**Final Report for Period:** 03/2005 - 02/2006       **Submitted on:** 06/20/2006

**Principal Investigator:** Stewart, Craig A.       **Award ID:** 0451237

**Organization:** Indiana University

**Title:**

SCI: ETF Early Operations - Indiana University

## Project Participants

**Senior Personnel**

    **Name:** Stewart, Craig

    **Worked for more than 160 Hours:** Yes

    **Contribution to Project:**

    **Name:** Voss, Brian

    **Worked for more than 160 Hours:** No

    **Contribution to Project:**

Brian Voss left Indiana University in March, 2005. Until his departure Brian worked extensively on TeraGrid related activities.

    **Name:** McRobbie, Michael

    **Worked for more than 160 Hours:** Yes

    **Contribution to Project:**

    **Name:** Shankar, Anurag

    **Worked for more than 160 Hours:** Yes

    **Contribution to Project:**

    **Name:** Simms, Stephen

    **Worked for more than 160 Hours:** Yes

    **Contribution to Project:**

    **Name:** McCaulay, David

    **Worked for more than 160 Hours:** Yes

    **Contribution to Project:**

D. Scott McCaulay is the IU TeraGrid site Lead, and has been responsible for operational project management on this project

**Post-doc**

**Graduate Student**

**Undergraduate Student**

**Technician, Programmer**

    **Name:** Lowe, Mike

    **Worked for more than 160 Hours:** Yes

    **Contribution to Project:**

    **Name:** Huffman, John

    **Worked for more than 160 Hours:** Yes

    **Contribution to Project:**

**Name:** Deximo, Christina
**Worked for more than 160 Hours:**     Yes
**Contribution to Project:**

**Name:** Grobe, Michael
**Worked for more than 160 Hours:**     Yes
**Contribution to Project:**

**Name:** Moore, Thomas
**Worked for more than 160 Hours:**     Yes
**Contribution to Project:**

**Other Participant**

**Research Experience for Undergraduates**

## Organizational Partners

**Purdue University**
IU and Purdue have collaborated throughout our involvement in the TeraGrid, particularly as regards putting resources on the TeraGrid and in developing portals

**University of California-San Diego**

**University of Texas at Austin**

**Pittsburgh Supercomputing Center**

**Oak Ridge National Laboratory**

**Argonne National Laboratory**

**UNIVERSITY OF ILLINOIS, NCSA**

## Other Collaborators or Contacts

We have worked with the Information Services and High Performance Computing Center of the Technische Universitaet Dresden as regards performance analysis of applications running on our TeraGrid-attached systems, using the software package Vampir-NG

## Activities and Findings

**Research and Education Activities:**
see attachment for desription of major research and education activities

**Findings:**
see attachment for desription of major findings from research and education activities

**Training and Development:**

see attachment for research and teaching skills and experience this project helped provide

**Outreach Activities:**

see attachment for outreach activities provided to the public to increase understanding and participation in science and technology

## Journal Publications

McCaulay, D.S., and M.R. Link, "Research data storage available to researchers throughout the US via the
TeraGrid.", Proceedings of SIGUCCS 2006, p. , vol. , (   ). Accepted

Li, H., Hart, D., Mueller, M., Markwardt, U., Repasky, R., Stewart, C.A., "GeneIndex: an open source system for exhaustive listing of words in very large genomes", Genome research, p. , vol. , (   ). Manuscript being revised, to be submitted within next 60 days

## Books or Other One-time Publications

## Web/Internet Site

**URL(s):**
www.ip-grid.org, www.teragrid.iu.edu, rac.uits.iu.edu, kb.iu.edu
**Description:**
These sites summarize IU's participation in the TeraGrid, in collaborations with Purdue University, and provide additional information and user support for HPC and Grid projects and Indiana.

## Other Specific Products

**Product Type:**

**Data or databases**

**Product Description:**

FlyBase - a database of genetic and molecular data for Drosophilla.  Includes data on all species from the family Drosophilla; the primary species represented is Drosophilla Melanogaster.

**Sharing Information:**

This database is available through to TeraGrid as a data collection.

**Product Type:**

**Data or databases**

**Product Description:**

GIS data for the state of Indiana, assembled and prepared by experts at Indiana University was added to the data resources provided to the TeraGrid by our partner site Purdue University.

**Sharing Information:**

Database is available to TeraGrid users as a resource.

