**PAPER • OPEN ACCESS**

# Storageless and caching Tier-2 models in the UK context

To cite this article: Samuel Cadellin Skipsey *et al* 2017 *J. Phys.: Conf. Ser.* **898** 062047

View the article online for updates and enhancements.

# Storageless and caching Tier-2 models in the UK context

**Samuel Cadellin Skipsey[1], Alastair Dewhurst[2], David Crooks[1], Ewan MacMahon[3], Gareth Roy[1], Oliver Smith[4], Kashif Mohammed[3], Chris Brew[4], David Britton[1]**

[1]University of Glasgow, Glasgow, G12 8QQ, UK
[2]STFC, RAL, Oxon, UK
[3]University of Oxford, Oxford, UK
[4]University of Durham, Durham, UK

E-mail: `samuel.skipsey@glasgow.ac.uk`

**Abstract.** Operational and other pressures have lead to WLCG experiments moving increasingly to a stratified model for Tier-2 resources, where "fat" Tier-2s ("T2Ds") and "thin" Tier-2s ("T2Cs") provide different levels of service. In the UK, this distinction is also encouraged by the terms of the current GridPP5 funding model. In anticipation of this, testing has been performed on the implications, and potential implementation, of such a distinction in our resources. In particular, this presentation presents the results of testing of storage T2Cs, where the "thin" nature is expressed by the site having either no local data storage, or only a thin caching layer; data is streamed or copied from a "nearby" T2D when needed by jobs.

In OSG, this model has been adopted successfully for CMS AAA sites; but the network topology and capacity in the USA is significantly different to that in the UK (and much of Europe).

We present the result of several operational tests: the in-production University College London (UCL) site, which runs ATLAS workloads using storage at the Queen Mary University of London (QMUL) site; the Oxford site, which has had scaling tests performed against T2Ds in various locations in the UK (to test network effects); and the Durham site, which has been testing the specific ATLAS caching solution of "Rucio Cache" integration with ARC's caching layer.

## 1. Introduction
Funding and political pressures, both within the UK, and in the wider WLCG, mandate a reduction in costs, and an increase in "efficiency". It is seen that this "efficiency" might be realised by a reduction in manpower costs at smaller Tier-2 sites. There have been various attempts at designing models for reducing the manpower cost for the compute component of WLCG operations at T2 sites; Additionally, the above funding and political pressures also mandate greater cooperation with and between research groups outwith the WLCG stakeholders. We are also mandated to develop means of operation which make it easier to support non-WLCG user groups; and changes to/alternatives to approaches which make our systems more flexible and useful to such groups. Technologically, trends for storage are also changing: spinning magnetic media is getting increasingly dense, with consequences for the optimal design of storage solutions; Flash and other solid-state storage technologies are becoming cheaper and denser

(but also less reliable, at the cheapest/densest end). As hard disk capacities grow, and thus the time needed to recover from hardware failures with RAID increases, the space of solutions has shifted increasingly to storage of data distributed across multiple servers - either explicit redundant copies, or as striped erasure coded blocks ("RAIN" or "RAIS", as opposed to the locally-resilient RAID). Our storage solutions have always supported RAIN-type redundancy - although only when users explicitly make multiple copies - and we must be careful that any solution we adopt is consistent with the increasing importance of this, for any part of the storage infrastructure which we consider "resilient". (Clearly, for non-resilient storage, this is not important.) For the purposes of this paper, we will discuss modifications of the "uniform" Tier-2 site into specialisations. In general "T2D" will represent a "fat" Tier-2 site, with large amounts of storage and compute, and a dedicated Storage Element (SE) service, similar to the original Tier-2 site design. "T2C" will represent a "slim" Tier-2 site, with reduced storage capability, but retaining compute provision. This can be seen as a weakening of the assumption that storage and compute will be co-located.

## 2. Ongoing Work
### 2.1. ARC Cache for T2Cs
The ARC[1] Compute Element, developed by NorduGrid, was designed from the start to support distributed Tier-2s, and thus has mechanisms to minimise data movement for jobs available. In particular, any ARC CE can be configured to prefetch data dependancies for a submitted workload in a site-local filesystem, only submitting the dependant job to the local batch system when the data is fully locally available.

