

1 Managing the ATLAS Grid through Harvester

2 *Fernando Harald Barreiro Megino*^{1,*}, *Aleksandr Alekseev*², *Frank Berghaus*³, *David*
3 *Cameron*⁴, *Kaushik De*¹, *Andrej Filipcic*⁵, *Ivan Glushkov*¹, *FaHui Lin*¹, *Tadashi Maeno*⁶,
4 and *Nicolo Magini*¹ on behalf of the ATLAS Collaboration

5 ¹University of Texas at Arlington, United States of America

6 ²Tomsk Polytechnic University, Russia

7 ³University of Victoria, Canada

8 ⁴University of Oslo, Norway

9 ⁵Jozef Stefan Institute, Slovenia

10 ⁶Brookhaven National Laboratory, United States of America

11 ⁷Iowa State University, United States of America

12 **Abstract.** ATLAS Computing Management has identified the migration
13 of all computing resources to Harvester, PanDA's new workload
14 submission engine, as a critical milestone for Run 3 and 4. This
15 contribution will focus on the Grid migration to Harvester. We have built a
16 redundant architecture based on CERN IT's common offerings (e.g.
17 Openstack Virtual Machines and Database on Demand) to run the
18 necessary Harvester and HTCondor services, capable of sustaining the load
19 of O(1M) workers on the grid per day. We have reviewed the ATLAS Grid
20 region by region and moved as much possible away from blind worker
21 submission, where multiple queues (e.g. single core, multi core, high
22 memory) compete for resources on a site. Instead we have migrated
23 towards more intelligent models that use information and priorities from
24 the central PanDA workload management system and stream the right
25 number of workers of each category to a unified queue while keeping late
26 binding to the jobs. We will also describe our enhanced monitoring and
27 analytics framework. Worker and job information is synchronized with
28 minimal delays to a CERN IT provided ElasticSearch repository, where we
29 can interact with dashboards to follow submission progress, discover site
30 issues (e.g. broken Compute Elements) or spot empty workers. The result
31 is a much more efficient usage of the Grid resources with smart, built-in
32 monitoring of resources.

33

34

* Corresponding author: barreiro [at] uta [dot] edu

35 **1 Introduction**

36 The Worldwide LHC Computing Grid (WLCG) [1] is a highly heterogeneous federation of
37 computing sites with different middleware and increasingly special resources, such as
38 Cloud or High Performance Computing (HPC) resources. PanDA [2] is the Workload
39 Management System for ATLAS [3], managing all production and user jobs across the
40 WLCG centers associated with the experiment. In order to exploit resources, PanDA is
41 based on the Pilot paradigm [4], where Pilot jobs are submitted to the batch systems at sites.
42 The Pilots retrieve the real payload from PanDA server and execute it.

43 Over the years numerous Pilot submission systems have been developed, frequently
44 specializing on a certain subset of resources and having independent code bases. The
45 Harvester project was born as an attempt to provide a universal Pilot submission system. At
46 the same time, in the case of Grid resources, some improvements were needed to increase
47 the stability and usage efficiency through a tighter integration with the PanDA Workload
48 Management System allowing a more informed decision taking.

49 This contribution will focus on core Harvester design decisions and other significant
50 aspects like new submission modes and monitoring. We will also show the process and
51 results of migrating all Grid resources to Harvester.

52 **2 Harvester design decisions**

53 **2.1 Lightweight vs High Performance execution modes**

54 In order to be a universal service for any type of resource, Harvester needs to provide
55 certain flexibility in terms of resource consumption and dependencies. In the case of HPCs,
56 Harvester frequently needs to run on edge nodes with strict restrictions on installation and
57 memory/CPU footprint of the application. In this case Harvester provides a lightweight
58 mode, using a SQLite database and limiting the number of internally managed threads.

59 On the contrary for the Grid, central Harvester instances manage very large, costly
60 infrastructures that need to be kept fully utilized. There are no important installation
61 restrictions and it is preferable to host high performance services. In this case, Harvester
62 typically uses a MySQL/MariaDB database and multiple Harvester processes can be started
63 through frameworks like uWSGI [5].

64 **2.2 Fast integration of new resources**

65 To reduce the lead time until a new resource is exploited successfully, the Harvester code
66 follows a plug-in approach: the code is split into a common core and a set of specific plug-
67 ins that can be configured for each resource. Typically, new resources only require the
68 implementation of resource-specific libraries for submission, monitoring and cleaning up
69 workers that have finished. For most of the cases, these are python modules with lengths of
70 a few hundred lines of code.

71 Up to date, plug-ins have been developed for HTCondor [6], ARC CE [7], Google
72 Compute Engine [8], Kubernetes [9], SAGA [10], PBS [11], Cobalt [12] and Slurm [13].

