

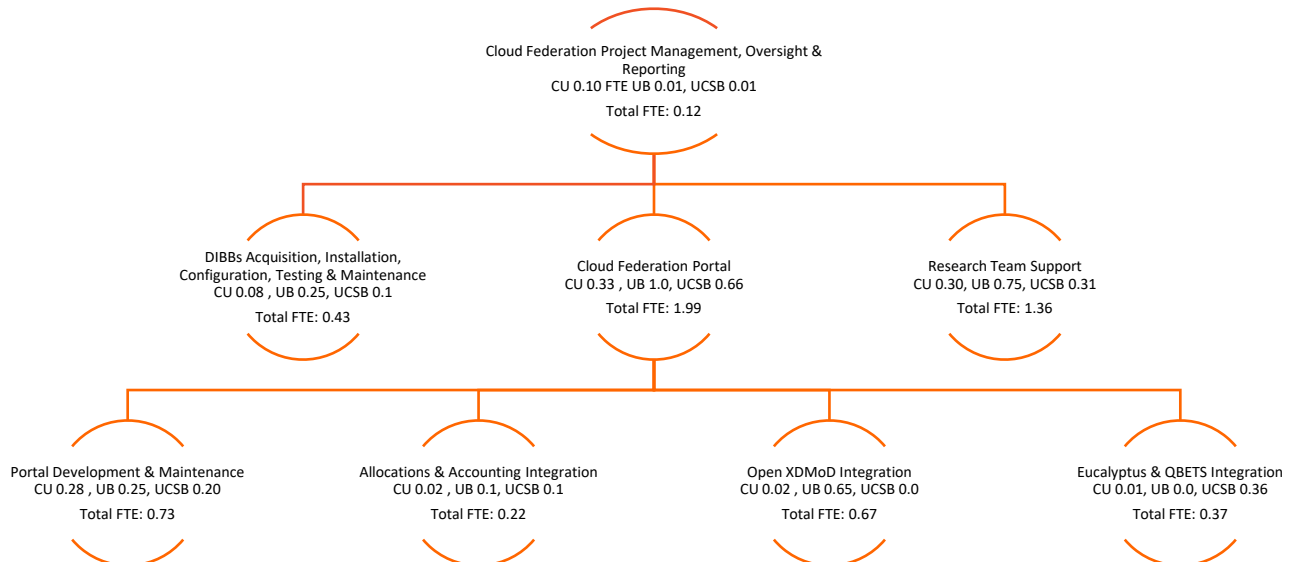
# CC\*DNI DIBBs: Data Analysis and Management Building Blocks for Multi-Campus Cyberinfrastructure through Cloud Federation

**Monthly Report 5/31/2016**

**Report 8 of 18: Part 1 of 2**

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

This is the eighth 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).





## Contents

<b>1.0 Cloud Federation Project Management, Oversight &amp; Reporting Report .....</b>	<b>3</b>
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 .....	3
1.6 Project Planning and Preparation .....	4
<b>2.0 DIBBs Acquisition, Installation, Configuration, Testing &amp; Maintenance Report .....</b>	<b>4</b>
2.1 Federation Resource Status Updates .....	4
2.2 Potential Tools: CloudLaunch & Supercloud .....	5
2.3 Industry Influence .....	5
<b>3.0 Cloud Federation Portal Report .....</b>	<b>6</b>
3.1 Software Requirements & Portal Platform .....	7
3.2 Integrating Open XDMoD and QBETs into the Portal .....	7
3.3 Allocations & Accounting .....	8
<b>4.0 Research Team Support .....</b>	<b>9</b>
4.1 General Update .....	9
4.2 Science Use Case Updates .....	9
Use Case 1: A Cloud-Based Framework for Visualization & Analysis of Big Geospatial Data ..	9
Use Case 2: Global Market Efficiency Impact .....	9
Use Case 3: High Fidelity Modeling and Analytics for Improved Understanding of Climate- Relevant Aerosol Properties .....	10
Use Case 4: Transient Detection in Radio Astronomy Search Data .....	10
Use Case 5: Water Resource Management Using OpenMORDM .....	10
Use Case 6: Mapping Transcriptome Data to Metabolic Models of Gut Microbiota .....	10
Use Case 7: Multi-Sourced Data Analytics to Improve Food Production .....	10
<b>5.0 Outreach Activities .....</b>	<b>10</b>
5.1 Presentation .....	10
5.2 Community Outreach .....	10