This "ARC Cache" is not widely used outside of the NorduGrid project itself. One reason for this is the dependance of the larger Virtual Organisations on "pilot jobs" for workload submission - the jobs seen by the site ARC will be the pilots (which have no data dependancies), and the actual workloads (which are pulled from remote sources by instantiated pilot jobs) are not visible to the ARC. In fact, generic pre-caches like the ARC Cache are usually difficult to integrate with late-binding abstractions like pilots.
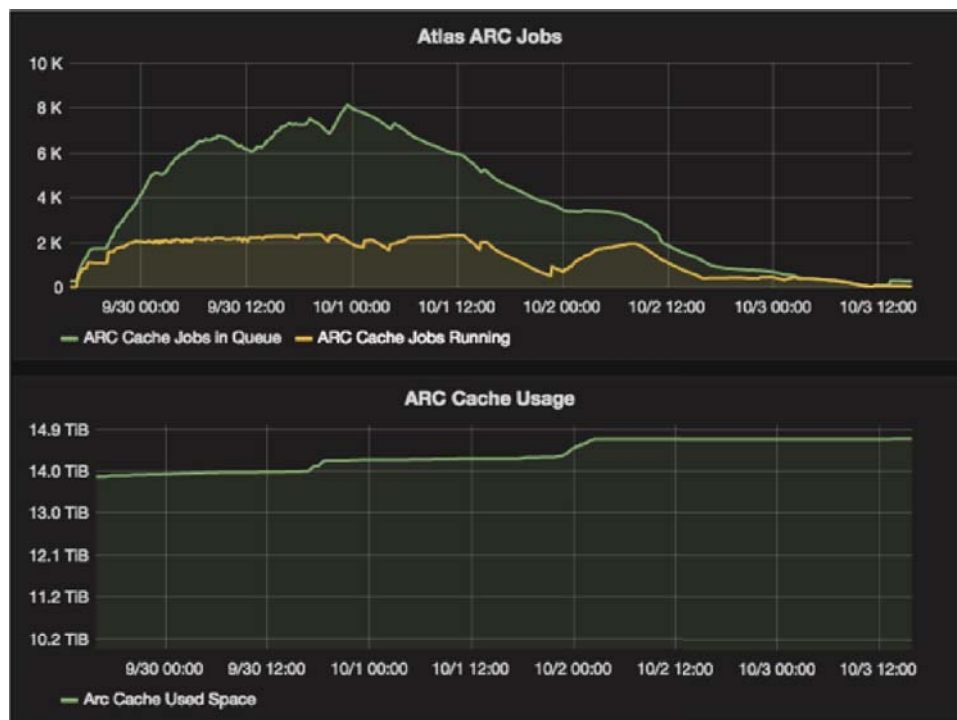
The ATLAS VO, however, has performed significant work on their "ARC Control Tower" (ACT)[2], essentially a conversion shim which interfaces with the ATLAS pilot system to pull real workloads, and then submit them directly to ARC sites which wish to use caching. This mechanism has been successfully used to submit to the NorduGrid project resources for many years.

Whilst the ScotGrid Tier-2 had been an early adopter of ARC CEs (and had attempted to promote ARC Caching at that time), the recent suggestions by the ATLAS experiment that further integration with ARC Caches might be investigated prompted another look at testing the ease of use of the caches in a UK context.

The Durham Tier-2 site, which has a particularly large CPU count in comparison to its storage, was an obvious candidate for testing this use case. Since July 2016, the site has configured their ARC CE to enable caching, against a shared filesystem (exported over NFS to the worker nodes), and has been added to the ATLAS ACT list so that they can receive jobs from that service.

Results (presented in the ATLAS Sites Jamboree, January 2017 [3]), show that workloads from the ACT are at least as efficient as conventional ATLAS pilots (which the site continues to receive). However, the total storage used by the cache is comparatively small, as figure 1 shows.

Despite the success of our trial, and the simplicity of implementing the caching in an existing ARC CE, this approach is limited by the lack of equivalents to the ACT for other pilot-dependant VOs. That said, the significant number of ATLAS-supporting sites in the UK means that this option is still of use for reducing the impact of ATLAS-related workloads.

