

# Cloud services for the Fermilab scientific stakeholders

S Timm<sup>1</sup>, G Garzoglio<sup>1</sup>, P Mhashilkar<sup>1\*</sup>, J Boyd<sup>1</sup>, G Bernabeu<sup>1</sup>, N Sharma<sup>1</sup>, N Peregonow<sup>1</sup>, H Kim<sup>1</sup>, S Noh<sup>2</sup>, S Palur<sup>3</sup>, and I Raicu<sup>3</sup>

<sup>1</sup>Scientific Computing Division, Fermi National Accelerator Laboratory

<sup>2</sup>Global Science experimental Data hub Center, Korea Institute of Science and Technology Information

<sup>3</sup>Department of Computer Science, Illinois Institute of Technology

E-mail: {timm, garzoglio, parag, boyd, gerard1, neha, njp, hyunwoo }@fnal.gov, rsyong@kisti.re.kr, psandeep@hawk.iit.edu, iraicu@cs.iit.edu

**Abstract.** As part of the Fermilab/KISTI cooperative research project, Fermilab has successfully run an experimental simulation workflow at scale on a federation of Amazon Web Services (AWS), FermiCloud, and local FermiGrid resources. We used the CernVM-FS (CVMFS) file system to deliver the application software. We established Squid caching servers in AWS as well, using the Shoal system to let each individual virtual machine find the closest squid server. We also developed an automatic virtual machine conversion system so that we could transition virtual machines made on FermiCloud to Amazon Web Services. We used this system to successfully run a cosmic ray simulation of the NOvA detector at Fermilab, making use of both AWS spot pricing and network bandwidth discounts to minimize the cost. On FermiCloud we also were able to run the workflow at the scale of 1000 virtual machines, using a private network routable inside of Fermilab. We present in detail the technological improvements that were used to make this work a reality.

## 1. Introduction

The Fermilab scientific program includes several running experiments, both the CMS experiment at the Energy Frontier, and the various neutrino and muon experiments on the Intensity Frontier. The ongoing data analysis and simulation for running experiments, combined with a large simulation load for future facilities and experiments, results in an unprecedented level of computing demand. This paper describes recent progress in the ongoing program of work to expand our computing to the distributed resources of grids and public clouds.

As part of the joint collaboration between Fermilab and KISTI, we have a program of work building towards distributed federated clouds. In the summer of 2014 the primary goal of this program was to demonstrate a federated cloud running at the 1000 Virtual Machine scale, using our local private cloud nodes and Amazon Web Services EC2. Our application of choice for this is the Cosmic Ray simulation of the NOvA experiment far detector at Ash River. This application requires

---

\* To whom any correspondence should be addressed.

negligible input and produces about 250MB of output per job, and is quite computationally intensive. The NOvA experimenters spent considerable effort in optimizing their code to run at various sites outside of Fermilab, including loading their code into the OSG OASIS CVMFS server. The NOvA experiment supplied us with a set of files and scripts, which would generate one full set of their cosmic ray Monte Carlo, 20000 input files in all, with one job per file.

## **2. Challenges towards large scale resources**

The main challenge that had to be addressed in expanding to larger scale on the public cloud was finding a scalable way to distribute the code to 1000 simultaneous jobs. Another challenge was the need to quickly convert a local virtual image to Amazon Web Services format so we could make changes in the setup as needed. On the private cloud we had to find a faster and more scalable way to deliver virtual machine images, and find a more scalable virtual machine instantiation API, as well as add a significant number of machines to our private cloud.

### *2.1. Description of Squid Caching and Shoal Discovery Services*

The CERN Virtual Machine File System (CVMFS) [1] is widely adopted by the High Energy Physics (HEP) community for the distribution of project software. CVMFS is a read-only network file system that provides access to files from a CVMFS Server over HTTP. When CVMFS is used on a cluster of worker nodes, a HTTP web proxy can be used to cache the file system contents, so that all subsequent requests for that file will be delivered from the local HTTP proxy server. Typically, a HEP computing site has a local or regional Squid HTTP web proxy [2], with the central CVMFS servers located at the main laboratory, such as CERN for the LHC experiments.

