

CC*DNI DIBBs: Data Analysis and Management Building Blocks for Multi-Campus Cyberinfrastructure through Cloud Federation

Monthly Report 3/30/2016

Report 6 of 18

Submitted by David Lifka (PI) lifka@cornell.edu

This is the sixth required monthly report of the Aristotle Cloud Federation team. We report on plans and activities for each area of the project Work Breakdown Structure (WBS).







1.0 Cloud Federation Project Management, Oversight & Reporting Report	3
1.1 Subcontracts	3
1.2 Project Change Request	3
1.3 Project Execution Plan	3
1.4 PI Meetings	3
1.5 Status Calls with ACI Director and Aristotle Program Manager	3
1.6 Project Planning and Preparation	4
2.0 DIBBs Acquisition, Installation, Configuration, Testing & Maintenance Report 2.1 Federation Resource Status Updates	4 4
2.2 Industry Influence: Eucalyptus	5
2.3 Potential Tools: CloudLaunch & Supercloud	5
3.0 Cloud Federation Portal Report	5
3.1 Software Requirements & Portal Platform	7
3.2 Integrating QBETS into Open XDMoD	7
3.3 Allocations & Accounting	8
4.0 Research Team Sunnort	9
4.1 General Undate	
4.2 Science Use Case Updates	9
Use Case 1: A Cloud-Based Framework for Visualization & Analysis of Big Geospatial	Data9
Use Case 2: Global Market Efficiency Impact	9
Use Case 3: High Fidelity Modeling and Analytics for Improved Understanding of Clin	nate-
Relevant Aerosol Properties	9
Use Case 4: Transient Detection in Radio Astronomy Search Data	9
Use Case 5: Water Resource Management Using OpenMORDM	9
Use Case 6: Mapping Transcriptome Data to Metabolic Models of Gut Microbiota	9
Use Case 7: Multi-Sourced Data Analytics to Improve Food Production	10
5.0 Outreach Activities	10
5.1 Media Outreach	
5.2 Industry Outreach	
•	





1.0 Cloud Federation Project Management, Oversight & Reporting Report

1.1 Subcontracts

All subcontracts are in place. Nothing new to report.

1.2 Project Change Request

No new project change requests were made this month.

1.3 Project Execution Plan

The Project Execution Plan (PEP) was approved by NSF on 12/18/2015. We are operating as planned and continuously updating our PEP on a monthly basis.

1.4 PI Meetings

Furlani and Wolski followed up with several External Advisory Committee members who inquired about Aristotle collaborations on QBETS and Open XDMoD. After Ian Foster put Rich Wolski (UCSB) in touch with the Globus Genomics team regarding QBETS capabilities, there was an effort to integrate QBETS into Globus Genomics to better predict the AWS spot market. The initial results were impressive; almost "too good to be true." Further testing seems to confirm the value of QBETS for predicting the AWS spot market. The XDMoD team took the standard QBETS distribution and tested it against the current NSF Service Provider system queues. Again, initial results were impressive and further interactions, particularly around its use for Comet at SDSC, will continue.

Discussions regarding allowing the UB and UCSB science teams to get started on Red Cloud at Cornell also took place. At UB, users are using their local UB cloud ("Lake Effect)" to develop and run scaled down versions of their use cases. Once the UB Aristotle cloud is available, these applications will be migrated. Cornell will provide the UB global finance science team access to Red Cloud in order to accommodate their high volume storage requirements.

Significant effort was put into scheduling the first Science Team Advisory Committee (STAC) meeting. This meeting is now scheduled for 4/1/2016. Lifka and Furlani will be in transit from the CASC meeting, so this meeting will be led by Adam Brazier, Susan Mehringer, and Resa Reynolds. Lifka will help set the agenda.

1.5 Status Calls with ACI Director and Aristotle Program Manager

On 3/7/2016, Lifka had a conference call with ACI director Irene Qualters. Irene was pleased to hear about the project's early results. The main points of discussions included cloud usage paradigms represented by the Aristotle science teams and interoperability between cloud platforms (Eucalyptus, OpenStack, AWS, Azure, and Google).

Lifka was invited by NSF to present an overview of Aristotle Q1/Q2 progress in Arlington on 4/1/2016.





