
The Virtual Grid Application Development Software (VGrADS) Project

Overview

Ken Kennedy
VGrADS Director
Rice University

<http://vgrads.rice.edu/>

The VGrADS Team

- VGrADS is an NSF-funded Information Technology Research project



THE UNIVERSITY
of NORTH CAROLINA
at CHAPEL HILL

Dan Reed



RICE

Keith Cooper
Ken Kennedy
Charles Koelbel
Richard Tapia
Linda Torczon



Jack Dongarra



Carl Kesselman



Fran Berman
Andrew Chien
Henri Casanova



Rich Wolski

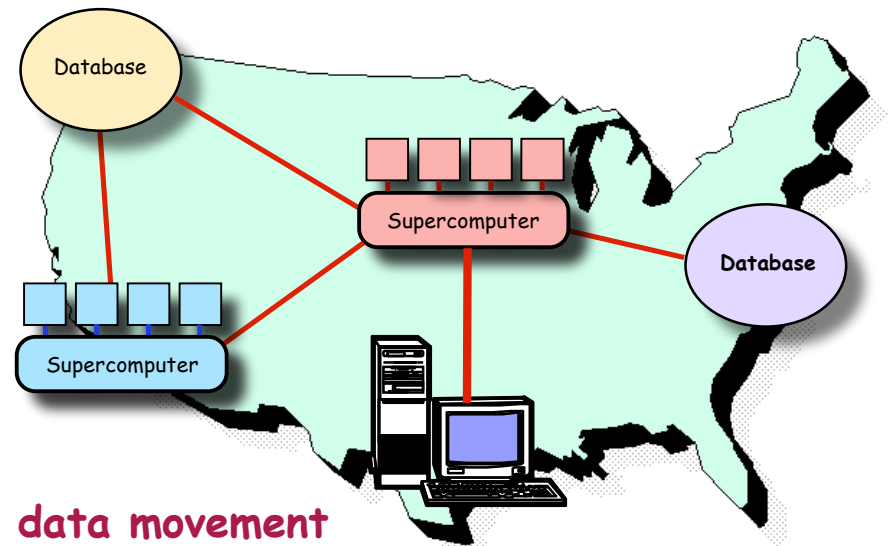


Lennart Johnson

- Plus many graduate students, postdocs, and technical staff!

Vision: Global Distributed Problem Solving

- **Where We Want To Be**
 - **Transparent Grid computing**
 - Submit job
 - Find & schedule resources
 - Execute efficiently
- **Where We Are**
 - **Low-level hand programming**
 - **Programmer must manage:**
 - Heterogeneous resources
 - Scheduling of computation and data movement
 - Fault tolerance and performance adaptation
- **What Do We Propose as A Solution?**
 - **Separate application development from resource management**
 - Through an abstraction called the **Virtual Grid**
 - **Provide tools to bridge the gap between conventional and Grid computation**
 - Scheduling, resource management, distributed launch, simple programming models, fault tolerance, grid economies

