

Position Paper: Exploring Science Gateway Use Cases for Cloud Computing

Mark Miller, Phil Papadopoulos, Amit Majumdar, Shava Smallen, Nancy Wilkins-Diehr,
Rick Wagner, Mahidhar Tatineni, Bob Sinkovits, Richard Moore, Mike Norman
San Diego Supercomputer Center, University of California San Diego

The popularity of CIPRES Science Gateway and the NeuroScience Gateway (NSG) show that there is a large demand for access to computational resources in the biology research community. The demand has been growing steadily in the past 4 years, driven by an increase in the amount of data available for analysis. For example, in the case of CIPRES, as the amount of DNA sequence is increasing exponentially, the size of the average submitted analytical job is increasing. NSG users are developing more complex neuronal models requiring HPC. No doubt as a result, the number of users of the resource is also increasing, as more users require larger computational resources than those available at their local institutions. The access provided by CIPRES has made it faster and easier for researchers to process the new wealth of DNBA data, turning it into scientific discoveries.

It is very clear that by provisioning large computational resources for public access, NSF has significantly accelerated the rate of scientific discovery in the global biology research community. While the impact of access is clear, there are reasons to ask whether the current model for Gateway compute engines can be further optimized by including cloud computing

First, CIPRES Gateway jobs are run on high end compute clusters that are designed for codes that scale to hundreds or thousands of cores. These resources are provisioned with costly advanced features such as high quality interconnect, parallel file systems, and solid state drives to push the boundary of large computational capabilities. However, the codes supported by CIPRES on these resources scale to no more than 64 cores, and a significant fraction of the jobs run by CIPRES are run on only 8 cores. For NSG there are neuronal simulations that scale from many hundreds to many thousands of cores but there are also parameter sweep type simulations that can take advantage of throughput possibly provided by cloud computing. In this context it will be interesting to explore what mix of highly scalable HPC resources and cloud resources could be utilized for science gateways for achieving the science impact needed while utilizing resources appropriately.

Second, the compute engines available through the XSEDE program are finite, time is awarded on these machines through a competitive allocations process, and the machines are always oversubscribed. With the recent retirement of some large clusters, the issue of oversubscription has become even greater. While the CIPRES Science Gateway remains in service, the impact of oversubscription and potential cuts to the allocation represent a clear threat to its continued operations. Similarly the NSG ran out of SUs within the first 10 months of a 12 month cycle of XSEDE allocation and required a supplemental allocation. Several questions are raised: What changes in policy should be made to preserve the scarce resources available to the popular Gateways like CIPRES? Should more appropriate high end computers be added to the NSF assets to manage the needs of these smaller jobs?

Third, when computational resources are provided at no cost to the user, limits must be placed on the maximum usage allowed for each individual user. Even in the best of times, it is clear that there is a population of users who have not attempted to use the CIPRES Science Gateway because their

computational needs are greater than can be accommodated by the policy limits that must be imposed by a free service. Unknown factors are: what is the size of this population? how are they meeting their needs currently? And, finally, can a fee-for service option be created to meet higher end needs that would be cost effective?

The question that arises logically from these issues is: Can a cloud provider offer a lower cost solution that is at once more fiscally scalable (users can pay for access and sustain the service beyond the limits of free service) and more computationally scalable (users can access as many computational cycles as is required).

Several steps would be taken to answer this question.

First, establish minimum requirements for provisioning for the set of phylogenetic, neuronal or other science gateway codes to be run.

Second, establish how these requirements mesh with the available provisioning options for a cloud provider.

Third, benchmark the codes on appropriate cloud resources to establish that networking and file system requirements are adequate, and do not cause issues with code runs, also, of course to ensure that the runs are achieving parallel efficiencies comparable to those achieved in a conventional cluster environment.

Fourth, create a distributable VM image that contains the requisite codes and that produces a robust work environment including a GUI that supports a domain user's workflow and job/results management tools within the cloud environment.

Potential advantages of a cloud strategy are:

- 1) Access is no longer confined to high end compute clusters, so those resources can be freed for other, highly scalable tasks
- 2) These low scalability calculations can be moved to less costly resources.
- 3) The easily portable work environment can be transferred among to user-owned or commercial cloud resources.
 - a. Assuming the user has access to modest clusters, rather than high end clusters, such provisioning would give them greater access to their own resources.
 - b. Assuming the user has sufficient financial resources the user has the option of purchasing access to cloud resources, user can access as many core hours of compute time as their budget allows.
 - c. The user could conceivably use their home environment for low end runs, and commercial resources for high end runs.

Potential risks of the cloud strategy:

- 1) Poor scalability of the code in the cloud environment
- 2) Poor file system or database performance in the cloud environment
- 3) Increased complexity: does the new user throughput really increase? Is the cost savings in computational equipment offset by millions of dollars in software development costs.
- 4) Poor portability of the platform. This would be devastating to the model.
- 5) Negative performance impact due to network performance, data management.