# DC24 Retrospective - INFN T1

Lucia Morganti on behalf of the Data management team @INFN-T1

# General overview and impressions from our site

- The DC24 is a very useful exercise to find bottlenecks within sites
  - Very much in favour of running preparatory tests to tune injection parameters, as done in the previous months for CMS, and of re-testing, as proposed by ATLAS and LHCb for the current week

- However, the DC  is a stress test greatly impacting sites and overloading storage endpoints
  - We had GGUS tickets and red SAM tests during the challenge
  - The DC24 time-range should be excluded from A/R computation, see e.g. https://ggus.eu/index.php?mode=ticket_info&ticket_id=165509 (thanks Stephan)
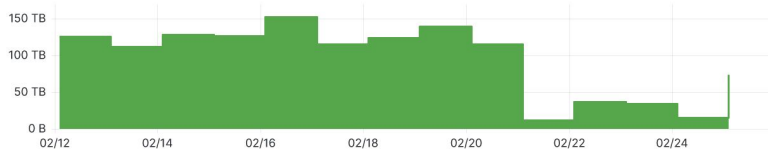
# General overview and impressions from our site

- The DC  is a stress test greatly impacting sites and overloading storage endpoints (continued)
  - Unfortunately, we do not have any way to regulate fluxes: once a StoRM WebDAV endpoint is overloaded and threads saturate, transfers fail; they are not queued or delayed. The more transfers are submitted, the worse it gets.
    - Is there a way for FTS to regulate injection based on success rate?

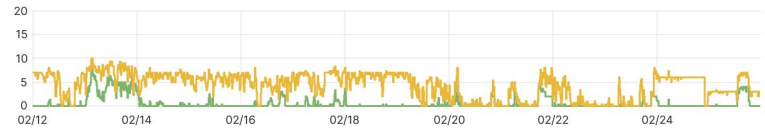# General overview and impressions from our site

- The DC is a stress test greatly impacting sites and overloading storage endpoints (continued)
  - On top of this, significant production load during the challenge, which in some cases had no impact (e.g. Alice) whereas in other cases heavily impacted the infrastructure

Recall bytes per day (stacked)

| | Mean | Last * | Max | Min | Total |
|---|---|---|---|---|---|
| tsm-hsm-6.cr.cnaf.infn.it | 93.8 TB | 74.0 TB | 153 TB | 11.8 TB | 1.31 PB |

Recall drives actually in use

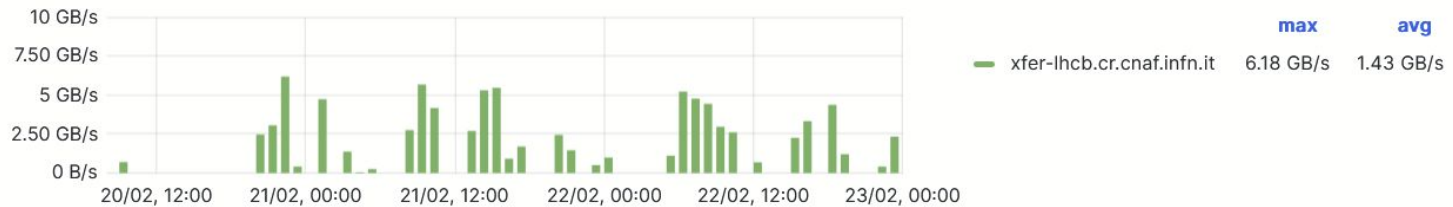| | min | max | avg | current | total |
|---|---|---|---|---|---|
| atlas t10000d | 0 | 9 | 1 | 0 | 765 |
| atlas ts1160 | 0 | 7 | 4 | 3 | 4662 |