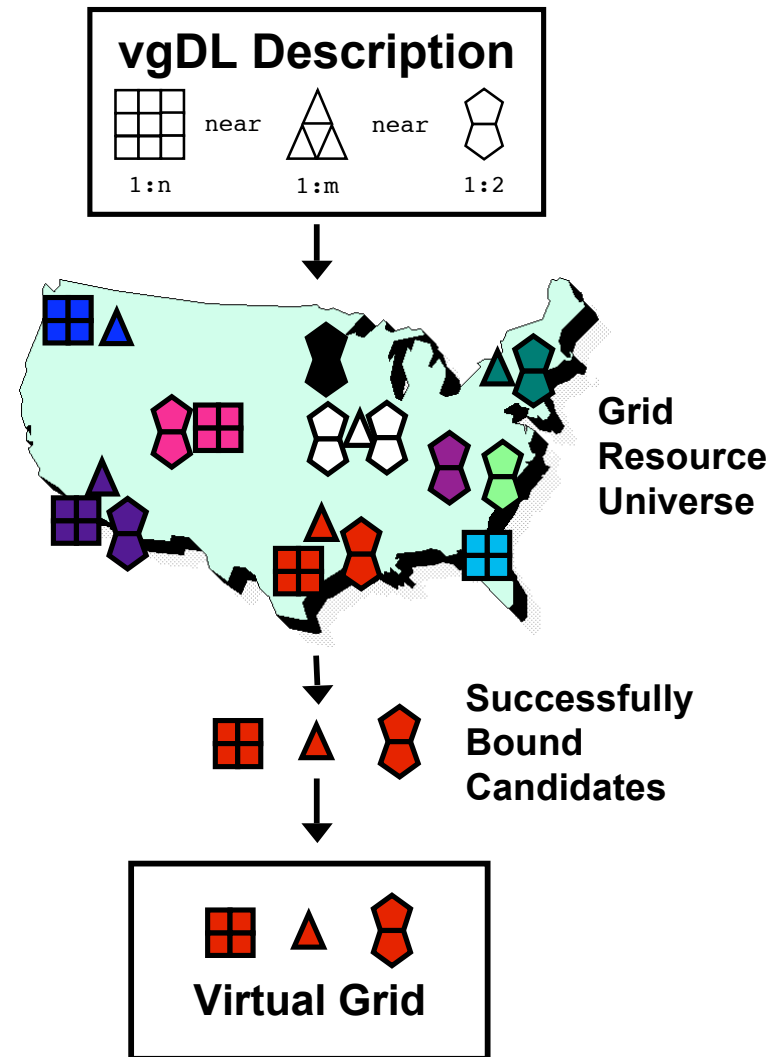


VGrADS Big Ideas

- **Virtualization of Resources**
 - Application specifies required resources in Virtual Grid Definition language (vgDL)
 - Give me a loose bag of 1000 processors, with 1 Gb memory per processor, with the fastest possible processors
 - Give me a tight bag of as many Opterons as possible
 - Virtual Grid Execution System (vgES) produces specific virtual grid matching specification
 - Avoids need for scheduling against the entire space of global resources
- **Generic In-Advance Scheduling of Application Workflows**
 - Application includes performance models for all workflow nodes
 - Performance models automatically constructed
 - Software schedules applications onto virtual Grid, minimizing total makespan
 - Including both computation and data movement times

Virtual Grids (VGs)

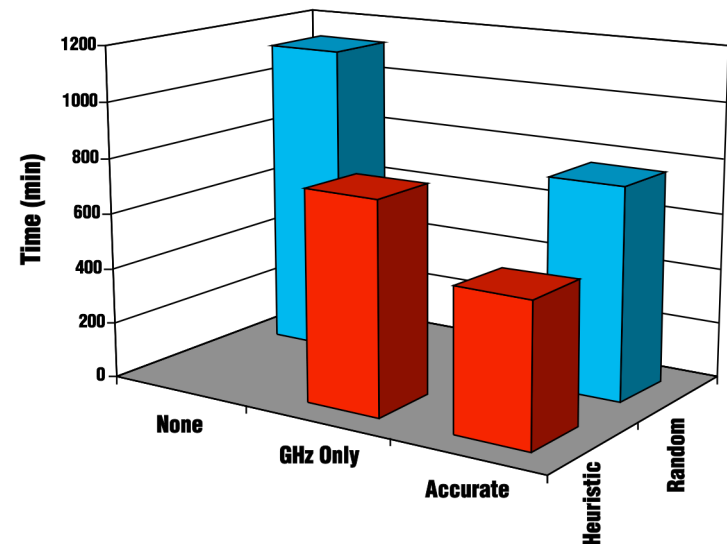
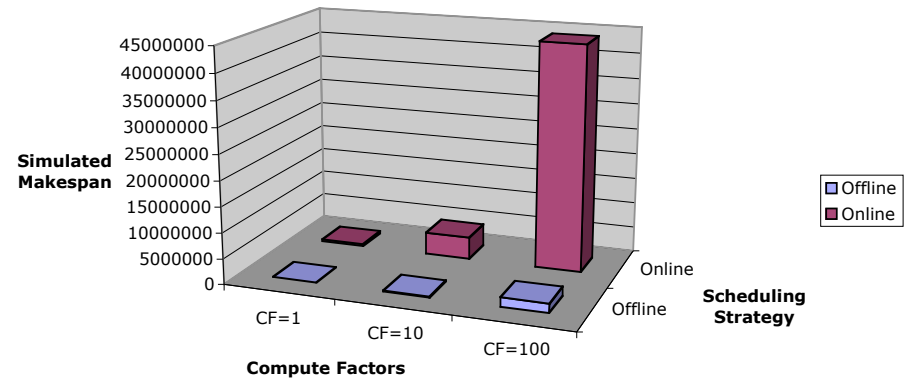
- A Virtual Grid (VG) takes
 - Shared heterogeneous resources
 - Scalable information service
- and provides
 - An hierarchy of application-defined aggregations (e.g. ClusterOf) with constraints (e.g. processor type) and rankings
- Virtual Grid Execution System (vgES) implements VG
 - VG Definition Language (vgDL)
 - VG Find And Bind (vgFAB)
 - VG Monitor (vgMON)
 - VG Application Launch (VgLAUNCH+DVCW)
 - VG Resource Info (vgAgent)



