# RAL-CMS summary

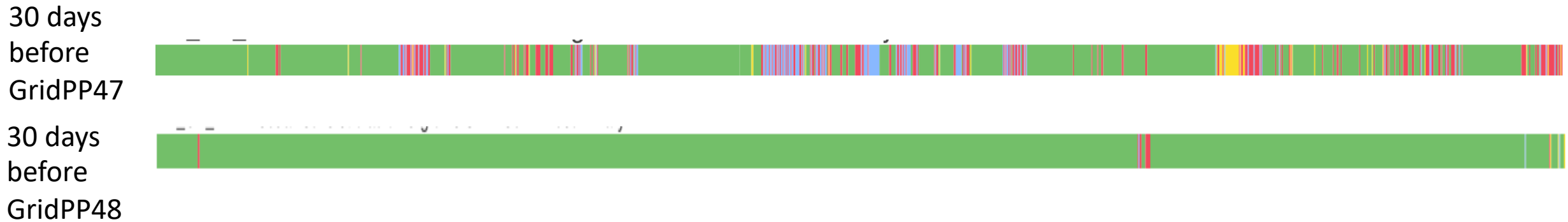Katy Ellis, CMS liaison at RAL Tier 1

02/09/22, GridPP48

# Content

- General talk about CMS at Tier 1 – status, problems and solutions – and some other stuff I know about in CMS world:
  - Operational issues
  - Job performance
  - Tape usage, tape families and automation
  - New monitoring
  - Rucio
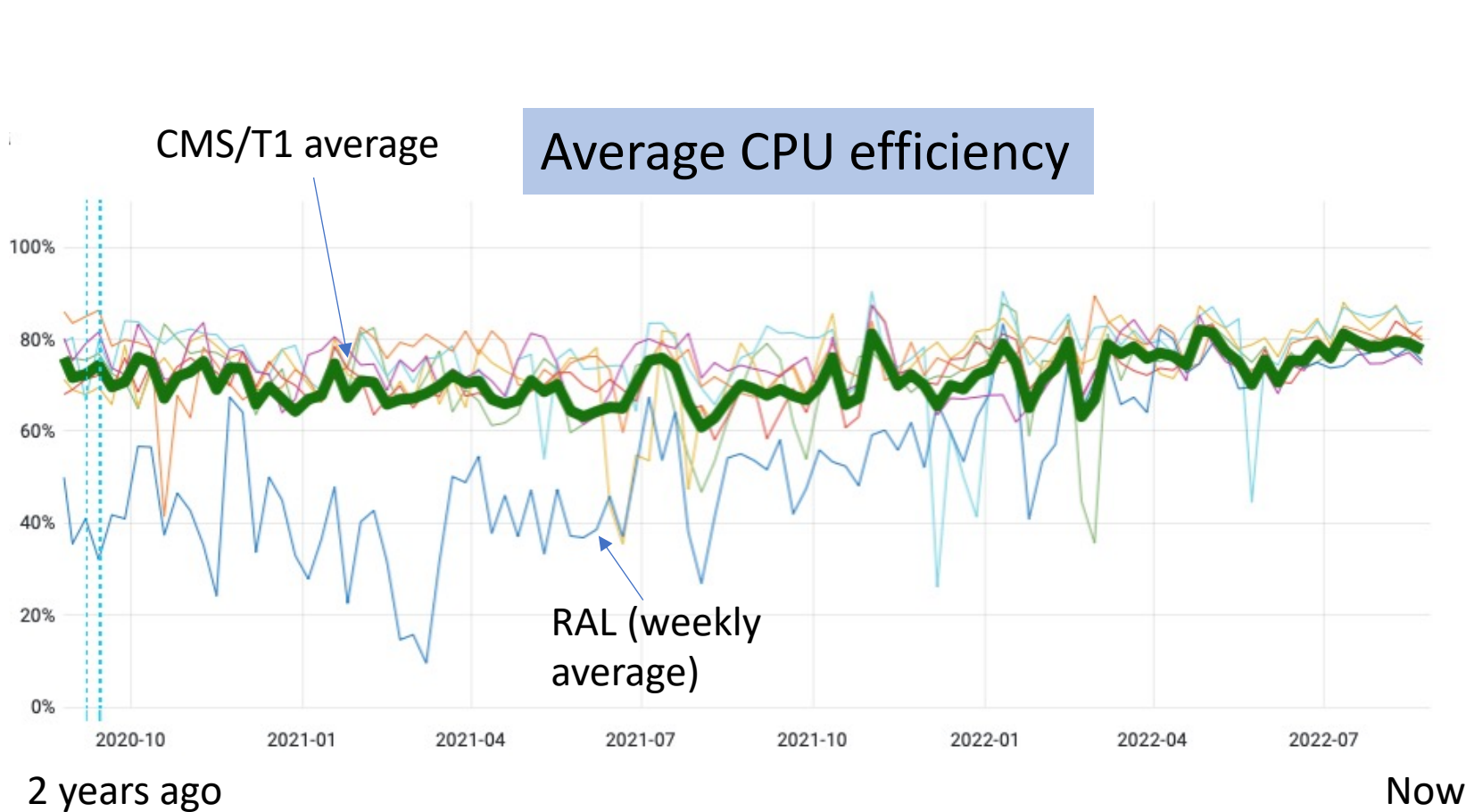  - CRAB and Support
  - Token status

