

---

# Numerical Libraries and the Grid: The GrADS Experiment

Antoine Petitet, Susan Blackford, Jack  
Dongarra, Brett Ellis, Graham Fagg,  
Kenneth Roche, and Sathish Vadhiyar

with lots of help from our collaborators  
(Rice, ANL, ISI, UCSB, UCSB, UH, UIUC)



Innovative Computing Laboratory  
UNIVERSITY OF TEXAS AT AUSTIN  
COMPUTER SCIENCE DEPARTMENT



Grid Application Development Software Project

---

## GrADS - Three Research and Technology Thrusts

- **GrADS - Grid Application Development Software**
  - NSF Next Generation Software (NGS) effort
- **Effort within the GrADS Project**
  - GrADS PIs: Berman, Chien, Cooper, Dongarra, Foster, Gannon, Johnsson, Kennedy, Kesselman, Mellor-Crummey, Reed, Torczon, Wolski
- **GrADSoft**
  - software infrastructure for programming and running on the Grid
  - Reconfigurable object programs
  - Performance contracts
  - Core Grid technologies
    - Globus, NetSolve, NWS, Autopilot, AppLeS, Portals, Cactus
- **MacroGrid**
  - Persistent multi-institution Grid testbed
- **MicroGrid**
  - Portable Grid emulator



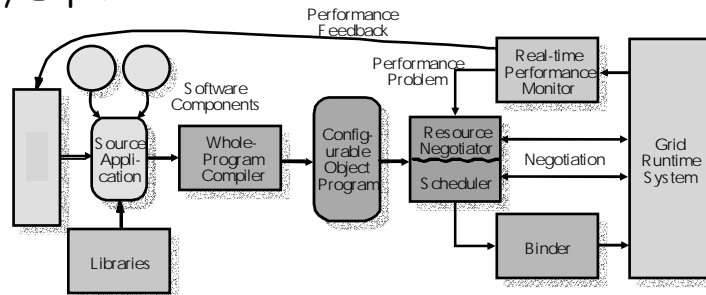
Innovative Computing Laboratory  
UNIVERSITY OF TEXAS AT AUSTIN  
COMPUTER SCIENCE DEPARTMENT



Grid Application Development Software Project

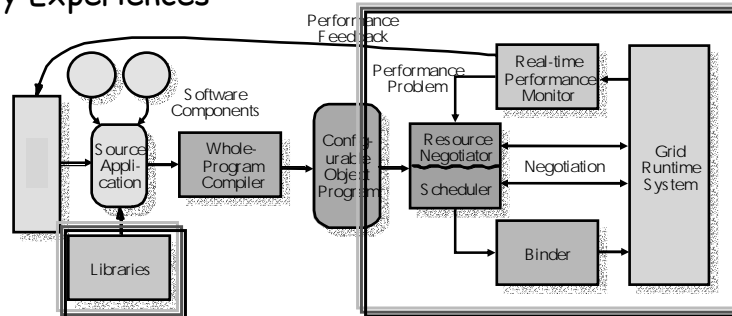
# Grid-Aware Numerical Libraries

- Using ScaLAPACK and PETSc on the Grid: Early Experiences



# Grid-Aware Numerical Libraries

- Using ScaLAPACK and PETSc on the Grid: Early Experiences



In some sense ScaLAPACK not an ideal application for the Grid.

Expanded our understand how various GrADS component fit together.

Key is managing dynamism.

# ScaLAPACK

**ScaLAPACK**  
A Software Library for Linear Algebra Computations on Distributed-Memory



- ScaLAPACK is a portable distributed memory numerical library
- Complete numerical library for dense matrix computations
- Designed for distributed parallel computing (MPP & Clusters) using MPI
- One of the first math software packages to do this
- Numerical software that will work on a heterogeneous platform
- Funding from DOE, NSF, and DARPA
- In use today by IBM, HP-Convex, Fujitsu, NEC, Sun, SGI, Cray, NAG, IMSL, ...
  - Tailor performance & provide support



Innovative Computing Laboratory  
UNIVERSITY OF MICHIGAN  
COMPUTER SCIENCE DEPARTMENT

Grids

Grid Applications Development Software Support

# ScaLAPACK Grid Enabled

- Implement a version of a ScaLAPACK library routine that runs on the Grid.
  - Make use of resources at the user's disposal
  - Provide the best time to solution
  - Proceed without the user's involvement
- Make as few changes as possible to the numerical software.
- Assumption is that the user is already "Grid enabled" and runs a program that contacts the execution environment to determine where the execution should take place.