VGrADS Tool Research

- Scheduling of workflow computations
 - Off-line look-ahead scheduling dramatically improves in total time
 - Accurate performance models significantly affect quality of scheduling
 - Batch queue behavior can be predicted accurately enough for scheduling decisions
- Fault tolerance
 - Diskless checkpointing for linear algebra computations (application-specific)
 - Temporal reasoning for fault prediction
 - Optimal checkpoint frequency for iterative applications

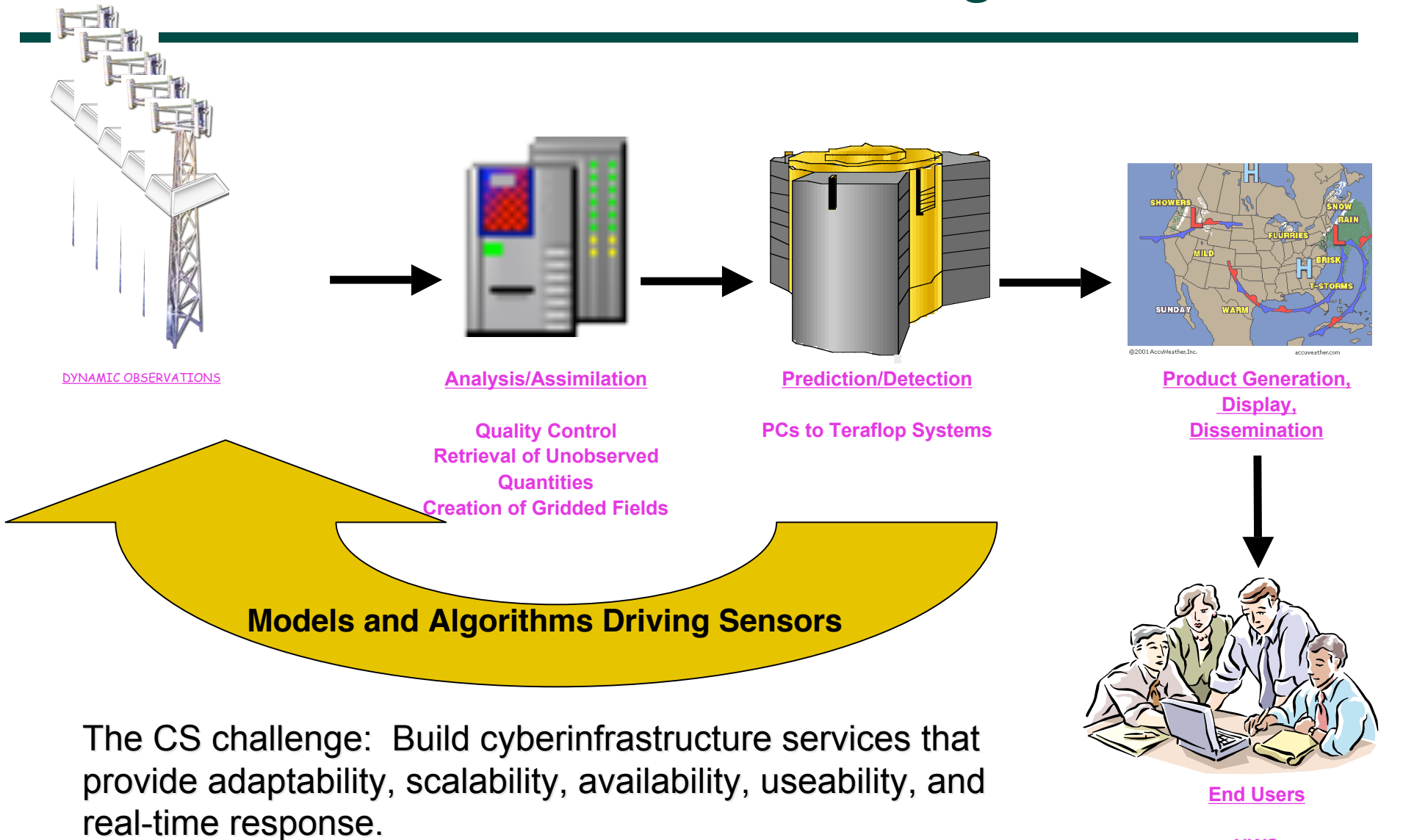
Online vs. Offline - Heterogeneous Platform (Compute Intensive Case)