On 3/11/2016, our monthly status call with NSF program manager Amy Walton occurred. Topics of discussion included:

- Updates on Supercloud testing on Jetstream.
- Further discussions regarding the possibility of holding a NSF DIBBs workshop in conjunction with the fall 2016 CASC meeting.
 - Awaiting direction from Amy Walton.
- Updates on XDMoD and QBETS early integration efforts, QBETs success to date with Globus Genomics, and potential use of QBETS by NSF Service Providers.
- Proof of concept activity moving basic science VMs from Cornell to Aristotle partner resources to ensure a consistent user experience. This will be followed by an all-to-all test where each partner tests a local VM at each of the partner sites. These tests are still being prepared.
- Continued progress with the science teams in getting their required software and workflows working.
- Preparing for the Science Team Advisory Board meeting now that the first quarterly External Advisory Committee meeting has occurred.
- Amy Walton pleased with quality of Aristotle project reports to date.

1.6 Project Planning and Preparation

Project planning and preparation by the Aristotle team continued in March 2016 with continued focus on the allocations and accounting system. UB is completing installation and testing of their first year hardware installment from Dell. UCSB has ordered their Dell hardware is starting to arrive. All of these efforts are described in more detail in this month's report.

2.0 DIBBs Acquisition, Installation, Configuration, Testing & Maintenance Report

2.1 Federation Resource Status Updates

• **CU**

The CU team deployed its year 1 hardware at Cornell. We added 168 cores (6GB/core) and 120TB SAN storage to our Red Cloud infrastructure.

• UB

UB's Central IT has acquired the public network space necessary for the Aristotle cloud and the UB team will be installing the Eucalyptus software for Aristotle in the coming weeks.

• UCSB

The UCSB team is waiting for their new hardware to arrive. They received first server the last week in March and more are on the way





	CU	UB	UCSB
		<u>ccr-cbls-</u>	
Cloud URL	euca4.cac.cornell.edu	2.ccr.buffalo.edu**	TBD**
HPE Helion			
Eucalyptus Version	4.2.1	4.2.1	4.1.2
Migrate to 4.2.1	1/1/2016	N.A.	TBD
Globus	Yes	Planned	Planned
InCommon	Yes	Planned	Planned
		Hardware installed.	
	Hardware deployed.	Working through	
	168 cores added to	network reqs. with	
	existing Red Cloud.	central IT. Install	Hardware
Hardware Quotes	376 total cores.	Euca stack next.	ordered.
Hardware Vendor	Dell	Dell	Dell
# Cores	168*	112-140 (target)	140 (target)
Ram/Core	4GB/6GB/8GB	6GB	8GB
10Gb Interconnect	Yes	Yes	Yes

The CU/UB/UCSB infrastructure planning table has been updated:

* 168 additional cores augmenting the existing Red Cloud

**UB and UCSB installing Aristotle as new cloud (not integrating with existing clouds)

2.2 Industry Influence: Eucalyptus

No updates this month.

2.3 Potential Tools: CloudLaunch & Supercloud

The Cornell team continues to massage the CloudLaunch code so that it can be run by outside clients. The Cornell CS Supercloud group was given access to Jetstream this month and is in early test phase. The first test will be to migrate a VM from Jetstream to Red Cloud and back. We plan to report results next month.

3.0 Cloud Federation Portal Report

No changes were made this month to the portal planning table below:

Portal Framework			
Phase 1	Phase 2	Phase 3	Phase 4
10/2015 - 3/2016	4/2016 - 10/2016	11/2016 - End	1/2017 - End
Gather portal requirements, including software requirements, metrics, allocations, and accounting. Install web site software.	Implement content/functionality as shown in following sections. Add page hit tracking with Google Analytics, as well as writing any site downloads to the	Implement content/functionality as shown in following sections. Add additional information/tools as needed, such as selecting where to run based on software/hardware needs	Release portal template via GitHub. Update periodically.
	database.	and availability.	



