

The production deployment of IPv6 on WLCG

J Bernier¹, S Campana², K Chadwick³, J Chudoba⁴, A Dewhurst⁵, M Eliáš⁴, S Fayer⁶, T Finnern⁷, C Grigoras², T Hartmann⁷, B Hoefft⁷, T Idiculla⁵, D P Kelsey⁵, F López Muñoz⁹, E Macmahon¹⁰, E Martelli², R Nandakumar⁵, K Ohrenberg⁷, F Prelz¹¹, D Rand⁶, A Sciabà², U Tigerstedt¹², R Voicu¹³, C J Walker¹⁴ and T Wildish¹⁵

¹ IN2P3 Computing Centre, Boulevard du 11 Novembre 1918, F-69622 Villeurbanne Cedex, France

² CERN, CH-1211 Genève 23, Switzerland

³ Fermi National Accelerator Laboratory, Batavia, IL 60510, U.S.A.

⁴ Institute of Physics, Academy of Sciences of the Czech Republic Na Slovance 2 182 21 Prague 8, Czech Republic

⁵ STFC Rutherford Appleton Laboratory, Harwell Oxford, Didcot, Oxfordshire OX11 0QX, United Kingdom

⁶ Imperial College London, South Kensington Campus, London SW7 2AZ, United Kingdom

⁷ Deutsches Elektronen-Synchrotron, Notkestraße 85, D-22607 Hamburg, Germany

⁸ Karlsruher Institut für Technologie, Hermann-von-Helmholtz-Platz 1, D-76344 Eggenstein-Leopoldshafen, Germany

⁹ Port d'Informació Científica, Campus UAB, Edifici D, E-08193 Bellaterra, Spain

¹⁰ The University of Oxford, Denys Wilkinson Building, Keble Road, Oxford OX1 3RH, United Kingdom

¹¹ INFN, Sezione di Milano, via G. Celoria 16, I-20133 Milano, Italy

¹² CSC Tieteen Tietotekniikan Keskus Oy, P.O. Box 405, FI-02101 Espoo

¹³ California Institute of Technology, Pasadena, Ca 91125, U.S.A.

¹⁴ Queen Mary University of London, Mile End Road, London E1 4NS, United Kingdom

¹⁵ Princeton University, Jadwin Hall, Princeton, NJ 08544, U.S.A.

E-mail: david.kelsey@stfc.ac.uk, ipv6@hepox.org

Abstract. The world is rapidly running out of IPv4 addresses; the number of IPv6 end systems connected to the internet is increasing; WLCG and the LHC experiments may soon have access to worker nodes and/or virtual machines (VMs) possessing only an IPv6 routable address. The HEPiX IPv6 Working Group (<http://hepox-ipv6.web.cern.ch/>) has been investigating, testing and planning for dual-stack services on WLCG for several years. Following feedback from our working group, many of the storage technologies in use on WLCG have recently been made IPv6-capable. The worldwide HEP computing community now needs to deploy dual-stack IPv6/IPv4 services on WLCG to allow such use of IPv6-only resources. This paper will present the IPv6 requirements, tests and plans of each of the four LHC experiments together with the tests performed both on the IPv6 test-bed and in targeted use of WLCG production services. This is primarily aimed at IPv6-only worker nodes or VMs accessing several different implementations of a global dual-stack federated storage service. The changes required to the operational infrastructure, including monitoring and security, will be addressed as will the implications for site management. The working group will present its deployment plan for dual-stack storage services, together with other essential central and monitoring services, to start during 2015.

1. Introduction

The much-heralded exhaustion of the IPv4 networking address is with us, etc. etc.

2. Status and hurdles of the worldwide IPv4→IPv6 transition

3. The 2014 survey of IPv6 readiness at WLCG sites

4. LHC Experiment requirements and main use case

5. LHC Experiment requirements and tests

The shortage of available ipv4 addresses implies that there is a significant possibility that new large computing facilities will not be able to give ipv4 addresses to all of the machines in their network. The most likely consequence is that for these sites, worker nodes – which constitute the largest fraction of independent computing nodes will have purely ipv6 network addresses. Hence, the main use case for the LHC experiments is to enable jobs to run on these machines, access their software areas and input data and upload their outputs to various grid storages or services as needed.

The LHC experiments generally assume [LHCassumption reference] that the storage on different sites and supporting middleware [middleware reference] like the LFC will either be directly dual-stack, or support dual stack operation in some way, enabling seamless access to the storage as needed for either downloading or saving. For example, it is expected that dual-stack squid proxies will be needed for CVMFS and xrootd will soon be dual-stack, to handle storage technologies like Castor which will not be ipv4 only. The servers that the LHC experiments use to handle the grid infrastructure are / will also be dual-stack.

