# CMS

Katy Ellis, GridPP51, 26 March 2024

# CMS news/objectives/O&C week (in my view)

**Reviewing 2023 data-taking**

**Resources for 2024 and beyond**

**Computing Design Report (CDR)**

**Alternative architectures**

**Migration to ALMA9**
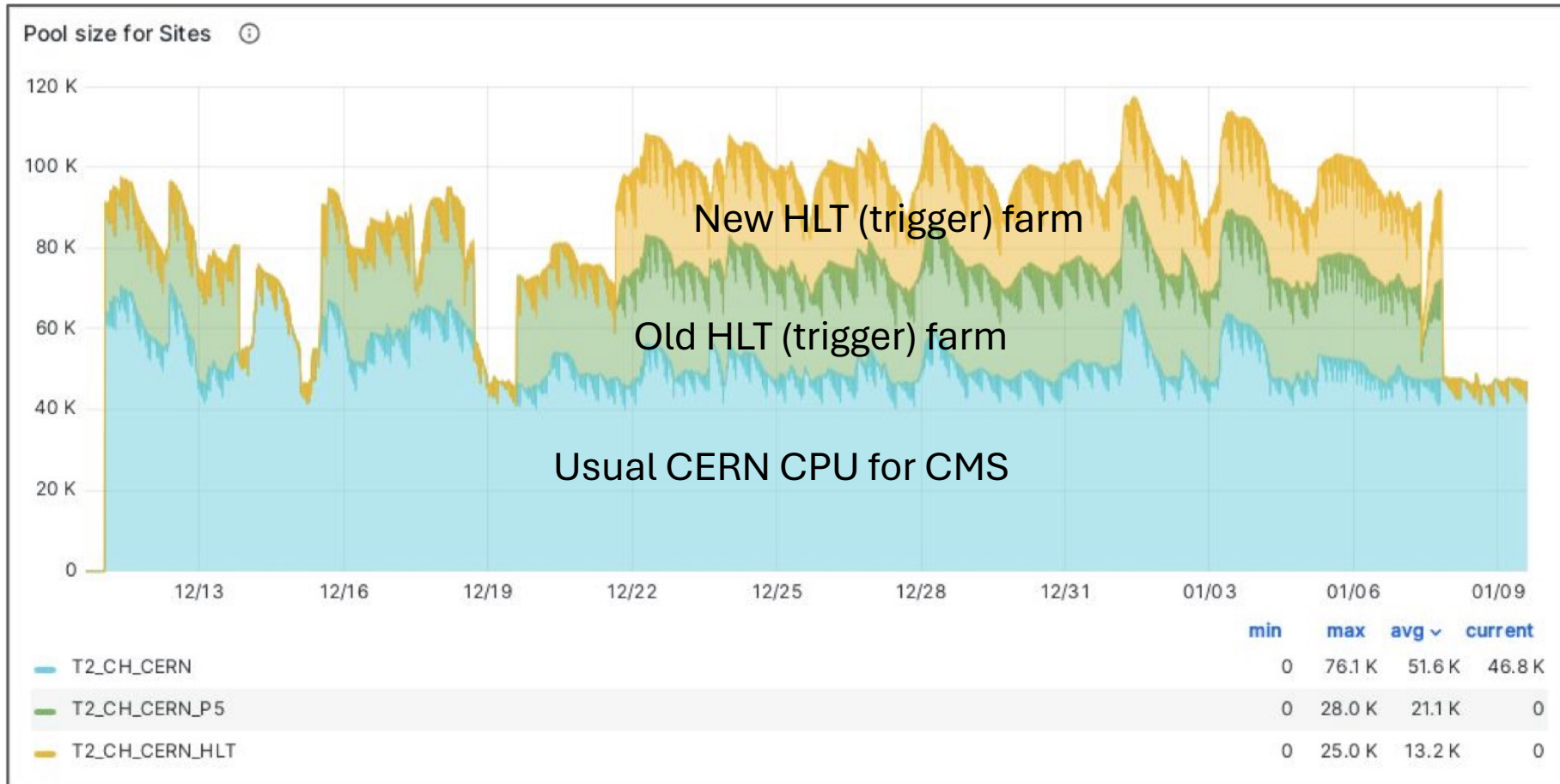
**Efficiency improvements**

Pilot overloading

**GPU latest**

+ DC24 + tokens

Pool size for Sites ⓘ

New HLT (trigger) farm

Old HLT (trigger) farm

Usual CERN CPU for CMS

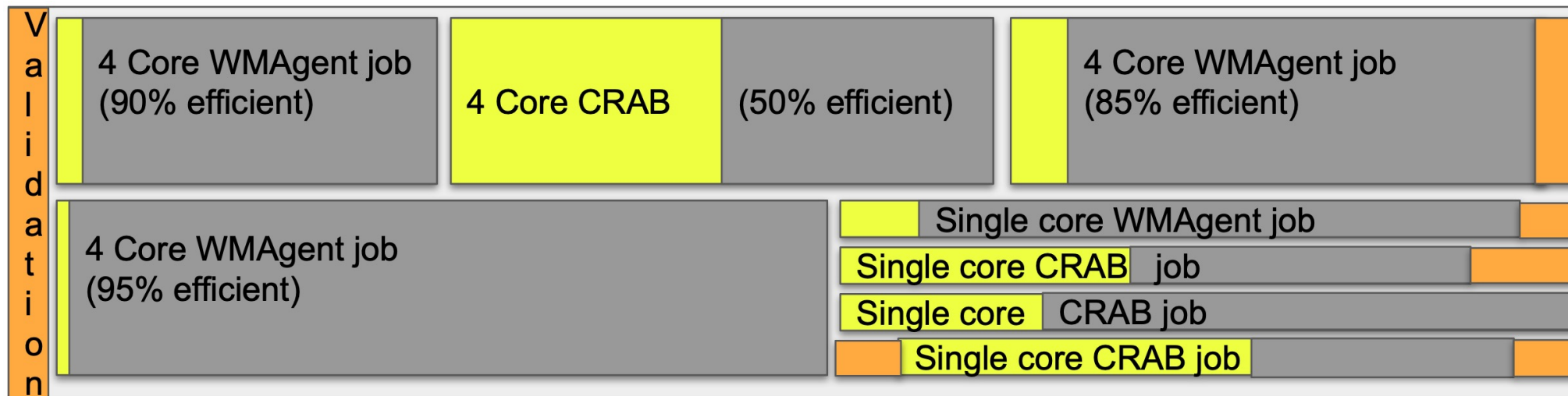| | | min | max | avg ⌄ | current |
|---|---|---|---|---|---|
| — | T2_CH_CERN | 0 | 76.1 K | 51.6 K | 46.8 K |
| — | T2_CH_CERN_P5 | 0 | 28.0 K | 21.1 K | 0 |
| — | T2_CH_CERN_HLT | 0 | 25.0 K | 13.2 K | 0 |

# Alternative architectures

- Discussion of removing support for PowerPC

- Tests with ARM continue
  - CMS tested on the Glasgow ARM farm before Christmas
  - Frustratingly - *The validation team observes some discrepancies that need to be understood and their impact on physics needs to be assessed. So at this point ARM is not yet validated for physics in CMS.*

# Overloaded pilots - motivation

- CMS wants to improve CPU efficiency (CPU time / wallclock time)
- Breaking down a typical CMS job (8 core, 48 hour duration)



- Inefficiencies are from scheduling (orange) and payload (yellow)

Images courtesy of CMS SI team

# Overloaded pilots – scheduling inefficiencies

**T1s have good scheduling efficiency** →

**For some T2s there is room for improvement – what can we do?** ↓

Detailed new monitoring is used to examine scheduling inefficiencies within the pilots

## Scheduling efficiency for T1 sites

| data.payload.GLIDEIN_ | Efficiency | Average idle CPUs | Average total CPUs |
|---|---|---|---|
| T1_US_FNAL | 0.973 | 897 | 33385 |
| T1_UK_RAL | 0.978 | 144 | 6562 |
| T1_RU_JINR | 0.970 | 609 | 20606 |
| T1_IT_CNAF | 0.957 | 370 | 8621 |
| T1_FR_CCIN2P3 | 0.956 | 350 | 7878 |
| T1_ES_PIC | 0.948 | 134 | 2581 |
| T1_DE_KIT | 0.941 | 584 | 9966 |