73 While it's generally not required for the Grid case, it's also possible to implement
74 data management plug-ins for more complex use cases.

75 **2.3 Queue unification**

76 ATLAS submits Pilots specifying the number of cores and memory for the jobs. In
77 the past, sites provided multiple separate queues per job type:

- 78 • Analysis: usually single core payloads running with a non-production proxy security
79 role
- 80 • Production: several different payloads summarized in Table 1.

81 **Table 1.** Requirements for the various production workloads

	CPUs	Memory
Single core	1	2 GB
Multi core	Depends on site, usually 8	2 GB/core
Single core, high memory	1	Depends on site, >2 GB
Multi core, high memory	Depends on site, usually 8	Depends on site, >2 GB/core

82
83 An average site for ATLAS would have an Analysis queue, a single core and a multi-core
84 production queue. Larger sites would also provide high memory queues. Pilots would be
85 submitted to these queues independently, causing random competition for the slots and
86 making it impossible to establish any control over the ratios.

87 In order to follow the Global Shares [14] priorities of ATLAS, it is more desirable to
88 unify all queues at a site into one single queue that can accept Pilots with different sizes.
89 These unified queues can be managed using either the Pull mode or the new Pull UPS mode
90 (see next subsection).

91 **2.4 Submission modes supported by Harvester**

92 Harvester supports the classical push and pull workflows, and has extended the pull
93 workflow for unified queues as follows:

- 94 • **Push (early binding):** the worker submitted to the batch queue is already assigned to a
95 particular job and requests the CPU, memory, and potentially other requirements, of
96 the specific job. The push workflow works natively with unified queues. The
97 disadvantage is the early binding: once the worker is submitted, you can't control the
98 queue time. If the queue time is long, a high priority job is stuck in the queue. Also,
99 during the queueing time, a more important job might have appeared and needs to wait
100 until the queue of previously submitted jobs has cleared.
- 101 • **Pull (late binding):** To circumvent the drawbacks of the push workflow, the pull
102 workflow submits workers without any pre-assigned job. The advantage is that the
103 WMS can choose the most important job right at the time the worker starts running.
104 However, all workers submitted to the queue are the same and thus this mode does not
105 support natively on unified queues. A dumb pilot submitter and multiple queues with
106 different requirements on the same site causes competition between them and ATLAS
107 can't control the ratio of workers across the queues.
- 108 • **Pull UPS (Unified Pilot Streaming):** This mode is an extension of the traditional pull
109 mode to support unified queues. All types of jobs are queued together in the same
110 queue and it's the WMS system deciding the fraction of workers of each type to
111 submit, depending on the global priorities. The WMS decision is published to
112 Harvester as a command.

115

116 **3 Harvester monitoring**

117 **3.1 Worker monitoring**

118 Harvester reports the worker information up to the WMS, in our case PanDA. PanDA
119 stores the worker information from the different Harvesters in one combined central Oracle
120 table and also mirrors it to an Elasticsearch repository managed by CERN IT [15]. Our
121 monitoring of choice follows the general trend to build dashboards on top of Elasticsearch,
122 because of its speed and flexibility. We have built dashboards in Kibana for expert users
123 and more user-friendly dashboards in Grafana for widespread use by site admins and
124 shifters.

125 The dashboards allow the user to interact with the Harvester data at a detailed worker
126 level and at high level overviews:

- 127 • The dashboards show a table with detailed worker information, which provide links to
128 the logs and to the job submission file, to provide easy debugging of worker failures.
- 129 • The dashboards also provide several plots to see the submission evolution broken down
130 by different criteria (Harvester instance, resource type, state, etc.). There are also plots
131 to see the worker distribution across the Computing Elements of a site.

132 **3.2 Service monitoring**

133 Our service monitoring collects important information about the Harvester instances:

- 134 • Disk, CPU and memory usage
- 135 • Worker submission, update and completion rates

136 These metrics allow to easily identify when an instance is in trouble or misbehaving.
137 We have an alerting system in place that sends a mail to the central Harvester team in case
138 a metric has exceeded a threshold. The alerting system has proven very useful and warns us
139 quickly about issues. These alerting systems are critical in times of low operational
140 manpower and shifters in ATLAS.

141 **3.2 Site monitoring**

142 We are also implementing monitoring to identify broken sites. There are views in place
143 grouping broken sites and broken Computing Elements, based on high ratios of workers in
144 bad states (failed or cancelled submission). We show also the error messages happening on
145 those sites.

146 We are working on a monitoring to identify inactive sites, where the submission rate
147 is lower than what would be expected. There can be different reasons for this, including
148 misbehavior of Harvester submission rate calculation for the site. In the future we would
149 like to automate actions, like issuing a reset of the worker submission calculation or
150 submitting tickets to the site.

