

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30

WLCG IPv6 deployment strategy

HEPiX IPv6 Working group

July 18, 2016
Version 0.1

Executive Summary

This document describes the Worldwide LHC Computing Grid's (WLCG) strategy, as proposed by the HEPiX IPv6 Working Group, to allow sites to provide IPv6 resources to the LHC experiments. In summary:

- From April 1st 2017 sites can provide IPv6-only CPU resources, if necessary.
- From April 1st 2017 sites can provide IPv6-only interfaces to their CPU resources, if necessary.
- Any site wishing to deploy IPv6-only CPU resources should contact the HEPiX IPv6 working group to discuss detailed plans.
- From April 1st 2017 central services must be accessible via IPv6 and by April 1st 2018 an equal quality of service via both IPv4 and IPv6 is required.
- By April 1st 2018 it should be possible for IPv6-only CPU resources to work within the WLCG with relative ease.
- All sites are encouraged to upgrade their storage to dual stack. CERN and the Tier 1s will be required to provide dual stack access to their storage from April 1st 2017.
- By the end of Run II enough sites should have upgraded their storage to dual stack to allow almost complete data availability over IPv6.
- ALICE use a federated storage model which requires all their data to be available via IPv6 before they can run jobs on IPv6-only CPU resources. ALICE will only be able to support IPv6-only CPU resources by the end of Run II.

31	Contents	
32	1 Introduction	3
33	2 Site requirements within the current computing model	4
34	2.1 Site Services	5
35	2.1.1 CPU	5
36	2.1.2 Disk	6
37	2.1.3 Tier-1 requirements	6
38	2.2 Shared services	7
39	2.2.1 CVMFS	7
40	2.2.2 FTS	7
41	2.2.3 PerfSonar	7
42	2.2.4 ETF test infrastructure	8
43	2.2.5 Frontier Service	8
44	2.2.6 Other Services	8
45	3 Experiment plans	9
46	3.1 ALICE	9
47	3.2 ATLAS	10
48	3.3 CMS	10
49	3.4 LHCb	11
50	4 Conclusion	12

51 1 Introduction

52 There are various motivations for WLCG sites to migrate services to IPv6.
53 The most obvious is the exhaustion of the IPv4 address space, which is al-
54 ready putting constraints on some countries and institutions. The WLCG is
55 expected to evolve under the assumption of flat cash funding for computing
56 resources and it is therefore important that sites are not hindered in their
57 procurement by unnecessary restrictions from the WLCG VOs. Hardware
58 procurements often have a significant lead time and will often be in produc-
59 tion for several years. Even if a site does not intend to switch to IPv6 any
60 time soon, they may well be making procurement decisions now which will
61 influence their decision to migrate.

62 Significant effort is also being put in by the WLCG community to in-
63 vestigate commercial cloud providers to see if they can provide resources
64 (normally CPU) more cost effectively than traditional Grid sites. Some
65 commercial providers charge more for machines with IPv4 connectivity over
66 those with IPv6 only connectivity[1]. The commercial sectors adoption of
67 IPv6 is significantly ahead of the WLCG. The rapid growth in the percent-
68 age of internet traffic going over IPv6 is expected to continue and large
69 companies such as Apple now mandate that software should be validated on
70 machines with IPv6 only connectivity[2].

71 The eventual goal with IPv6 deployment is for the entire internet to
72 migrate to IPv6 only. However as IPv6-only machines are (in general) unable
73 to talk to IPv4-only machines, during the migration, certain machines will
74 need to have both IPv4 and IPv6 addresses (dual stack machines). If the
75 WLCG wasn't so large and didn't have such a diversity of sites, it might
76 be possible to require all sites to upgrade to dual stack before allowing any
77 site to switch off IPv4. However in reality, this is completely unworkable
78 as some sites are already under pressure to migrate while others have not
79 started thinking about it. As there is an additional overhead in running dual
80 stack machines and the fact that the complete migration will take several
81 years, where possible machines should be migrated directly to IPv6.

