

Optimising network transfers to and from Queen Mary University of London, a large WLCG tier-2 grid site

C J Walker¹, D P Traynor¹, D T Rand², T S Froy¹ and S L Lloyd¹

¹Queen Mary University of London, Mile End Road, London E1 4NS, UK

²Imperial College London, Prince Consort Road, London SW7 2BW, UK

E-mail: C.J.Walker@qmul.ac.uk

Abstract. Optimising network performance is key to high bandwidth data transfers required for a Tier-2 site. We describe the techniques we have used to obtain good performance. Monitoring plays a key part, as does the elimination of bottlenecks and tuning TCP window sizes. Multiple parallel transfers allowed us to saturate a 1 Gbit/s link for 24 hours - whilst still achieving acceptable download speeds. Source based routing and multiple data transfer servers allowed us to use an otherwise unused “resilient” link.

1. Introduction

Analysis of the large quantities of data from the Large Hadron Collider (LHC) is performed using a distributed network of computing centres, the Worldwide LHC Computing Grid (WLCG). Queen Mary University of London (QMUL) hosts a Tier-2 WLCG site with 1.8 PB of storage. Making optimum use of Wide Area Network (WAN) links is key to filling the storage - which would take approximately six months at an average of 1 Gbit/s, our nominal WAN capacity until September 2012. This six months is commensurate with the reprocessing cycle of the experimental data.

As can be seen from figure 1, data transfers, and hence bandwidth requirements, are continually increasing. Also notable are the large fraction of international and intercontinental transfers - which are particularly sensitive to packet loss and TCP tuning. Increased use of technologies such as WebDAV[1] and xrootd[2] that enable remote interactive access to the storage, coupled with federated storage technologies such as FAX[3], may also increase bandwidth requirements.

A previous paper[4] describes the tuning of network and disk access within the cluster. This paper describes the steps we have taken to ensure high performance of data transfers over the WAN.

2. Network Monitoring

Monitoring is key to understanding, and therefore improving, network performance. We use a variety of monitoring and diagnostic tools to help us in this task.



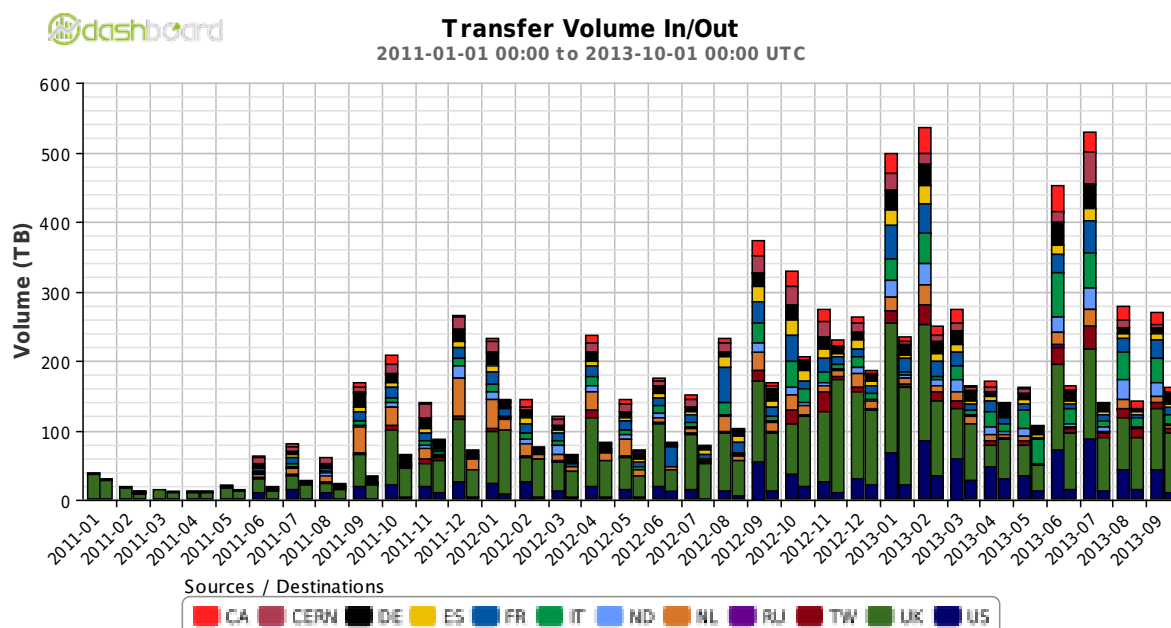


Figure 1. The increase in WAN data transfers, both into and out of QMUL, from January 2011 to September 2013. The network was upgraded to 10 Gbit/s in September 2012. Note that a significant proportion of traffic is from/to sites outside the UK.