## 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

The team held two meetings this month, the first covering progress across all Work Breakdown Structure areas and the second was a deeper dive into the capabilities and limitations of Supercloud. The Cornell Supercloud team from Computer Science demonstrated Supercloud's ability to migrate a running virtual machine (VM) from Cornell's Red Cloud (part of the Aristotle Federation) to Indiana's Jetstream, then to Microsoft's Azure, and finally, AWS. There was a great deal of discussion regarding how Supercloud works, how it performs network routing between different cloud providers, and what level of performance overhead is introduced by Supercloud. Two important questions were addressed:

1. What public cloud APIs does Supercloud support?  
*Answer: Supercloud supports OpenStack APIs. If you have a VM which makes API calls to a specific cloud provider (e.g., AWS) that VM will only work on AWS. It would, however, allow you to migrate between different AWS availability zones so if the cloud supports OpenStack APIs, then a VM using those APIs can migrate to that cloud without issue.*
2. Does Supercloud provide log data for accounting purposes to track the migration of a VM from one cloud provider to the next? (note: this is something Aristotle and other federations would require)  
*Answer: Not today, but it is something they know how to do and with funding would be made a priority.*

Overall Supercloud appears to have a lot of potential. Next step: we will next ask the Supercloud team to do a demonstration for the Aristotle External Advisory Committee. This is one of several technologies we are investigating that we believe can have a significant impact on the national community. We are working toward a supplemental proposal to move these technologies forward.

### 1.5 Status Calls

Project status calls were held on 5/10/2016 and 5/24/2016. Topics included:

- Aristotle team review of Supercloud with Supercloud team at Cornell to better understand its capabilities and limitations.
- CU and UB infrastructure team discussion on availability zones. Subsequently, UB decided to set up one availability zone for their Aristotle and Lake Effect clouds; CU also set up one zone. The CU and UB infrastructure teams will be meeting to discuss these implementations and the greater CI community need for a broadly applicable, open source accounting system.
- CU staff and partners discussed plans to implement the rest API and show simple graphs initially to be replaced eventually with Open XDMoD.

## 1.6 Project Planning and Preparation

Initial portal element design was completed and will be launch in mid-June.

A draft “XDMoD Requirements Document – Job Reporting for Cloud and Other Non-Traditional HPC Resources” was created this month. See Appendix A (pages 11-20).

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 and UCSB discussed Ceph configurations for providing EBS and S3 cloud storage. The HPE Helion Eucalyptus software was built against the “Hammer” version. Ceph documentation for the newer “Jewel” version shows improved RADOS gateway support. We agreed that the best practice is to go with “Hammer.”

CU, UB, and UCSB discussed cloud network configurations. CU explained why they are using “Edge” mode: it’s easy to deploy and there’s no shortage of IP space. UCSB talked about possibly using MidoNet (VPC).

- **CU**  
Cornell’s Red Cloud infrastructure was upgraded to Eucalyptus 4.2.2.
- **UB**  
UB’s Aristotle cloud is online with Eucalyptus version 4.2.2. UB staff are testing the system to flush out any configuration problems prior to adding Aristotle users.
- **UCSB**  
UCSB made significant progress this month in both services and software.

#### Services:

- Eucalyptus 4.2.2. running and verifying operation
- Ceph 9.94.7 (Hammer) for EBS and S3 storage running and verifying operation
- Planning UCSB’s website (console.aristotle.ucsb.edu)

#### Software:

- Integrating Ansible for provisioning and orchestration
- Planning integration with campus identity
- Configuring ELK Stack for monitoring and metrics.