Innovative Computing Laboratory  
UNIVERSITY OF MICHIGAN  
COMPUTER SCIENCE DEPARTMENT

Grids

Grid Applications Development Software Support

## To Use ScalAPACK a User Must:

- Download the package and auxiliary packages (like PBLAS, BLAS, BLACS, & MPI) to the machines.
- Write a SPMD program which
  - Sets up the logical 2-D process grid
  - Places the data on the logical process grid
  - Calls the numerical library routine in a SPMD fashion
  - Collects the solution after the library routine finishes
- The user must allocate the processors and decide the number of processes the application will run on
- The user must start the application
  - “mpirun -np  $N$  user\_app”
    - Note: the number of processors is fixed by the user before the run, if problem size changes dynamically ...
- Upon completion, return the processors to the pool of resources

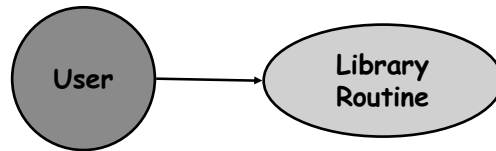


## GrADS Numerical Library

- Want to relieve the user of some of the tasks
- Make decisions on which machines to use based on the user's problem and the state of the system
  - Determine machines that can be used
  - Optimize for the best time to solution
  - Distribute the data on the processors and collections of results
  - Start the SPMD library routine on all the platforms
  - Check to see if the computation is proceeding as planned
    - If not perhaps migrate application



## GrADS Library Sequence



- Has “crafted code” to make things work correctly and together.

Assumptions:  
Autopilot Manager has been started  
and  
Globus is there.

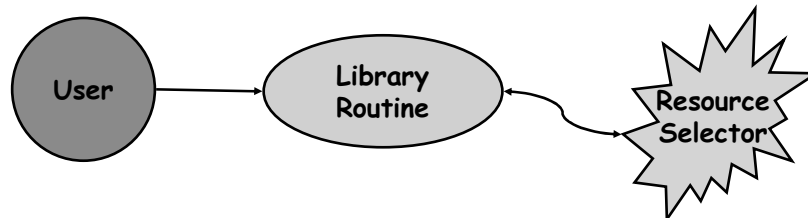


Innovative Computing Laboratory  
UNIVERSITY OF MICHIGAN  
COMPUTER SCIENCE DEPARTMENT

GrADS

Grid Application Resource Management Software Project

## Resource Selector



- Uses Globus'-MDS and Rich Wolski's NWS to build an array of values for the machines that are available for the user.
  - 2 matrices (bw,lat) 2 arrays (cpu, memory available)
  - Matrix information is clique based
- On return from RS, Crafted Code filters information to use only machines that have the necessary software and are really eligible to be used.



Innovative Computing Laboratory  
UNIVERSITY OF MICHIGAN  
COMPUTER SCIENCE DEPARTMENT

GrADS

Grid Application Resource Management Software Project

## Arrays of Values Generated by Resource Selector

- **Clique based**
  - 2 @ UT, UCSD, UIUC
    - Part of the MacroGrid
  - Full at the cluster level and the connections (clique leaders)
  - Bandwidth and Latency information looks like this.
  - Linear arrays for CPU and Memory
- **Matrix of values are filled out to generate a complete, dense, matrix of values.**
- **At this point have a workable coarse grid.**
  - Know what is available, the connections, and the power of the machines

xxxxxxx xxxxxxx xxxxxxx xxxxxxx xxxxxxx xxxxxxx	x	x	x
x	xxxxx xxxxx xxxxx xxxxx	x	x
x	x	xxxxx xxxxx xxxxx xxxxx xxxxx	x
x	x	x	xxxxxx xxxxxx xxxxxx xxxxxx xxxxxx xxxxxx



Innovative Computing Laboratory  
UNIVERSITY OF TEXAS AT AUSTIN  
COMPUTER SCIENCE DEPARTMENT

