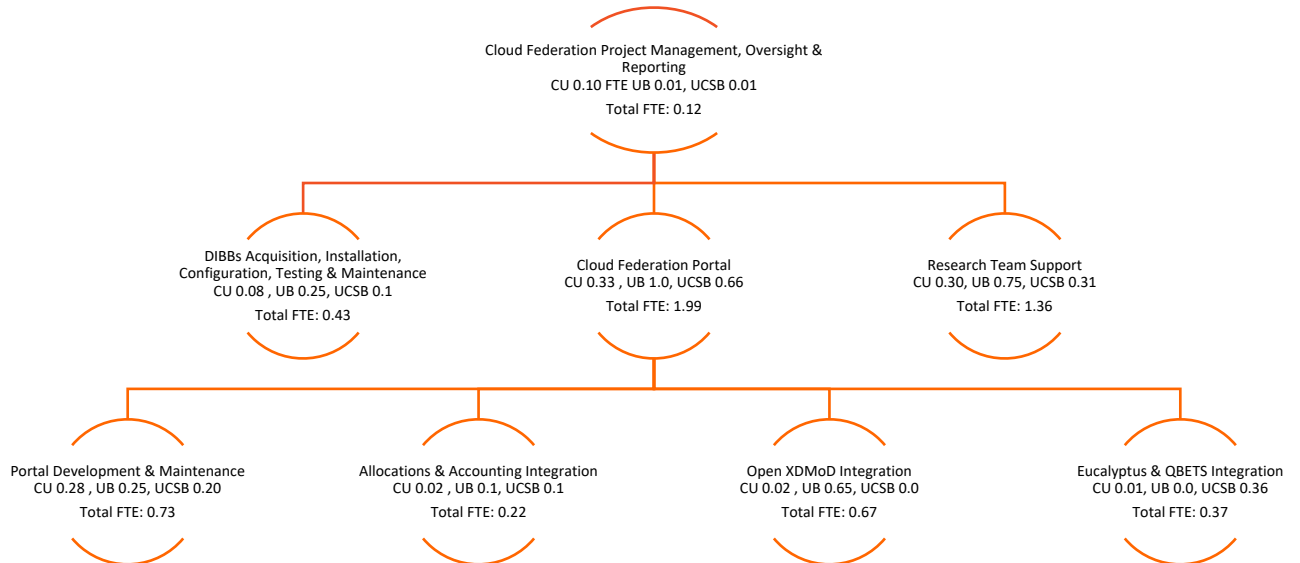# CC*DNI DIBBs: Data Analysis and Management Building Blocks for Multi-Campus Cyberinfrastructure through Cloud Federation

## Program Year 5: Quarterly Report 1

### 12/15/2019

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

This is the Program Year 5: Quarterly Report 1 of the Aristotle Cloud Federation team. We report on plans and activities for each area of the project Work Breakdown Structure (WBS).

Cloud Federation Project Management, Oversight & Reporting
CU 0.10 FTE UB 0.01, UCSB 0.01
Total FTE: 0.12

DIBBs Acquisition, Installation, Configuration, Testing & Maintenance
CU 0.08 , UB 0.25, UCSB 0.1
Total FTE: 0.43

Cloud Federation Portal
CU 0.33 , UB 1.0, UCSB 0.66
Total FTE: 1.99

Research Team Support
CU 0.30, UB 0.75, UCSB 0.31
Total FTE: 1.36

Portal Development & Maintenance
CU 0.28 , UB 0.25, UCSB 0.20
Total FTE: 0.73

Allocations & Accounting Integration
CU 0.02 , UB 0.1, UCSB 0.1
Total FTE: 0.22

Open XDMoD Integration
CU 0.02 , UB 0.65, UCSB 0.0
Total FTE: 0.67

Eucalyptus & QBETS Integration
CU 0.01, UB 0.0, UCSB 0.36
Total FTE: 0.37

## Contents

**1.0 Cloud Federation Project Management, Oversight & Reporting**

**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 quarter.

**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/Partner Meetings**

- Met with Weston Maggio, RightScale manager, to explain the intent of our open marketplace investigation (analyzing the cost and performance of running select benchmarks and applications on AWS, Azure, and Google using methods such as Terraform/Ansible automatic cluster creation). Secured additional technical team support from RightScale.
- Met with Sanjay Padhi, AWS head of research, and Stephen Fang, Google higher education manager. Secured public cloud credits from both providers to further the investigation.

**1.5 Project Planning Meetings/Status Calls**

10/8/2019 status call:
- Aristotle science use case lead Adam Brazier has been named technical lead for the new NSF Scalable Cyberinfrastructure Institute for Multi-Messenger Astrophysics (SCIMMA) grant. This 2-year $2.8M project is led by the University of Wisconsin, Milwaukee (UWM) and is supported by 10 collaborators, including Cornell. The code for SCIMMA will be developed and tested on Aristotle. This new use case is expected to have broad impact in the astrophysics community. See the Cornell news story at https://news.cornell.edu/stories/2019/10/cornell-partners-nsf-grant-astrophysics-institute and UWM announcement written by P. Redfern at https://uwm.edu/news/2-8-million-grant-boosts-effort-by-uwm-and-others-to-harness-cosmic-data.
- Dartmouth is currently debugging JuypterHub on Cornell's Red Cloud; this multi-user Hub will support classes at Dartmouth. It will be moved to Dartmouth's own cloud once that cloud is fully operational. This is an example of how a federation member can get an application up and running on another member's cloud, so that when their cloud is operational, they can simply move the app over. Dartmouth is also benefitting from cloud installation assistance from Cornell, UB, and UCSB which is saving Dartmouth systems administrators time and effort.
- The USCB science team redesigned and implemented new boards and is now preparing an edge cloud to monitor a California citrus test orchard for the upcoming frost season.
- Analysis of differential irrigation sensor data from nearly two seasons of almond grove growth is underway; interpretation of that data is challenging due to a very large number of sensors. The goal of this project is to use insights from this data to schedule and improve irrigation, taking into account specifics such as the difference between irrigating the sunny and shady sides of the trees.
- The UCSB "Edible Campus" program (growing food and establishing a food bank for students at food risk) has been initiated; the project team is interested in using the Aristotle cloud to process

farm sensor data. An Aristotle REU student is working with the project team to design the instrumentation.

