

Federating distributed storage for clouds in ATLAS

**F Berghaus¹, K Casteels¹, A Di Girolamo², C Driemel¹, M Ebert¹,
F Furano², F Galindo³, M Lassnig², C Leavett-Brown¹, M Paterson¹,
C Serfon², R Seuster¹, R Sobie¹, R Tafirout³ and R P Taylor¹**

¹ Department of Physics and Astronomy, University of Victoria, Finnerty Road, Victoria
V8P 5C2, Canada

² CERN, Geneva 1211, Switzerland

³ TRIUMF, Wesbrook Mall, Vancouver V6T 2A3 Canada

E-mail: frank.berghaus@cern.ch

Abstract. Input data for applications that run in cloud computing centres can be stored at distant repositories, often with multiple copies of the popular data stored at many sites. Locating and retrieving the remote data can be challenging, and we believe that federating the storage can address this problem. A federation would locate the closest copy of the data on the basis of GeoIP information. Currently we are using the dynamic data federation Dynafed, a software solution developed by CERN IT. Dynafed supports several industry standards for connection protocols like Amazon's S3, Microsoft's Azure, as well as WebDAV and HTTP. Dynafed functions as an abstraction layer under which protocol-dependent authentication details are hidden from the user, requiring the user to only provide an X509 certificate.

1. Introduction

We aim to run data-intensive applications on globally distributed opportunistic resources that have no local grid storage. The ATLAS experiment leverages a globally distributed system of infrastructure as a service (IaaS) clouds as part of its distributed computing system. These resources are integrated into the ATLAS distributed computing system using the Cloud Scheduler [1] technology developed at the University of Victoria. These IaaS resources are used opportunistically, and do not support any local grid infrastructure.

The workflows executed by high energy physics experiments often demand large volumes of input data or produce a significant volume of output data. We aim to use a data federation, such as Dynafed [2], to redirect the applications running on opportunistic resources to the optimal storage endpoint to retrieve input or deposit output data. We also aim to integrate storage solutions offered by cloud providers into the ATLAS distributed data management system using Dynafed.

In this paper we explain a system leveraging Cloud Scheduler and Dynafed which successfully executed functional test jobs as part of the ATLAS distributed computing system on the CERN OpenStack [3] cloud resource that read their input from and wrote their output to an object store implemented using Ceph [4] and exposing an S3 compatible gateway.

2. Conceptual Design

The ATLAS experiment leverages the resources of the Worldwide LHC Computing Grid, WLCG [5]. The computing centres that are part of the WLCG and support ATLAS provide

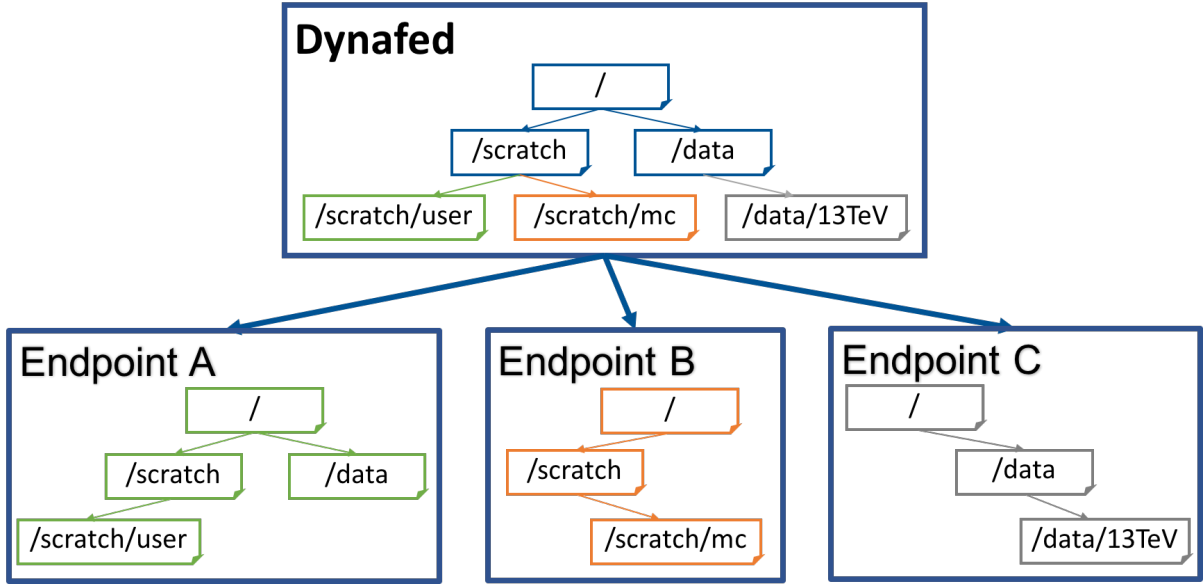


Figure 1. The dynamic federation is connected to multiple endpoints. Each endpoint may be a file system or an object store accessible using a protocol which allows redirection, such as HTTP. The dynamic federation provides a namespace that is a union of all the namespaces of the endpoints. That namespace is presented as a regular directory structure by the Dynafed. The content of a displayed directory is calculated when accessed.