82 This document describes the required steps to allow existing and future
83 sites and any opportunistic resource that may become available to provide
84 IPv6-only CPU resources. In order to provide CPU resources, a site also
85 needs to provide other services such as CEs, Squid caching proxies etc.; for
86 this reason, when this document refers to IPv6-only CPU resources it means
87 not only the WN but all related services can be IPv6-only.

88 In order not to penalise sites that choose to deploy IPv6-only CPU re-
89 sources, central services need to not only work with IPv6 but should provide

90 the same level of service (e.g. resilience and performance). Ideally the setup
91 would be identical and in cases where this is not possible the differences
92 should be clearly documented. Each of the WLCG VOs operate some central
93 services which are usually hosted at CERN. CERN is already able to
94 make these services dual stack. The WLCG also operates several central
95 services, which should be made dual stack.

96 **2 Site requirements within the current computing** 97 **model**

98 In the current computing model followed by the WLCG VOs a typical site
99 provides CPU and Disk resources. CERN and the Tier-1s also provide tape
100 resources. Data can be transferred to the site from many other sites. In
101 general, in the current computing model VO jobs are sent to sites where
102 the input datasets are present, although an increasing fraction of today's
103 jobs can read data remotely using the XRootD protocol. Jobs running on a
104 site's CPU resources normally only access the local storage. Some sites only
105 provide CPU resources in which case they normally make use of a nearby
106 site's storage by means of remote data access.

107 If a site was to upgrade their storage resources to dual stack and their
108 CPU resources to IPv6-only they would be able to run VO workflows that
109 require access to the local storage only. To allow workflows that require
110 remote access to data, the remote site would need to provide dual stack
111 access to their storage. While it is theoretically possible for the VOs to adapt
112 their workflow management systems to take this into account, it would be
113 a significant amount of effort for what is likely to be very little benefit. It
114 should therefore be assumed that jobs that require remote access to data will
115 not work until the vast majority of a VO's data is accessible via dual stack.
116 Note, due to the fact that VOs often make multiple copies of their data, it is
117 not necessarily true that all sites will need to upgrade their storage to dual
118 stack. The actual impact on a VO will depend on the number of jobs that
119 could potentially try to access data remotely. For VOs to take advantage
120 of [opportunistic] CPU only resources it will be necessary to have a nearby
121 dual stack storage for them to connect to.

122 The WLCG VOs, to different extents, all make use of federated XrootD
123 access. In the case of LHCb this is purely as a failover in case data access
124 at the local sites fails. For ATLAS (FAX) and CMS (AAA), as well as
125 making using of the failover mechanism, a small number of jobs make use
126 of the XRootD federation to access files remotely. This is normally to take

127 advantage of idle CPU at sites which lack the relevant data but have good
128 connectivity to the site that does. ALICE uses a fully federated storage
129 model.

130 The XrootD redirection mechanism has been designed to not direct re-
131 quests from an IPv4-only source to an IPv6-only destination or vice a versa.
132 Therefore if there are two copies of a file one on a dual stack storage and
133 one on a IPv4-only storage a job on an IPv6-only machine should be able
134 to access the file.

135 The Tier-1s play an important role in the WLCG and it is expected that
136 they lead the IPv6 integration and deployment as described in the document.

137 **2.1 Site Services**

138 **2.1.1 CPU**

139 A site that provides CPU resources to the WLCG is likely to have deployed
140 the following:

- 141 • Batch system: This manages the jobs running on the CPU resources
142 at a site.
- 143 • Computer Element (CE): VOs submit their jobs to site CEs. The role
144 of a CE is to convert this job submission over the Grid into something
145 that the local batch system understands.
- 146 • Worker Nodes (WN): These provide the job slots where jobs run. Out-
147 bound connectivity is normally assumed.
- 148 • Squid proxies: These are used to cache requests from jobs to CVMFS
149 and Frontier.
- 150 • Accounting: The number and usage statistics of jobs run is reported
151 on a monthly basis to APEL.
- 152 • Information Provider: Site information is provided by the BDII, while
153 usage of this service is dropping, it is still necessary for some functions.