## Contributions

**Contributions within Discipline:**

The principal discipline of this project is the development of large-scale cyberinfrastructure for enabling high-end computational research.  To that end, the ETF contributes by enabling researchers to address the most challenging computational problems by utilizing the integrated resources, data collections, instruments and visualization capabilities of the resource partners.

**Contributions to Other Disciplines:**

The ETF in general, and Indiana University's contributions, support research in a wide range of other disciplines, including the following:
-The development of research methods and tools that allow distributed data and computational resources to be utilized effectively
-Dynamic simulation studies in astronomy, particle physics and molecular simulations
-Real-time analysis of data for weather modeling, bioinformatics, astronomy and fusion energy simulations

**Contributions to Human Resource Development:**

IU has a strong and ongoing commitment to investing in people and to ensuring that the workforce of tomorrow represents the full richness of American society.

Indiana University coordinates and participates in IT research and education events at the regional and national levels:
-First annual TeraGrid conference at IU's Indianapolis campus
-Annual I-Light Symposium in conjunction with Purdue University
-The Annual SC conference, the premiere international even for supercomputing

TeraGrid membership enhances existing outreach efforts to interest and train people from traditionally underrepresented groups in the study and development of cyberinfrastructure:
-Sponsorship of undergraduate interns
-Participation in regional and national conferences, such as the Grace Hopper Celebration of Women in Computing, the Richard Tapia Celebration of Diversity in Computing, and the Indiana Women in Computing Conference

IU has also committed to community outreach in a variety of ways:
-Bringing grid computing information to the HPC community by way of portable stereoscopic visualization devices
-Sharing grid and high performance computing information with the general scientific community
-Encouraging an appreciation of the global value provided by our HPC and TeraGrid efforts in the lay public of Indiana
-Providing career encouragement in high performance computing to students from kindergarten to graduate school

**Contributions to Resources for Research and Education:**

The ETF is by definition an infrastructure project. All of the resources developed through this project contribute to the overall infrastructure available for research and education, including:
-Computing facilities
-Networking
-Data Storage and collections
-Visualization resources
-Development of applications and tools

**Contributions Beyond Science and Engineering:**

The cyberinfrastructure provided by the ETF has supported the development of applications in almost every area of public welfare. Some example uses include:
-Earthquake simulation
-Drug discovery
-Weather modeling
-Identifying brain disorders
-Modeling information processing
-Oil reservoir simulations
-Groundwater cleanup

<u>**Categories for which nothing is reported:**</u>

Any Book

# Activities and Findings

**1. Describe the major research and education activities of the project.**

In September of 2003, Indiana University, in partnership with Purdue University, received National Science Foundation (NSF) funding through award number 0338618 to join the TeraGrid, the NSF's flagship effort to build a national cyberinfrastructure. The TeraGrid brings together high performance computing resources, massive storage systems, visualization environments, instruments, data collections and people to create an integrated computational resource connected via high speed networks to support scientific discovery.

Subsequent awards 0451237 (ETF Early Operations) and 0504075 (SCI: TeraGrid Resource Partners: Indiana University) provided funding for Indiana to continue as an ongoing provider of resources to the TeraGrid beyond the initial construction phase and into the production phase of TeraGrid which began October 1, 2004.

This report summarizes the TeraGrid resource activities of Indiana University during the early production phase of the TeraGrid, through the end of the early operations grant, award number 0451237, from March 2005 through February 2006.

During the time of this award, Indiana was one of eight resource partners contributing to the TeraGrid, including: the University of Chicago/Argonne National Laboratory, the San Diego Supercomputer Center (SDSC) at UCSD, the Texas Avdanced Computing Center at UT-Austin, the National Center for Supercomputing Applications at UIUC, Indiana University, Purdue University, Oak Ridge National Laboratory and the Pittsburgh Supercomputing Center (PSC).

In addition to the computing and storage resources, data collections and visualization resources that Indiana has made available to TeraGrid users, IU has also provided direct support to researchers with TeraGrid allocations, and IU staff have been active in the internal operations and administration of the TeraGrid. IU personnel, including those funded from IU base funds, have participated in working groups dedicated to accounts management, portals, storage, architecture and security. IU has provided critical technical expertise and testing to the implementation and rollout of new technologies on the TeraGrid. This sharing of knowledge and experience among the supercomputing sites is one of the ways in which the TeraGrid becomes greater than the sum of its parts.

