

# The Dartmouth Green Grid

James E. Dobson<sup>1\*</sup>, Jeffrey B. Woodward<sup>1</sup>, Susan A. Schwarz<sup>3</sup>,  
John C. Marchesini<sup>2</sup>, Hany Farid<sup>2</sup>, and Sean W. Smith<sup>2</sup>

<sup>1</sup> Department of Psychological and Brain Sciences, Dartmouth College

<sup>2</sup> Department of Computer Science, Dartmouth College

<sup>3</sup> Research Computing, Dartmouth College

**Abstract.** The Green Grid is an ambitious project to create a shared high performance computing infrastructure for science and engineering at Dartmouth College. The Green Grid was created with the support of the Dean of the Faculty of Arts & Sciences to promote collaborative computing for the entire Dartmouth community. We will share our design for building campus grids and experiences in Grid-enabling applications from several academic departments.

## 1 Introduction

A Campus Grid enables the collaboration of multiple departments, labs, and centers within an institution. These groups are often in their own administrative domains but can share some common infrastructure. Dartmouth College has built a campus grid architecture called the “Green Grid” to leverage this shared infrastructure. Campus Grids can often be constructed faster than a larger multi-institutional Grid since there are common infrastructure services such as high speed networking, naming services (DNS), and certificate authorities already established. Dartmouth was an early adopter of Public Key Infrastructure (PKI) technology which is leveraged by the Globus software stack. The Green Grid would not have been possible without a close collaborative relationship between the departments providing these services including Network Services, the Dartmouth PKI Lab, and Research Computing. There are several existing campus grid projects at Virginia Tech [1], University of Michigan [2], University of Wisconsin, and University of Buffalo. In building the Green Grid we have attempted to follow the conventions established by the Global Grid Forum [3], the National Science Foundation’s Middleware Initiative (NMI), several large-scale national Grid projects [4, 5] and the work already done by existing campus grid projects.

## 2 Campus Grid Design

The Green Grid was designed to be a simple architecture which could expand to include the many scientific Linux-based clusters distributed around the Dartmouth campus. We followed two major design principles for the Green Grid:

---

\* Contact Author: HB 6162, Hanover, NH 03755 James.E.Dobson@Dartmouth.EDU

**Table 1.** Dartmouth's Campus-wide Grid Resources (Estimated)

Department	CPU <sub>s</sub>	CPU <sub>s</sub>	CPU <sub>s</sub>
	Phase I	Phase II	Phase III
Math	12	20	100
Research Computing	12	32	128
Tuck School of Business	12	12	128
Biology	12	12	128
Psychological and Brain Sciences	12	32	128
Computer Science	12	64	512
Physics	12	50	128
ISTS	12	16	128
Chemistry	12	12	60
Dartmouth Medical School	12	500	600
<b>Total</b>	60	750	1912

**No major centralized infrastructure** The Green Grid must exist as an informal collaboration between departments. No major systems or centralized infrastructure should be needed for the Grid to be operational. The groups providing resources should be able use their systems disconnected from the Green Grid. We did not want to create new authentication systems or stray too far from the standard technologies already deployed at the campus level.

**Local Control over Local Resources** For the Green Grid to successfully integrate the clusters owned by individual departments and PI's the architecture needed to enable the local cluster administrators to keep control over their own resources. The Green Grid doesn't require specific schedulers or resource management systems to be used. System administrators need to implement a small stack of software but not replace existing resource management systems.

## 2.1 Project Phases

**Phase I** The initial project phase included the purchase of 60 dual processor systems to bootstrap a Grid computing infrastructure. This system served as a reference architecture and immediately provided 120 processors (Table 1) available for running applications.

**Phase II** The second phase of the Green Grid project extends the Grid to include the dedicated Linux clusters housed within each department. The several labs of Linux desktops will be added to this infrastructure. Departments beyond the initial 10 are coming online with their own computer resources. We are currently in this phase of the project. The application requirements for the new users are being taken into effect and the final software stack is being defined for the cluster owners and system administrators to deploy.

**Phase III** We plan to look for solutions for extending the Green Grid beyond dedicated servers and clusters. The thousands of desktops on the campus could be integrated with the Green Grid infrastructure for the running of batch applications.

In the previous project phases we assume that the execution hosts will be running a Linux distribution on an x86 or x86 compatible system. In Phase III we will have to deal with true heterogeneity.