- Cornell is deploying Magnum on a test cloud with Puppet. Magnum makes container orchestration services such as Docker Swarm and Kubernetes available in OpenStack but is not officially part of the OpenStack platform. Dartmouth will test Magnum once the deployment is complete and provide feedback to the federation. Cornell is also looking at Red Hat Director, a toolset for installing and managing the OpenStack environment that is based primarily on TripleO.
- UB will look at data in RightScale to see what public cloud information the platform provides and to see if and how they can get that information into XDMoD.
- Steve Gallo is leaving the UB Aristotle team. Andrew Bruno will be the new Aristotle infrastructure lead at UB.

10/9/2019 progress meeting:
- We are working on improving JSON file security and the Aristotle dashboard, including adding a usage graph from the user perspective (in addition to the current PI perspective).
- Portal lead Mehringer and Cornell staff will create templates for the Aristotle portal and Aristotle database which will be made publicly available in PY5 on GitHub. The templates and database will be available so that other interested institutions can create their own cloud federations.
- All Cornell science use cases have been successfully containerized.
- While we know how to tag usage and run on credits with RightScale, the question is how to work with actual cloud bills and whether daily rather than minute-by-minute price information is the only price information available.
- All Aristotle instances of XDMoD are currently at v.8.1.2. V8.5.0 was released on 10/21/2019 but the federation module has not yet been completely tested to verify that it works with this new version. The XDMoD Aristotle federation can be viewed at https://aristotle-hub.ccr.xdmod.org/#main_tab_panel:about_xdmod?Federated. It appears that UCSB hasn't loaded data since 10/12/2019, so we need to investigate that.
- All federation sites are responsible for their own XDMoD updates rather than UB.

10/22/2019 status call:
- There is not currently a method via XDMoD to get information on how much storage a user has used. We don't have a storage mechanism to pull that information from the other sites. UB and Cornell will explore using XDMoD to capture this storage utilization information. We are collecting other cloud metrics in the meantime.
- UCSB will ramp up work on the DrAFTS 2.0 AWS pricing schema after a new hire is onboarded.
- Aristotle consultants helped get the Cornell REU student and Pat Reed water resource management team up and running on Azure.
- All instrumentation for the Citrus Under Protective Screening (CUPS) project (to protect CA citrus from the Asian citrus psyllid and citrus greening) will be hosted by Aristotle.
- Dartmouth continues to make progress on installing their first cloud. UCSB is helping them with configuration details. Cornell has completed the installation of Magnum on its test cloud for Dartmouth's use.
- Cornell developed a workaround to clean up instances that should not be running after a job is completed.
- UB is working with Red Hat on an OpenStack glitch, added GPU nodes, and purchased Globus.

**10/23/2019 planning meeting:**

- It's not clear Red Hat is committed to Magnum due to their own product initiatives.
- We'll be hooking up Azure to RightScale next to see what usage data can be generated.

**11/5/2019 status meeting:**

- Dartmouth may deploy a digital humanities use case on Aristotle's Red Cloud to support their Media Ecology Project: http://digitalhumanities.dartmouth.edu/projects/the-media-ecology-project/. This project uses ML to analyze significant amounts of video and archival film data.
- It was re-emphasized that each site is responsible for updating their XDMoD.
- UCSB is targeting next quarter for a prototype DrAFTS 2.0 and the following quarter for something fully operational.
- We confirmed allocations and accounting is working across all sites and is fully automated.
- Cornell upgraded their test cluster to Nautilus which still works with the test OpenStack cloud. They are also installing 4 new nodes and 1 new GPU node.
- UCSB is targeting this year's funds for more storage assets. The camera trap project has moved into Phase 2 which will require storage of a lot more image data and likely GPUs.
- Cornell is comparing runtimes on Open MORDM on classical HPC vs. cloud and also on a cluster at CAC. A doctoral student hopes to include these results in their dissertation and will credit Aristotle. Research on developing clusters of containers in the cloud continues and the idea of eventually developing a half day tutorial for a conference submission was discussed. No one really knows the answer on whether you can get full network performance in MPI running in containers. If there is performance degradation, we would need to show why and where.

**12/3/2019 status meeting:**

- Dartmouth has deployed their first OpenStack Cloud with the assistance of the Aristotle team. TripleO installation challenges and lessons learned are being documented for the benefit of others.
- UCSB scientists and collaborators are taking their first citrus frost prevention temperature measurements and ground has been broken on the Citrus Under Protective Screening (CUPS) project whose infrastructure will be supported by Aristotle. A final paper on differential irrigation of California almond trees is being written. In addition, a new camera trap project is underway which will provide real-time images through a solar-powered network pipeline.
- UCSB is hiring a new programmer to work on DrAFTs 2.0 which will be accessible through the Aristotle portal.
- UB met with Red Hat Engineering at SC19. We are currently running a TripleO supported version of Red Hat (v.13). We will need to move to a Director-based version in 2020. Going from v.13 to v.16 will be a complete reinstall (no direct path). eGPU support is available right now in v.13 with proper patches. If we update to the latest packages in v.13 repos, we should be able to reconfigure eGPU support. Red Hat offers help installing Director 16 if needed; they can do custom consulting onsite for a day or so to get us going. We're thinking about this.

**12/11/2019 planning meeting:**

- We plan to work on a paper (tentatively titled "Cost Effective Workflows for Scientific Computing and HPC in the Cloud") that compares costs and performance of running Linear Algebra Benchmarks, Lake_Problem (a multi-VM MPI app with low communication), WRF (a multi-VM MPI with data management considerations) and FRB_Pipeline (an embarrassingly parallel app that can use the AWS Spot market to process lots of data with no communications). These apps will be

tested on multiple clouds (AWS, Azure, and/or Google Cloud). Methods will be discussed, including our use of Terraform/Ansible for automated cluster creation.