Indiana University has continued to partner closely with Purdue University in their TeraGrid efforts, continuing a history of research and technology partnerships between these two institutions which includes the creation of the I-light network ([www.i-light.org](www.i-light.org)) which connects the Purdue campus in West Lafayette, the Indiana University campus in Bloomington, and the joint Indiana University Purdue University campus in Indianapolis to each other and to the Abilene network and the commodity internet. It was this

advanced fiber infrastructure that created the foundation on which both Indiana and Purdue Universities proposed to join the TeraGrid as "IP-Grid" (Indiana Purdue Grid).

The partnership between Purdue and Indiana Universities as part of the TeraGrid has been extremely productive. Indiana and Purdue leveraged the existing I-light network and created a highly cost-effective, joint 20 GB/sec connection to the TeraGrid backplane, joining at Chicago. Perhaps more importantly has been the intellectual collaboration between the two universities. Technical teams from IU and Purdue meet regularly – meetings are now at least monthly, sometimes more frequently. This has permitted sharing of expertise and development activities. For example as Purdue develops the Nanohub as a TeraGrid portal for nanotechnology and IU develops the Hydra portal for bioinformatics as a TeraGrid portal, IU and Purdue have shared expertise, and shared a common voice in explaining within the TeraGrid project overall what our needs were in order to be able to add such unique resources to the TeraGrid.

Indiana University and Purdue University have also collaborated extensively on outreach and education activities. This ranges from tutorials presented at several conferences, to the education program at SC05, to symposia designed to attract new researchers to the TeraGrid. IU and Purdue periodically put on a joint conference about network-based and grid computing activities called the I-Light Symposium. The 2005 I-light symposium was focused on grid computing in general, and the TeraGrid in particular. The 2005 symposium also marked the beginning of a transformation in the conference itself. It is no longer a "State of Indiana" event. It is becoming very much a regional conference, drawing attendees from several neighboring states.

The overall collaboration between Indiana and Purdue has indeed proven to be of significant and lasting value to both institutions, the overall TeraGrid project, and most importantly, the national science community.


## 1a) TeraGrid Computing Resources Provided by Indiana University

At the end of the construction phase of the TeraGrid. Indiana had one computing resource available to TeraGrid users, the AVIDD IA-64 Linux cluster which consists of 32 1.3 GHz Itanium 2 processors.  This system provides approximately 275,000 TeraGrid Service Units annually.  Like the other sites added to the TeraGrid as part of the TeraGrid Expansion Program (TEP), Indiana's initial contribution of computing resources was modest in capacity compared to the systems provided to the TeraGrid by NCSA, PSC and SDSC.  By the end of the period covered by this grant, IU was well into the planning process for a major new computing system which was announced in March of 2006.  This new IBM e1350 cluster has more processing power than any system currently connected to the TeraGrid.  The new "Big Red" system will be available to TeraGrid users for the September 2006 allocations, and with a peak capacity of 20.48 Teraflops, it will allow IU to exceed their TeraGrid SU commitment for years to come.
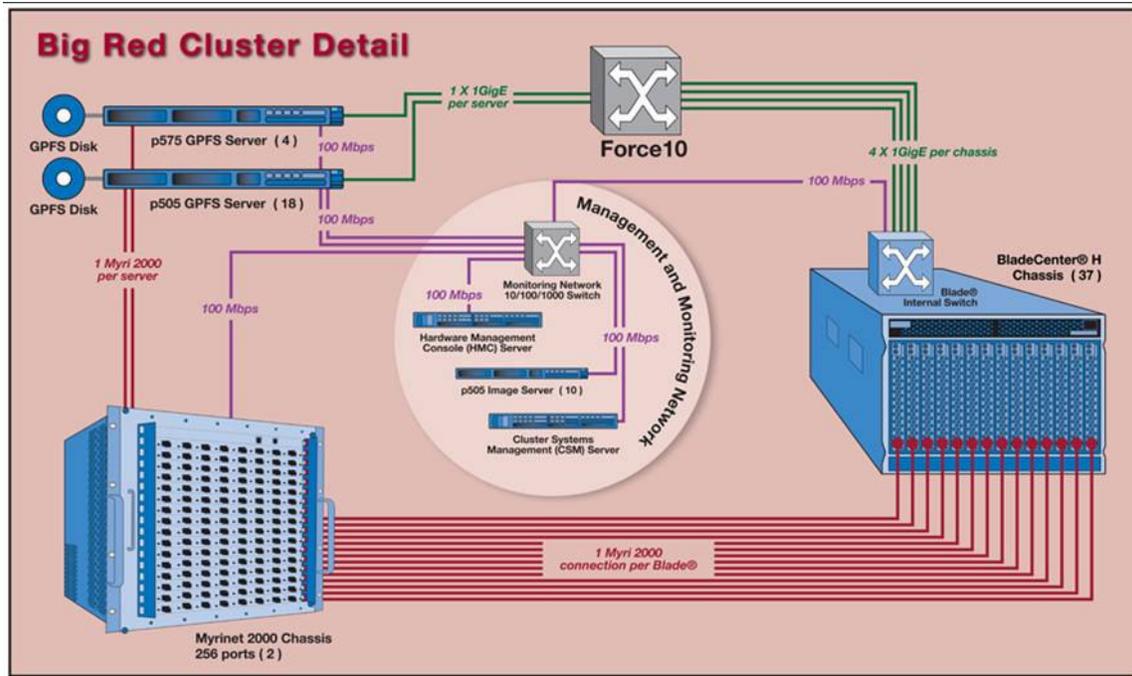
**Figure 1: IU's New 20.48 TF Supercomputer**