VGrADS: What's New

- **SC'04**
 - Scheduling EMAN application
 - Aware of performance models
- **SC'05**
 - Find and Bind (FAB) for resource selection
 - Scheduling EMAN application
 - Aware of batch queue predictions (and performance models)
- **SC'06**
 - Virtual Grid "slots" for resource availability
 - Start time + duration
 - Uses advance reservations where available
 - Uses batch queue prediction elsewhere
 - Scheduling LEAD application
 - Aware of reservations and batch queue predictions (and performance models)

The LEAD Vision: A Paradigm Shift



LEAD Portal – Experiment Builder

The screenshot displays the LEAD Portal Experiment Builder interface within a Microsoft Internet Explorer browser window. The browser address bar shows the URL <https://portal-dev.leadproject.org>. The page header includes the LEADPORTAL logo (LINKED ENVIRONMENTS FOR ATMOSPHERIC DISCOVERY) and the NSF logo (SPONSORED BY THE NATIONAL SCIENCE FOUNDATION). The navigation menu includes HOME, MY WORKSPACE, ABOUT LEAD, DATA SEARCH, EXPERIMENT (selected), VISUALIZE, EDUCATION, RESOURCES, and HELP. The current page is titled "Experiment Builder" and contains an "Experiment Wizard" section.

The Experiment Wizard section is titled "Specify a name, description, and select workflow". It shows the following fields:

- Name: VGrADSTest
- Description: (empty text area)
- Workflow: (dropdown menu)

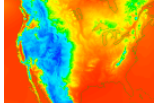
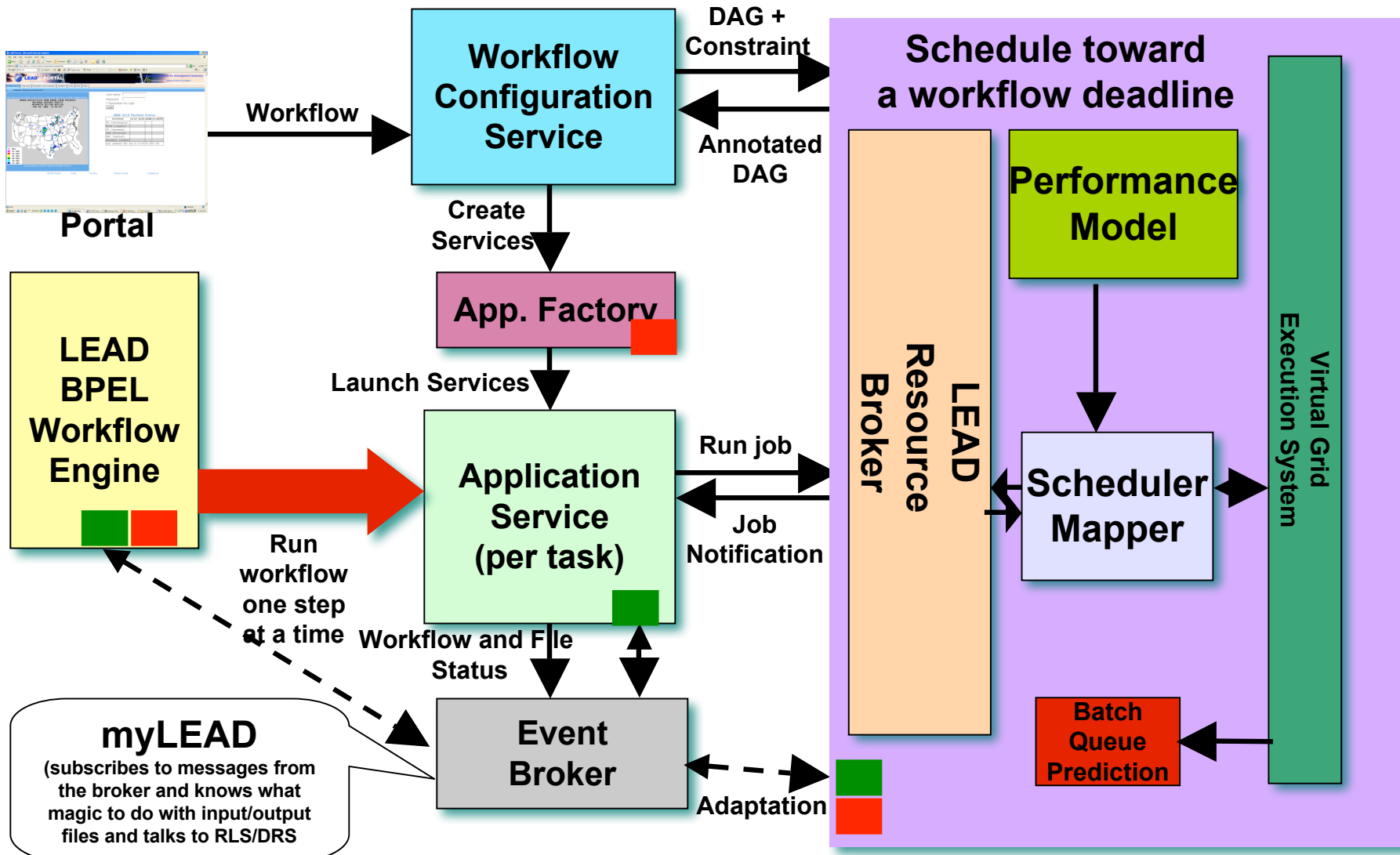
The selected workflow is "ADAS-Initialized-WRF-Forecast". The description for this workflow is: "A workflow to run WRF Forecast with ADAS initialized data".