**Figure 1.** Snapshot of live monitoring of UKI-SCOTGRID-DURHAM site's ARC Cache, with number of running and queued jobs submitted via the ATLAS ACT (and thus dependant on the ARC Cache). Changes in the cache utilisation are usually followed by the jobs dependant on the newly cached data switching to a running state. The cache grows monotonically in this snapshot as it is still much smaller than the available space on the storage system.
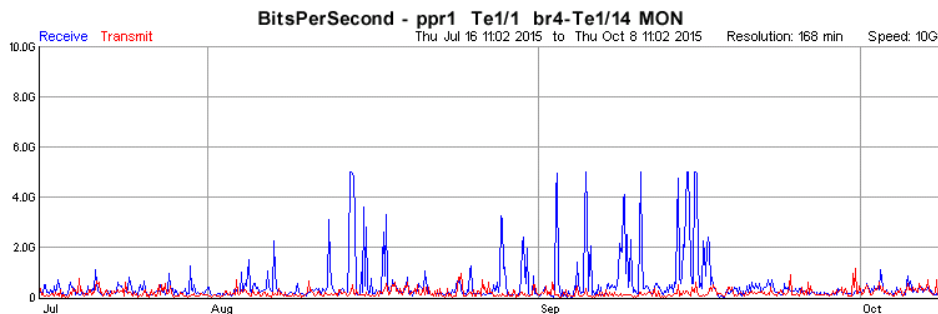
*2.2. Federation with "close" T2D*
Both the ATLAS and CMS VOs have expressed some interest in smaller sites simplifying by removing their SEs entirely, in favour of using a larger, nearby, sites' storage. The simplification on the VO's side is in the reduction of the number of storage locations that must be managed by their tools. On the smaller sites' side, having no local storage at all is a clear reduction in the number of managed systems, software and hardware wise.

In the UK, there are multiple test examples of such setups: within the LondonGrid Regional Tier-2, the University College London site has been running ATLAS workloads for several years without any local SE - instead treating the storage at Queen Mary University of London (QMUL) as local. This works well, for the subset of workloads which the UCL site receives from ATLAS. However, the job mix is mostly production work with low data requirements, and the networking environment is ideal, as the London area is the best served location for academic network links in the country.

The SouthGrid Bristol site reported results, in 2015, of a local user attempting to run CMS analysis workloads with wholly remote, direct, data access. The figures produced, figure 2, caused some concern, as they suggest a clear scaling limit for efficiency at the site, around 250 jobs (filling the site's 5 Gbit/s share of the University's 10 Gbit/s external link).

Work has been attempted, since then, to replicate these tests with ATLAS workflows, but has been significantly inhibited by staff movement, and by issues with the ATLAS orchestration systems themselves. Initial tests were planned to involve completely replicating the CMS approach - running jobs which were completely agnostic to the location of the data they needed

**Figure 2.** Network load at the UKI-SOUTHGRID-BRISTOL site during tests by a local CMS researcher of job execution using the CMS remote file access facility on the local cluster. The peaks at 5Gbit/s represent the maximum bandwidth of the site, and occurred at around 250 jobs concurrently executing.

- using the ATLAS HammerCloud framework. However, HammerCloud was not at the time (and is still not now, to our knowledge) capable of operating in such a manner; additionally, the ATLAS data management and job orchestration systems were not capable of associating jobs with a site completely ignoring data location.

More recently, work was done at the Oxford and RALPP (RAL Tier-2) sites to partly replicate the Bristol test with CMS workflows. The Oxford site was associated with the RALPP SE, for all CMS work, via modifications at the CMS Site configuration.