- Next steps are: (1) measure the cloud cost of a small scripted WRF run on AWS (later on Goggle, use rates times runtime, figure out Azure billing), (2) work with the Pryor Group or M. Sullivan to run and capture for scripting a significant WRF simulation, (3) share public NCAR WRF Docker Hub link and dockerfile, (4) prototype the AWS Spot workflow with an FRB sample or stub Python code running an embarrassingly parallel workflow, (5) develop a Docker image for the Linear Algebra Benchmark (OSU MPI micro benchmarks) and an optional Linpack, (6) pursue Azure EA resource labeling for RightScale billing ingest.

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

## 2.1 Hardware Acquisition

- UB purchased two nodes, each with 2x NVIDIA V100 GPUs; they also purchased a Globus license.

## 2.2 Installation, Configuration, and Testing

- Cornell updated their test Ceph Cluster to Nautilus in anticipation of updating their production Ceph. They deployed 5 new cloud nodes, 4 standard compute nodes, and 1 node with 4x NVIDIA V100 GPUs (all 5 purchased with non-DIBBs funds). PY5 funds will be used to upgrade the Red Cloud network with more bandwidth, and to add Ceph storage assets.
- UCSB ran into snags while updating their OpenStack cloud. Issues were raised with LDAP (a known OpenStack-keystone bug (BZ#1727976) and networking container image versions; they are working with Red Hat Support to resolve. Their Globus data transfer implementation is in progress and they will work with Cornell to stand up their Globus endpoint. UCSB also spent time working with the Dartmouth team; their knowledge and experience were integral to Dartmouth's successful OpenStack build. PY4 funds are being used to add storage assets; PY5 funds will likely be used to purchase additional GPU nodes.
- UB is in the process of upgrading all of their OpenStack components to the latest RHOSP 13 packages so that they can implement support for vGPU. They completed the purchase of Globus and Cornell will work on setting up their endpoint soon. They added 2 GPU nodes (2x V100s) to their OpenStack cloud and will use these to provide vGPUs. UB had a problem adding new storage assets (OSDs) to Ceph and had to rollback. They will test their procedures and scripts in their development environment. They hope to have the assets added to their production cloud in the next few weeks. UB has not finalized plans for PY5 hardware funds yet.
- Dartmouth worked with the UCSB and Cornell teams on building their OpenStack cloud with Red Hat's Director-based installation method (TripleO). They had trouble with their node's IDRACS going bad, but with UCSB's assistance they worked through the errors and successfully deployed a test OpenStack cloud using TripleO. They will now re-deploy to include the Horizon web interface and allow Globus authentication to provide single sign-on and officially join the Aristotle Cloud Federation. We expect this to be completed next quarter.

## 2.3 Federated Identity Management

Researchers use single sign-on at any member site.

## 2.4 Cloud Status by Site

The chart below shows each site's production cloud status. Dartmouth's cloud is in test mode.

| | Cornell | Buffalo | UCSB |
|---|---|---|---|
| **Cloud URL** | https://redcloud.cac.cornell.edu | https://lakeeffect.ccr.buffalo.edu/ (access only to federation) | https://openstack.aristotle.ucsb.edu/ |
| **Status** | Production | Production | Production |
| **Software Stack** | OpenStack | OpenStack | OpenStack |
| **Hardware Vendors** | Dell | Dell, Ace | Dell, HPE, DXC |
| **DIBBs Purchased Cores** | *616 | **256 | 356 |
| **RAM/Core** | 8GB | up to 8GB | 9GB Dell, 10GB HPE |
| **Storage** | Ceph (1244TB) | Ceph (768TB) | Ceph (528TB) |
| **10gb Interconnect** | Yes | Yes | Yes |
| **Largest instance type** | 28core/240GB RAM | 24core/192GB RAM | 48core/119GB RAM |
| **Globus File Transfer** | Yes | In Progress | In Progress |
| **Globus OAuth 2.0** | Yes | Yes | Yes |
| **Total Cores (DIBBs purchased cores + existing cores) = 2424** | * 616 additional cores augmenting the existing Red Cloud (1252 total cores). | ** 256 additional cores augmenting the existing Lake Effect Cloud (600 total cores). | ***356 cores in UCSB Aristotle cloud (572 total cores, Aristotle is separate from UCSB campus cloud) |

## 2.5 Tools

- Red Hat OpenStack – Cornell, UB, and UCSB all have production OpenStack clouds. Dartmouth has a test OpenStack cloud. UB spoke with Red Hat engineering at SC19 about support for non-Director-based installations like the ones UB and Cornell are running. The engineer gave UB details for how to obtain support. Cornell contacted Red Hat with the information and Red Hat is hoping to have the support path fixed this quarter.

## 3.0 Cloud Federation Portal Report

Content updates to the project portal are ongoing (https://federatedcloud.org).

Open XDMoD continues to monitor data ingestion from all sites, as well as provide the utilization data (https://federatedcloud.org/using/federationstatus.php).

Scripts for collecting data have been automated to generate core hours used per month; this information is available on the project web page.

InCommon certificate domain verification was completed.

The portal planning table was not updated this quarter:

| Portal Framework | | | |
|---|---|---|---|
| **Phase 1** | **Phase 2** | **Phase 3** | **Phase 4** |
| **10/2015 – 3/2016** | **4/2016 – 12/2016** | **1/2017 - 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. Available through CU doc pages. | Update materials to be federation-specific and move to portal access. | Add more advanced topics as needed and after implementation in Science Use Cases, 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 – 12/2017** | **4/2017 – 12/2017** | **1/2018 - End** |
| Cross-training expertise across the Aristotle team via calls and science group visits. | Hold training for local researchers.  Offer Webinar for remote researchers.  Use recording/materials to provide asynchronous training 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 – 3/2017** | **4/2017 – End** |
| Plan how to achieve seamless login and key transfer from portal to Euca dashboard. | Login to the portal using InCommon. | Beta testing Euca 4.4 with Euca console supporting Globus Auth. Will deploy and transition to Euca 4.4 on new Ceph-based cloud. | Transition to OpenStack console with Globus Auth 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. | No longer relevant since Globus Auth will let us interface with Euca web console | N/A | N/A |
| Allocations and Accounting | | | |
| Phase 1 | Phase 2 | Phase 3 | Phase 4 |
| 10/2015 – 3/2017 | 3/2017 –5/2018 | 6/2017 – 10/2018 | 6/2017 – 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. | Display usage and CPU hours by account or project on the portal. Integration hooks for user and project creation/deletion and synchronization across sites. Note: due to OpenStack move, account creation across sites is delayed. | 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. |

### 3.1 Software Requirements & Portal Platform

No software changes were made to the portal platform this quarter.