The workflow diagram consists of the following components and connections:

- ConfPropertiesFile Config** (yellow box) is the starting point, with arrows pointing to **Wrf_Static_Preprocessor**, **ADASDataFiles Config**, **ADAS_Interpolator**, **Terrain_Preprocessor**, and **Lateral_Boundary_Interpolator**.
- ADASDataFiles Config** (yellow box) has an arrow pointing to **ADAS_Interpolator**.
- Terrain_Preprocessor** (yellow box) has an arrow pointing to **ADAS_Interpolator**.
- Lateral_Boundary_Interpolator** (yellow box) has an arrow pointing to **ADAS_Interpolator**.
- ADAS_Interpolator** (yellow box) has an arrow pointing to **ARPS2WRF_Interpolator**.
- Wrf_Static_Preprocessor** (yellow box) has an arrow pointing to **ARPS2WRF_Interpolator**.
- ARPS2WRF_Interpolator** (yellow box) has an arrow pointing to **WRF_Forecasting_Model**.
- WRF_Forecasting_Model** (yellow box) has an arrow pointing to **WRF_Output_Files Config**.
- WRF_Output_Files Config** (yellow box) is the final output of the workflow.

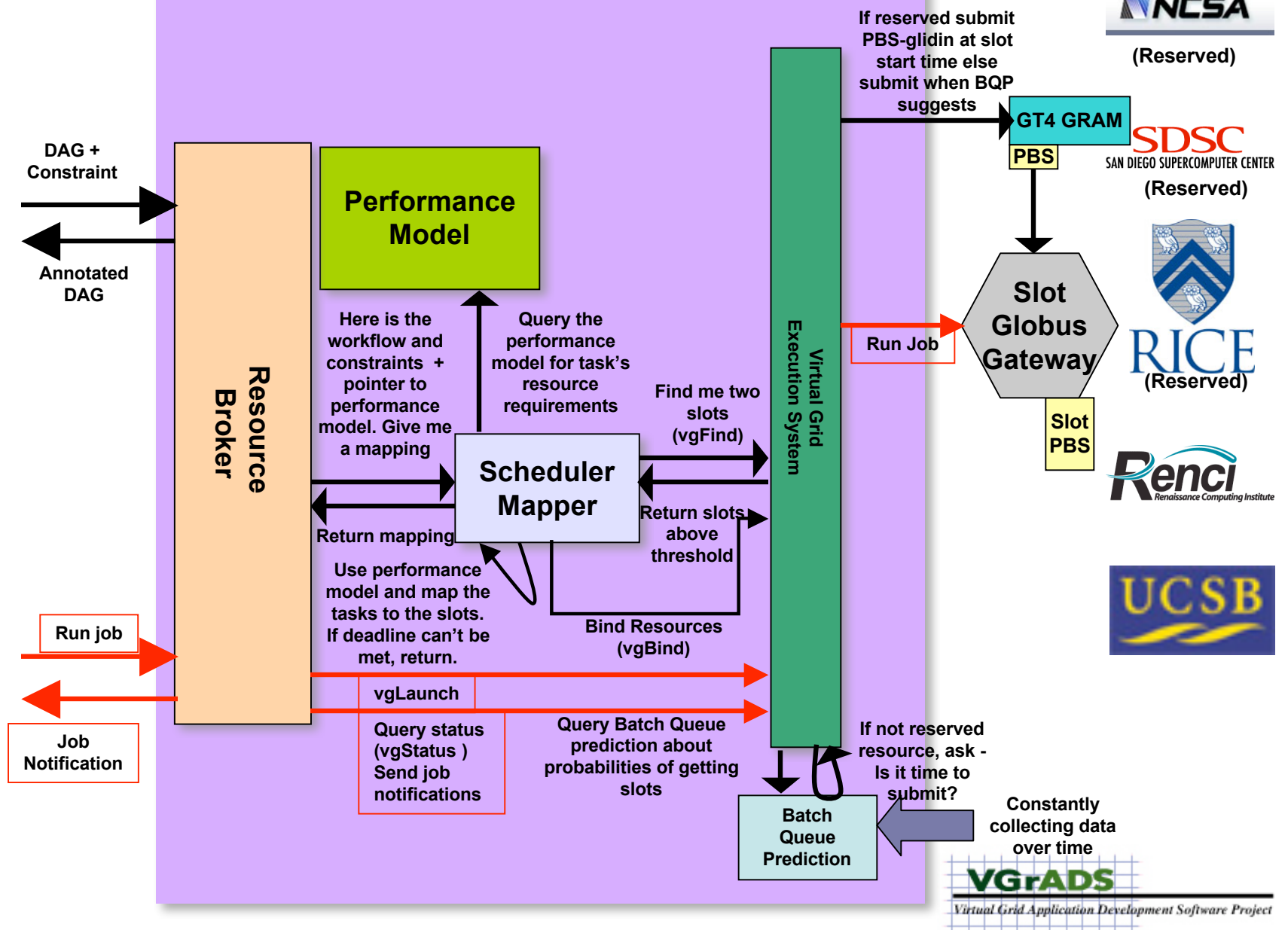
The Windows taskbar at the bottom shows the Start button, several application icons, and the system tray with the time 11:53 AM and date 11/13/2007.

VGrADS Application Collaboration



LEAD
Linked Environments for Atmospheric Discovery

Schedule toward a workflow deadline



Some Future Challenges

- Parallelism in the LEAD workflow manager