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

	CU	UB	UCSB
<b>Cloud URL</b>	euca4.cac.cornell.edu	ccr-cbls-2.ccr.buffalo.edu	console.aristotle.ucsb.edu
<b>HPE Helion Eucalyptus Version</b>	4.2.2	4.2.2	4.2.2
<b>Migrate to 4.2.1</b>	1/1/2016	Completed	5/1/2016
<b>Upgrade to 4.2.2</b>	5/25/2016	5/2016	5/25/2016
<b>Globus</b>	Yes	Planned	Planned
<b>InCommon</b>	Yes	Yes	Yes
<b>Hardware Status</b>	Deployed 3/2016	Internal user testing	Installation testing
<b>Hardware Vendor</b>	Dell	Dell	Dell
<b># Cores</b>	168*	144**	140**
<b>Ram/Core</b>	4GB/6GB/8GB	8GB	9GB
<b>Storage</b>	Dell SAN	Dell SAN	Ceph
<b>10Gb Interconnect</b>	Yes	Yes	Yes

\* 168 additional cores augmenting the existing CU Red Cloud

\*\*UB and UCSB Aristotle clouds are separate from existing campus cloud resources

## 2.2 Potential Tools: CloudLaunch & Supercloud

CloudLaunch is on hold; development efforts will resume during summer 2016.

Supercloud capabilities were demonstrated by Cornell CS to CU, UB, and UCSB on 5.24.2016 via WebEx. Supercloud source code will be made available to the Aristotle team so we can investigate how Supercloud's features might be incorporated into the Federation. The next step is for the Supercloud team to give a demo to the Aristotle External Advisory Board.

## 2.3 Industry Influence

A conference call was held between the CU and HPE Helion Eucalyptus teams on 5.28.2016 to discuss adding Security Assertion Markup Language (SAML) support to the HPE Helion Eucalyptus console which would allow the Aristotle Federation to use InCommon authentication. The discussion was very positive and the HPE team was enthusiastic to help. Further conversations amongst the CU team regarding CILogon and Globus Auth (the identity and access management infrastructure being adopted by XSEDE) led to another discussion with the HPE team. We changed our request from SAML support to OAuth2 support. Implementing OAuth2 will allow the Federation to not only use InCommon but also use any identity providers supported by Globus Auth. CU will make arrangements with Globus to register a Globus Auth resource provider for development and testing by HPE.

HPE is discussing this feature with its architectural team and will provide Aristotle team with plans for implementing it.



### 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 database.	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 and availability.	Release portal template via GitHub. Update periodically.
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 Keys			
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 Users, Projects and collections of CPU usage and Storage Usage of the federated cloud.	Implement project (account) creation in the database and display on the portal. Integration hooks for user and project creation/deletion and synchronization across sites.	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.
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

The new portal design was implemented in May and will be publicly available in mid-June.