The use of IaaS cloud resources is becoming a realistic solution for HEP workloads [3, 4], and CVMFS is an effective means of providing the software to the virtual machines (VMs). Historically each VM has a list of the available Squid servers and CVMFS servers. In most cases, the Squids have been remote, which is not optimal for security or network bandwidth purposes. The optimal Squid may be different depending on the location of the cloud. Further, one can imagine dynamically instantiating Squid servers in an opportunistic cloud environment to meet application demand. This requires a discovery service to locate the optimal squid server. However, there is currently no mechanism other than Shoal for locating the optimal Squid server. As a result, we use Shoal as a service that can dynamically publish and advertise the available Squid servers. Shoal is ideal for an environment using both static and dynamic Squid servers.

The CernVM File System (CernVM-FS) provides a scalable, reliable and low maintenance software distribution service. It was developed to assist High Energy Physics (HEP) collaborations to deploy software on the worldwide-distributed computing infrastructure used to run data processing applications. CernVM-FS is implemented as a POSIX read-only file system in user space (a FUSE module). Files and directories are hosted on standard web servers and mounted in the universal namespace /cvmfs. Internally, it uses content-addressable storage and Merkle trees in order to maintain file data and meta-data. CernVM-FS uses outgoing HTTP connections only; thereby it avoids most of the firewall issues of other network file systems. It is actively used by small and large HEP collaborations. In many cases, it replaces package managers and shared software areas on cluster file systems as means to distribute the software used to process experiment data.

Squid is a caching proxy for the Web supporting HTTP, HTTPS, FTP, and more. It reduces bandwidth and improves response times by caching and reusing frequently requested web pages.

Squid has extensive access controls and makes a great server accelerator. It runs on most available operating systems, including Windows and is licensed under the GNU GPL.

Shoal is divided into three logical modules, a server, an agent, and a client. Each package is uploaded to the Python Package Index [5] (the standard method of distributing new components in the Python language). The Shoal Server maintains a list of active Squid servers in volatile memory and receives AMQP messages from them. It provides a RESTful interface for Shoal Clients to retrieve a list of geographically closest Squid servers. It provides a basic support for Web Proxy Auto-Discovery Protocol (WPAD), and provides a web user interface to easily view Squid servers being tracked. The Shoal Agent is a daemon that runs on Squid servers to send an Advanced Message Query Protocol (AMQP) [6] message to Shoal Server on a set interval. Every Squid server wishing to publish its existence runs Shoal Agent on boot. Shoal Agent sends periodic heartbeat messages to the Shoal Server (typically every 30 seconds).

The Shoal Client is used by worker nodes to query Shoal Server to retrieve a list of geographically nearest Squid servers via the REST interface. Shoal Client is designed to be simple (less than 100 lines of Python) with no dependencies beyond a standard Python installation. The Shoal Server runs at a centralized location at Fermilab with a public and fixed IP address. For agents (i.e. Squid servers), Shoal Server will consume the heartbeat messages sent and maintain an up-to-date list of active Squids. For clients, Shoal Server will return a list of Squids organized by geographical distance and load. For regular users of Shoal Server, a web server is provided. The web server generates dynamic web pages that display an overview of Shoal. All of the tracked Squid servers are displayed and updated periodically on Shoal Server's web user interface, and all client requests are available in the access logs.

AMQP forms the communications backbone of Shoal Server. All information exchanges between Shoal Agent (Squid Servers) and Shoal Server are done using this protocol, and all messages are routed through a RabbitMQ [7] Server.

## *2.2. Design and Implementation of Squid Service in the Cloud*

This project aims to automate installation and deployment of CVMFS (network file system), Squid (on-demand caching service) and Shoal (squid cache publishing and advertising tool designed to work in fast changing environments) on dynamically instantiated large scale batch of virtual machines on Fermi Cloud (private cloud) and on Amazon Web Services (Public cloud) using Serverless Puppet (puppet apply). As a part of this work, we developed puppet modules and scripts for installing and deploying shoal client, agent and server. We also fixed potential bugs in Shoal Server and made it suitable to publish Squid Servers running on EC2 instances that do not have a static public IP.

The architecture of large scale batch of dynamically instantiated FermiCloud and EC2 worker nodes provided with network file system, on-demand caching service and a cache publishing and advertising tool is shown in Figure 1.

When a new Shoal Server node is instantiated, the Shoal Server and its dependent components including Apache and RabbitMQ server are installed on startup of the machine. For this purpose we constructed a puppet script which can be run at boot time from a normal puppet server, or as a standalone script with puppet apply, using the cloud-init features of AWS or the contextualization features of OpenNebula respectively.

When a Squid Server is instantiated dynamically, the Squid service and the Shoal Agent are installed on the startup of the machine, using puppet apply. The shell script to execute the puppet apply command is included as part of the launch of the virtual machine and uses the cloud-init features of AWS, or the contextualization features of OpenNebula respectively.