**GrADS**

Grid Applications Development Software Support

## ScaLAPACK Performance Model

$$T(n, p) = C_f t_f + C_v t_v + C_m t_m$$

$$C_f = \frac{2n^3}{3p} \text{ — Total number of floating-point operations per processor}$$

$$C_v = \left(3 + \frac{1}{4} \log_2 p\right) \frac{n^2}{\sqrt{p}} \text{ — Total number of data items communicated per processor}$$

$$C_m = n(6 + \log_2 p) \text{ — Total number of messages}$$

$t_f$  — Time per floating point operation

$t_v$  — Time per data item communicated

$t_m$  — Time per message

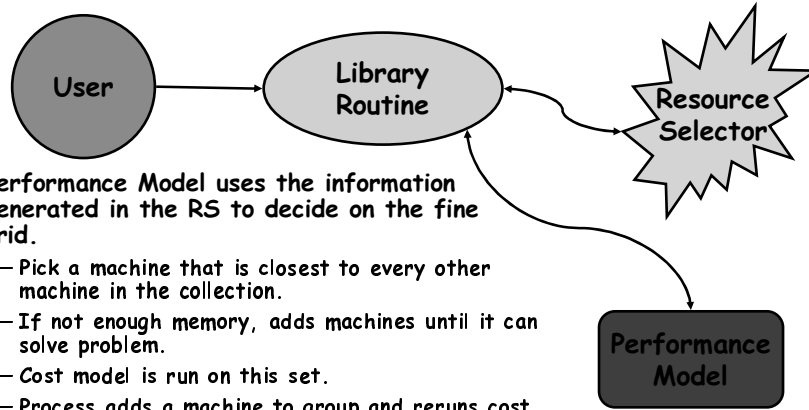


Innovative Computing Laboratory  
UNIVERSITY OF TEXAS AT AUSTIN  
COMPUTER SCIENCE DEPARTMENT

**GrADS**

Grid Applications Development Software Support

## Performance Model



- Performance Model uses the information generated in the RS to decide on the fine grid.
  - Pick a machine that is closest to every other machine in the collection.
  - If not enough memory, adds machines until it can solve problem.
  - Cost model is run on this set.
  - Process adds a machine to group and reruns cost model.
  - If “better”, iterate last step, if not stop.



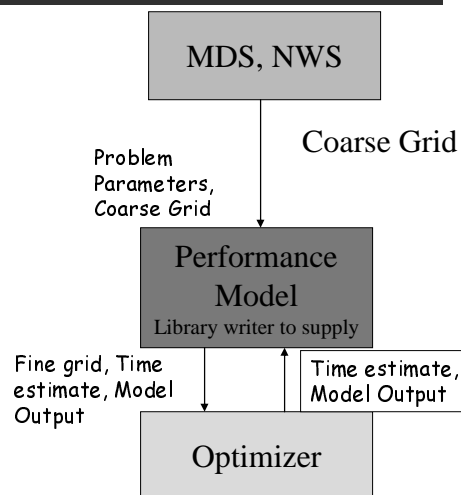
Innovative Computing Laboratory  
UNIVERSITY OF MICHIGAN  
COMPUTER SCIENCE DEPARTMENT

GrADS

Grid Application Resource Management Software Project

## Resource Selector/Performance Modeler

- Refines the course grid by determining the process set that will provide the best time to solution.
- This is based on dynamic information from the grid and the routines performance model.
- The PM does a simulation of the actual application using the information from the RS.
  - It literally runs the program without doing the computation or data movement.
- There is no backtracking in the Optimizer.
  - This is an area for enhancement and experimentation.

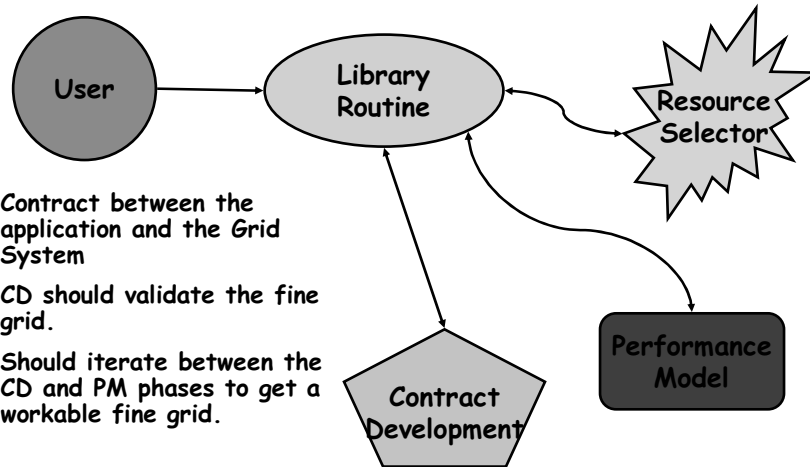


Innovative Computing Laboratory  
UNIVERSITY OF MICHIGAN  
COMPUTER SCIENCE DEPARTMENT

GrADS

Grid Application Resource Management Software Project

## Contract Development



- Contract between the application and the Grid System
- CD should validate the fine grid.
- Should iterate between the CD and PM phases to get a workable fine grid.