## Scheduling efficiency for T2 sites

| data.payload.GLIDEIN_ | Efficiency ↑ | Average idle CPUs | Average total CPUs |
|---|---|---|---|
| T2_BR_UERJ | 0.0601 | 20.6 | 21.9 |
| T2_CN_Beijing | 0.600 | 107 | 268 |
| T2_IN_TIFR | 0.680 | 238 | 743 |
| T2_UK_SGrid_Bristol | 0.718 | 13.6 | 48.2 |
| T2_RU_INR | 0.730 | 20.9 | 77.4 |
| T2_UK_London_Brunel | 0.737 | 67.1 | 255 |
| T2_FR_IPHC | 0.746 | 369 | 1451 |
| T2_HU_Budapest | 0.755 | 165 | 675 |

**Similar for T3s** ↓

## Scheduling efficiency for T3 sites

| data.payload.GLIDEIN_ | Efficiency ↑ | Average idle CPUs | Average total CPUs |
|---|---|---|---|
| T3_US_NotreDame | 0.367 | 14.0 | 22.2 |
| T3_US_FNALLPC | 0.569 | 234 | 543 |
| T3_TW_TIDC | 0.587 | 10.7 | 25.9 |
| T3_US_PuertoRico | 0.588 | 14.8 | 35.9 |
| T3_US_Rutgers | 0.668 | 13.5 | 40.8 |
| T3_UK_ScotGrid_GLA | 0.677 | 188 | 584 |
| T3_KR_KNU | 0.733 | 35.0 | 131 |
| T3_UK_London_QMUL | 0.790 | 164 | 781 |

# Overloaded pilots – unused memory



Average memory allocation to payloads for fully utilized pilots ⓘ

1 - Sum data.payload.Memory / Sum data.payload.TotalSlotMemory

| | min | max | avg |
|---|---|---|---|
| | 0.424 | 0.759 | 0.692 |

# Overloaded pilots - testing

- During the last months IC and RALPP have been part of a test
- Overloaded pilots config looks like this:

  - Site CE/**batch system** parameters (i.e. what we request from a CE) →

    ```
    <submit_attr name="RequestCpus" value="8"/>
    <submit_attr name="+maxMemory" value="20240"/>
    ```

    } 25% overloading

  - CMS **scheduling** parameters (i.e. how we tell HTCondor to use it)
    **Overloading pilots by only increasing the internal scheduling parameters**

    →
    ```
    <attr name="GLIDEIN_CPUS" value="10"/>
    <attr name="GLIDEIN_MaxMemMBs" value="25000"/>
    ```

- One CE had the new config, the other(s) had the old config
- Results on next slides
- CMS is extending the testing to more sites
- CMS is also finding scheduling efficiencies via whole-node pilot extension at selected sites

# Overloaded pilots – test results

- Results can be monitored(?) through EGI accounting ("Row variable = Submit Host") but only monthly granularity

- At RALPP, increase in efficiency is clear (March 2024 numbers):

| | | |
|---|---|---|
| https://heplnx206.pp.rl.ac.uk:60000/arex | 96.91% | Overloaded CE |
| https://heplnx207.pp.rl.ac.uk:60000/arex | 81.17% | Un-changed |

- At IC,

| | |
|---|---|
| ceprod00.grid.hep.ph.ic.ac.uk:9619/ceprod00.grid.hep.ph.ic.ac.uk-condor | 59.29% |
| ceprod01.grid.hep.ph.ic.ac.uk:9619/ceprod01.grid.hep.ph.ic.ac.uk-condor | 58.37% |
| ceprod02.grid.hep.ph.ic.ac.uk:9619/ceprod02.grid.hep.ph.ic.ac.uk-condor | 62.92% |
| ceprod03.grid.hep.ph.ic.ac.uk:9619/ceprod03.grid.hep.ph.ic.ac.uk-condor | 62.03% |

These numbers are wrong!