# Operational issues since GridPP47

- Webdav transfers and **SAM tests** much more stable

30 days before GridPP47
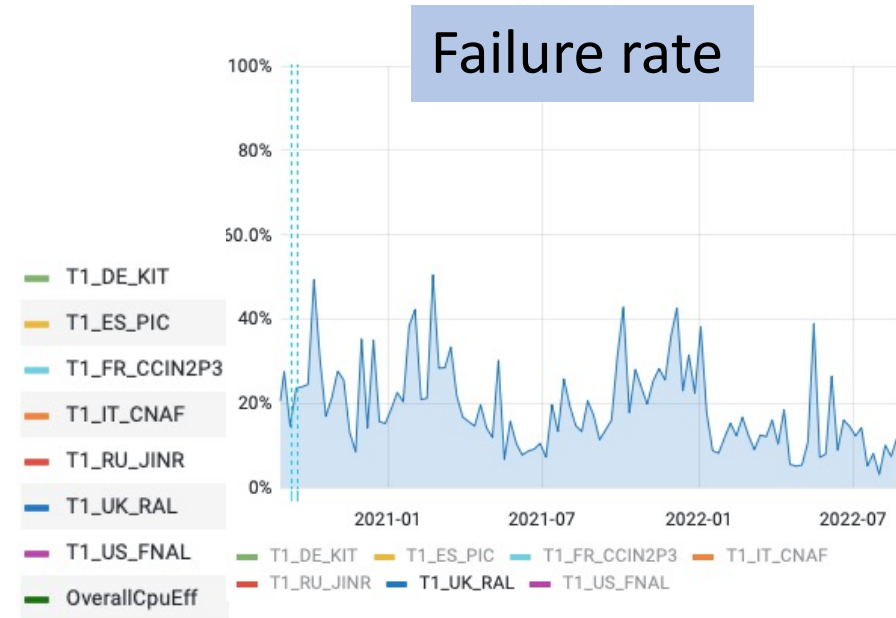


30 days before GridPP48



- SAM test failures on all CEs during draining of one ARC-CE
  - We don't know what causes this, as CEs are supposed to be independent!
  - Non-test jobs are typically not affected
  - Grid-services team now tries to avoid draining CEs (instead just doing a straight reboot when updates are required)