154 There are some CPU resources (HPC, commercial clouds) that have different
155 setups. However these normally place constraints on the VOs that would ac-
156 tually make it easier to be IPv6 compliant (e.g. no outbound connectivity).
157 WN normally make up the majority of machines (and hence IP addresses)
158 run by a site. Not all batch systems used by WLCG sites are IPv6 compli-
159 ant, a site thinking about upgrading to IPv6-only CPU resources may first
160 want to migrate their batch system.

161 From April 1st 2017 sites will be allowed to deploy IPv6-only CPU re-
162 sources (and all related services). It is expected that the first few sites to
163 migrate will choose a gradual upgrade, which will hopefully avoid problems
164 that could significantly affect the sites availability.

165 **2.1.2 Disk**

166 WLCG sites deploy a range of storage solutions. In general data can ei-
167 ther be accessed directly from the storage node or via gateway machines
168 (sometimes known as a doors or proxy machines). Sites can use a variety of
169 transfers protocols internally; however the LHC VOs rely on the XRootD
170 and GridFTP protocols, both of which have been shown to be IPv6 com-
171 pliant. There is also a push within the WLCG to use http and this is also
172 IPv6 compliant. Both dCache and DPM, the most popular storage services
173 run by WLCG sites are already being run by a small number of sites as
174 dual stack in production. Other storage service such as StoRM have been
175 shown to be IPv6 complaint. For storage services which aren't IPv6 com-
176 pliant (e.g. Castor) it is still possible to provide dual stack access via an
177 XRootD/GridFTP dual stack gateway service.

178 All sites are encouraged to upgrade their storage to dual stack. Even if
179 a site does not intend to migrate to IPv6 soon, if it provides external access
180 to its services via dual stack gateways these will help the VOs with data
181 access.

182 **2.1.3 Tier-1 requirements**

183 In addition to CPU and Disk resources, Tier-1s also provide Tape backed
184 storage. Tape backed data is not used by jobs running on the Grid and there
185 is no requirement at this time to make this service dual stack. Having said
186 that most Tier-1 use dCache which provides a common interface to access
187 both disk and tape back files so it is expected that tape backed service will
188 become dual stack at the same time the disk is.

189 Tier 1s will be required to provide dual stack access to their storage with
190 the following requirements:

- 191 • At least 1Gb/s and 90% reliability by April 1st 2017.
- 192 • At least 10Gb/s and 95% reliability by April 1st 2018.

193 Even if there are Tier-1s, that haven't started to think about IPv6, it should
194 be possible to fulfil the April 2017 goal with a testbed setup which should

195 be easily achievable. Any central services a Tier-1 provides will also need to
196 be made dual stack with a similar timeline and reliability as for the storage.

197 **2.2 Shared services**

198 **2.2.1 CVMFS**

199 All the WLCG VOs as well as many others distribute their software across
200 the Grid using CVMFS. The software is uploaded to a Stratum-0 server
201 (located at CERN for the WLCG VOs) which then mirrors the data to
202 several Stratum-1 servers[4]. Jobs will access the VO software from a cache
203 on the local disk; if the file is not available, it will be looked for in the
204 site Squid server, which in turn, will contact a Stratum-1 if needed. Squid
205 3.x is IPv6 compliant¹ and is being used in production by some sites. It
206 is essential that the Stratum-1 service at CERN is upgraded to dual stack
207 by April 2017. When possible the Tier 1 should upgrade their service to
208 dual stack and all Tier-1s should be upgraded by April 1st 2018 at the very
209 latest.