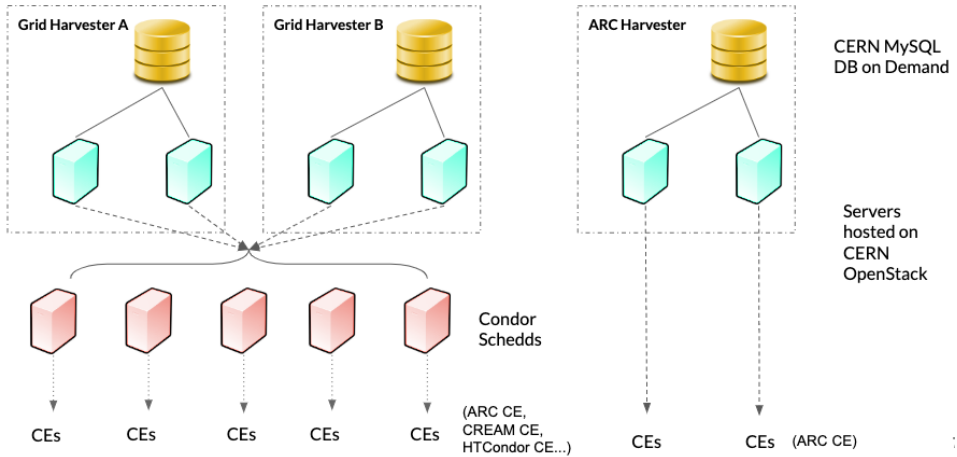
151 **4 Migration of the Grid to Harvester**

152 **4.1 Central infrastructure**

153 We have setup a central high-performance Harvester infrastructure for the whole Grid
 154 migration. It is based on CERN IT provided services. The databases are hosted by the DB
 155 on Demand [16] project, while the servers are virtual machines from the OpenStack service.
 156 We have distributed the Grid sites across three Harvester instances:
 157

- 158 • Harvester A: US, CERN Tier 0, France, Russia, Canada, Netherlands
- 159 • Harvester B: CERN, Germany, UK, Taiwan, Italy, Spain
- 160 • Harvester ACTA: ARC CE sites

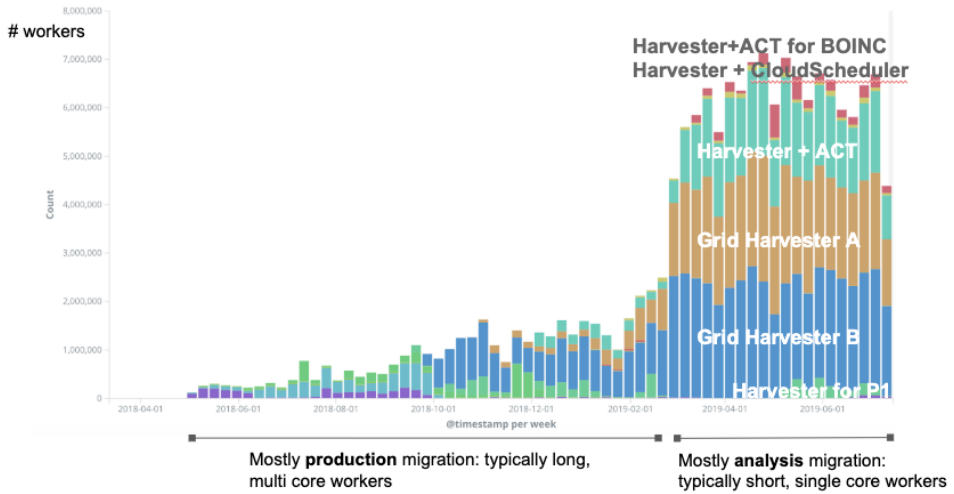
161 Each instance consists of two active machines and one shared database for each
 162 instance. The redundant machines are in different zones of the CERN computing center. If
 163 one machine is unavailable, the other machine can take over the full load.
 164 Non-ARC computing elements rely on HTCondor as interface. There are five load-
 165 balanced HTCondor machines for the whole grid.



166
 167 **Fig. 1.** Harvester infrastructure overview

168 **4.2 Migration process**

169 The whole grid was migrated June 2018 to March 2019. During the migration process we
 170 watched stability and scaling of the services, improved the efficiency of the HTCondor
 171 interfaces, simplified the configuration of queues through the ATLAS Grid Information
 172 System (AGIS) [2017] and improved some interactions with the database.



173

174 **Fig. 2.** Overview of the Grid migration progress to Harvester, in number of workers. This histogram
 175 is extracted from our Harvester worker monitoring.