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

UCSB is completing the process of hiring a programmer to make the DrAFTS 2.0 prototype (developed by a graduate student working on the project).  The student prototype is anticipated shortly after the first of the year and the plan is to hand it off to the programmer to make it ready for general purpose usage.

### 3.3 Application Kernels Containerization in the Cloud

The XDMoD team continues to develop application kernels (AK), i.e., short benchmarks based on the actual real-world applications and benchmarks, which are used for performance monitoring as well as benchmarking. One of the difficulties in AK use is the need to compile them on each monitored resource. The AK building can be a lengthy process due to dependencies on multiple libraries. Containerization would simplify the installation process and allow rapid deployment of AK for performance monitoring purposes.

During this quarter, we made production versions of HPCC and NAMD containers, improved the OpenStack integration, and added app kernels execution and performance analysis to the Aristotle hub federated XDMoD portal. Previously AKs were compiled for generic CPU architectures. For the majority of compute-intensive applications, an architecture-specific optimization can give a significant increase in performance.  HPCC and NAMD containers automatically detect the number of cores and the highest supported vectorized instructions set. This allows running the executable optimized specifically for that instruction set.  The binaries were generated for SSE2, AVX, AVX2 and AVX512 (separate versions for

Intel Skylake-X and Knights Landing CPU) instructions set. The performance analysis showed that in the case of HPCC there is no significant difference between the optimizations. This is due to the usage of MKL libraries, which internally select routings with proper optimization. In the case of NAMD, utilization of the AVX512 binary gives a 50% performance boost in comparison with a generic binary (SSE2). In clouds, the physical CPU architecture is often exposed to a virtual environment as a common denominator between all physical nodes, and it currently corresponds to AVX or AVX2 instruction sets on Aristotle systems.

HPCC and NAMD app kernel execution were added to Aristotle's XDMoD. It allows app kernel performance monitoring and visualization as well as a comparison between different sites (Figure 1). XDMoD also provides service for periodic performance reports and anomaly detection (Figure 2). It should be noted that we use a previously developed XDMoD-AppKernel module for HPC resources performance monitoring and only added support for execution on OpenStack with Docker container, as well as installed it on Aristotle's XDMoD.
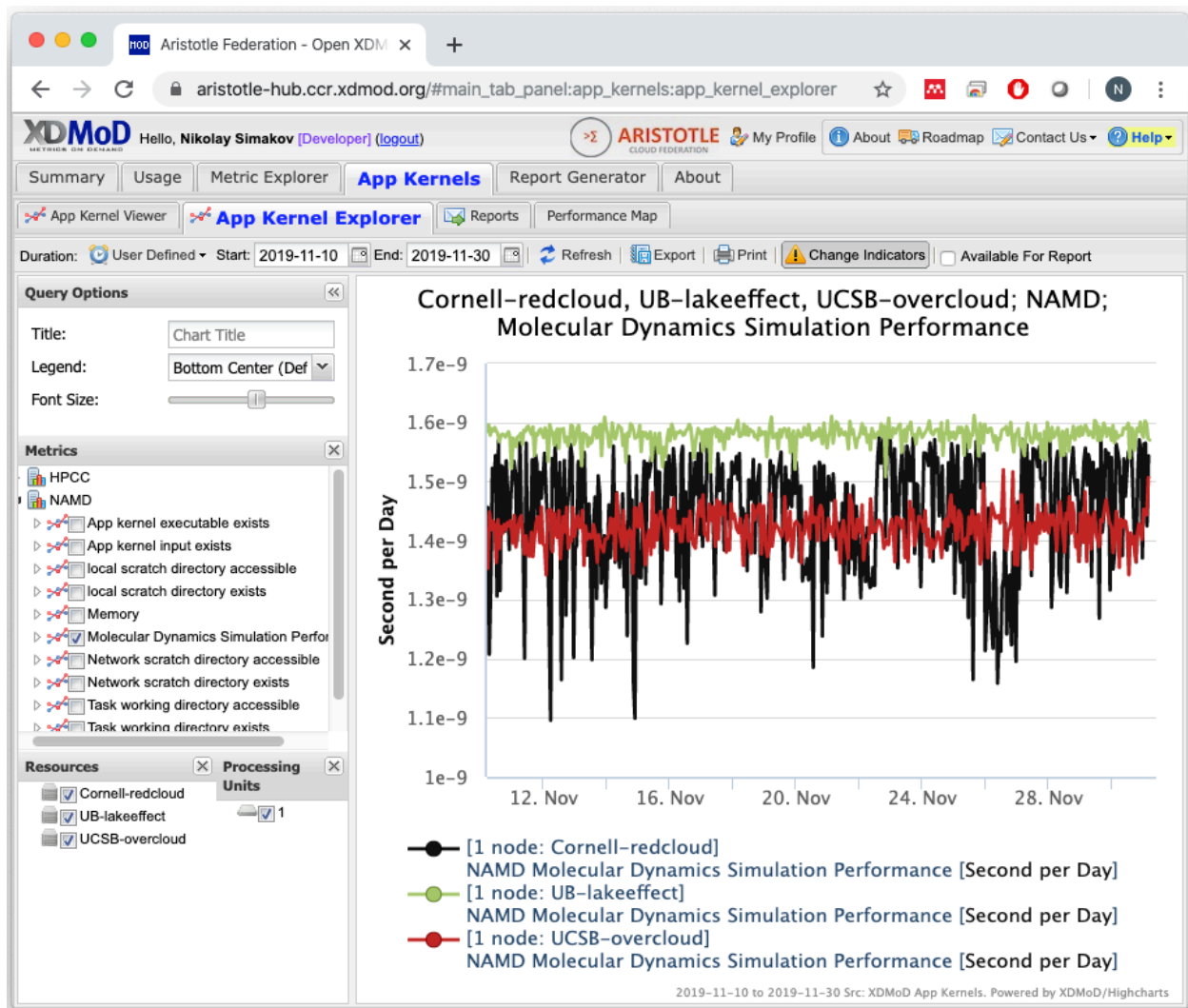