Preliminary graphs of data transfer rates from RALPP, and job numbers and efficiencies at Oxford, (figure 3) suggest that there is a complex relationship between the number of jobs, the data flows induced, and the resulting efficiency impact. Work is ongoing to more finely characterise the mix of jobs running on the Oxford site at any point in time, for correlation with the other observed data.
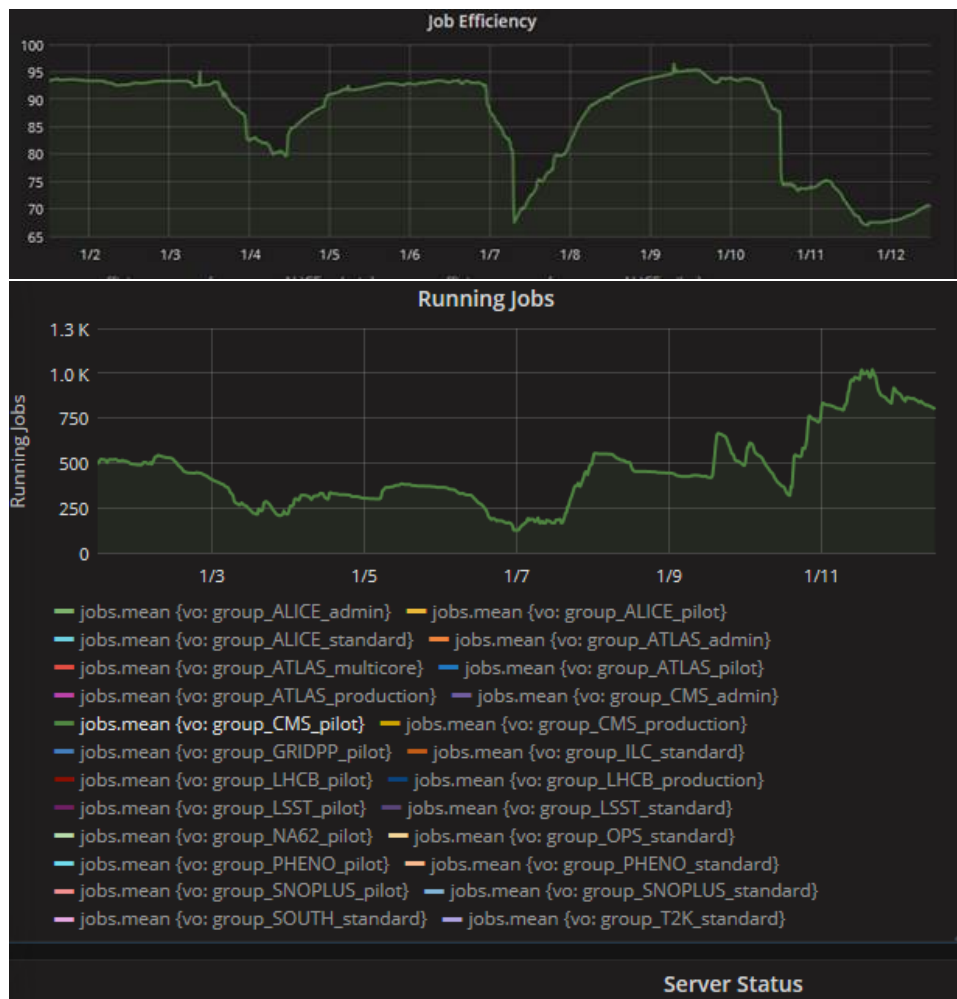
## 3. Future Work
In addition to the presented statistics, we have other projects in development to investigate other modes of simplification.

### 3.1. Xrootd Federation
The Xrootd [4] protocol is designed naturally to support hierarchical operation, supporting a tiered tree of redirector services for consolidating access to multiple underlying repositories. Indeed, this is in use by both the ATLAS and CMS VOs as part of their failover and remote access approach for sites: the ATLAS 'FAX' [5] and CMS 'AAA' [6] initiatives both being built upon hierarchies of Xrootd services. Whilst the intent for CMS is to allow access of any data, any time, any where; the ATLAS use-cases have always been more focussed on failure recovery (if the file is not available locally, the FAX redirectors allow for location of a remote copy, at the cost of additional latency).

However, another use for Xrootd Federation would be for provision of "opaque" Tier-2 storage federations. By presenting a single XRootd endpoint, itself redirecting across the combined storage elements of the Tier-2 sites, the view of external VOs would be considerably simplified from the current plethora of endpoints visible. This is especially true for the ATLAS VO, which has traditionally had more Tier-2 supporting sites than CMS within the UK. As sites within a UK Tier-2 are assumed to be relatively network-close, redirection between the spokes of such a federation should be near best-case for latency and bandwidth. Care, however, must be taken

**Figure 3.** Concurrent CMS jobs at the UKI-SOUTHGRID-OXFORD site, running with their "local storage" set to RALPP's storage endpoint, alongside the measured efficiencies for those jobs. Note the reduction in efficiency when running around 1000 jobs, which may be due to network starvation.

here: the UK academic network is not a full mesh, and connectivity between physically close locations is not always as direct in network terms. A good example is the environment between the two Scottish Universities of Glasgow and Edinburgh.

A pilot of such a framework, connecting the ScotGrid Tier-2 sites (Durham, Edinburgh and Glasgow) is currently under development, with services installed at Edinburgh and Durham at present. With VOs increasingly moving away from requiring a complete SRM at a site, our belief is that a light-weight, Xrootd only, service should be sufficient for access for many job types.