# Overloaded pilots – test results

- Results can be monitored through EGI accounting ("Row variable = Submit Host") but only monthly granularity

- At RALPP, increase in efficiency is clear (March 2024 numbers):

| | |
|---|---|
| https://heplnx206.pp.rl.ac.uk:60000/arex | 96.91% |
| https://heplnx207.pp.rl.ac.uk:60000/arex | 81.17% |

Overloaded CE
Un-changed

- At IC

| | Normal | Overloaded | Enabled |
|---|---|---|---|
| T2_ES_PIC | 82.01% | 90,88% | April 2023 |
| T2_ES_CIEMAT | 82.14% | 97.49% | April 2023 |
| T2_UK_London_IC | 74.82% | 83.37% | Dec 6th 2023 |
| T2_UK_SGrid_RALPP | 80.11% | 95.31% | Dec 6th 2023 |
| T2_IT_Bari | 81.27% | 84.97% | Feb 20th 2024 |
| T2_IT_Legnaro | 81.71% | 95.51% | Feb 20th 2024 |
| T2_IT_Pisa | N/A | N/A | Feb 20th 2024 |
| T2_IT_Roma | N/A | N/A | Feb 20th 2024 |
| T1_FR_CCIN2P3 | 87.41% | 97.61% | March 4th 2024 |
| T1_IT_CNAF | 82.46% | 84.71% | March 4th 2024 |
| T2_BE_IIHE, T2_CH_CSCS, T2_EE_Estonia, T2_US_Caltech, T2_US_Vanderbilt, T2_US_Wisconsin | | | March 14th 2024 |

ceprod00.grid.hep.p
condor
ceprod01.grid.hep.p
condor
ceprod02.grid.hep.p
condor
ceprod03.grid.hep.p
condor

# Tokens - storage

- SAM tests include *storage* tests with tokens
  - Many sites are passing these
  - RAL T1 needed a new XRootD version but tests are green since 4th March
- CMS found that not every site passing SAM tests were ready for production transfers with Rucio
- 'Good' sites did use tokens for DC24
  - 19 sites from the start, 25 sites by the second week
  - However, this was a minimal implementation for both CMS and FTS
  - Now reverted
- Concern about interactions with IAM

# Tokens - compute

- CMS Submission Infrastructure and Workload Management System are ready

- Some issues with ARC-CEs
  - Rollout has been delayed

# Tokens - compute

- CMS Submission Infrastructure and Workload Management System are ready

- Some issues with ARC-CEs
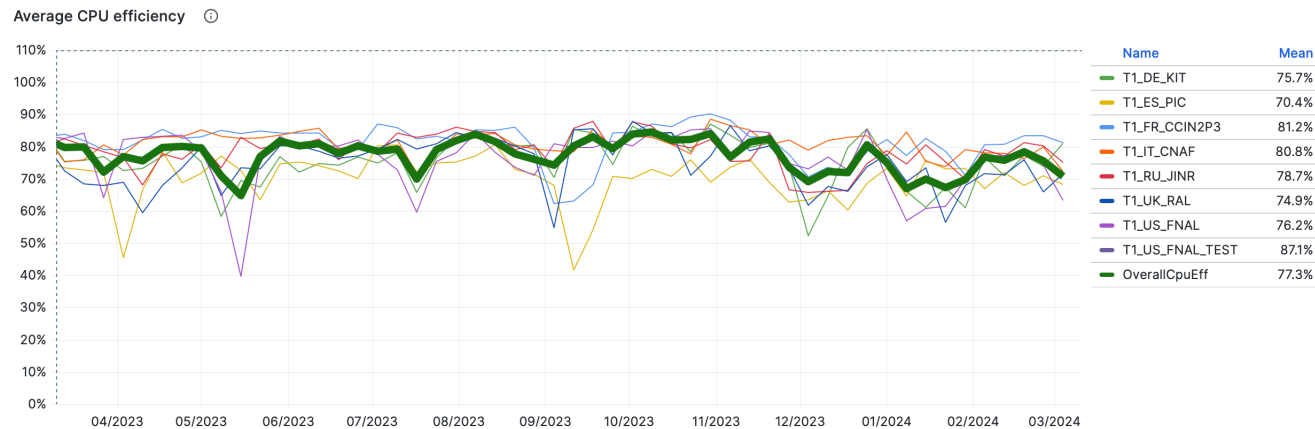  - Rollout has been delayed