Documentation			
Phase 1	Phase 2	Phase 3	Phase 4
10/2015 – 3/2016	4/2016 - 10/2016	11/2016 – End	1/2017 - End
Basic user docs, focused on getting started. Draw from existing materials.	Update materials to be federation-specific.	Add more advanced topics as needed, including documents on "Best Practices" and "Lessons Learned." Check and update docs periodically, based on ongoing collection of user feedback.	Release documents via GitHub. Update periodically.
Training			
Phase 1	Phase 2	Phase 3	Phase 4
10/2015 - 3/2016	4/2016 - 10/2016	11/2016 - 3/2017	4/2017 - End
Cross-training expertise across the Aristotle team via calls and 1-2 day visits.	Hold 1 day training for local researchers. Offer Webinar for remote researchers. Use recording and materials to provide training asynchronously on the portal.	Add more advanced topics as needed. Check and update materials periodically, based on training feedback and new functionality.	Release training materials via GitHub. Update periodically.
User Authorization and Ke	ys		
Phase 1	Phase 2	Phase 3	Phase 4
10/2015 - 1/2016	2/2016 - 5/2016	6/2016 - 9/2016	10/2016 – End
Plan how to achieve seamless login and key transfer from portal to Euca dashboard.	Login to the portal using InCommon.	Get 4.2.1 federated key after InCommon login.	Move seamlessly to Euca console after portal InCommon login.
Euca Tools		-	
Phase 1	Phase 2	Phase 3	Phase 4
10/2015 - 3/2016	4/2016 - 12/2016	1/2017 – End	1/2017 – End
Establish requirements, plan implementation.	Implement minimal set of Euca Tools to bridge portal to Euca console.	Add/refine/update, based on ongoing collection of user feedback.	Release via GitHub. Update periodically.
Allocations and Accounting			
Phase 1	Phase 2	Phase 3	Phase 4
10/2015 - 3/2016	3/2016 - 5/2016	6/2016 - 9/2016	10/2016 – End
Plan requirements and use cases for allocations and account data collection across the federation. Design database schema for	Implement project (account) creation in the database and display on the portal. Integration hooks for user and project creation/deletion	Automate project (account) creation by researcher, via the portal.	Report on usage by account, if the researcher has multiple funding sources. Release database schema via GitHub.
osers, i rojects anu		1	1



collections of CPU usage and Storage Usage of the federated cloud.	and synchronization across sites.		
Metrics and Usage			
Phase 1	Phase 2	Phase 3	Phase 4
10/2015 - 3/2016	2/2016 - 5/2016	6/2016 - 10/2016	11/2016 - End
Implement graphs of basic usage data, including % utilization, available resources, and user balance, using scripts from Cornell CAC for basic data collection. Provide documentation for installing XDMoD and SUPReMM at individual sites.	Install Open XDMoD/SUPReMM at individual sites and begin data collection. This includes the installation of SUPReMM and the data collection piece at the federation sites. Begin integration with federated authentication providers.	Federated data collection across sites. Ship data from the individual sites to UB. We can summarize data remotely and send the summarized data or collect all raw data and summarize locally. Other job information will be federated as well using the prototype model under development with OSG. Display federated metrics in Open XDMoD at UB.	Release materials via GitHub. Update periodically.

3.1 Software Requirements & Portal Platform

Planning continued on the portal, currently populated with a "Coming Soon" page at <u>http://www.federatedcloud.org</u>. The SSL certificate has been installed. Next, we will add InCommon authentication. In order to make the portal easy and free to recreate, a Bootstrap web framework was chosen to handle formatting. The first draft implementing the framework and the portal design has been completed. Next month, we expect to implement the framework and start populating content.

3.2 Integrating QBETS into Open XDMoD

Wolski revived a clustering version of QBETS which allows for tighter batch job wait time predictions based on using clustering algorithms to predict restricted classes of job wait times. The code can cluster jobs based on requested maximum wait time, number of processors, or a product of the two. The clustering version of QBETS was downloaded from GitHub and compiled by UB Center for Computational Research (CCR) personnel. We ran the code in all of the different clustering modes using CCR Rush batch job data. Typically, 1-3 clusters were found for either clustering by requested maximum run time or by processor. In some cases the bounds on the predicted wait times within the various clusters were tightened compared to the non-cluster single QBETS prediction. We also ran it successfully using XSEDE TACC Stampede batch job wait data. The next step will be to run clustering QBETS on other XSEDE resource data sets and compare the predicted wait times between the resources. Unfortunately, most other XSEDE resources do not currently report requested wait time to the TGCDB. We are working on obtaining the requested wait time directly through the SUPReMM data pipelines. When this data has been obtained, we will run clustering QBETS on the wait time data and compare between XSEDE resources to see if clustering QBETS would be a viable batch job wait prediction tool.