2.1. Active network monitoring

We have deployed both PerfSONAR[5] and RIPE[6] probes for active network monitoring. Three PerfSONAR machines are deployed, one for bandwidth monitoring, one for latency, and a third to test IPv6 [7] and jumbo frames performance (see section 4). The RIPE probe performs a similar job to the PerfSONAR latency instance, but is deployed at many more sites.

2.2. Passive network monitoring

Figure 2 shows monitoring of the WAN link provided by JANET, the UK National Research and Education Network (NREN). This monitoring allows us to monitor traffic over the WAN link, and therefore see whether this is the limiting factor. The figure shows saturation of a 1 Gbit/s link in March 2012 (top), so this clearly is the limiting factor here, but in February 2013 (bottom), the average rate inbound is 6 Gbit/s on a 10 Gbit/s link, so the bottleneck lies elsewhere - and may be outside QMUL.

The ATLAS experiment monitors data transfer rates for individual files between grid sites. Figure 3 shows rates for files larger than 1 GB from Taiwan. There is a clear reduction in transfer rates for individual files between 7 and 19 September. The reduction coincided with the upgrade of the WAN link from 1 Gbit/s to 10 Gbit/s, but only affected transfers from some sites, while others saw increased transfer rates. A PerfSONAR host deployed at the Taiwan Tier-1 enabled us to establish that traffic to QMUL was not being routed via the preferred route, but was taking a congested fall back route instead. Further debugging using RIPE probes established that this was because route advertisements from QMUL were not being propagated correctly. Once this problem was fixed, transfer rates were higher than before the network upgrade.

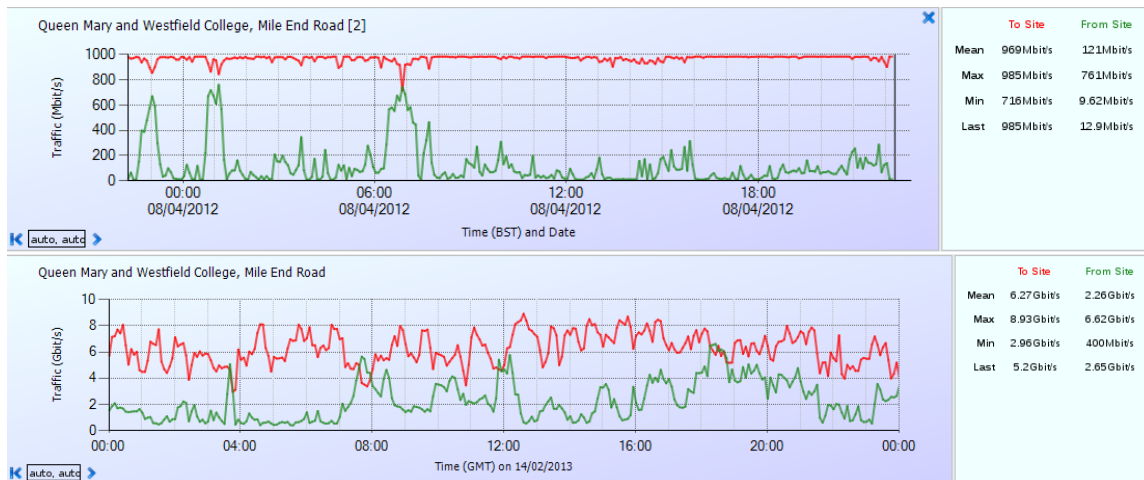


Figure 2. Network traffic: March 2012 saturating a 1 Gbit/s network (top), and February 2013 averaging over 60% of a 10 Gbit/s link(bottom).