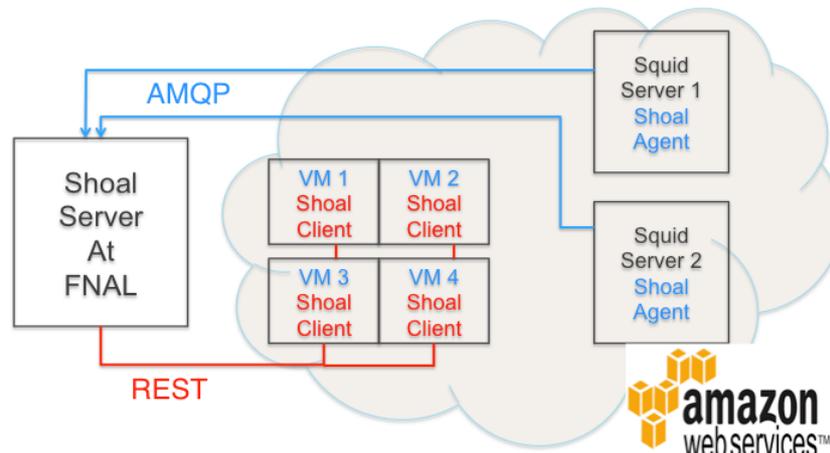


Figure 1: Architecture of Dynamically Instantiated EC2 Instances with On-Demand Caching Service and a Cache Publishing Service

When a worker node is instantiated, the Shoal Client, the CVMFS client, and other unrelated software can be installed at the startup of the machine via a startup script that uses puppet apply. This is done on FermiCloud. It could be done on AWS as well but in practice we pre-install it and send it along with our virtual machine worker node image. The Shoal Client is a cron job that queries the Shoal Server using the REST interface to get the closest Squid Server and is configured to run every 2 hours. Shoal Client updates the proxy address in the CVMFS configuration file. So that when CVMFS client tries to download any software from CVMFS server, the request passes through the Squid Server, whose IP address is configured in CVMFS configuration file. The Shoal Client was later modified to modify the system configuration files to also send other non-CVMFS related traffic such as fetching certificate revocation lists to the dynamically detected nearest squid server.

Our goal was to have the FermiCloud VM's on the private net use the on-site Squid servers and the AWS VM's use Squid servers that were instantiated on AWS. In this way we did not have to open the Fermilab Squid servers to the outside internet and we also delivered a much faster Squid service due to decreased latency. We found that one m3.large Squid server in AWS was sufficient to serve the load for caching code for 1000 running jobs. Any element code would only be fetched once to the AWS squid server and all other AWS virtual machines would access it from there.

### 2.3. Virtual machine conversion system

We also required a faster way to upload specialized virtual machines to Amazon Web Services. It was our goal to run a very similar virtual machine image on Amazon as we run on our private cloud, based on Scientific Linux. We leveraged our existing mechanism that manufactures our stock image for the private cloud and added extra steps to it that strip out a few Fermilab-specific configurations and add a few Amazon-specific configurations. The standard image is uncompressed from qcow2 to raw format. Then it is copied to another running image on AWS which has an extra disk mounted. The extra disk is then saved as a snapshot. Once the virtual machine image is stored on Amazon then all copies of the virtual machines are launched from that. We settled on the Amazon m3.large instance

type because it had two available cores and enough default scratch space (30GB SSD) to host two jobs at once.

#### *2.4. Private cloud scalability improvements*

Our private cloud, FermiCloud, was running OpenNebula 3.2 at the time. The “econe” emulation of the AWS EC2 API was functional for launching small numbers of virtual machines and had been used at times when an experiment needed a batch slot with memory larger than the then-default 2GB per batch slot on FermiGrid. There were known problems if you tried to launch 20 or more virtual machines. The file system of GFS2 on SAN was adequate for the initial 23 hosts that it served but could not be expanded to a large scale. In addition we did not have public IP address space readily available to launch 1000 more virtual machines on the public network.

140 old Dell PowerEdge 1950 machines with dual quad-core Xeon processors and 16 GB of RAM were made available for this test. A private virtual LAN was made available to these nodes for use by the virtual machines. This private LAN was assigned a routable private network which could be reached from anywhere inside Fermilab. We installed a new test head node based on OpenNebula 4.8. For an image repository we used space on our BlueArc NAS server. This proved to be quite scalable.