**Figure 1:** Snapshot of Aristotle's XDMoD website with visualization of NAMD app kernel execution results on all three Aristotle Cloud sites.
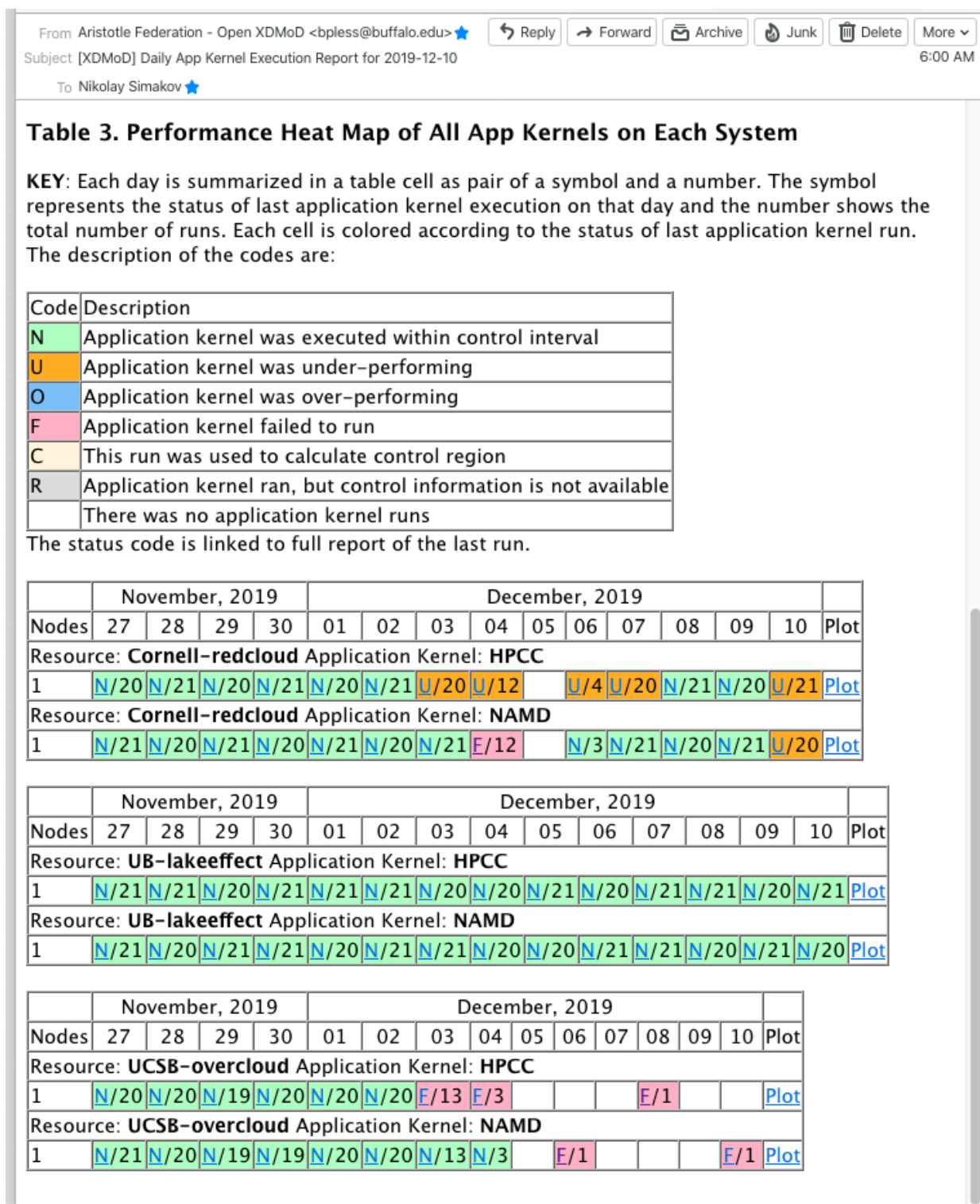
**Figure 2:** Snapshot of e-mail report from Aristotle's XDMoD on app kernels performance
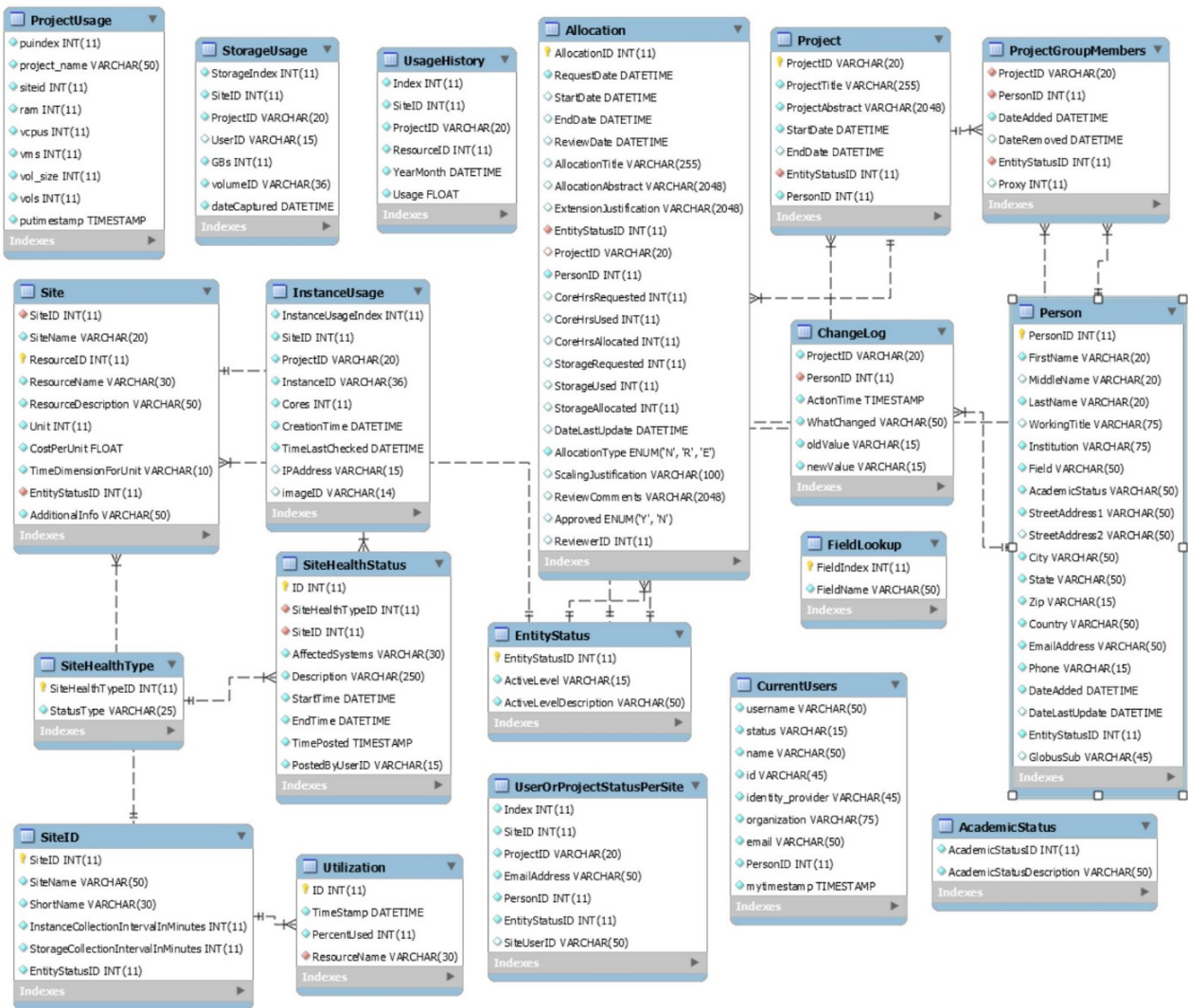
Future work includes:

- finalizing containers for HPCG and NWChem AKs
- identifying ways to reliably measure file system performance
- continuation of AK deployment on Aristotle Cloud for performance monitoring
- comparison between different sites and addition of operational metrics.

During the deployment of AK on Aristotle Cloud we noticed that time from instance start to access can exceed ten minutes. This can be undesirable for certain Cloud use cases (for example, spin-off instances only when they are needed). Thus, we want to add OpenStack operational metrics to monitor this and similar operations.

### 3.4 Allocations & Accounting

No changes were made to the database schema this quarter:

**4.0 Research Team Support**

**4.1 Science Use Case Team Updates**

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

UB Professor Varun Chandola and collaborators continued to develop a community infrastructure on Aristotle Cloud to support the Lower Great Lakes Resiliency Effort that is spearheaded by the faculty members at UB, and includes many regional partners, including academic, government and non-profit organizations. As a first step, we created an interaction platform, based on an open-source collaborative platform (Rocket.Chat) that is hosted on Aristotle Cloud. Additionally, we have a version of our webGlobe system, tailored for this collaboration, also running on Aristotle Cloud. By the start of next year, we will be heavily using this infrastructure to engage the partners into research and partnership building, all culminating into a large NSF proposal, that will be submitted in Spring 2020.

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

UB Professor Dominick Roesch, 3 PhD students, and collaborators are using the financial framework and underlying data hosted on Aristotle Cloud. They continue to investigate how price deviations (market inefficiencies) affect liquidity (the ease at which you can buy or sell). Their comparative analysis of night and day asset pricing, a collaboration with Terrence Hendershott and Dmitry Livdan of UC Berkeley, has been accepted for publication by the *Journal* of *Financial Economics*. Other papers are under development.

**Use Case 3: Application of the Weather Research and Forecasting (WRF) Model for Climate-Relevant Simulations on the Cloud**

*Précis objectives of Cornell Professor Sara C. Pryor's and Post Doctorate Associate Tristan J. Shepherd's current suite of simulations:*

1. Quantify impact of resolution (to convective permitting scales) on near-surface flow (i.e., wind speed) regime fidelity
2. Examine scales of coherence in wind fields. Specifically, spatial scales of calms (i.e., wind speeds $< 4$ m/s), and spatial scales of intense wind speeds (i.e., wind speeds $>$ the local $90^{th}$ percentile value)
3. Quantify the platform dependence of wind simulations (i.e., quantify the differences in near-surface wind regimes from simulations conducted on conventional HPC and the cloud)
4. Examine inter-annual variability in near-surface wind speeds (can we simulate it, what is the source?)
5. Evaluate impact of large wind turbine (WT) developments on downstream climate (local to mesoscale).

*Objective 4:* We are assessing the possible impact of climate variability and change in both the average (expected) wind resource, the inter-annual variability in power production, and the conditions in which wind turbines will operate. These analyses were conducted with the Weather Research and Forecasting (WRF) model at 12 km grid-spacing (resolution) over the eastern USA. We used these simulations to quantify the spatiotemporal variability in one aspect of wind turbine operating conditions (extreme wind speeds) and possible changes in the magnitude and interannual variability of expected wind power generation. We also have begun to elaborate approaches that can be applied to assess the differential

credibility of model-derived assessment of these properties at different locations using examples drawn from the eastern U.S.