STOP PRESS!
Ticket requesting the enabling of
access received yesterday morning!
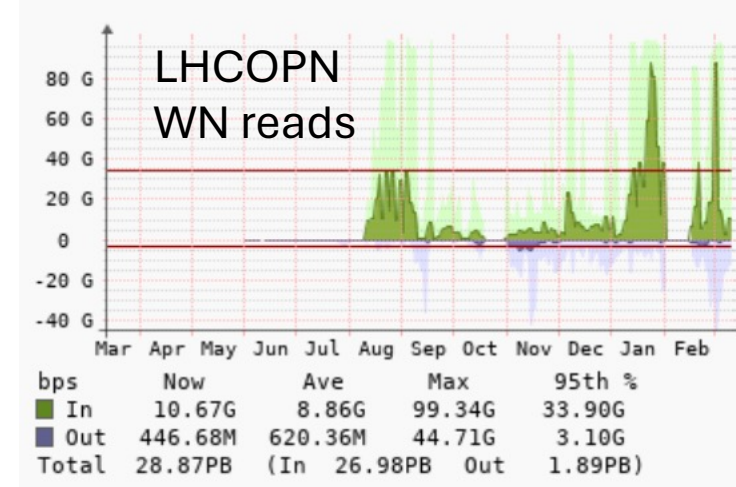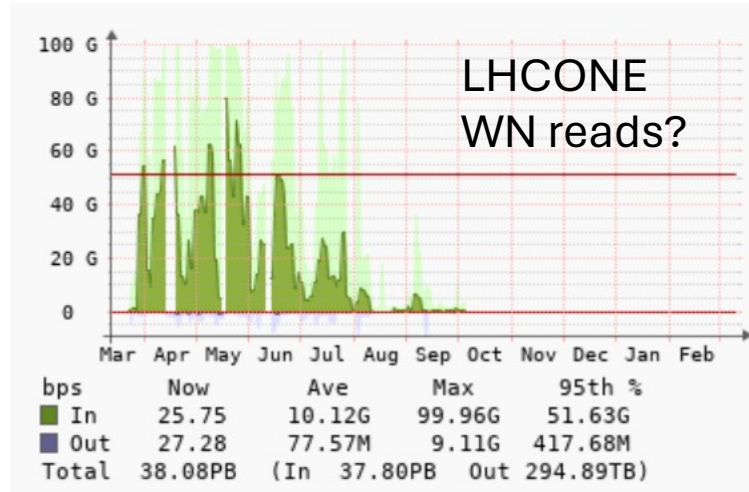
# CMS UK computing

How's it going?