The new Big Red system will be the fifth computing system that IU will make available to the TeraGrid. Following the initial AVIDD-IA-64 cluster, IU made a portion of the 32-bit AVIDD-I cluster available, roughly doubling the service units provided.

The Hydra portal, which had already been available to local users at IU, was made available to TeraGrid users in the Fall of 2005. This system runs on a Condor pool that makes use of unused cycles on several thousand desktop machines running Microsoft Windows. While this system makes a tremendous amount of computing power available, it is currently limited in how it can make use of its resources. The Hydra portal provides researchers in BioInformatics access to a suite of applications including Blast, MEME and FastDNAml. While these are popular applications, the fact that they are readily available on other systems with better performance has limited the widespread use of the Hydra portal. The Windows condor pool back-end is successful as a high-throughput system, but is not competitive in processing speed with other HPC systems available for running BioInformatics applications. IU, along with Purdue, is looking at other alternatives to make use of these unused cycles on Windows machines. One solution under development at IU is the use of this Windows condor pool for distributed rendering for visualization.
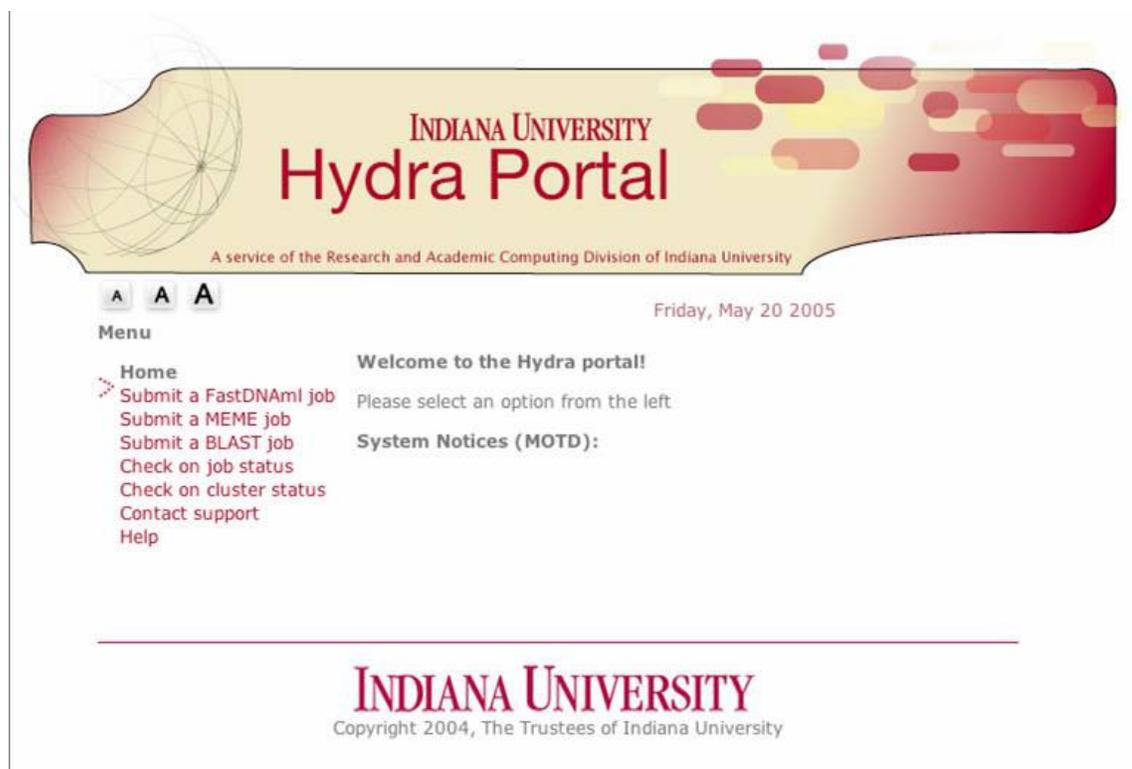
**Figure 2: The Indiana University Hydra Portal**