*Objective 1 & 5:* We are evaluating the sensitivity of wind farm wake effects and power production from two wind farm parametrizations (the commonly used Fitch scheme and the more recently developed Explicit Wake Parameterization (EWP) to the resolution at which the model is applied. These simulations were conducted for a 9-month period for a domain encompassing much of the U.S. Midwest but centered on Iowa. The two horizontal resolutions considered at 4 by 4 km and 2 by 2 km grid cells, and the two vertical discretizations employ either 41 or 57 vertical layers (with the latter having double the number in the lowest 1km). Higher wind speeds are observed close to the wind turbine hub-height when a larger number of vertical layers are employed (12 in the lowest 200 m, versus six), which contributes to higher power production from both wind farm schemes. Differences in gross capacity factors for wind turbine power production from the two wind farm parameterizations and with resolution are most strongly manifest under stable conditions (i.e. at night). The spatial extent of wind farm wakes when defined as the area affected by velocity deficits near to wind turbine hub-heights in excess of 2% of the simulation without wind turbines are considerably larger in simulations with the Fitch scheme. This spatial extent is generally reduced by increasing the horizontal resolution and/or increasing the number of vertical levels. These results have important applications to projections of expected annual energy production from new wind turbine arrays constructed in the wind shadow from existing wind farms.

We have recently begun new simulations for a domain centered over the Southern Great Plains in order to explore whether the inferences we have drawn regarding simulation sensitivity to resolution are fully generalizable irrespective of the base climate. These simulations are particularly interesting in terms of exploring what is possible on a single VM. The inner-most domain comprises 247 by 247 grid cells with 57 vertical levels and is repeated three times; once without wind turbines operating, once with one wind farm parameterization, and then for a third time with the second wind farm parameterization. Thus, it is very computationally expensive and RAM demanding. After some initial difficulties with the node failing, the simulation is now proceeding with the following timing statistics; 120 hours of simulation are requiring approximately 79 hours of compute time. This case would thus make an exceptional candidate for a trial of simulations across multiple VMs.

*Activities planned for next quarter:*
- Our activities will focus on additional WRF simulations and analyses of WRF output generated to data in support of Objectives 1, 4, and 5.

*Journal manuscripts:*
- Pryor S.C., Shepherd T.J., Volker P., Hahmann A.N. and Barthelmie R.J. 'Wind theft' from onshore wind turbine arrays: Sensitivity to wind farm parameterization and resolution. *Journal of Applied Meteorology and Climatology* (JAMC-D-19-0235, in review).
- Shepherd T.J., Barthelmie R.J. and Pryor S.C. Sensitivity of wind turbine array downstream effects to the parameterization used in WRF. *Journal of Applied Meteorology and Climatology* (JAMC-D-19-0135, in review).
- Pryor S.C., Shepherd T.J., Bukovsky M. and Barthelmie R.J. (2019): Assessing the stability of wind resource and operating conditions. *Journal of Physics: Conference Series* (in press).
- Letson F., Shepherd T.J., Barthelmie R.J. and Pryor S.C. (2019): Modelling hail and convective storms with WRF for wind energy applications. *Journal of Physics: Conference Series* (in press).

- Letson F.W., Barthelmie R.J., and Pryor S.C. (2019): RADAR-derived precipitation climatology for wind turbine blade leading edge erosion. *Wind Energy Science* (in press, https://doi.org/10.5194/wes-2019-43

*Presentations:*
- Pryor S.C., Shepherd T.J., Bukovsky M. and Barthelmie R.J. (2019): Assessing the stability of wind resource and operating conditions. *North American Wind Energy Academy WindTech Conference*, Amherst, USA, October 2019 (*oral presentation*).
- Letson F., Shepherd T.J., Barthelmie R.J. and Pryor S.C. (2019): Modelling hail and convective storms with WRF for wind energy applications. *North American Wind Energy Academy WindTech Conference*, Amherst, USA, October 2019 (*oral presentation*).
- Shepherd T.J., Barthelmie R.J. and Pryor S.C. (2019): Assessment of wind turbine impact on future climate in GCM-driven WRF simulations. *North American Wind Energy Academy WindTech Conference*, Amherst, USA, October 2019 (*oral presentation*).
- Shepherd T.J., Barthelmie R.J. and Pryor S.C. (2019): Quantifying array-array effects using WRF model simulations: A sensitivity analysis. *North American Wind Energy Academy WindTech Conference*, Amherst, USA, October 2019 (*poster presentation*).

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

Led by Cornell Professor Jim Cordes and CAC Computational Scientist Adam Brazier, the main focus of the Radio Transient Detection Use Case this quarter has been:

- implementing the workflow for the Friends-Of-Friends (FOF) algorithm in parallel form
- retaining the ability to generalize this parallelization for other algorithms
- preparing to deploy this workflow in containerized form to a cluster of Spot instances in an automated cloud deployment.

Using the previously developed method to split up large portions of data into smaller chunks, we can run processing steps on several chunks at once using multiprocessing. We are currently developing the final features for the pipeline to automate installation in the Docker container, as well as scripting deployment of a cluster of AWS Spot Instances with Terraform and Ansible. We expect the general framework of an automated deployment of a cluster for Spot instances will be useful to the larger scientific community. We are testing on data with known FRB detections for development runs, and intend to do production runs on terabytes of new unsearched data when we deploy the cluster soon. The runs will support the RightScale investigation, as all of the runs will be profiled for performance and the costs will be measured and analyzed.

We have continued to maintain and update our Radio Astronomy containers that are publicly available on Docker Hub, especially with the recent major release of PRESTO 3.0, which is a dependency of our pipeline. The PRESTO 3.0 release includes full compatibility with Python 3, which makes our pipeline more useful to the community searching for Pulsars and Fast Radio Bursts (FRBs) because it is written in Python 3 and leverages PRESTO's new features. Further work has also been done to add more general Radio Astronomy methods, such as Pulsar detection, to the pipeline. Several standard methods for Pulsar detection are being implemented within our pipeline to be easily customized and run, similar to our implementations of FRB detection methods. For our aforementioned cloud runs, and any future runs on new data, we will have the opportunity to run searches for candidates of many interesting types of astrophysical phenomena within the same software, making the best use of cloud resources. Finally, we are

exploring the possibility of adding GPU-based (Neural Net) search and analysis methods for FRBs developed by our collaborators to the pipeline in the near future.