Staging activity from ATLAS during the DC

# Interpreting DC monitoring

Does it make sense to reference "average transfer rate" for a period when FTS does not submit transfer requests?
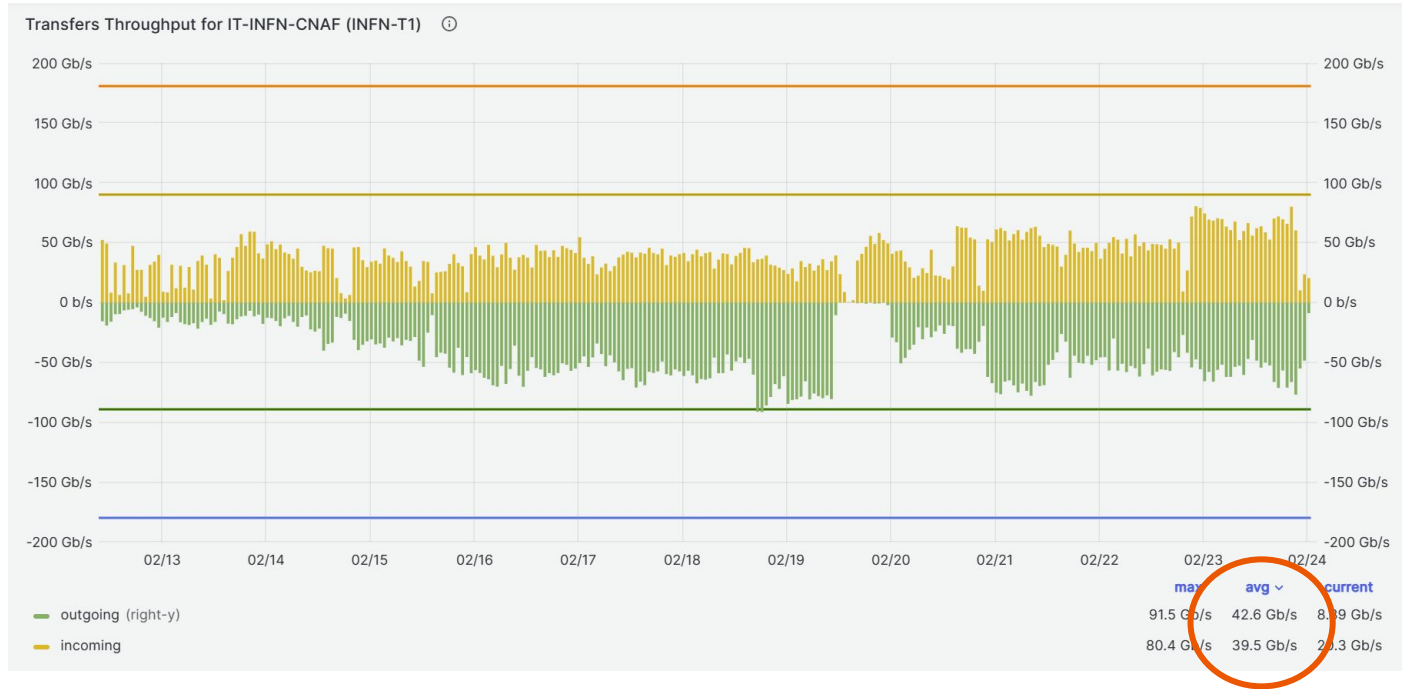


LHCb, Tape-Disk, which is actually Disk_buffer-Disk
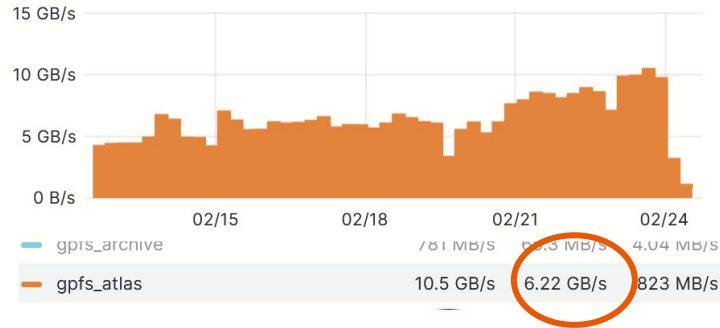(FTS plot provided by A. Rogovskiy)

# Interpreting DC monitoring

- Throughputs reported by FTS monitoring for our site are much lower than what we observe
  - Are we measuring an important contribution from production load?
    - Unfortunately, we cannot disentangle.
    - How are other sites dealing with this?
  - Is FTS throughput computed and reported only for successful transfers?
    - Again, unfortunately we cannot disentangle in the traffic we measure.
    - Shouldn't success rate be reported together with throughput?
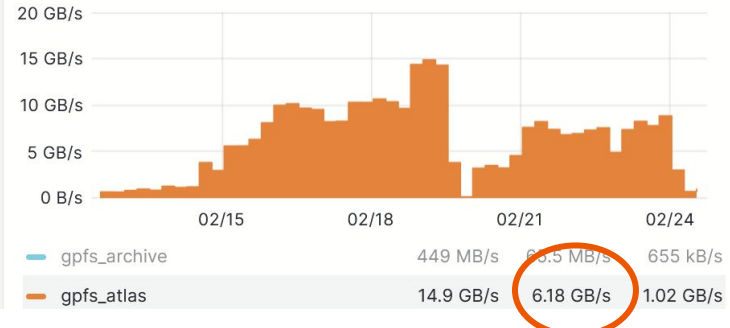- How about using/ comparing different metrics, e.g. transferred TB per day?