# Job performance – improvement sustained

CMS/T1 average

Average CPU efficiency

Failure rate

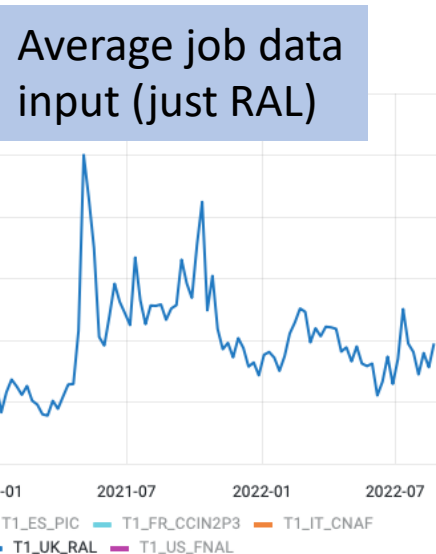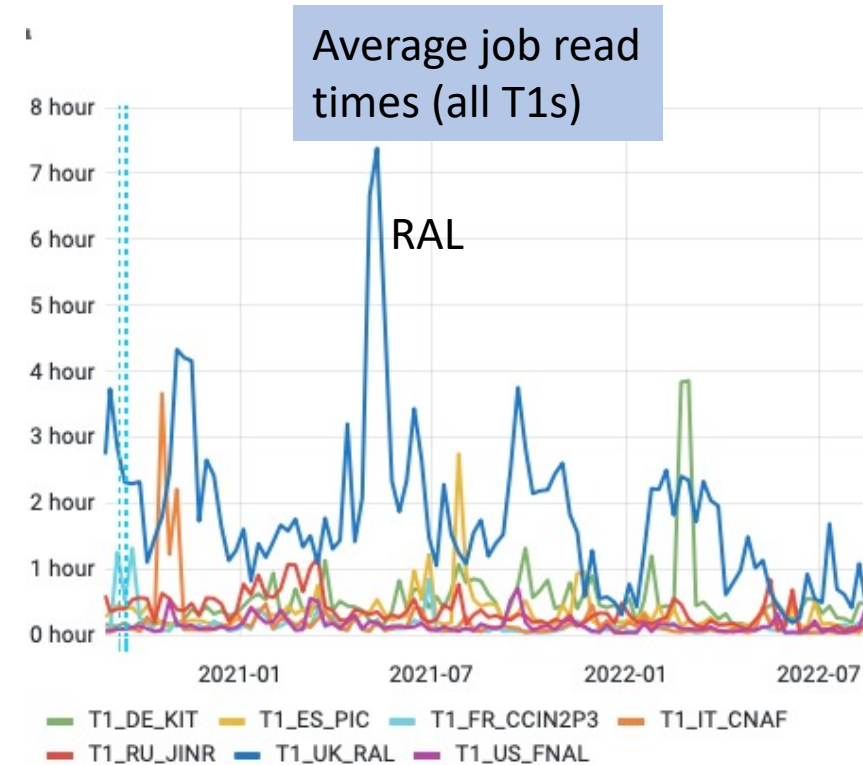RAL (weekly average)

2 years ago

Now

(CPU efficiency = CPU time / Core time)

The changes made to the network and batch farm caused continuous improvement in 2021/2022 and this has been sustained over the last ~6 months

4

# Job read volumes and times


Average job data input (just RAL)


Average job data input (all T1s)

RAL

KIT and RAL sit above the other T1s most of the time due to use of the 'Lazy-Download' setting


Average job read times (all T1s)

RAL

However, the read time at RAL is always longer…but signs of improvement are good…

# Different job types at RAL T1



Average CPU efficiency

Running cores

2 years ago

It will be interesting to see if performance is maintained when CMS switches back to Processing-dominated workflows

# Comparison between WN tranches



- Consistent behaviour over last month.
- New tranche has joined in recent weeks
- Slightly worse efficiencies from the oldest tranches – but not bad considering!