# RAL – job efficiency

- Job efficiency has been around the T1 average in the last year
- Occasional drops generally attributed to periods of "drain" (i.e. when the summary of SAM tests have failure rates>10% for 2 days out of 3)
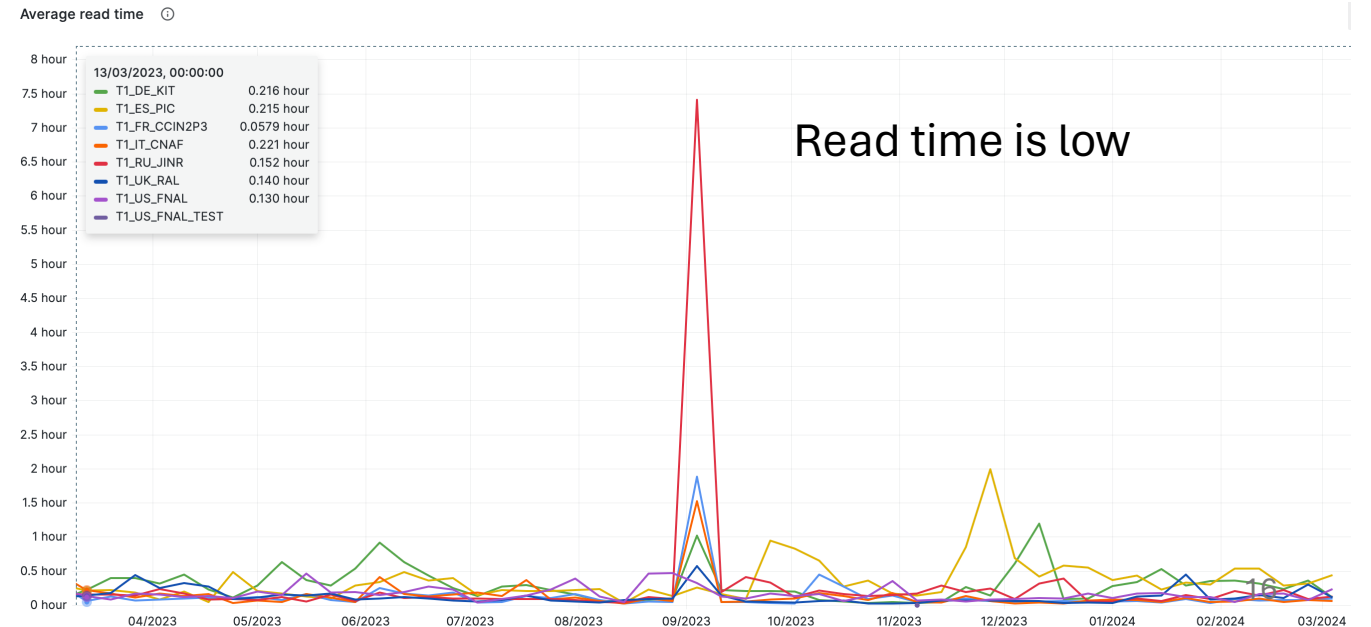
Average CPU efficiency
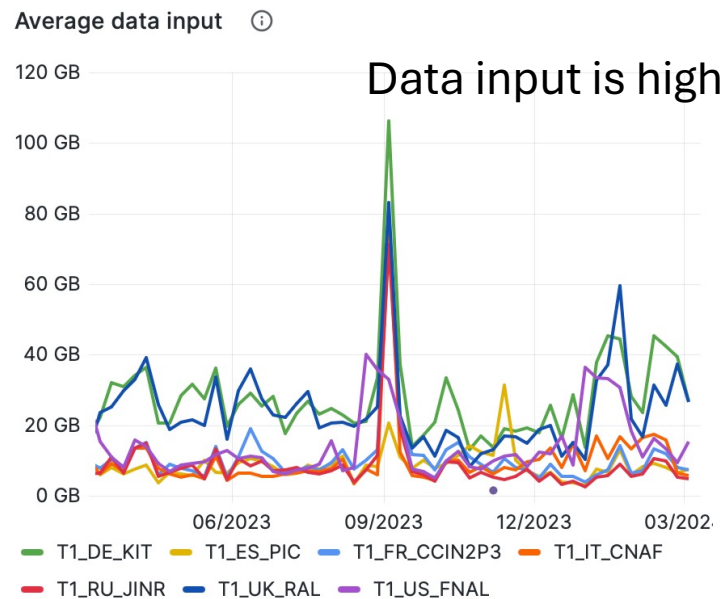
| Name | Mean |
| --- | --- |
| T1_DE_KIT | 75.7% |
| T1_ES_PIC | 70.4% |
| T1_FR_CCIN2P3 | 81.2% |
| T1_IT_CNAF | 80.8% |
| T1_RU_JINR | 78.7% |
| T1_UK_RAL | 74.9% |
| T1_US_FNAL | 76.2% |
| T1_US_FNAL_TEST | 87.1% |
| OverallCpuEff | 77.3% |

# Offsite reads?

How much data is being read?



LHCONE
WN reads?



LHCOPN
WN reads

Bandwidth near saturation

Lazy-download

Data input is high



Read time is low

# RAL disk transfers via FTS

- Although improved on previous years, we were aware that disk transfers via Echo gateways did not give the expected performance according to our infrastructure capabilities. Pre-DC24 tests in 2023 confirmed this.

- Liaison tests in Jan 2024 showed unexpectedly that newer gateways on the new network showed worse performance for single file write speed (~5MB/s) and in iPerf tests than older gateways on the old network (~15MB/s).