3.3 Allocations & Accounting

Last month we finalized our plans for data collection and reporting for compute time, and drafted our plans for storage usage. Defined use cases include researchers who have an allocation, who have purchased time, who want to run across the federation, and all combinations thereof. We established a convention for federation ProjectIDs that are unique across the federation, and that drawing from multiple accounts must be done under separate ProjectIDs. This will allow us to accurately account and report usage. The database schema for allocations and usage data has been updated (see new diagram below). The database has been created. This month, we began implementing this plan; four project groups have been created at Cornell, with project, user, and usage data, enabling the researchers to begin setting up instances.

UB has developed a REST API that will expose Eucalyptus accounting data based on the scripts provided by Cornell. This month Cornell began testing and implementing this API for use within the federation.







4.0 Research Team Support

4.1 General Update

4.2 Science Use Case Updates

Use Case 1: A Cloud-Based Framework for Visualization & Analysis of Big Geospatial Data We deployed the base webglobe server on <u>http://128.205.11.214/wglobe/Index.html</u>. The server is currently running on the UB Lake Effect cloud and will be eventually migrated to the UB Aristotle cloud. The analysis stack currently setup in the cloud consists of:

- An HDFS with climate data stored as distributed NetCDF files.
- iGlobe server which analyzes the distributed NetCDF data using the SciSpark API which runs on top of a Spark cluster. The iGlobe server creates the Spark cluster "on-demand" and destroys the cluster after the analysis is done.
- Web server that hosts the webglobe client. The client allows users to interact with the underlying climate data through their browsers.
- For the next step, we will deploy more advanced distributed analytic codes to interact with the data. The web interface will also be developed further to allow users better interaction with the climate data.

Use Case 2: Global Market Efficiency Impact

We developed and tested the capability to create and destroy VMs with SAS installed on the machines. The primary researcher has been able to execute the basic codes as part of the use case. Currently, more extensive jobs are being executed.

Use Case 3: High Fidelity Modeling and Analytics for Improved Understanding of Climate-Relevant Aerosol Properties

Instances (single-process and distributed-memory) were migrated from a development project to a newlycreated Aristotle project for testing and development.

Use Case 4: Transient Detection in Radio Astronomy Search Data

Brazier met with Cordes (PI) and Chatterjee and representatives of two other radio astronomy projects: the Long Wavelength Array (LWA-New Mexico) and the Murchison Widefield Array (MWA-Australia) to discuss the possibility of adding additional data sets. Cordes, Chatterjee and Brazier planned basic software architecture and functional requirements.

Use Case 5: Water Resource Management Using OpenMORDM

No updates this month. We plan to build instances for OpenMORDM in April 2016.

Use Case 6: Mapping Transcriptome Data to Metabolic Models of Gut Microbiota

Barker (CU) communicated with PI Douglas on possible additional requirements, as well project setup on Aristotle. He is investigating the use of Supercloud as a possible means to scale out to other clouds. Initial versions of libSBML and COBRApy packages for Nix were built and software documented.

A Windows instance is preferred by this science team for MATLAB. This was created and user documentation is being added, including edits pertaining to how to create a Windows instance. We set up





an NFS server for serving application files and for investigating MATLAB Linux network client installation to NFS volume. A Samba server was set up, tested and documented to allow file sharing across Windows and Linux. We also documented scalable file access for cross-cloud or many-instance file/app access and for starting MATLAB from a shared filesystem.

Use Case 7: Multi-Sourced Data Analytics to Improve Food Production

The soil moisture sensing project experienced a hardware failure. Sedgwick is meeting with the vendor in early April to get new equipment. We have about a month's worth of data, but it is intermittent due to hardware instability in the sensors. A new site will come online as soon as the hardware issues are resolved.

The image processing survey is proceeding. Currently, the camera trap data is downloaded manually from each camera once a quarter or so. The latest data will be uploaded in early April.

Soil electrical conductivity (EC) survey is also scheduled for early April although it may still be too wet. A consultant will make the determination.

The Citizen Science deer survey is being rescheduled.

Sedgwick found a tenant for a test farm. Planting of peppers will be scheduled after the first set of soil moisture sensors are sited.

All computer science components are up and functioning. We are waiting on data from the various science projects which are subject to pretty typical agricultural logistics.

5.0 Outreach Activities

5.1 Media Outreach

The "Will the cloud change scientific computing" article was featured by *InsideHPC*: <u>http://insidehpc.com/2016/02/will-the-cloud-change-scientific-computing/</u>.

5.2 Industry Outreach

Cornell wrote/edited an HPE Helion Eucalyptus case study on Red Cloud and Aristotle that will be available online in a month or two.

