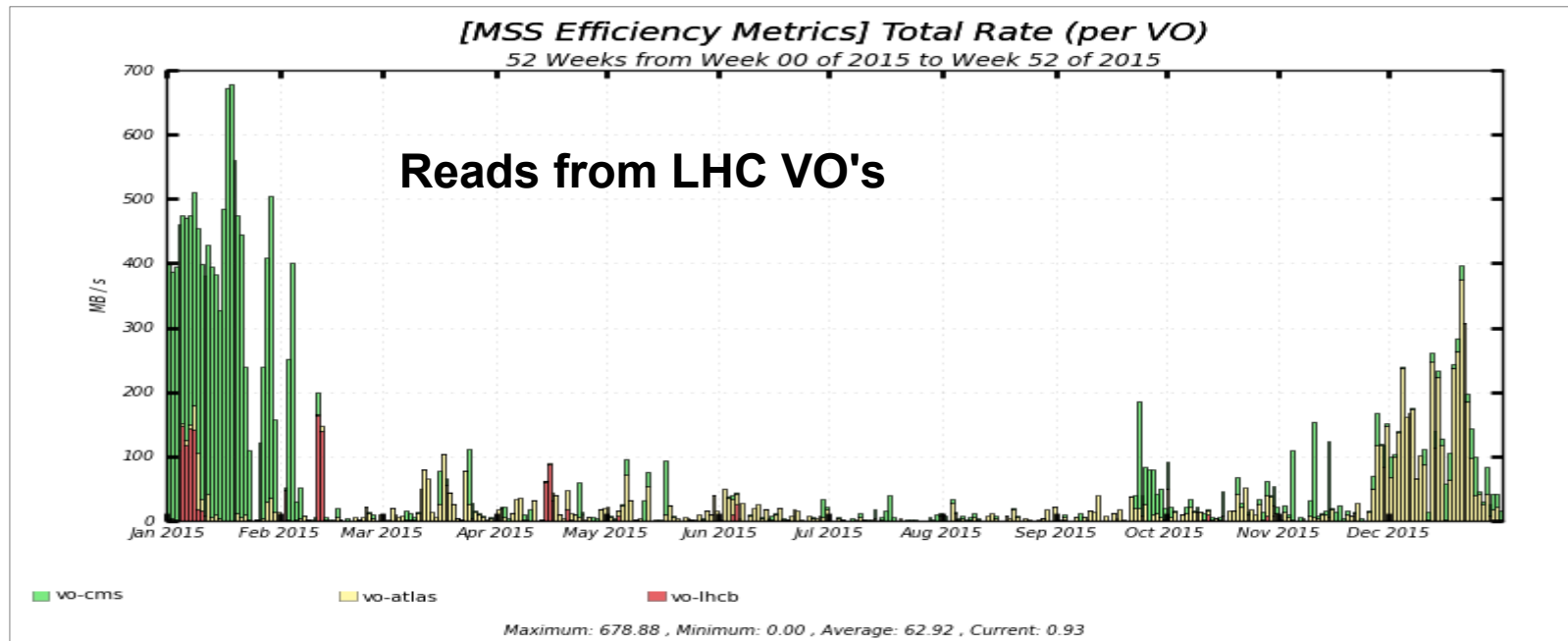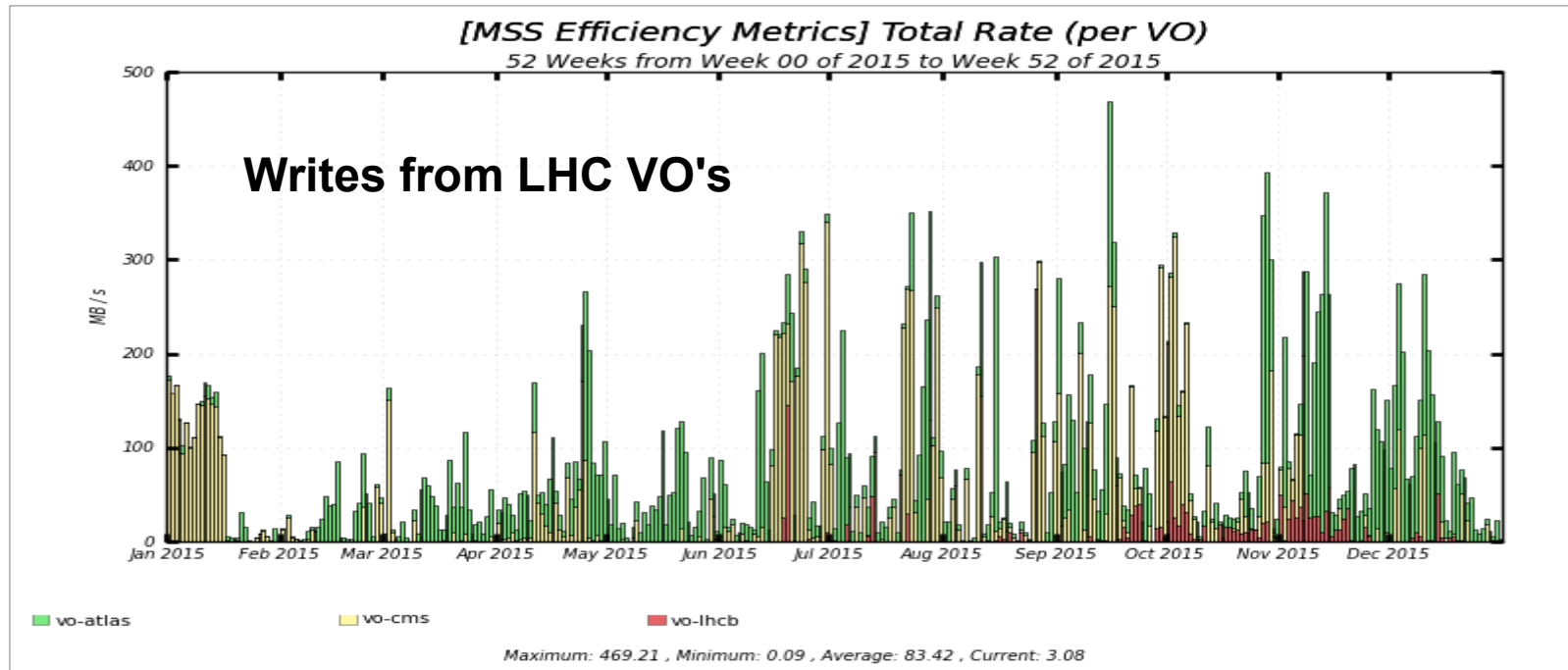
## ATLAS vs. GridKa tape **worst case**

- ~600MB average file size
- ↝ 2500 files fit on one LTO-5 tape
- ≤10 files queued per tape on average recall
- every second wrap a file ↝ 90s winding to reach a file
- **+4s** to read 600MB + fraction to mount tape
  ≈ 100s total read time/file
- ≈ 6MB/s throughput in the worst case

## files queued per tape



atlas Avg. #Files queued per Tape last 100_da

Files per Tape

55
50
45
40
35
30
25
20
15
10
5
0

Week 14  Week 16  Week 18  Week 20  Week 22  Week 24  Week 26

f01-DG5-117 TSS average queue length     now  0.01
f01-DG5-119 TSS average queue length     now  0.00

# Tape activity monitoring (PIC)

# Tape activity monitoring (SARA)

- Write and read pool requests for the last year
  (see http://web.grid.sara.nl/dcache.php?r=year)