# Echo gateway tests

- Transferred the same (small) file from Prague to Echo via each gateway 5 times (3rd Jan)
- Gateways were loaded with regular traffic at the time, unless otherwise specified
- Calculated the average rate for each gateway
- Arbitrarily coloured the rate green if it managed higher than 10MB/s

| Dest GW | Comment | Network | Average rate (MB/s) |
|---|---|---|---|
| xrootd.echo.stfc.ac.uk:1094 | Alias | Mix | 9.43 |
| SVC16 | No production work | New | 17.09 |
| GW14 | | Legacy | 14.82 |
| GW15 | | Legacy | 15.64 |
| GW16 | | Legacy | 17.04 |
| GW4 | | Legacy | 15.40 |
| GW5 | | Legacy | 14.76 |
| GW6 | | Legacy | 14.60 |
| GW7 | | Legacy | 11.79 |
| SVC01 | | New | 14.51 |
| SVC02 | | New | 7.03 |
| SVC03 | | New | 8.60 |
| SVC05 | | New | 9.38 |
| SVC11 | | New | 8.41 |
| SVC13 | | New | 7.46 |
| SVC14 | | New | 5.11 |
| SVC15 | | New | 5.68 |
| SVC17 | | New | 9.27 |
| SVC18 | | New | 5.40 |
| SVC21 | | New | 8.07 |
| SVC22 | | New | 10.04 |
| SVC23 | | New | 6.94 |
| SVC24 | | New | 4.20 |
| SVC25 | Didn't work | New | |
| SVC26 | No production work | New | 21.54 |
| SVC97 | | Legacy | 13.88 |
| SVC98 | | Legacy | 12.59 |
| SVC99 | | Legacy | 10.71 |

# Echo gateway tests (II)

- Subsequent tests with larger files coming from CERN showed broadly consistent results (with higher rates as expected of larger files)

- iPerf tests also agreed, showing the performance is asymmetric (better rates when reading from RAL) and no significant difference between IPv4 and IPv6

- Tests with several different sources (UK, Europe, US) to two gateways (new and legacy network) again demonstrated consistent results with the legacy network faster

# Echo gateway tests (III)

- During the data challenge Echo gateways were network tuned, ECN value changed, and load-balancing algorithm changed

19th March, after the above changes, before gateway moves

| Source | Rate (MB/s) average over 5 transfers | |
|---|---|---|
| | gw16 | svc14 |
| KIT | 32.46 | 46.21 |
| FNAL | 19.71 | 31.33 |
| RALPP | 46.92 | 46.46 |
| IC | 47.29 | 49.10 |
| DESY | 40.97 | 49.55 |
| Purdue (US) | 46.56 | 43.98 |
| CSCS (CH) | 43.80 | 50.02 |

Still legacy network

Always new network

Single, new network gateway improvement in time

| | 5/01/24, SVC14 | 19/03/24, SVC14 |
|---|---|---|
| T2_CH_CERN | 16.43 | |
| T1_DE_KIT_Disk | 14.48 | 46.21 |
| T1_US_FNAL_Disk | 13.66 | 31.33 |
| T2_UK_SGrid_RALPP | 13.80 | 46.46 |
| T2_UK_London_IC | 32.97 | 49.10 |
| T2_DE_DESY | 8.29 | 49.55 |
| T2_US_Purdue | 14.92 | 43.98 |
| T2_CH_CSCS | 14.56 | 50.02 |

# IC

- Included in the pilot overloading tests
- Interested in further network tests
- Alex Richards now working for CMS, spending part of his time developing core Rucio code

# Brunel

- Long-running storage migration from DPM to Ceph
- [Ticket](#)
- Site is now out of the 'waiting room' and fully back in production as of 5[th] March
- Some variable job performance to investigate

# Shoveler - Monitoring software specific to XRootD transfers

- Has been tested at RAL Tier 1 for...ages.
- Observed issues for some time
  - Stopping randomly and failing to restart
  - No throughput monitoring provided / no consistent comparison with internal monitoring
- Now being rolled out to more sites
- Important for CMS as a large fraction of our data movement is via AAA (which uses XRootD)
- (In my opinion) requires "official" validation by CMS before being used
  - Possibly a task for Katy...
- Then we can use the monitoring information to make improvements

# Summary

- CMS computing is ready for a bumper year of data!
- Working hard on the challenges of HL-LHC although person-power is a limiting factor
- Trying to get everything we can out of the resources we have
- Further progress for I/O at RAL
- Brunel migrated disk storage away from DPM