The fourth computing system made available by IU to TeraGrid users is the MD-GRAPE system, which has been set up for TeraGrid job submission through the AVIDD cluster. The MD-GRAPE system is a unique computing resource which is tuned for a specific type of computing problem, molecular dynamics. These MD-GRAPE cards offer a tremendous amount of computing power relative to their cost, power consumption and space requirements. Our benchmarks indicate that 1 MD-GRAPE card performs as well as 36 of the Itanium 2 processors in our AVIDD-IA64 cluster. So given the right problem, our humble 4 CPU MD-GRAPE system would be capable of producing over 1.2 million TeraGrid Service Units annually. To date, although our MD-GRAPE system is available to TeraGrid users nationally, we have done little to promote it and usage has been limited to local IU researchers already familiar with the system. Being such a specialized entity, more training and support is required to use this system effectively than would be the case for a more general purpose supercomputer. Our plan is to continue and to expand our use of MD-GRAPE systems. We are interested in acquiring the next generation MD-GRAPE3 systems when they become widely available. We are also working on making packages such as CHARMM and Amber available on our existing system. Our conversations with researchers expressing interest in these systems indicate that we would be able to generate usage on these systems if these packages were supported.

**1b) Data Storage Resources and Data Collections Provided by Indiana University**

One of Indiana's original goals as a TeraGrid resource provider was to make long term storage available to researchers using IU's HPSS-based Massive Data Storage System (MDSS).  Currently IU provides 1 TB of storage upon request to researchers with TeraGrid allocations.  Currently access is available through a custom command line interface, but later in 2006 use of this storage will be more convenient when HPSS supports access via GridFTP.  The TeraGrid has begun to list storage resources in the POPS system to advertise their availability to researchers, and has begun on an experimental basis to accept requests for large amounts of storage as part of the allocation process, with an exchange rate of 10,000 Service Units = 1 TB of annual storage.  Indiana is interested in making its storage resources available in this fashion.

By the end of the early operations grant period, Indiana University had announced and begun work on their Data Capacitor project.  The Data Capacitor is a high speed and high volume storage system to provide for short term storage needs for applications that need to deal with bursts of high volumes of data, for example instruments that produce a large amount of data that must be stored before it can be summarized or analyzed, or applications that produce large term intermediate data sets necessary for calculations.  The Data Capacitor is designed to meet the needs of these applications for fast short term storage and retrieval of massive amounts of data, to allow for types of processing that would otherwise have been impossible.  The Data Capacitor is planned as a TeraGrid resource, and as part of the development and testing of this system, we are working with other TeraGrid sites on tuning the performance of moving large amounts of data efficiently between our sites.

Part of IU's commitment is to provide data sources to TeraGrid users.  Some are already available, such as the FlyBase fruit fly genome database.  IU's Centralized Life Sciences Data (CLSD) is scheduled to be made available in the 2$^{nd}$ half of calendar 2006.  The CLSD service provides a single, SQL-based interface for querying a variety of public Life Sciences Data, including BLASTable sequence databanks and non-relational datasets that have been transformed into relational tables.


**1c) Other Resources and Contributions Provided by Indiana University**

IU's visualization resources are available to TeraGrid users who travel to our Indianapolis or Bloomington campuses.  These resources include a high-resolution display wall and configurable virtual reality theater in Indianapolis and an immersive CAVE in Bloomington.  Plans are underway to make these resources remotely accessible in 2006, including providing tools for batch rendering.

