

## **MAGIC Meeting April 4, 2012**

### **Attendees**

Gabrielle Allen	NSF
Shane Canon	Berkeley Labs
Rich Carlson	DOE/SC
Chip Elliot	BBN/NSF
Geoffrey Fox	Indiana University
Bob Grossman	Univ. of Chicago
Shantenu Jha	Rutgers Univ.
Chen Kan	
Dan Katz	NSF
Miron Livny	Univ. of Wisconsin
Bryan Lyles	NSF
David Martin	
Grant Miller	NCO
Mike Nelson	Georgetown Univ.
Manish Parashar	Rutgers Univ.
Barry Schneider	NSF
Alan Sill	TTU
Michaela Taufer,	Univ. of Delaware
John Towns	NCSA
Von Welch	Indiana Univ.

### **Action Items**

### **Proceedings**

This MAGIC meeting was chaired by Gabrielle Allen of the NSF and Rich Carlson of DOE/SC. The meeting focused on Challenges and Opportunities of Integrating Clouds with CyberInfrastructure. This review was organized by Manish Parashar.

### **Overview: Manish Parashar**

Clouds join traditional CyberInfrastructure as viable platforms for scientific exploration and discovery. Primary questions for this use of clouds include: What formulations, applications and standards are appropriate? How can workflows effectively utilize the cloud infrastructure?

Possible usage models for cloud infrastructure include:

- Managing complexity by outsourcing
- Democratization
- Applications are driven by the science, not by available resources through using cloud abstractions for science

There are many challenges including identifying application types and capabilities that can be supported by clouds. Are there new science applications that might be enabled through the use

of cloud resources? Running applications on community clouds can facilitate quick start-up and provide good performance but they are not suitable for many uses. Costs add up quickly for use of cloud resources.

NERSC carried out research on a cloud testbed and determined that:

- The basic cloud is 8 x as expensive as dedicated resources.
- New application formulations can be supported that are asynchronous, resilient, and elastic
- New delivery modes are enabled such as client plus cloud accelerators
  - o HPC in the cloud
  - o HPC plus the clouds where clouds complement HPC and Grid resources
  - o HPC as the cloud: Expose HPC/Grid resources using elastic on-demand cloud abstractions

### **XSEDE and Cloud Technologies: John Towns**

The XSEDE Program is developing its strategy for cloud technologies. They are in the process of defining approaches and possibilities. Currently, they need to identify what they need to know to incorporate complimentary cloud capabilities to current XSEDE resources. They will develop pilot capabilities based on their findings. Initially use cases will be used to identify cloud capabilities that are important for that use case.

### **FutureGrid: Geoffrey Fox**

FutureGrid currently supports a wide range of cloud applications. It has 4500 cores at 5 sites. It supports distributed, flexible software and IaaS, PaaS, and HPC. Activities on FutureGrid include Nimbus, Eucalyptus, HPC, Hadoop, MapReduce, XSEDE, Twister, OpenStack and others. There are about 190 current projects using clouds. The clouds enable:

- Multiple users (long tail of science)
- Internet of things (sensornets)
- Iterative applications such as MapReduce and including most data analysis
- Exploiting elasticity and platforms

It has been observed that good strategies include building an application as a service, building on existing cloud deployment such as Hadoop, use PaaS if possible, design for failure, and addressing the fact that data is moving. Clouds are not suitable for all applications. They leverage commercial software investments. Academic cloud software needs investment in core capabilities plus the “Platform”.

FutureGrid currently is compatible with Google, Amazon, and other commercial cloud resources. They have experience with Microsoft’s Azure. Current example use cases on FutureGrid include HubZero, XSEDE. FutureGrid supports collaboration on team Wikis. For domain-specific computing environments it is their experience that it is better to use the computing environment through direct interaction rather than using a portal, for example see [hadoop.apache.org/mapreduce](http://hadoop.apache.org/mapreduce).

For access to burst resources, we need to implement the ability for many users to submit actual and potential use cases. An information gathering activity is needed to identify what works and what does not work.