As an example of the above, we look at LHCb [LHCb reference] which is the experiment on the LHC, optimised for studying beauty and charm physics. LHCb uses the DIRAC [DIRAC reference] interware to manage its grid operations. The DIRAC software was coded to be able to handle both ipv4 and ipv6 addresses in late 2014, with the modifications being easy enough to make by non-expert programmers. While testing the processes on a dual-stack machine, it was found that there was a significant number of connections which were not going through to the servers which was finally traced back to a missing `enable_ipv6` option to compile python by an external provider causing errors in identifying ipv6 addresses. Using a version of the library, the problem with dropped connections went away and testing DIRAC will restart soon.

In general, testing of the grid middleware by the different experiments is going ahead as fast as possible given the manpower and time constraints of the LHC startup and the immediate issues with handling purely ipv6 worker nodes are expected to be sorted by sometime in 2016.

[LHCassumption reference] : WLCG pre-GDB - IPv6 Workshop (A. Dewhurst et al.), <https://indico.cern.ch/event/313194/session/1/contribution/3/0/material/slides/0.pdf> [middleware reference] : https://wiki.ege.eu/wiki/Middleware_products_verified_for_the_support_of_IPv6, <http://hepiv6.web.cern.ch/wlcg-applications> [LHCb reference] : LHCb Collab. (A. Alves et al.), JINST 3, S08005 (2008), <http://dx.doi.org/10.1088/1748-0221/3/08/S08005> [DIRAC reference] : DIRAC: a community grid solution (A. Tsaregorodtsev et al.), 2008 J. Phys.: Conf. Ser. 119 062048, <http://iopscience.iop.org/1742-6596/119/6/062048>

6. Testbed operation: testing FTS3/dCache

6.1. The Transfer Testbed

The transfer testbed was upgraded in March 2015. Until then, it operated with gridFTP transfers between all sites, providing a low-level test of connectivity and functionality for the almost two years that it ran. Since March 2015 the testbed uses FTS3 to initiate the transfers, moving up the middleware stack. Since FTS3 is used for the vast majority of experiment transfers in WLCG this provides an important full-stack test.

At the present time, the testbed consists of 7 storage elements at sites distributed around Europe. One is IPv6-only, the rest are all dual-stack. All the SEs are running dCache. Most

are stable installations, but one (DESY) is rebuilt every morning with the latest patches from dCache, providing a valuable regression-test for both the dCache and IPv6 teams.

As before, each site serves as both a source and a destination, with each source sending a 1 GB file to each destination. The file-size is validated at the destination using **gfal-ls**, then the destination is cleaned with **gfal-rm** and the transfer duration is recorded. Then the cycle is repeated after a short delay, to avoid abusing the hardware/network with too much traffic. Physical file names are specified using the SRM protocol.

Two FTS3 servers are deployed for the testbed, one at Imperial College and one at KIT, though currently only the one at KIT is used.

Figure. 1 shows the transfers in the FTS3 testbed so far. Most sites transfer efficiently in both directions, but the effect of the firewall at KIT on inbound traffic can be clearly seen.