## 2.2 Certificate Management

The Dartmouth PKI Lab, initially chartered by Internet2, has been pursuing both development and deployment of campus PKI as well as experimental research. Both aspects of this work are being integrated into the Green Grid.

On the deployment end, we've done a phased roll-out of X.509 identity certificates to the campus population, and retrofitted principal campus information services to use PKI-based authentication and authorization; for some sensitive administrative services, PKI is the only option now permitted. Initially, Dartmouth Certification Authority (CA) issued low-assurance certificates online, by confirming the identity of the remote user via the college's legacy campus-wide userid/password system. Recently, the Dartmouth CA started offering high-assurance certificates, required for higher-assurance services. In addition to the standard identity lookup, the CA requires a form of in-person authentication before it will issue a high-assurance certificate; in some versions, the user's private key is protected inside a USB dongle.

The Green Grid bootstraps on this PKI. When a new client installation is done, the user can have the software enroll him or her in the Dartmouth PKI: we obtain the user's username and password, and post these, over HTTPS, to the CA's Web-based enrollment system, and receive a low-assurance certificate. If the client already has a Dartmouth-certified keypair and the keystore permits export, the client can opt to export and use that keypair instead of getting a fresh one. Our current plan is to offer a MyProxy[8] service for users to store their Green Grid credentials.

The Green grid is also bootstrapping on the PKI Lab research work. For example, our *Secure Hardware Enhanced MyProxy (SHEMP)* [9] project takes advantage of TCPA/TCG secure hardware and the eXtensible Access Control Markup Language (XACML) [10] to harden a MyProxy credential repository. This system allows a user to specify the conditions under which the repository should use her private key, and for what purposes. We plan on piloting this within the Green Grid.

## 3 Applications

The Green Grid is an important research tool which is bringing computational scientists, students, and researchers across the campus together on a shared computer platform.

### – Bioinformatics

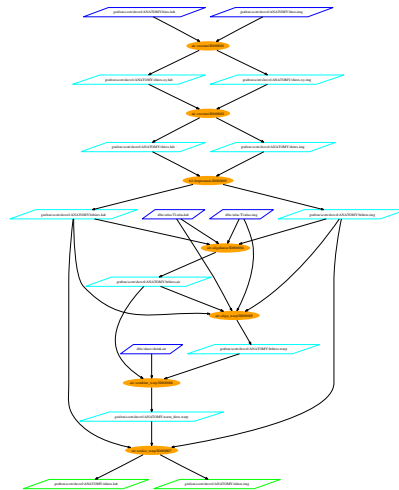
A number of research groups at Dartmouth are using gene sequencing. One researcher plans to sequence >10,000 random Damsselfly genes to construct a list of genes that may be involved in genetic networks mediating the responses. This application will access data stored in a network database. Another Bioinformatics application takes a sample dataset and randomizes a column 1,000 times to create new input datasets which are then run through the Green Grid.

– Math

One of the first applications to be run on the Green Grid was C code from the Math department which was graduate student project to search for special Proth primes with 1,500 digits. This code was compiled for the x86-64 architecture using the GNU C compiler and the GNU GMP library. This application was a single static binary which was able to read data and write from standard UNIX facilities.

– Virtual Data System

The Virtual Data System [6] (VDS) from the IVDGL and GriPhyN [7] projects is used, in part, for some of the applications running on the Green Grid. Virtual Data is being integrated into applications and methods used in the research labs within Psychological and Brain Sciences. Site selection for VDS is done using a random selection from an array of sites (departmental Globus Gatekeeper nodes) which can run the requested application. We have run a spatial normalization workflow (Fig 1) on four sites during the initial test runs.



**Fig. 1.** Example fMRI Workfbw

### 3.1 Software

The Globus Gatekeeper nodes on the Green Grid are all using the Globus Toolkit version 3.2.0. We have made a few simple modifications to both the Gatekeeper and the GridFTP server to obtain multiple AFS tokens. The Green Grid uses pre-WS Globus services such as GRAM. We are using MDS-2 for a Grid information service. Each

departmental Globus Gatekeeper node that is also reporting into a GRIS server. web-interface (mdsweb) is used to display data from each of the departmental Gatekeepers. In addition to the standard Globus services we are also using the the GSI-SSH package for remote shell access and small file transfers. Our current distribution of the grid-mapfile is through a http server. In the future we would like to use a relational database to store our authorization data using both the SAZ[12] and VOMS services.