 - Parallel steps in different slots or within one slot
- Accurate Slot Requests Through Preliminary Scheduling
 - Minimization of wasted slot time
 - Accurate scheduling, better queue prediction
 - Dynamic adaptation of slot reservations
 - Requires some form of resource equivalence:
 - For step B, I need the equivalent of 200 Opterons, where 1 Opteron = 3 Itanium = 1.3 Power 5 (from perf models)
- Increased Schedule Robustness
 - Minimizing variation along the critical path
- Scheduling to Minimize Cost
 - In the presence of cycle exchange rates
 - Get the minimum-cost resources to solve the problem by the given deadline

VGrADS at SC'06

- **Booth Talks and Demos**

- Tuesday, noon - GCAS booth (1825)
- Tuesday, 2:30 - USC booth (2246) [Not live]
- Wednesday, 1:00 - SDSC booth (1915)
- Thursday, 10:30 - RENC I booth (1143)
- What you'll see
 - LEAD running on several clusters
 - Scheduler mapping LEAD components to slots
 - vgES managing slots via batch queue prediction

- **Papers**

- "Improving Grid Resource Allocation via Integrated Selection and Binding" by Kee, et al. - Wednesday, 10:30
- "Toward a Doctrine of Containment: Grid Hosting with Adaptive Resource Control" by Ramakrishnan, et al. - Wednesday, 11:00
- "Evaluation of a Workflow Scheduler Using Integrated Performance Modeling and Batch Queue Wait Time Prediction" by Nurmi, et al. - Thursday, 2:00