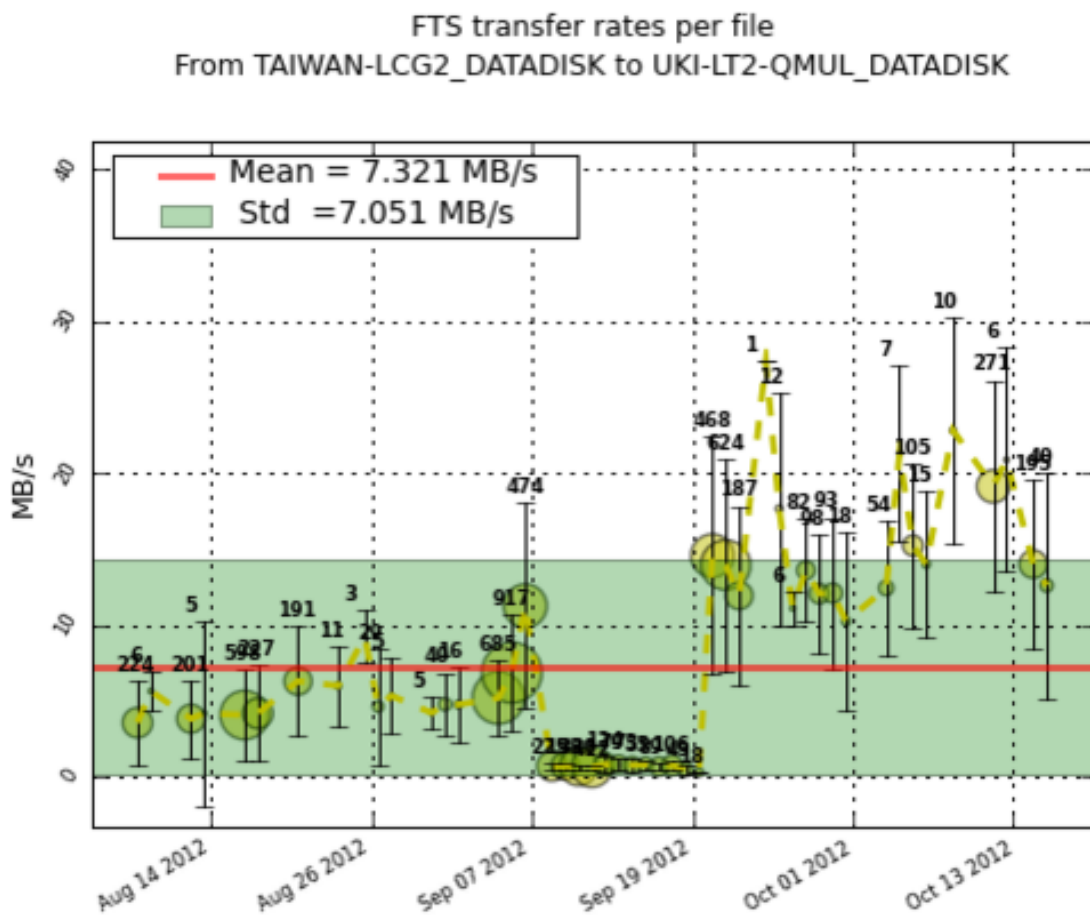


Figure 3. Transfer rates from Taiwan to QMUL for files larger than 1GB. Numbers by each data point give number of files. Note the decrease between 7th September and 19th September when traffic took a suboptimal route.

3. Network tuning

The monitoring described previously allowed us to measure actual network performance. Comparing actual with expected performance enabled us to find and fix bottlenecks. Our strategy was to obtain high performance on individual transfers, then to increase the number of parallel transfers to increase aggregate data rates.

3.1. QMUL Tier-2 Cluster

The Tier-2 has 3500 CPU cores, 1.8PB of Lustre[8] disk storage, and a 10 Gbit/s connection to the internet. The cluster has been optimised for parallel IO using the Lustre filesystem, and should not present a bottleneck to data transfers. Details of the IO performance of the cluster are described in the previous paper[4]. Two data transfer nodes are used, with GridFTP[9] the main protocol used to transfer data over the WAN.