## Use Case 5: Water Resource Management Using OpenMORDM

The Patrick Reed Group has successfully performed OpenMORDM WaterPaths production runs at scale on Aristotle Red Cloud resources using up to 224 CPU cores on 8 worker VMs organized into an on-demand cloud MPI cluster with Docker container-based software deployment. Job runtime and software performance step timing was captured and will be published in forthcoming work authored by B.C. Trindade. The general result was that cloud resources leveraging the convenience of virtualization and portable Docker container deployment can be easily scaled out to engage many physical CPU cores to reduce overall job step timing relative to small hardware resources, yet room is still left for platform and application specific optimizations that might attempt to approach the performance of optimized MPI resources such as large XSEDE clusters.

Cornell has developed Terraform and Ansible scripts for on demand creation of Red Cloud OpenStack cluster VMs with Docker installed, provisioned a cluster-shared NFS storage server, and assisted with OpenMORDM WaterPaths MPI Docker container image creation with documented running instructions. MPI debugging tools have been installed and provided in the produced Docker image and the tooling should be broadly portable to other platforms that can run Docker containers on network connected VMs, including public cloud resources where other on demand MPI cluster creation has been demonstrated.

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

The Angela Douglas Lab at Cornell focused on adjusting parameters for numerical robustness within SteadyCom and polished scripts that generated tables and figures for an upcoming manuscript. We modified nutrient uptake values to be more biologically relevant for two of the three media used in our study. Prior to these adjustments, the simulated media allowed growth rates beyond what was biologically feasible.

We split all analyses into two variations: a first group that included inorganic ions and a second group that did not. With the new parameters, we investigated the impact of nutrient depletion on *Drosophila* gut microbiota growth dynamics by simulating three growth media with variable nutrient composition: (1) a nutrient replete medium with all nutrients provided in excess, (2) a base medium comprised of the minimal combination nutrients required for growth by all microbiota, (3) a nutrient depleted medium with only glucose, glycerol, ammonia, sulfate and phosphate as primary sources of carbon, nitrogen, sulfur and phosphorus respectively.

Our simulations show switching between nutrient-replete and nutrient-depleted diets results in large shifts in the growth dynamics of gut microbiota. Acetic acid producing bacteria, in particular *Acetobacter tropicalis*, displayed the highest growth yields in nutrient-replete and base media and lactic acid producing bacteria, in particular *Lactobacillus plantarum*, displayed the highest growth yields in the nutrient depleted medium. Microbial community composition influenced the growth dynamics of *Drosophila* gut microbiota, co-culture of acetic acid and lactic acid producing bacteria enhanced *Acetobacter* growth in nutrient-replete and base media and *Lactobacilli* growth in the nutrient depleted medium.

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

*Citrus Frost Prevention (Lindcove Research and Extension Center, Exeter, CA):*
Sensor installations are now producing frost prevention data at the citrus orchard at Visalia, California. There has yet to be a freeze, but the temperature data being gathered and analyzed by Aristotle are now available to the farm operations managers. We have done an initial scrub with them and the operators are asking for a better user interface, which we are developing. There will be a meeting before the end of the year in Visalia to discuss the validity of the analysis, the new interface requirements, and the observed chill hours so far.

*Citrus Under Protective Screening (pest protection):*
The CUPS facility broke ground mid-November. The science team visited the facility to conduct range tests for the new XBEE and 950 MHz radio installation. The team also completed the first phase of a new edge-based computational infrastructure that will be required to monitor the facility. Finally, the team developed a coordination plan with the automatic irrigation vendor to ingress sensor data from the irrigation system. This system will be installed shortly after the first of the year.

*Differential Irrigation*
The science team is now in the process of analyzing a full year's worth of differential irrigation data. UCSB data scientists on the team have produced the necessary data sets and error analysis. The research is now in the hands of the agricultural scientists on the project who are based at Fresno State.

*New Publication*
Golubovic, N., Krintz, C. & Wolski, R. (2019). A scalable system for executing and scoring K-means clustering techniques and its impact on applications in agriculture. *International Journal of Big Data Intelligence.*
https://sites.cs.ucsb.edu/~ckrintz/papers/centaurus-journal18.pdf


## 5.0 Community Outreach and Education

## 5.1 Community Outreach

- Rich Knepper gave an invited presentation at *PEARC19* on "Red Cloud and Aristotle: campus clouds and federations." The corresponding paper authored by the Aristotle team of R. Knepper, S. Mehringer, A. Brazier, B. Barker, and R. Reynolds has now been published: https://www.cac.cornell.edu/about/pubs/RedCloudAndAristotle.pdf. It discusses lessons learned from helping researchers leverage Red Cloud and Aristotle, leverage other research cloud infrastructure, and transition to public cloud.
- NSF REU student experiences on the Aristotle project were featured in an October 1 news release: https://www.cac.cornell.edu/about/news/191001.aspx
- Aristotle infrastructure team lead Resa Reynolds gave a presentation on "Federated Clouds: the Aristotle Project" at the 2019 DellXL Meeting that was held in Scottsdale, AZ: https://federatedcloud.org/papers/Federated%20Clouds%20-%20the%20Aristotle%20Project%20-%20Reynolds%20DellXL.pdf
- Cornell featured the Aristotle project at their SC19 Conference exhibit and "SC19 Cornell News Highlights" featured several Aristotle stories, including "OpenStack Cloud

Implementation Toolkit Under Development by XSEDE and Aristotle" and "Dartmouth Joins Aristotle to Explore the Federated Cloud Computing Model."
https://federatedcloud.org/science/CornellNewsHighlightsSC19.pdf

## 5.2 Education

- Cornell hosted the 2019 Higher Education Cloud Computing Forum on November 6-8. The Forum provided an opportunity for CIOs and dev ops evangelists to focus on strategy and direction rather than deep dives into technical solutions: https://blogs.cornell.edu/cloudforum/home/agenda/.
- Cornell also hosted AWS training workshops on November 12-13 focusing on containers and AI/ML services, and delivered a November 13 seminar on how to use Globus for data transfers.
- The Aristotle team is currently developing three new seminars to be offered in Spring 2020:
    1. Moving Researchers to the Cloud
    2. Containerization Techniques for Cloud & HPC
    3. GPUs in the Cloud.