IU's Knowledge Base (KB) has been chosen by the GIG as the framework for an online help and user support environment for TeraGrid users.  The TeraGrid KB will be supported and maintained by IU resources.
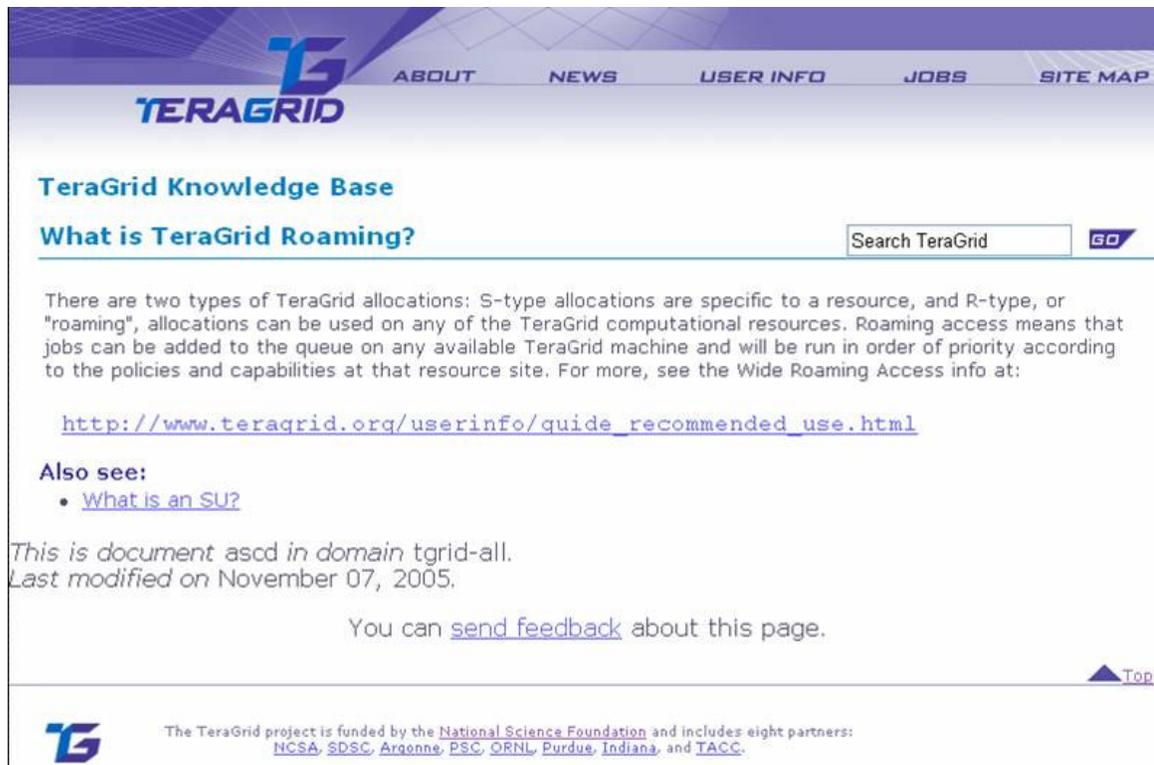
**Figure 3: The TeraGrid Knowledge Base**

As a resource provider, Indiana's original goals were to focus both on the life sciences and on usability. We are achieving those goals through the delivery of specialized tools such as the Hydra portal and CLSD, through providing user support in the form of the KB, and by supporting the full life cycle of analysis through our focus on storage, data collections and visualization in addition to computing resources.

2. **Describe the major findings resulting from these activities.**

Indiana University, along with the other TeraGrid sites, has contributed to the success of the TeraGrid Expansion Project (TEP). In less than two years, Indiana has progressed from being a new partner in the project to being a major contributor in a number of areas, from computational resources to grid portal development, and has provided useful contributions to the overall governance of the TeraGrid.

Indiana's success in the TeraGrid can be seen as evidence that the TeraGrid can function as something greater than the sum of its parts. The acceleration of growth experienced at Indiana in technical expertise, support for research, and utilization of our resources would not have been achieved without the daily interaction and cooperation we experience with our colleagues at the other centers in the TeraGrid. Likewise the growth that the TeraGrid continues to enjoy could not be achieved without the contributions of

new partner sites and the significant new resources and projects they bring.