210 **2.2.2 FTS**

211 ATLAS, CMS and LHCb all use the FTS service extensive for data move-
212 ment around the Grid. All VOs are encouraging sites to make their storage
213 dual stack. Transfers via two dual stack service should go via IPv6, however
214 it is the FTS server which initiates the negotiation and sends a PASV (on
215 IPv4) or an EPSV (on IPv6) to the destination and sends the IP (for the
216 corresponding protocol) and port to the source. Therefore all FTS services
217 should be upgraded to allow transfers between dual stack sites to go over
218 IPv6.

219 Currently the FTS service at CERN is dual stack. There are IPv4 only
220 FTS services at RAL, BNL and Fermilab that are used by the LHC VOs.
221 While it is possible to work around this all FTS services should be upgraded
222 to dual stack when possible and by April 1st 2018 at the very latest.

223 **2.2.3 PerfSonar**

224 PerfSonar instances are required at all WLCG sites to implement the net-
225 work monitoring infrastructure. All Tier-1s were requested to provide a
226 dual stack perfSonar instance and GGUS tickets have now been submitted

¹There is a bug in the handling of HTTP caching headers, whose resolution is expected for July 2016.

227 to those that have not. PerfSonar is a very good way of checking that the
228 migration to IPv6 won't caused any network/routing problems. All sites are
229 requested to provide a dual stack PerfSonar instance by April 2018 at the
230 latest. While it is not essential for all Tier 2s to migrate, it would be con-
231 cerning if they are unable to provide a PerfSonar instance by this time. Any
232 site unable to provide a PerfSonar instance by April 2018 will be requested
233 to provide a clear description of their IPv6 plans.

234 **2.2.4 ETF test infrastructure**

235 A separate ETF test infrastructure will need to be set up to monitor IPv6-
236 ready sites. This must be done by April 2017. This will be run in parallel to
237 the production ETF test infrastructure. This service will provide sites with
238 low level monitoring to help them identify problems with their IPv6 migra-
239 tion and not used for official availability metrics unless the site is providing
240 some resources on IPv6-only. From April 2018 the official ETF infrastruc-
241 ture will be migrated to dual stack. From this point on production work
242 going over IPv6 should be considered entirely normal. This will hopefully
243 encourage sites to investigate IPv6 before April 2018.

244 **2.2.5 Frontier Service**

245 ATLAS and CMS both use the Frontier Service[5] to access conditions data
246 across the Grid. The Frontier service has three components:

- 247 • Frontier client: This software is run by ATLAS and CMS jobs. It
248 converts a conditions database query into an HTTP request. The
249 Frontier Client was made IPv6 compliant in January 2016.
- 250 • Squid proxy: Sites are expected to deploy squid servers to cache the
251 conditions data requests. Squid 3.x is IPv6 compliant.
- 252 • Frontier Launchpad: This converts the HTTP requests back into database
253 queries which are then submitted to the conditions database. Frontier
254 launchpads use squid proxies to cache requests and therefore should
255 be IPv6 complaint.

256 **2.2.6 Other Services**

257 There are several other services such as certificate authorities, software
258 repositories, the GOCDB/OIM, GGUS, VOMS and the BDii. These are
259 not used directly by jobs but are needed when configuring the site. These

260 services should be made dual stack when possible and ideally by April 2018
261 (although some services might not fall under the WLCG banner). It will
262 depend heavily on the site setup as to whether the lack of IPv6 connectivity
263 will cause problems. Problems will have to be followed up by the HEPiX
264 working group as they appear.

265 Some sites also provide VO boxes which may need to be made dual stack.
266 This is covered in the next section.

267 **3 Experiment plans**

268 The WLCG VO plans to allow IPv6 only CPU sites are detailed below.
269 In general the motivation for VOs to support IPv6 is to be able to take
270 advantage of any opportunistic resources that maybe IPv6 only. The VOs
271 agree that sites should be able to migrate to IPv6 if it gives performance,
272 cost or operational advantages. The VOs recommend all sites to upgrade
273 their storage to dual stack and would expect a large number to have by the
274 end of Run II. During Run II, the VOs still expect good site availability and
275 reliability and where possible sites should retain their IPv4 connectivity until
276 the end of Run II, even in a degraded form, as a precaution (e.g. don't hand
277 back IPv4 addresses or completely decommission NATs if not necessary).