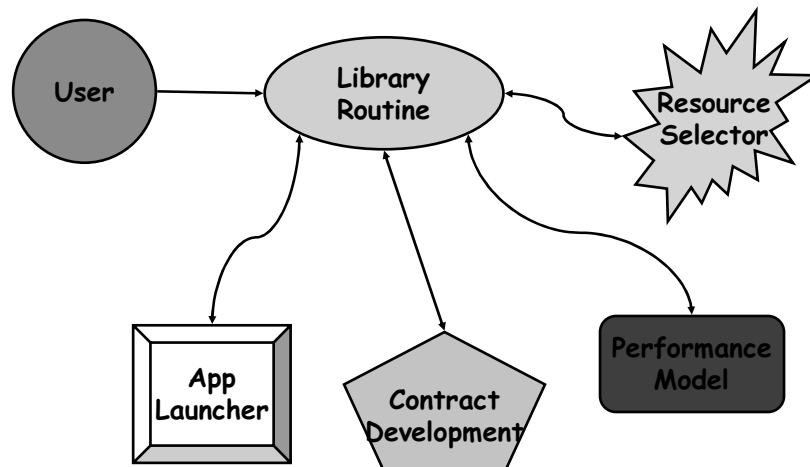


Innovative Computing Laboratory  
UNIVERSITY OF TORONTO  
COMPUTER SCIENCE DEPARTMENT

GrADS

Grid Application Development Software Support

## Application Launcher



`"mpirun -machinefile -globusrl fine_grid grid_linear_solve"`



Innovative Computing Laboratory  
UNIVERSITY OF TORONTO  
COMPUTER SCIENCE DEPARTMENT

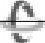
GrADS

Grid Application Development Software Support




## Experimental Hardware / Software Grid

MacroGrid Testbed	TORC	CYPHER	OPUS
Type	Cluster 8 Dual Pentium III	Cluster 16 Dual Pentium III	Cluster 8 Pentium II
OS	Red Hat Linux 2.2.15 SMP	Debian Linux 2.2.17 SMP	Red Hat Linux 2.2.16
Memory	512 MB	512 MB	128 or 256 MB
CPU speed	550 MHz	500 MHz	265 – 448 MHz
Network	Fast Ethernet (100 Mbit/s) (3Com 3C905B) and switch (BayStack 350T) with 16 ports	Gigabit Ethernet (SK-9843) and switch (Foundry FastIron II) with 24 ports	Myrinet (LANai 4.3) with 16 ports each



Advanced Computing Laboratory  
UNIVERSITY OF TORONTO  
COMPUTER SCIENCE DEPARTMENT



GrADS  
Global Applications Research Environment Software Support

- Globus version 1.1.3
- Autopilot version 2.3
- NWS version 2.0.pre2
- MPICH-G version 1.1.2
- ScaLAPACK version 1.6
- ATLAS/BLAS version 3.0.2
- BLACS version 1.1
- PAPI version 1.1.5
- GrADS' "Crafted code"

Independent components being put together and interacting

## Performance Model Validation

	Opus14	Opus13	Opus16	Opus15	Torc4	Torc6	Torc7
mem(MB)	215	214	227	215	233	479	479
speed	270	270	270	270	330	330	330
load	1	0.99	1	0.99	1	1.04	0.87

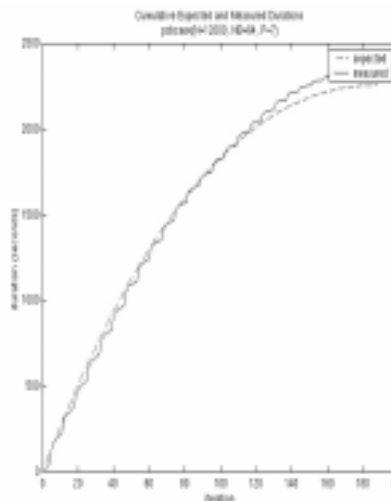
Speed = 60% of the peak

Latency	Opus14	Opus13	Opus16	Opus15	Torc4	Torc6	Torc7
Opus14	-1	0.24	0.29	0.26	83.78	83.78	83.78
Opus13	0.24	-1	0.24	0.23	83.78	83.78	83.78
Opus16	0.29	0.24	-1	0.23	83.78	83.78	83.78
Opus15	0.26	0.23	0.23	-1	83.78	83.78	83.78
Torc4	83.78	83.78	83.78	83.78	-1	0.31	0.31
Torc6	83.78	83.78	83.78	83.78	0.31	-1	0.31
Torc7	83.78	83.78	83.78	83.78	0.31	0.31	-1