### 3.2 Integrating Open XDMoD and QBETs into the Portal

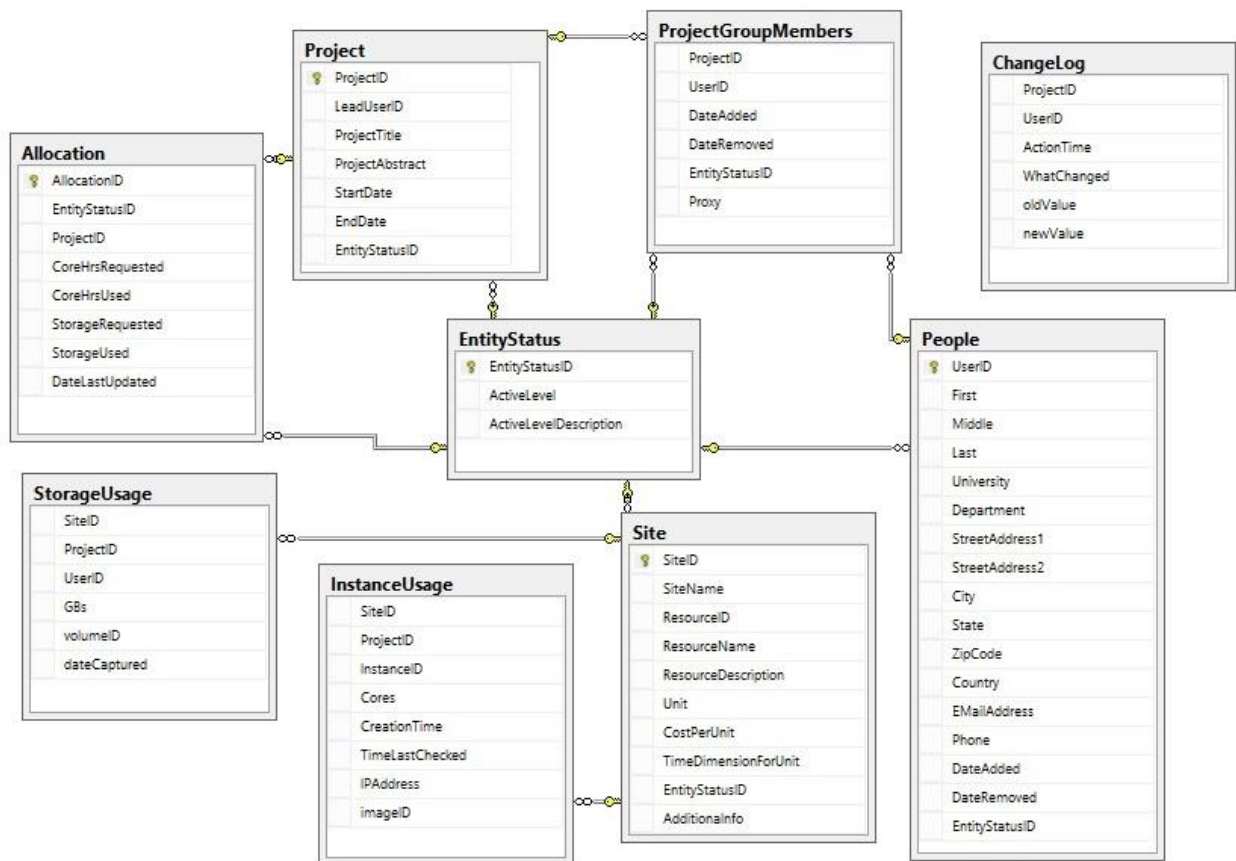
The existing XDMoD data warehouse was developed primarily to report on data generated by individual HPC jobs run on traditional HPC resources. With the advent of open source cloud solutions such as Eucalyptus and Open Stack, as well as non-traditional HPC resources such as Hadoop running alongside

traditional HPC clusters at many centers, we must re-examine the infrastructure used to store and report on center utilization as well as the definition of an HPC job within XDMoD. The capabilities of the XDMoD data warehouse will be updated to support these new resource types and become more flexible to better manage new types developed in the future. This includes cloud accounting and low-level job performance, job reservations, and job arrays.

A 1.0 draft of the “XDMoD Requirements Document: Job Reporting for Cloud and Other Non-Traditional HPC Resources” was created on 5.23.2016 and is available in the Appendix (section 6.0) of the report.