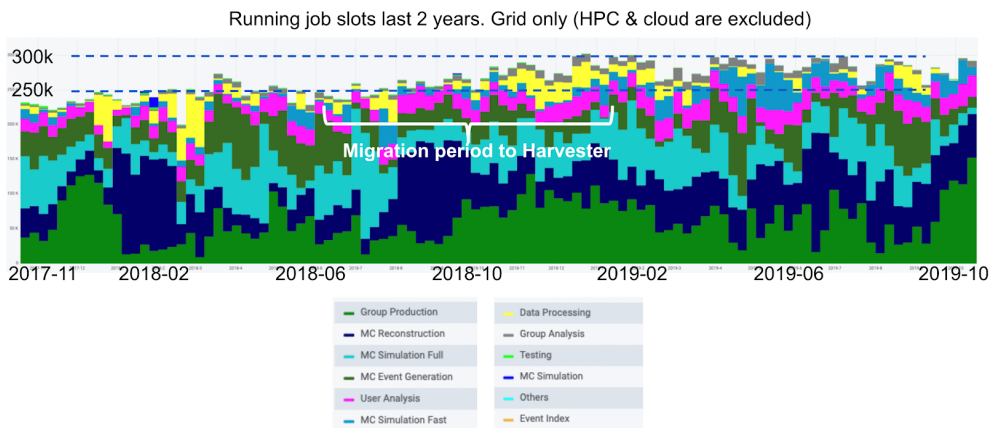
176 During the first period we migrated production queues to Harvester. Although
 177 production jobs occupy the better part of the grid slots, they are rather easy to handle from a
 178 Harvester load point of view: they run for multiple hours and typically occupy multiple
 179 cores. During the migration a major effort was devoted to unify production queues to the
 180 pull UPS mode explained in section 2.4.

181 Once the infrastructure was proven to be scalable and stable, we also migrated
 182 analysis queues to Harvester. Analysis jobs represent less than 20% of the grid slots, but
 183 they are short and run on a single core. From a Harvester load point of view, they are
 184 heavier and require a much higher worker submission rate.

185 5 Results

186 Harvester is a project that took 2 years from initial discussion to full roll out. It is
 187 contributing to a better usage of the ATLAS Grid. During the roll out of Harvester, we
 188 observed a significant usage of job slots. While there are multiple potential contributions
 189 happening in parallel (e.g. 3-4% of pledge increase or increased usage of the ATLAS Tier
 190 0), we believe that Harvester has made a significant contribution to the increase for two
 191 reasons:

- 192 • Harvester is more aggressive than previous Pilot factories in worker submission
 193 and competes better on shared sites across multiple experiments.
- 194 • Through the queue unification and the pull UPS submission, the worker
 195 submission is also more intelligent. It avoids uninformed competition between
 196 ATLAS queues at the same site and follows better the ATLAS job priorities.



197

198 **Fig. 3.** Running job slots on Grid resources before and after the migration to Harvester

199 Unified queues for production are working very well and we are currently working on
 200 further unifying production and analysis queues, so that a standard Grid site only needs to
 201 provide one queue for all job types.

202 Lastly, we want to enhance automation of issues on the Grid to reduce the need of
 203 human operational effort.

204 Acknowledgements

205 This research was enabled in part by support provided by National Science Foundation
 206 (www.nsf.gov) and US Department of Energy (www.energy.gov/).

207 References

- 208 1. LHC Computing Grid: Technical Design Report, document LCG-TDR-001, CERN-
 209 LHCC-2005-024 (The LCG TDR Editorial Board) (2005)
- 210 2. T. Maeno et al. *J. Phys. Conf. Ser.* **898** 052002 (2017)
- 211 3. ATLAS Collaboration 2008 The ATLAS Experiment at the CERN Large Hadron
 212 Collider *J. Inst.* **3** S08003
- 213 4. P. Nilsson et al. *J. Phys. Conf. Ser.* **513** 032071 (2014)
- 214 5. uWSGI <https://uwsgi-docs.readthedocs.io/en/latest/>
- 215 6. HTCondor <https://research.cs.wisc.edu/htcondor/>
- 216 7. ARC CE <http://www.nordugrid.org/arc/ce/>
- 217 8. Google Compute Engine <https://cloud.google.com/compute/>
- 218 9. Kubernetes <https://kubernetes.io/>
- 219 10. SAGA <https://saga-python.readthedocs.io/en/latest/>
- 220 11. PBS <https://www.nas.nasa.gov/hecc/support/kb/121/>
- 221 12. Cobalt <https://trac.mcs.anl.gov/projects/cobalt>
- 222 13. Slurm <https://www.slurm.schedmd.com/>
- 223 14. F. Barreiro Megino et al. *EPJ Web Conf.*, **214** (2019) 03025

- 224 15. P. Saiz et al. EPJ Web Conf., (to be published)
- 225 16. R. Aparicio et al. J. Phys. Conf. Ser. **664** 052021 (2015)
- 226 17. A. Anisenkov et al. J. Phys. Conf. Ser. **664** 062001 (2015)