# Tape families

- How should we group tape data so it can be efficiently written, read, stored and deleted?
  - Also consider how many tapes we should write to simultaneously? (e.g. if CMS has access to 8 drives)

Write continuously to any tape with 8 drives?

Put every kind of data on a different tape?

**Maximum** write and space efficiency!

**Minimum** read and deletion efficiency ☹

**Maximum** read and deletion efficiency!

**Minimum** write and space efficiency. ☹

# Tape families

- How should we group tape data so it can be efficiently written, read, stored and deleted?
    - Also consider how many tapes we should write to simultaneously? (e.g. if CMS has access to 8 drives)

Write continuously to any tape with 8 drives?

Optimum

Put every kind of data on a different tape?

**Maximum** write and space efficiency!

**Minimum** read and deletion efficiency ☹

**Maximum** read and deletion efficiency!

**Minimum** write and space efficiency. ☹

# Tape families

- How should we group tape data so it can be efficiently written, read, stored and deleted?
  - Also consider how many tapes we should write to simultaneously? (e.g. if CMS has access to 8 drives)

Write continuously to any tape with 8 drives?

Optimum

Put every kind of data on a different tape?

**Maximum** write and space efficiency!

**Minimum** read and deletion efficiency ☹

**Maximum** read and deletion efficiency!

**Minimum** write and space efficiency. ☹

- I consulted with CMS colleagues on how data is read from tape and determined that 'data tier' (data type), e.g. raw, partly-processed, fully-processed, was the best way to split 2022 experiment data
  - It would be good to verify this via monitoring and analysis

# Shoveler (new XRootD monitoring)

- https://github.com/opensciencegrid/xrootd-monitoring-shoveler
- Useful for monitoring of CMS' 'AAA' service, which allows CMS jobs to access Any Data, Any Place, Any Time
  - Previously not well-monitored (but high on CMS's 'wanted' list)
- Running as a test at RAL Tier 1 since end of May on the AAA proxy host machines
- Monitoring only available in Kibana so far

# Shoveler plot?

MB/s

Shoveler

2022-08-03 23:00   2022-08-06 23:00   2022-08-09 23:00   2022-08-12 23:00   2022-08-15 23:00   2022-08-18 23:00   2022-08-21 23:00   2022-08-24 23:00   2022-08-27 23:00   2022-08-30 23:00

per 60 minutes

● ceph-gw10 **0.167**     ● ceph-gw11 **0.145**

RAL internal monitoring of ceph-gw10

Data rate gw10

Data from Echo to proxy     Data from proxy to client

12

# List of recent Rucio/FTS fixes/improvements

- Immediately failure of an entire FTS batch job if one file was missing from the source

- 'Destination file exists' transfer errors   Fix partially in use
  - File write is attempted. Maybe successful, maybe not...
  - If successful, no message passed back to tell FTS
  - FTS checks the file
    - Fix Part 1 – If file has correct checksum then next transfer attempt is marked as success.
      IT WORKS!
    - Fix Part 2 – If the file does not have correct checksum then allow overwrite
      NOT (YET) APPROVED