Latency in msec

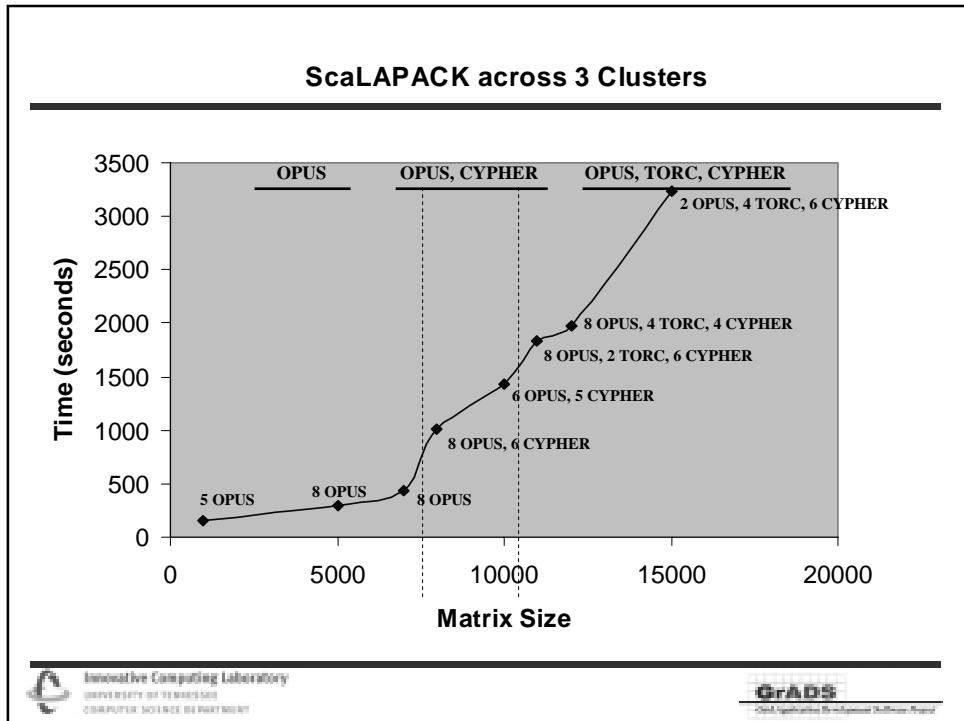
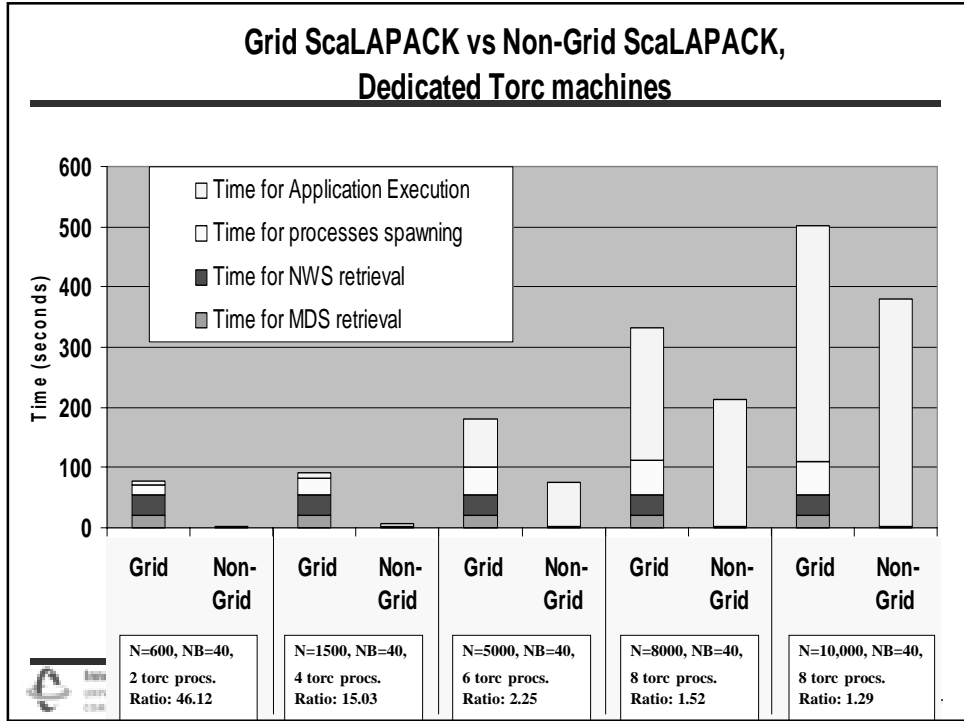
Bandwidth	Opus14	Opus13	Opus16	Opus15	Torc4	Torc6	Torc7
Opus14	-1	248.83	247.31	246.38	2.83	2.83	2.83
Opus13	248.83	-1	244.54	240.94	2.83	2.83	2.83
Opus16	247.31	244.54	-1	247.54	2.83	2.83	2.83
Opus15	246.38	240.94	247.54	-1	2.83	2.83	2.83
Torc4	2.83	2.83	2.83	2.83	-1	81.96	56.47
Torc6	2.83	2.83	2.83	2.83	81.96	-1	50.9
Torc7	2.83	2.83	2.83	2.83	56.47	50.9	-1