Configuration of the image was done by a script that ran at boot time and executed a “puppet apply” script to add the minimum configuration necessary to make the node able to run a grid job.

#### *2.5. Performance Evaluation*

Figures 1a and 1b show the results of the final and largest trial, in which 1000 jobs ran simultaneously both on our local private cloud and on Amazon Web Services. The completion time of the jobs was relatively the same. The ramp-up factor on AWS was limited only by the submission rate that was configured into our GlideinWMS factory. The slower ramp-up time for jobs on our local cloud was more due to the virtual machine launching pattern we were using at the time, which caused all 8 of the virtual machines to launch on the same node simultaneously, significantly stressing the local disk. We have since switched to a scheduling algorithm, which distributes the launch more equally across the cloud.

We ran a total of 3300 jobs on AWS in the largest trial, shown in Figure 1A. Each job generated on average 250MB of output, the total output was 467GB for which we incurred \$51 in data transfer charges. We incurred \$398 in virtual machine charges, with a peak of 525 virtual machines running.

### **3. On-demand scalable services**

The summer of 2014 testing demonstrated the successful operation of a classical service discovery model of locating a squid server. Since that time we have demonstrated an automated launch of a squid service in Amazon Web Services. This service makes use of the Elastic Load Balancer service to provide a single entry point to the service, which can have multiple squid servers backing it up, and an Auto Scaling Group to automatically add more squid servers when the load goes higher, and remove them when the load goes down. A CloudFormation service can be used to start the whole stack of services. We plan to use native cloud scaling mechanisms like this wherever possible in future as we address cloud workflows of 100-1000 times larger than the ones described in this work.