Launching from the LEAD Portal

- Work in Progress

VGrADS



[Slots Only](#) | [Slots + Map](#) | [Archives](#)

Slots



time →

UNSUBMITTED - RUNNING - DONE - FAILURE

```

Terrain_Preprocessor ARPS2_WRF_Interpolator
Wrf_Static_Preprocessor

```

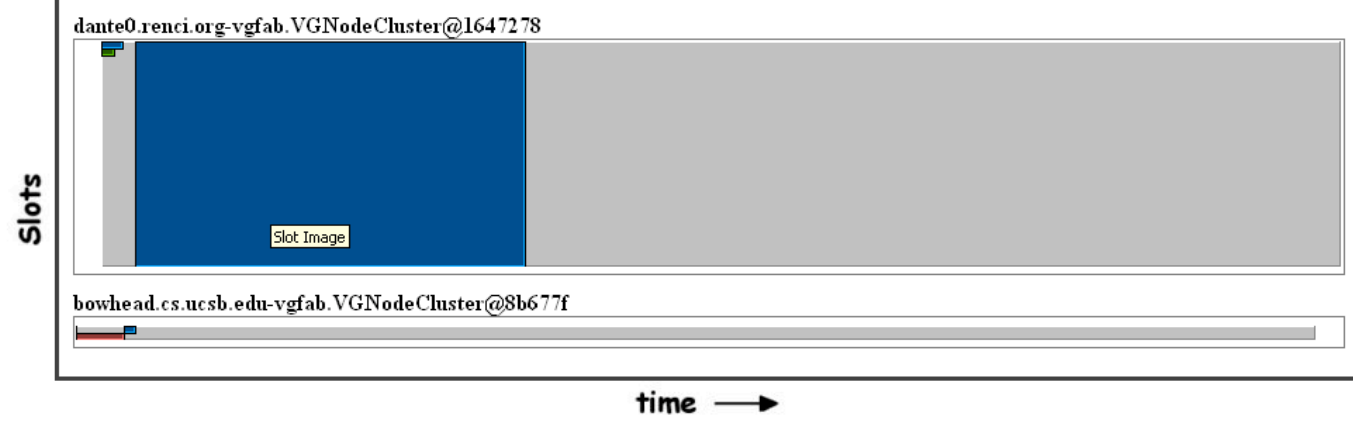
Message Log

Project

VGrADS



Slots Only | Slots + Map | Archives



UNSUBMITTED - RUNNING - DONE - FAILURE

Lateral_Boundary_WRF_Forecasting_Model
_D_Model_Data_

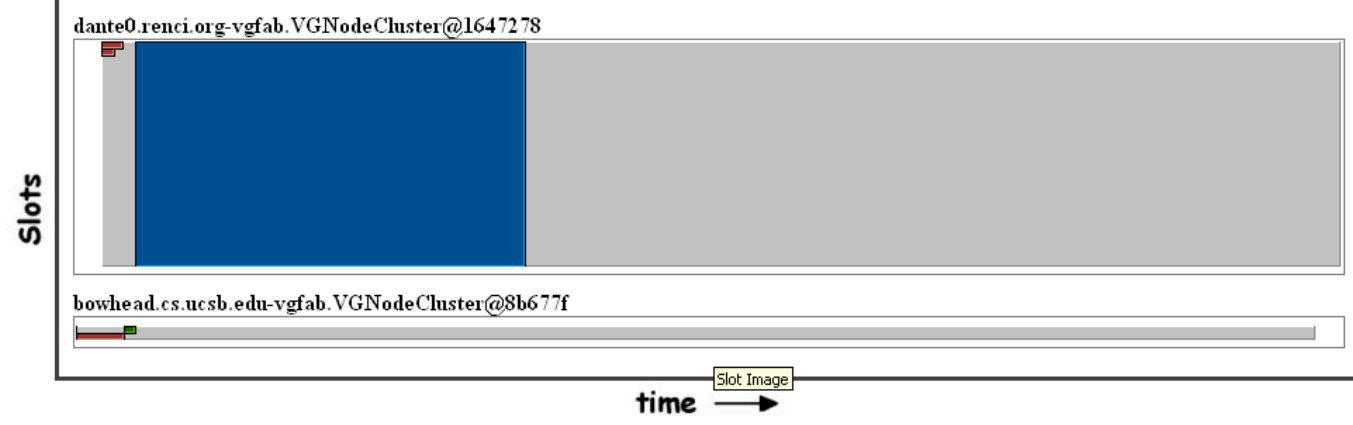
Message Log

Project

VGGrADS



[Slots Only](#) | [Slots + Map](#) | [Archives](#)