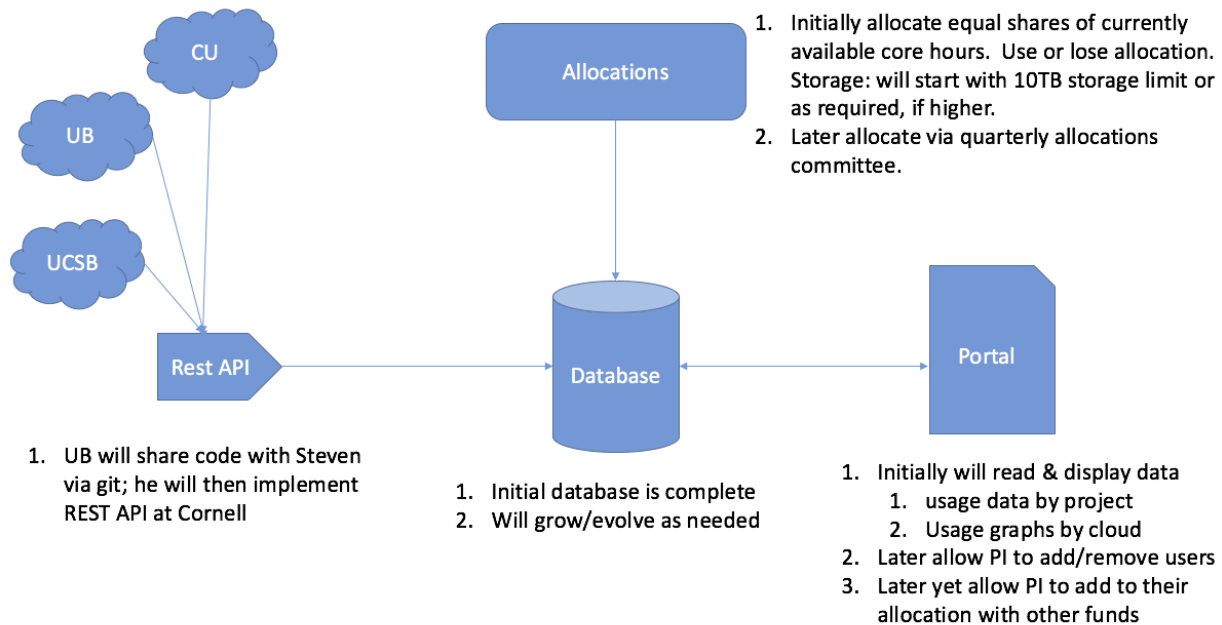
### 3.3 Allocations & Accounting

The database schema for allocations and usage data, included below, has been expanded to include data on Site Health. This capability was added so that staff can post patch days, scheduled down times, unscheduled down times, updates, etc.





Current plans for allocations and tracking have been updated as shown in the figure below:



## 4.0 Research Team Support

### 4.1 General Update

- First quarter allocations agreed for use on CU's Aristotle node with fixed use-or-lose core-hours available for all, pending full deployment of the accounting system. Announcement to researchers is imminent.
- All Cornell use cases (#3-#6) now have working base images, with more software being installed by and in conjunction with the science teams.
- Use case #6, "Mapping Transcriptome Data to Metabolic Models of Gut Microbiota," has been selected for trial Dockerization.

### 4.2 Science Use Case Updates

#### Use Case 1: A Cloud-Based Framework for Visualization & Analysis of Big Geospatial Data

The interactive visualization interface and the associated backend system is now complete. We are currently doing extensive testing of the interface.

#### Use Case 2: Global Market Efficiency Impact

The migration of the OpenNebula image to the Aristotle cloud is now complete and will be documented for the portal. This advance will enable the PI to seamlessly migrate their application to Aristotle. Agreement for purchase of the Thomson Reuters Tick History (TRTH) database was secured for the

second finance use case and that database will be available to all Federation members at the Aristotle Cloud Federation UB node.

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

Initial user testing of the image is underway with an MPI version of WRF Chem.

### **Use Case 4: Transient Detection in Radio Astronomy Search Data**

PRESTO, a large suite of pulsar search and analysis software used for transients detection, and all its dependencies, was installed. The research team is working on porting single-pulse detection software to the image. Planning of data decimation (factor of 4 in time resolution and 2 in frequency) is underway.

### **Use Case 5: Water Resource Management Using OpenMORDM**

An external collaborator, David Hadka (The Pennsylvania State University's Applied Research Laboratory), was added to this project. Hadka developed the Borg Multiobjective Evolutionary Algorithm with Patrick Reed which is particularly effective at solving complex engineered systems in civil engineering and other fields.

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

This research project will be containerized. Bessem Chouaia, a CU Postdoc in the Douglas Lab, will define requirements for a Docker image. Demonstrating containerization will be a key technical advancement.

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

UCSB campus networking is installing long-range WiFi at the Sedgwick Reserve. No currently-known ETA on its availability but when it goes in, the plan is to automate the acquisition of the animal imagery data. Currently the collection is done manually via USB thumb drive followed by an upload.

Research teams are now starting to get good data from the soil moisture sensors (the ones where we had previous installation problems). Rain last week registered properly and they went through an irrigation cycle and found that they were putting on way too much water. This is excellent news. There is also a student signed up to work on the animal imagery application this summer. Work is needed on the data management architecture, however, because currently the upload to box.com (UCSB's preferred data preservation system) will take 49 days.

## **5.0 Outreach Activities**

### **5.1 Presentation**

Tom Furlani provided a presentation on the Aristotle DIBBs project at the Best Practices in Data Infrastructure Workshop held at the Pittsburgh Supercomputing Center, 5.17 - 5.18.2016:  
<https://www.psc.edu/index.php/bpdi-workshop>.

### **5.2 Community Outreach**

MIT inquired about the Aristotle project and Lifka briefed Chris Hill on 5.20.2016. Hill develops software and computing infrastructure for the simulation of atmospheric, oceanic, and geophysical flows.