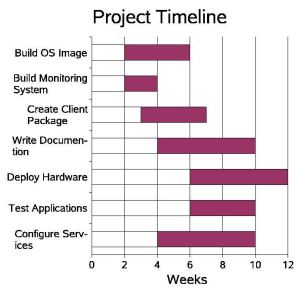
### **3.2 Standards**

We are using standards developed by the Grid community and the National Science Foundation's Middleware Initiative (NMI). These standards create a software and protocol stack that is required to be on the Green Grid. Local to each node is a temporary scratch location which is published in a catalog. Nodes that have a shared file system (such as NFS or AFS) publish the location of a cluster-wide temporary file location. Some applications are installed into the Green Grid AFS Cell (grid.dartmouth.edu). This cell is currently available on all Green Grid systems. It is expected, however, that this will not be the case in the future. Users can submit jobs through the Globus GRAM interface using simple scripts or a DRM system such as PBS or Condor-G. Using wrapper scripts on top of these standards students are able to run many existing applications at a larger scale and with site independence.

## **4 Experiences**

The concept of a campus-wide grid was conceived in May of 2004 and operational four months later (Fig 2). The Green Grid project started with the requirements from seven professors in four departments. With a demonstrated need for additional computational capacity we moved quickly to organize a proposal. We presented this proposal to the administration to seek support and the funds to purchase the hardware to bootstrap our Grid. We selected a platform which would provide participating groups compatibility with existing applications and high performance. The AMD Opteron architecture was selected for the 64-bit capabilities and price point. We solicited bids from four major system vendors and placed the order. Once the hardware began to arrive we started work on creating a single system image environment based on the Linux 2.4 kernel We constructed a website (<http://grid.dartmouth.edu>) for the distribution of Grid information, project descriptions, and procedures for using the Grid.

We have a small number of users from our academic departments utilizing the Green Grid for computational applications. These users, for the most part, have embarrassingly parallel or serial jobs to run. We have started computer science research into topics on Grid security, authorization, policy, and workflow. It is expected that the Green Grid could be used in undergraduate distributed systems classes and in laboratory science such as a fMRI data analysis offered by the Psychological and Brain Science Department. We have found some of the tools for establishing Grid credentials to not be either flexible enough or user friendly. Our current client software package includes a tool for



**Fig. 2.** Green Grid Implementation Schedule

certificate enrollment which has made this process easier. We have started to explore the use of a Grid Portal with the Open Grid Computing Environment (OGCE). We plan to have the portal available for users to manage their PKI credentials and Grid jobs (Fig 3).

**Grid Job Submission Portlet**

[Monitor submitted jobs](#)

**Job name**

**Host name**

**Port**

**Certificate to use**   
Note: If you don't see all your certificates above, ti

**Executable**

**Arguments**

**Directory**

**Standard Output File**

**Standard Error File**

**Output** Job submitted successfully.

**Fig. 3.** OCGE Grid Interface

## 5 The Intergrid

The Green Grid was designed to follow the model used on National Grids such as Grid3, the TeraGrid, and the EU-Data Grid. This design should provide for a trivial connection of Green Grid resources to larger Grid communities. A single department could, for

example, join a science domain-specific Grid. The Green Grid can, as a whole, become a resource provider for a project such as the Open Science Grid[13].

The Green Grid's role in a federated Grid project such as the Open Science Grid is that of a resource provider. Dartmouth will have multiple virtual organizations (VO's) that will want to use resources available through OSG. We are participating in discussions on site and service agreements [16] that will provide policy statements on the integration of large campus grid with OSG. In addition to policy language and agreements there are technical issues around such topics as PKI cross-certification that will need to be worked out.

Recent work in the area of connecting Grids (e.g., [14, 15]) indicates that Bridge CAs can be used for Grid authentication across organizations. Jokl et al. found that two widely used Grid toolkits could be modified to support authentication with a Bridge CA [15]. Their experiments used a testbed Bridge CA at the University of Virginia [14] with five large universities cross-certified to the Bridge CA.

## 6 Futures

The Green Grid currently has 120 processors on-line. It is shortly expected to grow to several hundred with the addition of Linux clusters in Computer Science, Research Computing, Dartmouth Medical School, and Psychological and Brain Sciences. The initial bootstrap infrastructure deployed in Phase I will be replaced as the systems are integrated in each department's local infrastructure. We expect several additional Grid services to appear on the Green Grid shortly including Replica Location Service (RLS) servers, high volume GridFTP servers, Virtual Data Catalogs (VDC), and the SEMP proxy certificate management system. The OSG is due to come online in the spring of 2005 with Green Grid resources. Our work on this will provide an example for other institutions who wish to participate in this project.