**2a) Research Supported by IU TeraGrid Resources (selections)**

   Linked Environments for Atmospheric Discovery (LEAD) brings advances in cyberinfrastructure tools and techniques to the meteorology community with the goal of enabling more accurate and timely forecasts through on-demand execution of forecast models.  The LEAD portal, hosted at Indiana University, allows scientists on the TeraGrid to easily access atmospheric data. Workflow and myLEAD cooperate to provide automatic organization of the user's space, provenance collection, and performance information recording. These tools together promise significant enhancements in the quality of the experience for experts and novice alike in working with the products of computational mesoscale meteorology research.
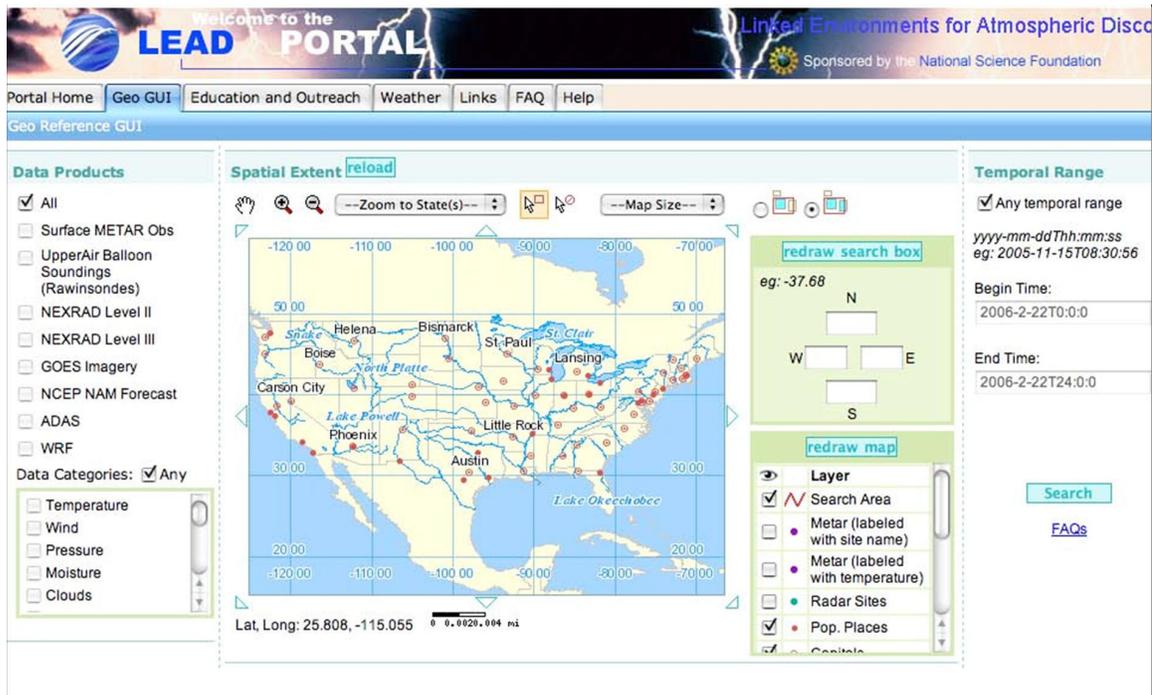


**Figure 4: The LEAD Portal**

   Two key tools are at the core of the cyberinfrastructure of Linked Environments for Atmospheric Discovery (LEAD): the workflow design and execution tool, and the user's personal workspace (i.e., "myLEAD"). The workflow tool enables experts and non-experts alike to run complex data $\Rightarrow$ model $\Rightarrow$ analysis $\Rightarrow$ visualization workflows using any of the resources available to the community. The user workspace provides a private space for a user to store model results and a host of other information related to the user's

investigations. Access to both private and community data and experiment products are accessible through a visual query interface.

The MD-GRAPE is a specialized VLSI system useful for any sort of dynamical simulation studies. Indiana University has successfully enabled use of GRAPE systems in astronomy, particle physics, and molecular simulations. IU physicist Chuck Horowitz and colleagues (including one of the UITS staff members involved in the TeraGrid) published a paper regarding the surface structure of neutron stars based on simulations performed using IU's GRAPE-6 system. These calculations would have taken impractically long without the special capabilities of the GRAPE-6 systems.

Phylogenetic estimation – the process of inferring evolutionary histories by comparing DNA sequence information – is tremendously valuable in enabling scientists to uncover how evolution brought life on earth to its current state. Phylogenetic estimation is well suited to grid computing, because the tasks involved can be run well in parallel and there is a high computation to communication ratio. Indiana University has worked for several years to create parallel and grid-enabled versions of one of the most popular software programs for phylogenetic estimation, fastDNAml.

## 3. Describe the opportunities for training and development provided by your project.

In 2006 Indiana University hosted the first TeraGrid Conference, attended by over 450 people representing 98 institutions in 24 states.  The conference featured a full day of tutorials, from entry level training to help new users get started using the TeraGrid in either a command line environment or through a science gateway, through more complex hands-on presentations on web services and visualization.  The tutorials were attended by over 200 people.