278 **3.1 ALICE**

279 Unlike the other LHC VOs, ALICE uses fully federated storage, any site
280 can access the storage element of another site if needed (reading, writing
281 and data transfers). Therefore in order to ensure all job types can run on
282 IPv6-only CPU all data needs to be accessible over IPv6. Some data is
283 stored on multiple sites and therefore it does not necessarily mean all sites
284 will need to be dual stack. To support IPv6, the site storage elements need
285 to run xrootd v.4. The central ALICE Grid services have been tested to
286 run on IPv6 and are running in dual stack mode for over a year. For sites
287 supporting ALICE the current situation is:

- 288 • One third of the sites are still running SEs with xrootd v.3.
- 289 • 5% of the SEs are running in dual stack mode, while the remaining
290 are IPv4.

291 ALICE request that all sites that provide them with disk resource provide
292 dual stack storage by the end of Run II.

293 **3.2 ATLAS**

294 The ATLAS workload management system is called PanDA [3]. Pilot facto-
295 ries generate pilot jobs which are sent directly to CEs at sites. Once these
296 pilots are started by the batch system, they will contact a central PanDA
297 Server to pull in a job (done via http). They will also contact the Rucio
298 server for file lookup (done via http) and the local storage. Some ATLAS
299 jobs access conditions data using the Frontier service. At the end of the job
300 the pilot will write the output files to a local SE. Every 30 minutes while the
301 job is running the pilot will report to the Panda server (via http). It will
302 also contact the PanDA server at the end of the job. ATLAS jobs running
303 on IPv6 WN will need access to the following resources:

- 304 • The production PanDA server nodes.
- 305 • The Rucio authentication nodes.
- 306 • The Rucio production nodes.
- 307 • The Frontier servers at CERN, IN2P3, RAL and Triumf.

308 The pilot factories that submit jobs to CEs have been made dual stack.
309 ATLAS also use the ARC Control Tower (aCT) to submit jobs primarily to
310 NorduGrid but potentially any sites running an ARC CE. This will also need
311 to be made dual stack. ATLAS are working on making all these services
312 dual stack by April 2017.

313 **3.3 CMS**

314 The job submission middleware, glideinWMS, is used to launch HTCCondor
315 worker nodes and its major components (frontend and factory). These have
316 been validated as IPv6 compliant. Some of the glidein factories are already
317 dual stack. HTCCondor itself is fully IPv6-compliant, but the collectors and
318 schedds still need to be all dual-stack in production in order to support
319 IPv6-only worker nodes.

320 The central services hub, cmsweb.cern.ch, has been validated for dual
321 stack operation. The CMS-specific job management systems (WMAgent for
322 production and CRAB3 for analysis) have not yet been fully tested on IPv6,
323 but they are expected to work with little effort needed. In any case, they
324 do not need to be in dual stack for the foreseeable future.

325 The data management system, PhEDEx, uses the Oracle client for com-
326 munication between local site agents and the central service. Tests have not

327 yet been done, but Oracle 12c fully supports IPv6. The global and regional
328 XRootD redirectors for AAA are not fully dual stack.

329 CMS plans to immediately start upgrading all services to dual-stack.
330 Upgrades will be coordinated to minimise operational disruption and will
331 be completed by the end of Run II. For services contacted by worker nodes
332 (like HTCondor) these will be given priority and the aim is to have them
333 done by April 1st 2017.

334 At the time of writing, only eleven CMS sites expose IPv6 addresses for
335 their services. No problems are observed, either in the ETF tests or for real
336 production or analysis jobs.