3.2. Buffer sizes

High bandwidth links to international and intercontinental sites have high latency, and hence large amounts of data need to be in transit unacknowledged to obtain maximum performance. We have therefore followed the ESnet[10] recommendations for high bandwidth delay product links and increased TCP buffer sizes to obtain good performance.

3.3. Router technology

The commercial router initially used proved incapable of advertising routes while handling 1 Gbit/s of traffic. The resulting “route flapping” caused considerable performance degradation. A replacement Linux PC based router running Quagga[11] was deployed. This has recently been upgraded to a Xeon X5670 CPU with Intel X520 10 Gbit/s network cards and is capable of handling our 10 Gbit/s WAN link.

3.4. Source based routing

We were able to make use of a backup 1 Gbit/s link by using two data transfer servers on different IP addresses. Inbound traffic was routed separately by advertising a route with higher priority via Border Gateway Protocol (BGP) for a subset (a “/27”) of our address space. Source based routing was then used to ensure that outbound traffic used both links.

One problem we encountered was that the link to our machine room was provisioned as a bonded pair of 1 Gbit/s links, and all the traffic was going down one of these links resulting in performance still being capped at 1 Gbit/s. To fix this, we changed the MAC address of one of the data transfer servers. This allowed us to use both halves of the bonded link.

Since upgrading to a resilient pair of 10 Gbit/s links, this technique has been used to route our bulk data traffic over a different link to general university traffic. This ensures that any additional latency introduced by our traffic filling the link doesn’t affect interactive applications used by general university users.

4. Future work

4.1. Jumbo frames

In principle, jumbo frames (a Maximum Transmission Unit (MTU) greater than the default of 1500 bytes - typically 9000), reduce overhead and permit increased performance. Performance improvements are most likely on latency limited links, but only if jumbo frames are enabled end to end. Performance measurements are ongoing, but preliminary results do not indicate an increase in transfer failures.

4.2. IPv6

We currently have two monitoring machines - a RIPE probe, and a PerfSONAR host running dual stack IPv4/IPv6. Performance has generally been good. We have, however, found poor performance over IPv6 to one site due to a routing problem. This has now been fixed, but illustrates the importance of testing. We plan to deploy some production machines as dual stack in the near future.

5. Conclusions

We have carefully monitored our network performance and found and fixed bottlenecks causing packet loss. We have also tuned the TCP stack of our data transfer nodes and run multiple transfers in parallel. We are now able to transfer large amount of data at high speeds over long distances.

References

- [1] Furano F, Alvarez A, Devresse A, Hellmich M P and Manzi A, 2013 An HTTP Ecosystem for HEP Data Management *CHEP 2013*
- [2] Xrootd (<http://xrootd.org/>)
- [3] Vukotic I, On behalf of the ATLAS Collaboration, 2013 Data Federation Strategies for ATLAS using XRootD *CHEP 2013*
- [4] Walker C J, Traynor D P and Martin A J, 2012 Scalable Petascale Storage for HEP using Lustre *J. Phys.: Conf. Ser.* **396** 042063
- [5] Hanemann A, Boote J W, Boyd E L, Durand J, Kudarimoti L, Lapacz R, Swany D M, Zurawski J and Trocha S, 2005 PerfSONAR: A Service Oriented Architecture for Multi-Domain Network Monitoring *Proceedings of the Third International Conference on Service Oriented Computing*, **3826** 241–4 (<http://www.perfsonar.net/>)
- [6] RIPE atlas probe <http://atlas.ripe.net/>
- [7] Deering S and Hinden R, 1998 Internet Protocol, Version 6 (IPv6) *RFC 2460* <http://www.ietf.org/rfc/rfc2460.txt>
- [8] *The Lustre Filesystem* - <http://www.lustre.org>
- [9] Bresnahan J, Link M, Khanna G, Imani Z, Kettimuthu R and Foster I 2007 *Proceedings of the First International Conference on Networks for Grid Applications (GridNets 2007)*
- [10] ESnet Fasterdata Knowledge base <http://fasterdata.es.net/>
- [11] Boggis SA 2010 *JANET Networkshop 38 Manchester*