## **XSEDE**

XSEDE is investigating the possibility of two categories of pilot projects:

- Use cases already using clouds
- Use cases with potential for important capabilities using clouds

A plan needs to be developed for implementing pilot projects. The pilots need to evaluate:

- Authentication, authorization, accounting (AAA). How do we bridge XSEDE AAA with a cloud service provider?
- What are the specifics for integration with commercial sector clouds versus academic-based clouds?
- How do we calculate/implement chargeback for commercial providers?
- How do we purchase short-term cloud services and account for this in the project budget?

## **Open Science Data Cloud: Robert Grossman**

The Open Cloud Consortium is a not-for-profit consortium that includes NASA, several Federal labs, and university users. The Open Science Data Cloud is run at 3 sites with 3000-4000 nodes. This is expected to double this year. It is focused on big-data science and medium-sized collaborations that include: 3 biology, one satellite, genomics, social science, and other projects. The implemented clouds include utility, data, storage, and HPC applications.

New opportunities for cloud capabilities need:

- Simplifying the management, analysis, and sharing of data
- New opportunities for discovery by integrating data at scale within and across disciplines
- Think of the data center as an instrument to produce a new science of data

Challenges include:

- Building the infrastructure at the required scale,
- New software tools for exploring and integrating data at scale,
- Bringing together enough data to change how we make discoveries
- Security compliance
- Paying for cloud resources

Next steps needed include:

- Pilot projects to address the challenges
- Creating data hubs
- Creating large-scale services between data hubs
- Creating open-source blueprints for container-based facilities

## **Magellan Project: Shane Canon**

Clouds provide easier access to resources needed to meet a surge. They enable real-time analysis from instruments and detectors. They provide redundancy for critical services. Workloads

should be shifted to the cloud where it makes economic sense. A marketplace is needed for software supporting research applications; SaaS in a cloud might provide this capability.

Challenges for current clouds include missing software to seamlessly integrate resources; user administration of cloud resources; unpredictable or poor performance for many applications; data transfers to commercial providers; surprising costs for data storage and data transfers, particularly with commercial clouds; an evolving security model.

Immediate steps that are needed include establishing peerings with commercial providers such as Amazon, Google, and Microsoft. Reduced costs for data storage and transfers to commercial providers need to be realized. We need to identify applications to act as demonstration projects and we need to investigate operational changes in CyberInfrastructure to support more workloads. The

### **GENI: Chip Elliot**

Clouds are becoming a planetary-scale information utility (SaaS). They provide an abstraction of infrastructure, virtualization with multiple tenancy, and elasticity and dynamics. GENI is an infrastructure project providing at-scale infrastructure. It enables researchers to develop a slice across many different resources to support their research with dedicated resources that will not be interfered with by other research being concurrently carried out. This enables researchers to try out different architectures over the GENI infrastructure. GENI is building out resources at 14 campuses implementing:

- OpenFlow over usual infrastructures
- GENI Racks with OpenFlow to provide storage and computational power
- WiMax or LTE access

The regional science networks are deploying OpenFlow and GENI Racks, e.g., in I2, NLR and regional resources. Campuses and dorms are being enabled to allow students to experiment almost at scale. GENI is planning to expand from 14 campuses to 100-200 campuses. University CIOs are being engaged. They are interested in interfederating with other cloud projects such as Probe.

### **Discussion**

Discussion among the MAGIC members identified that next steps include pilot projects for clouds. The Federal agencies are particularly interested in what research needs to be done to make sure these infrastructures are useful to the science communities and users.

### **AI**

Manish Parashar will organize a follow-up session on Integrating Clouds with CyberInfrastructure for the May MAGIC meeting to discuss how do we address the research challenges in this area?

### **Upcoming Meetings of Interest**

April 10-12: Globus World Meeting  
October: E-Science conference of IEEE, Chicago

IEEE Computing in Science and Engineering Magazine – Special Issue on Cloud Computing in Science and Engineering - <http://www.computer.org/portal/web/computingnow/cscfp4>

**Future MAGIC Meetings**

May 2, 2:00-4:00, NSF, Room II-415

June 6, 2:00-4:00, NSF, Room II-415