337 **3.4 LHCb**

338 LHCb uses the DIRAC framework to submit jobs to the grid. DIRAC
339 officially supports IPv6 and some other VOs, who use DIRAC, are already
340 using a dual stack service in production. LHCb submits generic pilot jobs
341 to CEs as needed. When these pilots start on a WN, they contact the
342 LHCb DIRAC central services for available tasks (via dips) which are then
343 executed. If input data is needed, they contact the relevant storages using
344 the sites SRM² to access the data. Production jobs typically download the
345 data to the worker node, as they know exactly how much data is needed.
346 User jobs stream data from the storage directly.

347 Once the job is done, the pilot will upload the output to a storage lo-
348 cation. If the default preferred location is not available, all other possible
349 locations (available for LHCb) are tried in turn until successful and a request
350 is set in the central services of LHCb to transfer the file to the preferred lo-
351 cation when possible. If no location is available, the job ends up in status
352 "failed", and could be resubmitted depending on the conditions.

353 LHCb jobs running on an IPv6 only WN will need access to the following
354 resources :

- 355 • LHCb's DIRAC central services
- 356 • Storage services supporting LHCb
- 357 • Optionally, one of six VO-boxes at LHCb Tier-1 sites

358 Currently there is one Tier-1 storage and one Tier-2D storage that support
359 LHCb in a dual-stack configuration. The LHCb central services are being

²The job is given a list of locations of the input files by DIRAC. It currently contacts the site SRMs in turn to retrieve the data. This will in future be updated to bypass the SRM and construct the file location automatically using the information available.

360 moved to dual-stack machines which will aim to be complete by April 1st
361 2017. There is one outstanding issue with the gLite software which has
362 problems submitting to dual-stack cream CEs which needs to be fixed [6].

363 4 Conclusion

364 The LHC VOs are committed to being able to work on the Grid over IPv6.
365 Much work still remains to be done to make this a reality. The HEPiX
366 IPv6 working group is validating that all essential software is IPv6 com-
367 pliant. Software developers should consider IPv6 compliance a standard
368 requirement and the emphasis should be on them to test this. All the VOs
369 have analysed their workflows on the grid and have provided a list of ser-
370 vices which they will need to make dual stack. While exact time lines have
371 not been agreed the amount of work required is sufficiently small that it
372 should be achievable by April 1st 2017 without significantly disrupting nor-
373 mal WLCG operations.

374 From April 1st 2017 sites will be allowed to deploy IPv6-only CPU re-
375 sources. Sites wishing to deploy IPv6-only CPU resources must deploy dual
376 stack storage if they provide it. All sites are encouraged to upgrade their
377 storage to dual stack. From the contact the HEPiX IPv6 working group
378 has with sites, we believe that there are at most one or two sites that wish
379 to urgently upgrade making up less than 2% of the pledged WLCG CPU
380 resources. Even though the initial migration may be small, it is important
381 that the deadlines are adhered to as this will allow site admins to plan the
382 long term evolution of their site. Any site wishing to upgrade should be in
383 contact with the HEPiX IPv6 working group to ensure that the inevitable
384 teething problems are resolved promptly. By April 1st 2018 it should be
385 possible to deploy IPv6-only CPU resources with relative ease and by the
386 end of Run II enough sites should have upgraded their storage to dual stack
387 to allow almost complete data availability via federated XrootD over IPv6.

388 References

- 389 [1] <https://www.mythic-beasts.com/servers/virtual>
390 [2] <https://developer.apple.com/news/?id=05042016a>
391 [3] <https://twiki.cern.ch/twiki/bin/view/PanDA/PanDA>
392 [4] <http://cernvm-monitor.cern.ch/cvmfs-monitor/atlas.cern.ch/>

393 [5] <http://frontier.cern.ch/>

394 [6] https://ggus.eu/index.php?mode=ticket_info&ticket_id=120586