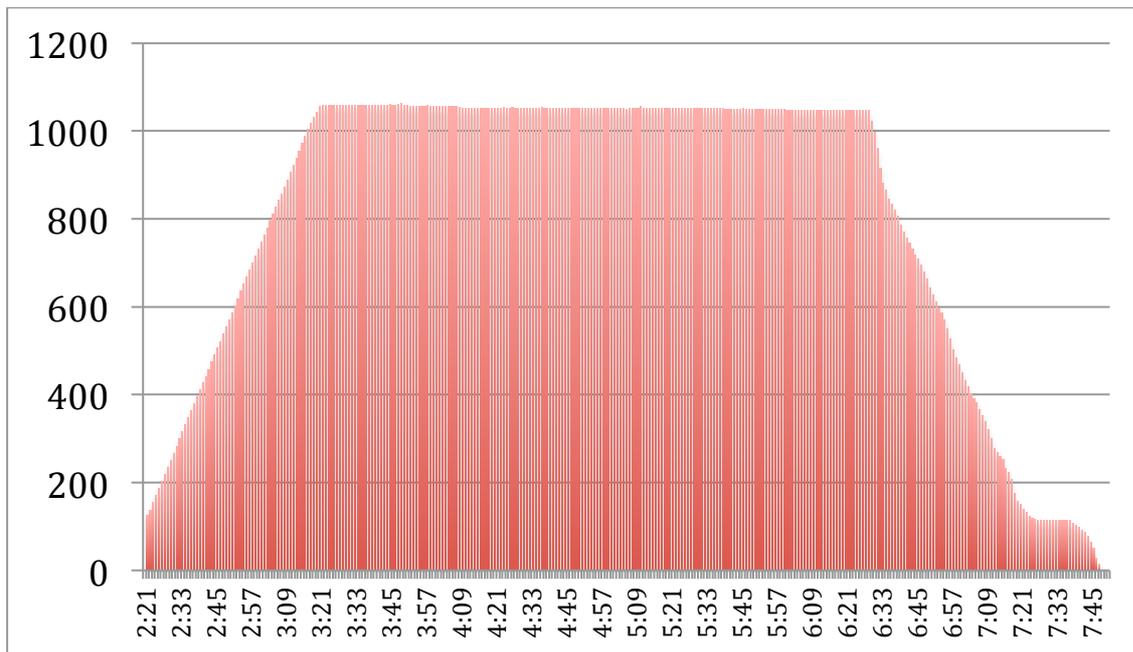


Figure 2a. Number of jobs running as a function of time of day, Amazon AWS

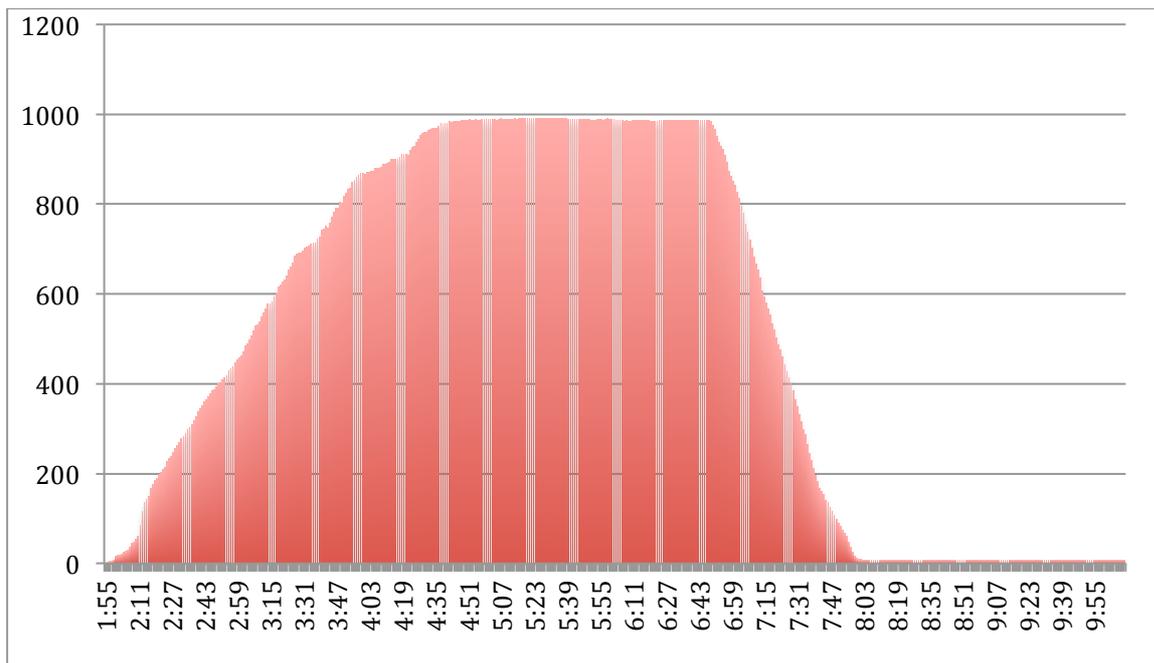


Figure 2b. Number of jobs running as a function of time of day, Fermilab private cloud

#### 4. Conclusions and Future Work

We have demonstrated the capacity to successfully run an Intensity Frontier application at scale both in the public and the private cloud. The cloud submission capacities remain available to our production users. Our future program of work will now focus on integrating the commercial cloud providers more closely into our facility operations. The overall goal is to make the presence of the

cloud resources be transparent to our end users, running even the most data-intensive applications on the cloud if necessary. We are grateful especially for all the work that was done by the NOvA experimenters to help make this happen.

### **Acknowledgements**

Work supported by the U.S. Department of Energy under contract No. DE-AC02-07CH11359, and by CRADA FRA 2014-0002 / KISTI-C14014.

### **References**

- [1] Sfiligoi, I., Bradley, D. C., Holzman, B., Mhashilkar, P., Padhi, S. and Wurthwein, F. (2009). "The Pilot Way to Grid Resources Using glideinWMS", 2009 WRI World Congress on Computer Science and Information Engineering, Vol. 2, pp. 428–432. [doi:10.1109/CSIE.2009.950](https://doi.org/10.1109/CSIE.2009.950).
- [2] J. Blomer et al, Status and future perspectives of CernVM-FS J. Phys.: Conf. Ser. 396052013, doi:10.1088/1742-6596/396/5/052013
- [3] Squid - HTTP proxy server <http://www.squid-cache.org>, 2015
- [4] Gable, Ian, et al. Dynamic web cache publishing for IaaS clouds using Shoal. *Journal of Physics: Conference Series*. Vol. 513. No. 3. IOP Publishing, 2014.
- [5] I. Gable et al, A batch system for HEP applications on a distributed IaaS cloud J. Phys.: Conf. Ser. 331062010, doi:10.1088/1742-6596/331/6/062010
- [6] Python Package Index <https://pypi.python.org/>, 2015
- [7] RabbitMQ - AMQP Messaging software, <http://www.rabbitmq.com>, 2015
- [8] S.Vinoski, Advanced Message Queuing Protocol, IEEE Internet Computing 1087, doi:10.1109/MIC.2006.116