In partnership with Purdue University, IU hosts the annual I-Light Symposium (named after the high-speed network connecting the major research campuses in Indiana) to share research and educational applications of advanced networking, computing and visualization with the broader academic audience across the state.  The 2005 I-Light Symposium had a theme of grid computing, and featured an afternoon of TeraGrid-related presentations. (http://www.iupui.edu/~ilight/symposium05/)

Indiana University has also had a presence at the annual SC conference, the premiere international event for supercomputing, since 1997. In 2003, IU reported on its AVIDD facility, and the formation of the IP-grid. In 2004, IU delivered presentations on the Visualization tools developed for the TeraGrid, the Massive Data Storage System, computer simulations using the MD-GRAPE 2, to name a few. In conjunction with Purdue, IU provided infrastructure support, interactive content, in-depth tutorials, and support staff for the Education Program at Supercomputing 2005. One specific activity

provided an opportunity for Randy Heiland, Associate Director of the Scientific Data Analysis (SDA) Lab, one of the Pervasive Technology Labs at IU, to contribute to the SC Education Program. In a show of collaboration, Purdue University researchers, affiliated with the Education Program, invited Heiland to prepare and deliver a Flash-animated, voiced-over Powerpoint presentation. The presentation, "Introduction to Distributed Computing", along with accompanying open source software used by the SDA Lab, was made available on CDs to all teachers who participated in the SC Education Program. This and other online modules were also made freely available through the Sakai project (collab.sakaiproject.org).  The SDA Lab actively promotes K-12 science education and outreach. (http://sda.iu.edu/K-12)

**4. Describe outreach activities your project has undertaken.**

   IU has strong, ongoing commitments to investing in people and to ensuring that the workforce of tomorrow represents the full richness of American society. Through the IP-grid partnership with Purdue and by joining the TeraGrid, we are enhancing existing outreach efforts to interest and train in cyberinfrastructure people from traditionally underrepresented groups. Results of the research enabled by this proposal will be rapidly disseminated. And our inclusion in the TeraGrid will certainly accelerate nationally funded research underway at our university and magnify the value of NSF efforts in our institution and its two major research campuses.

   IU received an REU supplement to the early operations grant, which went to provide work experience and educational opportunities to two undergraduates.

   One IU undergraduate internship was awarded this year to Rishi Verma who worked on a portal to data, including image data and videos, on development of embryos.  Mr. Verma has developed a prototype OGCE-compliant portal called "EmbryoGrid" designed to meet the needs of embryologists. Our key contact within the embryology community is Dr. Charles Little, of Kansas State University. Dr. Richard Repasky is the local supervisor of Mr. Verma. We are continuing development of EmbryoGrid as a TeraGrid resource, and hope to have it in general use by 2006.

   Another undergraduate intern, Matthew Burks, performed laboratory work in inorganic chemistry, accompanied the IU team to the SC2005 conference, and attended several talks on supercomputing applications in chemistry - thus helping interest and prepare him for use of high performance computing applications in graduate school.

   IU has committed to significant outreach to many communities in a variety of ways: bringing grid computing information to the HPC community, sharing grid and high performance computing information with the general scientific community, encouraging an appreciation of the global value provided by our HPC and TeraGrid efforts in the lay public of Indiana, and providing career encouragement in high performance computing to students from kindergarten to graduate school.

   For the local and regional communities, IU has utilized the portable stereoscopic visualization devices provided as part of an NSF MRI project (AVIDD, award #

0116050) to present scientific research and educational content to over 2,500 individuals through campus outreach days, IT awareness fairs, and on-site school demonstrations.  Of particular note was the summer 2004 partnership with the Indianapolis-Marion County Public Libraries that toured eight branch libraries in different demographic regions as part of the "NASA in Your Library" program.  One portion of the distributed AVIDD cluster, and one of AVIDD's 3D visualization systems, was installed on the IU Northwest campus located in Gary, Indiana. As a result, many students from traditionally underserved groups were able to use advanced scientific computing and visualization systems.

   IU has also sponsored and actively participated in several national conferences focused on traditionally underrepresented groups, including the Grace Hopper Celebration of Women in Computing and the Richard Tapia Celebration of Diversity in Computing.   Thanks in part to this sponsorship and a very active Women in Computing organization (WIC@IU), IU had 14 representatives at Grace Hopper 2004. IU researchers also presented results at Tapia 2005. The WIC@IU group has also developed an interactive experience called "Just Be" that seeks to break common stereotypes about people in computing.  With support from UITS, this program has been presented at a number of middle schools and high schools in Indiana, Kentucky, and Missouri.