# Reminder of multihop in Rucio/FTS transfers

Rucio knows data
is here

Rucio rule wants
to have data here

Source site ⟶ 'Middle hop' site ⟶ Destination site

'Multihop' transfer:
1 FTS job

(Source or Destination site may be isolated from all other sites, and only connected to the world via the multihop/middle hop site)
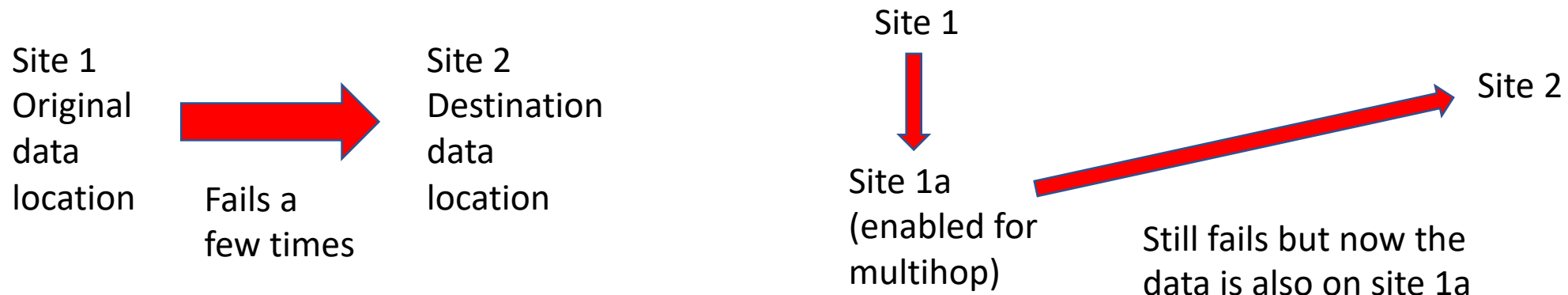
# List of recent *multihop* fixes/improvements

- After initial submission Rucio 'forgets' transfer is multihop <span style="color:green">FIXED in Rucio</span>

- No overwrite allowed on 'middle hop'
  - Overwrites typically not permitted for CMS tape systems, but middle hop locations are typically disks. <span style="color:green">FIXED in FTS+Rucio</span>

- Final destination failures causing unnecessary additional hops
  - Files try to transfer through _Test sites. <span style="color:green">FIXED by simple config change</span>
  - Additional files being stored – we cannot afford the space, particularly on CERN-EOS. <span style="color:red">NOT FIXED</span>

# List of recent multihop fixes/improvements

- After initial submission Rucio 'forgets' transfer is multihop FIXED in Rucio

- No overwrite allowed on 'middle hop'
  - Overwrites typically not permitted for CMS tape systems, but middle hop locations are typically disks.      FIXED in FTS+Rucio

- Final destination failures causing unnecessary additional hops
  - Files try to transfer through _Test sites.      FIXED by simple config change
  - Additional files being stored – we cannot afford the space, particularly on CERN-EOS.      NOT FIXED

Site 1
Original
data
location

Fails a
few times

Site 2
Destination
data
location

Site 1

Site 1a
(enabled for
multihop)

Site 2

Still fails but now the
data is also on site 1a

# Consistency checking

- CMS scripts compare the Rucio database of files with a file list from the sites

- Compiles two lists:
  - 'Missing' files – those that should be at the site but are not
  - 'Dark' files – those that should not be at the site, but are

- For missing files:
  - Asks Rucio to copy the file using an alternative replica. If none exists then remove the file from Rucio

- For dark files:
  - If files appear repeatedly (e.g. 4 consecutive weeks) in the list, then delete the files from site

# CRAB and support

- 'CMS Remote Analysis Builder'  a.k.a., User analysis software
- Development: Small features to improve user experience
- Maintenance: Responding to a changing environment
- Documentation and user support
- (Katy) new feature testing
  - GPU tests soon??

# Token support

- Of course, WLCG has milestones in mind which apply to CMS
  - https://zenodo.org/record/7014668
- On the CMS side, there is a Roadmap – some small changes made
- There are currently 'brain-storming' type meetings
  - Discussion on how to open this up to more people, but then focus it down to the right people
  - Several computing coordination groups need to take part, and we will require cross-coordination between groups
  - Estimation of person-power is required

# Conclusions

- Tier 1 has been broadly operationally stable since GridPP47
- The reported performance improvement has been sustained

# Conclusions

- Tier 1 has been broadly operationally stable since GridPP47
- The reported performance improvement has been sustained
- Run 3 data is arriving and assigned to appropriate tape families (automated system)
- XRootD monitoring is lacking, but efforts to improve this
- CMS continues to observe Rucio and respond to undesirable behaviour
- CMS token support in early stages