Dartmouth is in the unique position of also operating the EDUCAUSE-chartered *Higher Education Bridge CA (HEBCA)*, which is intended to connect the PKIs of Higher Education institutions. Since the Grid community is about sharing resources, and HEBCA is positioned to enable PKI trust relationships between academic institutions, it seems like a natural evolution to use HEBCA to connect Grids in higher education.

## 7 Acknowledgments

We would like to thank the following people for their hard work, consultation, and patience in getting the Dartmouth Green Grid project off the ground:

**Research Computing** John Wallace, David Jewell, Gurcharan Khanna,  
and Richard Brittain

**PKI Lab** Kevin Mitcham and Robert J. Bentrup

**Network Services** Sean S. Dunten, Jason Jeffords, and Robert Johnson

**Physics** Bill Hamblen and Brian Chaboyer

**Math** Francois G. Dorais and Sarunas Burdulis

**Computer Science** Tim Tregubov and Wayne Cripps  
**Biology** Mark McPeck  
**Tuck School of Business** Geoff Bronner and Stan D. Pyc  
**Thayer School of Engineering** Edmond Cooley  
**Dartmouth Medical School** Jason Moore, Bill White, and Nate Barney  
**Dean of the Faculty of Arts & Sciences** Michael S. Gazzaniga  
and Harini Mallikarach

We would like to also extend our thanks to Distributed Systems Lab at the University of Chicago and Argonne National Lab: Catalin L. Dumitrescu, Jens-S. Voeckler, Luiz Meyer, Yong Zhao, Mike Wilde, and Ian Foster.

James Dobson is supported in part by grants from the National Institutes of Health, NIH NS37470 and NS44393

## References

- [1] Ribbens, C.J., Kafura, D., Karnik, A., Lorch, M.: The Virginia Tech Computational Grid: A Research Agenda. Tr-02-31,, Virginia Tech (December 2002)
- [2] The University of Michigan: MGRID (2004) <http://www.mgrid.umich.edu>.
- [3] Global Grid Forum: The Global Grid Forum (2004) <http://www.ggf.org>.
- [4] Pordes, R., Gardner, R.: The Grid2003 Production Grid: Principles and Practice. In: Thirteenth IEEE International Symposium on High-Performance Distributed Computing (HPDC13). (2004)
- [5] Johnston, W.E., Brooke, J.M., Butler, R., Foster, D.: Implementing Production Grids for Science and Engineering. In Foster, I., Kesselman, C., eds.: The Grid: Blueprint for a New Computing Infrastructure. Morgan Kaufmann (2003)
- [6] Foster, I., Voeckler, J., Wilde, M., Zhao, Y.: The Virtual Data Grid: A New Model and Architecture for Data-Intensive Collaboration. In: First Biennial Conference on Innovative Data Systems Research. (2004)
- [7] Avery, P., Foster, I.: The GriPhyN Project: Towards Petascale Virtual Data Grids (2001) <http://www.griphyn.org>.
- [8] Novotny, J., Tuecke, S., Welch, V.: An online credential repository for the grid: Myproxy (2001)
- [9] Marchesini, J.C., Smith, S.W.: SEMP: Secure Hardware Enhanced MyProxy. Technical Report TR-2004-525, Computer Science Department, Dartmouth College (2004) <http://www.cs.dartmouth.edu/~carlo/research/shemp/tr2004-525.pdf>.
- [10] OASIS: XACML 1.1 Specification Set. <http://www.oasis-open.org> (2003)
- [11] von Laszewski, G., Foster, I.T., Gawor, J.: Cog kits: a bridge between commodity distributed computing and high-performance grids. In: Java Grande. (2000) 97–106
- [12] Sehki, V., Mandrichenko, I., Skow, D.: Site authorization service (SAZ). CoRR **cs.DC/0306100** (2003)
- [13] The Open Science Grid Consortium: The Open Science Grid (2004) <http://www.opensciencegrid.org>.
- [14] SURA: SURA NMI Testbed Grid PKI Bridge CA. <https://www.pki.virginia.edu/nmi-bridge/> (2004)
- [15] Jokl, J., Basney, J., Humphrey, M.: Experiences Using Bridge CAs for Grids. In: UK Workshop on Grid Security Experiences. (2004)
- [16] Open Science Grid Technical Policy Group: Open Science Grid Service Agreement Policy (2004)