a global storage infrastructure for the experiment data and simulated events. They may be accessed using standard protocols, such as HTTP with WebDAV extensions. Dynafed supports storage backends that offer HTTP and WebDAV access and promises sufficient scalability to create the appearance of a single virtual namespace for the entire ATLAS data catalogue. Figure 1 shows how Dynafed could unify the namespaces of attached storage elements into a single namespace.

Dynafed allows the usage of cloud storage systems such as S3, SWIFT, and Azure. On the user-facing side Dynafed still presents an HTTP interface implementing authentication and authorization through an X509 public key infrastructure. Dynafed supports grid security infrastructure extensions of X509 with VOMS attributes [6]. Credentials may be presented as certificate and key or as a proxy, which allows the additional use of VOMS attributes. When Dynafed forwards clients to cloud storage systems, it translates their X509 credentials to pre-signed URL that permit, for a limited time, access to the cloud storage system.

3. Data Access

The dynamic federation used for this work was configured to use three endpoints: one at CERN, one at TRIUMF, and one at the University of Victoria. Each endpoint was a CephS3 object store. Table 1 illustrates the task division within the dynamic federation to handle client requests.

To access data through Dynafed a client makes a request for a file using HTTP optionally with WebDAV extensions. We will focus on WebDAV from here on. In our configuration the Dynafed only allows access to members of the ATLAS Virtual Organization. The client must provide X509 credentials with the request. The credential must be signed by a trusted certificate authority. If the credential contains VOMS extensions certifying the user to be a member of the ATLAS collaboration, access is granted. Without VOMS extensions the Dynafed checks the credential against all current members of the ATLAS collaboration and grants access if a match

Table 1. The dynamic web federation is an Apache server running the LCGDM implementation of WebDAV. The namespace usually managed by LCGDM has been replaced by the uniform general redirector (UGR) which translates the requests to the web file system to the connected endpoints. The endpoint modules handle the communication with the configured endpoints. All requests are cached in memory on the server as well as in a second-level cache which may be shared across multiple load-balanced servers.

Component	Purpose
Apache	Load the <code>lcgdm_dav</code> module and start up a WebDAV server
<code>lcgdm_dav</code>	Configure dmlite and load the uniform general redirector
dmlite	as namespace plugin
UGR	Configure authentication and endpoints
<div style="display: flex; align-items: center;"> <div style="writing-mode: vertical-rl; transform: rotate(180deg); margin-right: 5px;">Plugins</div> <div style="border-left: 1px solid black; padding-left: 5px;"> dmlite WebDAV/HTTP S3 Azure </div> </div>	Communicate with endpoints on request
Memcached	Cache recent redirects to distributed object caching system

is found. Administrators and data management services have privileged accounts. Authorized clients are redirected to a signed URL on the closest CephS3 endpoint. Authorization is granted explicitly for reading, writing, listing, and/or deleting operations.

When a file is requested, the dynamic federation checks whether the locations of the file are already in its cache. If so, the cached entry is used, otherwise each endpoint is queried for the file after name translation to that endpoint. Dynafed waits for responses¹ to collect and cache. The resulting endpoints are evaluated for proximity to the client and the client is redirected to the closest copy. The Dynafed regularly polls all connected endpoints to determine if they are reachable. Should an endpoint be unresponsive, requests will not be forwarded to it.

4. Application Workflow

In order to integrate the dynamic federation into the ATLAS distributed computing and data management system, it was defined as a storage element associated with the **CERN-EXTENSION**² ATLAS site. It was configured to be accessible using WebDAV and flagged as special in the ATLAS grid information system to allow Rucio [7] to select a copy tool implementation which does not move or rename files.

The input datasets for analysis and production functional tests were transferred to the dynamic federation using the file transfer service at CERN. Once the transfers completed successfully the data was registered manually in the Rucio data catalogue.

With the input data registered in the data catalogue it was possible to run grid jobs against the data in the dynamic federation. The jobs were executed on virtual machines hosted on the CERN OpenStack using the Cloud Scheduler technology as illustrated in Figure 2. The resulting data and logs were uploaded to the CephS3 storage via Dynafed upon job completion.

Some additional development is required for full integration of cloud storage into the ATLAS distributed data management: bulk transfers negotiated between storage endpoints using the HTTP protocol must be fully supported, and the data management system must be able to parse the checksums of files on cloud storage.