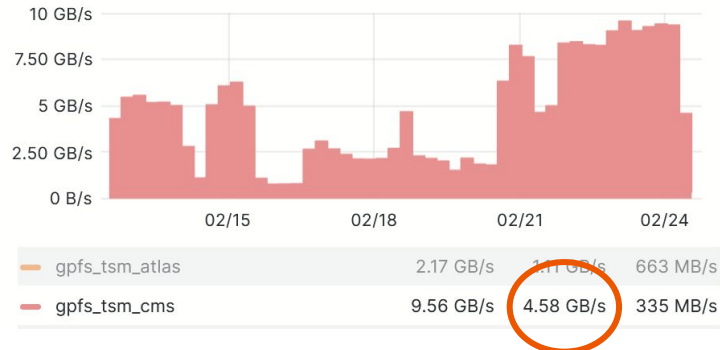# Monit plot provided by A. Forti for ATLAS+CMS
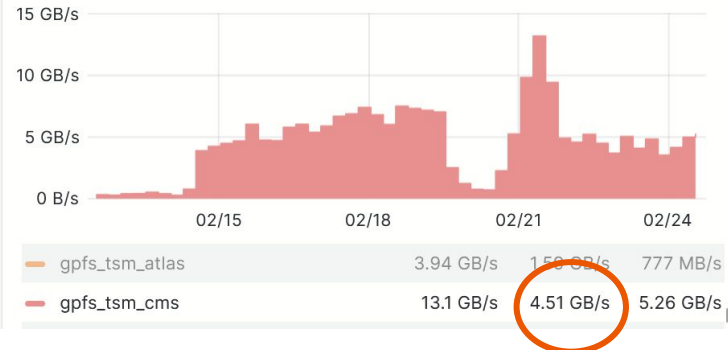
Gateway traffic in (non POSIX writing)

| | | | |
|---|---|---|---|
| gpfs_archive | 781 MB/s | 65.3 MB/s | 4.04 MB/s |
| gpfs_atlas | 10.5 GB/s | **6.22 GB/s** | 823 MB/s |

Gateway traffic out (non POSIX reading)

| | | | |
|---|---|---|---|
| gpfs_archive | 449 MB/s | 65.5 MB/s | 655 kB/s |
| gpfs_atlas | 14.9 GB/s | **6.18 GB/s** | 1.02 GB/s |

Gateway traffic in (non POSIX writing)

| | | | |
|---|---|---|---|
| gpfs_tsm_atlas | 2.17 GB/s | 1.11 GB/s | 663 MB/s |
| gpfs_tsm_cms | 9.56 GB/s | **4.58 GB/s** | 335 MB/s |

Gateway traffic out (non POSIX reading)

| | | | |
|---|---|---|---|
| gpfs_tsm_atlas | 3.94 GB/s | 1.50 GB/s | 777 MB/s |
| gpfs_tsm_cms | 13.1 GB/s | **4.51 GB/s** | 5.26 GB/s |

We measure 85 Gb/s OUT and 86 Gb/s IN   (FTS says 42 Gb/s OUT and 39.5 Gb/s IN)

# Interpreting DC monitoring

In the mixed scenario, it seems there was an imbalance in the requests T0-T1 vs T1-T1 that affected the metric "model vs reality" independently from the site performance.

| T1 Site | Minimal (T0→T1) | | | Flexible (T0→T1) | | | Flexible (T0+T1→T1) | | |
|---|---|---|---|---|---|---|---|---|---|
| | model | reality | [%] | model | reality | [%] | model | reality | [%] |
| BNL-ATLAS | 60.0 | 25.9 | 43 | 68.4 | 21.2 | 31 | 82.1 | 57.1 | 70 |
| FZK-LCG2 | 32.0 | 34.1 | 107 | 39.0 | 13.2 | 34 | 59.4 | 43.2 | 73 |
| IN2P3-CC | 38.0 | 36.4 | 96 | 44.2 | 1.4 | 3 | 59.1 | 21.4 | 36 |
| INFN-T1 | 23.0 | 22.0 | 96 | 28.3 | 8.9 | 31 | 39.4 | 47.6 | 121 |

From ATLAS presentation, last Retrospective

# Things we will investigate/improve

- CMS
  - We are planning to align the StoRM WebDAV instances dedicated to CMS since we observed higher load and higher failures in those servers having lower number of CPU cores
- LHCb
  - We'll re-think LHCb hardware configuration so to accommodate their workflow, given @INFN-T1 tape buffer and disk are on the same filesystem, managed by the same endpoints
- StoRM developers are working at improving efficiency (e.g. https://github.com/italiangrid/storm-webdav/pull/40) and introducing performance markers in StoRM WebDAV