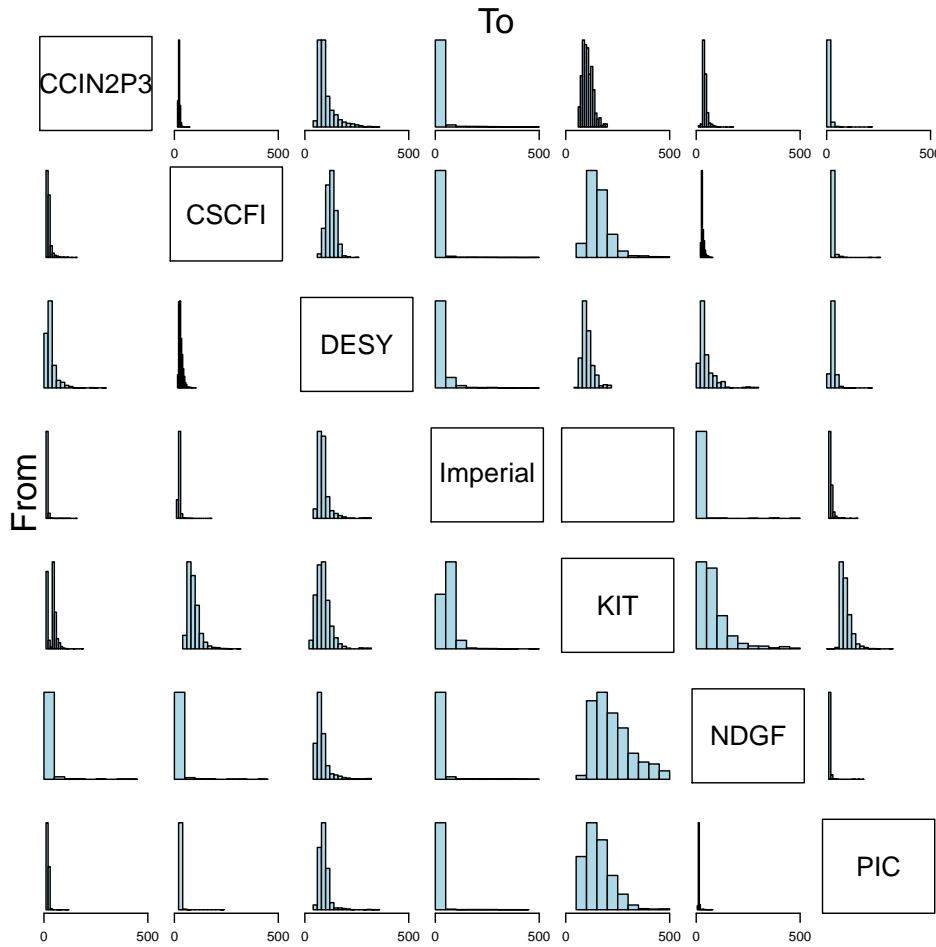


Figure 1. The FTS3 transfer testbed. Rows show transfers from the named site, columns show transfers to the destination. The horizontal axis for all plots is fixed at 500 seconds, i.e. transfers that proceed at less than 2 MB/sec will overflow.

6.2. FTS3 server and dCache SE at KIT

For managing file transfers between sites, an FTS3 instance is setup at KIT. Furthermore, a storage element based on dCache 2.10 on Scientific Linux is created. Both instances are rolled out on physical machines.

6.2.1. FTS3 FTS3 supports IPv6 in its baseline version. The service has to bind to IPv6 locally in the FTS3 conguration (IP=::) and to be enabled explicately for gfal.2 (IPV6=true). The host is available in the DNS as default dual-host, with IP4 and IPv6 announced in the A and AAAA records, and also with IPv4-only or IPv6-only names with an -ipv4 or -ipv6 appendix, respectively. Thus, all aliases have to included in the host certicate.

File transfers are successfully brokered by the FTS3 instance via IPv4 and IPv6 between the sites. Since most FTS3 instances in production use a separated database instance for performance and failsafe reasons, moving the database to a dedicated machine was tested as well. For the SQL db backends supported by FTS3, IPv6 support had been implemented in MySQL v5.5.3 and MariaDB v5.5.35, which are not available in the baseline SL6 repositories. MariaDB was installed on a dedicated host with version 5.5.42. After binding mysql locally to IPv6 as well ([mysqld] bind-address = ::), the database could be connected remotely with an IPv6-ready mysql client. For the FTS3 service to connect to the remote database via IPv6, the address has to be escaped explicitly, i.e., encapsulating the IP as [IP6]:PORT/fts3 and may depend on the version of the database library used by the FTS3 service.

6.2.2. dCache based SE A dCache instance is setup on a dedicated host. The main hurdle for file transfer access was a reverse lookup by the instance when receiving a file request. After explicately setting the dual-stack, IPv4-only, and IPv6-only host names in the dCache configuration (srm.net.local-hosts=hostname-ipv4,ipv6) and the general hosts file, the storage element is accessible via IPv4 and IPv6 as well.

7. IPv6 readiness of storage technology for WLCG

8. Status of LHCOPNv6/LHCONEv6

9. LHCONE/LHCOPN

Based on the actions initiated at Grid Deployment Board in November and December 2014, to request tier-1s to join the HEPiX-IPv6 working group and to encourage sites moving their production endpoints to dual stack even if this requires concessions of the quotation of their site relatability and site availability. The proposal of the LHC experiment Atlas was to - request that all Tier-1s has to provide ? besides an IPv6 peering to LHCOPN, ? a dual stack PerfSONAR machine by first of April 2015 - request that Tier-2s provide ? Besides an IPv6 peering to their LHCONE connection, ? a dual stack PerfSONAR machine by August 2015. At the last LHC[OPN/ONE] meeting a proposal was stated that: LHCOPN connecting Tier-1 sites to CERN gets ipv6 ready by 1. April 2015 and LHCONE connecting Tier-[123] sites geting ipv6 ready in August 2015. No objections to this proposal were presented. The url: <http://maddash.aglt2.org/maddash-webui/index.cgi?dashboard=Dual-Stack>

10. IPv6 perfSONAR measurements

11. Outlook and future plans

References

- [1] <http://hepex-ipv6.web.cern.ch>
- [2] All Internet Engineering Task Force Requests For Comments (RFC) documents are available from URLs such as <http://www.ietf.org/rfc/rfcNNNN.txt> where NNNN is the RFC number, for example <http://www.ietf.org/rfc/rfc2460.txt>
- [3] See for instance <http://www.google.com/ipv6/statistics.html>. The 5% global connectivity threshold was crossed in January 2015.
- [4] Salichos M, Keeble O, Alvarez Ayllon A, Kamil Simon M 2013 FTS3 - Robust, simplified and high-performance data movement service for WLCG *also presented at CHEP 2013*