Bandwidth in Mb/s



This is for a refined grid

GrADS  
Global Applications Research Environment Software Support

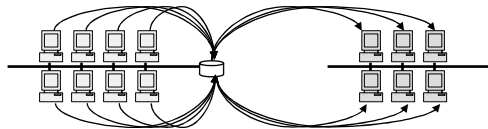


## Largest Problem Solved

- Matrix of size 30,000
  - 7.2 GB for the data
  - 32 processors to choose from UIUC and UT
    - Not all machines have 512 MBs, some little as 128 MBs
  - PM chose 17 machines in 2 clusters from UT
  - Computation took 84 minutes
    - 3.6 Gflop/s total
    - 210 Mflop/s per processor
    - ScaLAPACK on a cluster of 17 processors would get about 50% of peak
    - Processors are 500 MHz or 500 Mflop/s peak
    - For this grid computation 20% less than ScaLAPACK

## Contracts, Checkpointing, Migration

- We are using University of Illinois' Autopilot to monitor the progress of the execution.
- The application's software has the ability to perform a checkpoint and can be restarted.
  - We manually inserted the checkpointing code.
- If the application is not progressing as the contract specifies we want to take some corrective action.
  - Go back and figure out where the application can be run "optimally".
  - Restart the process from the last checkpoint, perhaps rearranging the data to fit the new set of processors.



## General Library Interface

```

graph TD
    User((User)) --> LR((Library Routine))
    LR --> RS[Resource Selection]
    LR --> AL[App Launcher]
    LR --> CD{{Contract Development}}
    LR --> PM[Performance Model]
  
```

- We have a start on a general interface for numerical libraries.
  - It's can be a "simple" operation to plug in other numerical routines/libraries.
  - Developing migration mechanisms for contract violations.
- Today a library writer needs to supply
  - Numerical Routine
  - Performance Model
- The rest of the framework can remain the same.

Innovative Computing Laboratory  
UNIVERSITY OF TENNESSEE  
COMPUTER SCIENCE DEPARTMENT
 **GrADS**  
Grid Applications Development Software Project

## Conclusions

---

- Experiments are driving GrADS development
  - Handcrafted developed leading to an automated design.
  - Exposes a number of areas for improvement
  - Very positive feed back to component developers with each experiment.
- GrADS will automate much of the decisions in the Grid environment to provide best time to solution.
  - Adaptivity to the dynamic environment.
  - As the complexities of the Grid increase need to develop strategies for self adaptability.
- Developing a basic infrastructure for computational science applications and software in the Grid environment.
  - Lack of tools is hampering development today.

Web pages:

- <http://icl.cs.utk.edu/grads/>
- <http://www.hipersoft.rice.edu/grads>

- Thanks to other GrADS researchers.

Innovative Computing Laboratory  
UNIVERSITY OF TENNESSEE  
COMPUTER SCIENCE DEPARTMENT
 **GrADS**  
Grid Applications Development Software Project