The additional benefit of a pure-Xrootd approach is that caching layers can also be introduced, via intermediate Xrootd services. Whilst caching data on transfer is a low-cost optimisation, it is not clear how advantageous it will be. Results in table 1 from analysis of the number of accesses on files from various VOs at the Glasgow site suggest that the majority of files are accessed only once or twice, minimising the benefits of a cache. For some VOs, for example CMS, the benefits appear stronger, so this decision may depend on the mix of VOs supported

by a site.

**Table 1.** Histogram of file popularity for 50 Virtual Organisations at UKI-SCOTGRID-GLASGOW Tier-2, derived from access logs in DPM database, for data up to August 2016. Note that test files used for common functional tests are the most common component of the 10+ bin, and are not removed from this data.

| | Number of Accesses per file | | | | | | | | | |
| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10+ |
|---|---|---|---|---|---|---|---|---|---|---|
| ATLAS | 2784430 | 810069 | 126566 | 38919 | 18358 | 15855 | 6205 | 1687 | 1086 | 9476 |
| CMS | 7 | 9600 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 252 |
| Biomed | 37953 | 541 | 1056 | 287 | 203 | 149 | 82 | 44 | 20 | 87 |
| Mice | 14363 | 4954 | 313 | 137 | 3 | 0 | 0 | 0 | 0 | 55 |
| Pheno | 11027 | 2807 | 1344 | 454 | 16 | 6 | 2 | 0 | 0 | 1 |

*3.2. Data CVMFS*

The Cern Virtual Machine File System[7], CVMFS, has been a very successful innovation in software delivery across the WLCG. However, it has previously had limitations: the use of unsecured HTTP, whilst allowing easy caching in tier via Squid services, also means that it is impossible to distribute non-open data with it; the assumptions made in HTTP, and most web browsers, also make it much less efficient to distribute larger files.

Work done at Universities of Nebraska and Fermilab[8] has resulted in a series of commits to CVMFS, allowing a new 'Data CVMFS' operational mode. This was developed initially to support researchers at LIGO in their use of WLCG resources in analysis, where the datasets are relatively small, and relatively static, but are exceptionally valuable (and need to be secured). Data CVMFS essentially adds HTTPS support, including support for X509 certificates used on the Grid, as well as providing data movement over Xrootd/HTTP.

While we do not anticipate widespread applicability of Data CVMFS to the majority of customers of GridPP, we are planning a series of small tests of the service, potentially with UK LIGO. As mentioned earlier in the paper, the network environment of the UK (and Europe) differs significantly from that of the academic USA - with less bandwidth, and more complex hierarchies. Discussion with Brian Bockelman suggested strongly that testing is needed before assuming that the same model would work as well [9]. Potential developments to support operation in a more constrained network environment might include reintroduction of local cache tiers (which would need to be implemented separately to Squids, as HTTPS traffic is opaque by design), or other smarter data placement approaches.

**4. Summary**

Whilst some work is ongoing, the work in the UK on developing and exploring ways of running T2C storage has already produced useful results. We're confident that caching models, even if not ARC CE based, are suitable solutions to simplify at least some smaller Tier-2 sites in the UK, and a "taskforce" to engage with smaller sites (which also have less manpower to enact changes) has been assembled to progress with this in the near future.

**references**
[1] Ellert M et al. 2007 *Future Gener. Comput. Syst.*, **23**(1):219240 . http://doi.org/10.1016/j.future.2006.05.008

[2] Nilsen J K et al 2015 *J. Phys.: Conf. Ser.* 664 062042 doi:10.1088/1742-6596/664/6/062042
[3] ATLAS Sites Jamboree, Indico https://indico.cern.ch/event/440821/
[4] XRootd project page: http://www.xrootd.org/
[5] Gardner R et al. 2014 *J. Phys.: Conf. Ser.* 513 042049 doi:10.1088/1742-6596/513/4/042049
[6] Bloom K et al. 2014 *J. Phys.: Conf. Ser.* 513 042005 doi:10.1088/1742-6596/513/4/042005
[7] Blomer J et al. 2012 *J. Phys.: Conf. Ser.* 396 052013 doi:10.1088/1742-6596/396/5/052013
[8] Dykstra D et al. 2015 *J. Phys.: Conf. Ser.* 664 042012. doi:10.1088/1742-6596/664/4/042012
[9] *Private conversation*, B Bockelman and the author, November 2016