¹ Up to a given timeout set to 3 seconds here

² **CERN-EXTENSION** is an ATLAS site defined as a part of the **CERN-PROD** WLCG and GOCDB site.

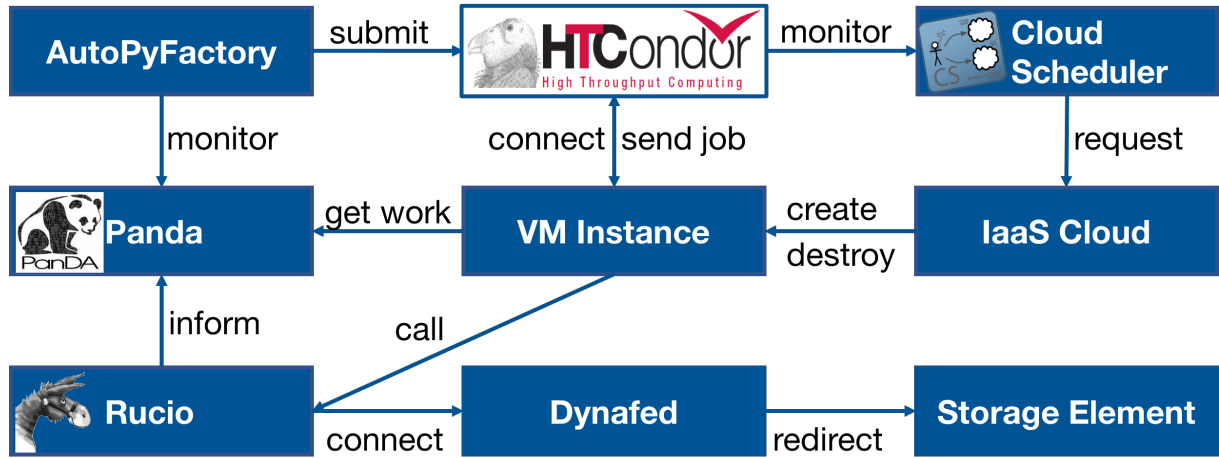


Figure 2. A client (AutoPyFactory or Harvester) submits pilot wrapper scripts to an HTCondor queue. The queue is monitored by a Cloud Scheduler. The Cloud Scheduler makes requests to connected cloud interfaces in response to queued jobs. The cloud infrastructures create virtual machine instances and provide user-data to cloud-init running in the virtual machines for configuration. The VM instances are configured to connect to HTCondor and start consuming jobs from the queue. The jobs are pilot wrapper scripts which download the pilot. The pilot gets tasks from a PanDA [8] queue and uses Rucio to download input data and upload results. Rucio is configured to contact the dynamic federation using the HTTP protocol with WebDAV extensions. The federation forwards the Rucio client to the closest available storage element.

5. Summary

It was shown that ATLAS jobs can retrieve and deposit their data on a cloud storage system accessed via a dynamic federation using the HTTP protocol with WebDAV extensions. The jobs ran on virtual machine instances in a cloud and could be scheduled anywhere in the distributed cloud system currently running as part of the ATLAS production system. Further development is necessary to allow the execution of production or analysis jobs against the dynamic federation. Work is ongoing to integrate the dynamic federation with the Belle-II experiment and the DIRAC workload management system. While this development is being pursued against opportunistic cloud resources it should also be useful in the context of volunteer resources [9].

References

- [1] Cloudscheduler project, "Cloudscheduler" [software], version 1.12, 2017, Available from <http://cloudscheduler.org/> [accessed 2017-10-20]
- [2] Dynafed project, "Dynafed" [software], version 1.3.1, 2017, Available from <http://cern.ch/lcgdm/dynafed-dynamic-federation-project> [accessed 2017-10-20]
- [3] Openstack project, "OpenStack" [software], various versions, Available from <https://www.openstack.org/> [accessed 2017-10-20]
- [4] Ceph project, "Ceph" [software], various versions, Available from <http://ceph.com> [accessed 2017-10-20]
- [5] Bird I 2011 Computing for the Large Hadron Collider *Ann. Rev. Nucl. Part. Sci.* **61** 99
- [6] Foster I, Kesselman C and Tuecke S 2001 The Anatomy of the Grid: Enabling Scalable Virtual Organizations *International Journal of Supercomputer Applications* **15** 3 200 - 222
- [7] Rucio project, "Rucio" [software], version 1.2.5, Available from <http://rucio.cern.ch> [accessed 2017-10-21]
- [8] Maeno T 2008 PanDA : Distributed Production and Distributed Analysis System for ATLAS *Journal for Physics* **119** 6
- [9] LHC@home project, "ATLAS@home" [software], version BOINC 7.8, Available from <http://lhathome.web.cern.ch> [accessed 2017-10-21]