[UNSUBMITTED](#) - [RUNNING](#) - [DONE](#) - [FAILURE](#)

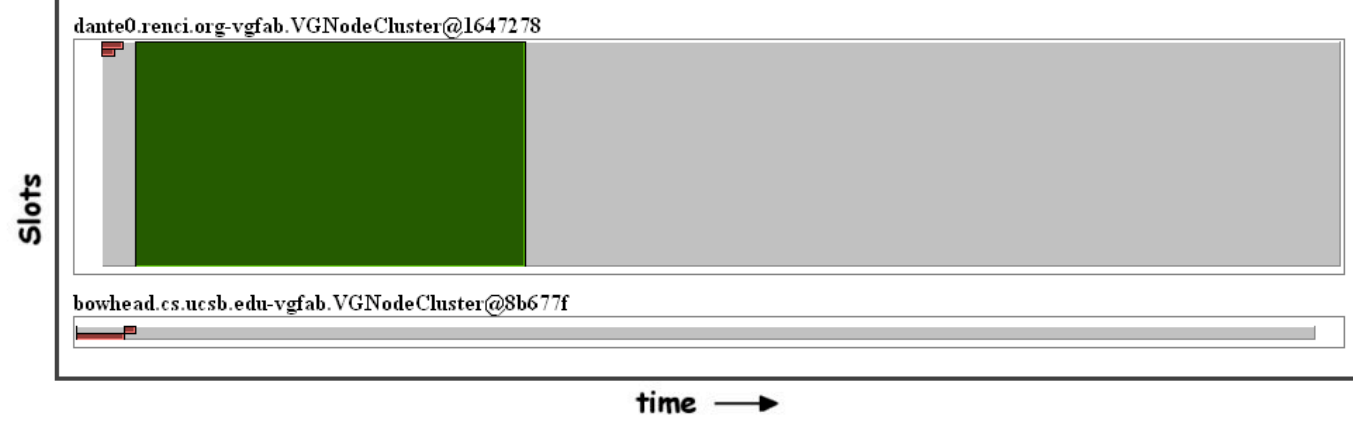
```
Terrain_Preprocessor ARPS2_WRF_Interpolator  
Wrf_Static_Preprocessor
```

Message Log

VGGrADS



[Slots Only](#) | [Slots + Map](#) | [Archives](#)



[UNSUBMITTED](#) - [RUNNING](#) - [DONE](#) - [FAILURE](#)

```
Lateral_Boundary WRF_Forecasting_Model  
_D_Model_Data_
```

Message Log

Project



WELCOME TO THE LEAD PORTAL



Linked Environments for Atmospheric Discovery (LEAD) makes meteorological data, forecast models, and analysis and visualization tools available to anyone who wants to interactively explore the weather as it evolves. The LEAD Portal brings together all the necessary resources at one convenient access point ... [read more](#)

Remember my login

[Forgot your password?](#)
[Create new account](#)

FEATURES FOR ANYONE INTERESTED IN THE WEATHER

Researchers	With university, government, or industry affiliations	<input type="button" value="GET FEATURES"/>
Educators	At college and university level, high school, or middle schools	<input type="button" value="GET FEATURES"/>
Students	At graduate, undergraduate, middle and high school levels	<input type="button" value="GET FEATURES"/>
Visitors	Newcomers and the curious	<input type="button" value="GET FEATURES"/>

QUICK LINKS

- [Live Weather](#)
- [LEAD Grid](#)
- [Glossary](#)
- [Website Help](#)
- [Frequently Asked Questions](#)

POPULAR TOOLS

<p>Visualize Weather Data</p> <p>Integrated Data Viewer MORE ></p> 	<p>Make a Forecast or Analysis</p> <p>Experiment Builder MORE ></p> 	<p>Access Weather Data</p> <p>Geographic Region Search MORE ></p> 
--	---	--

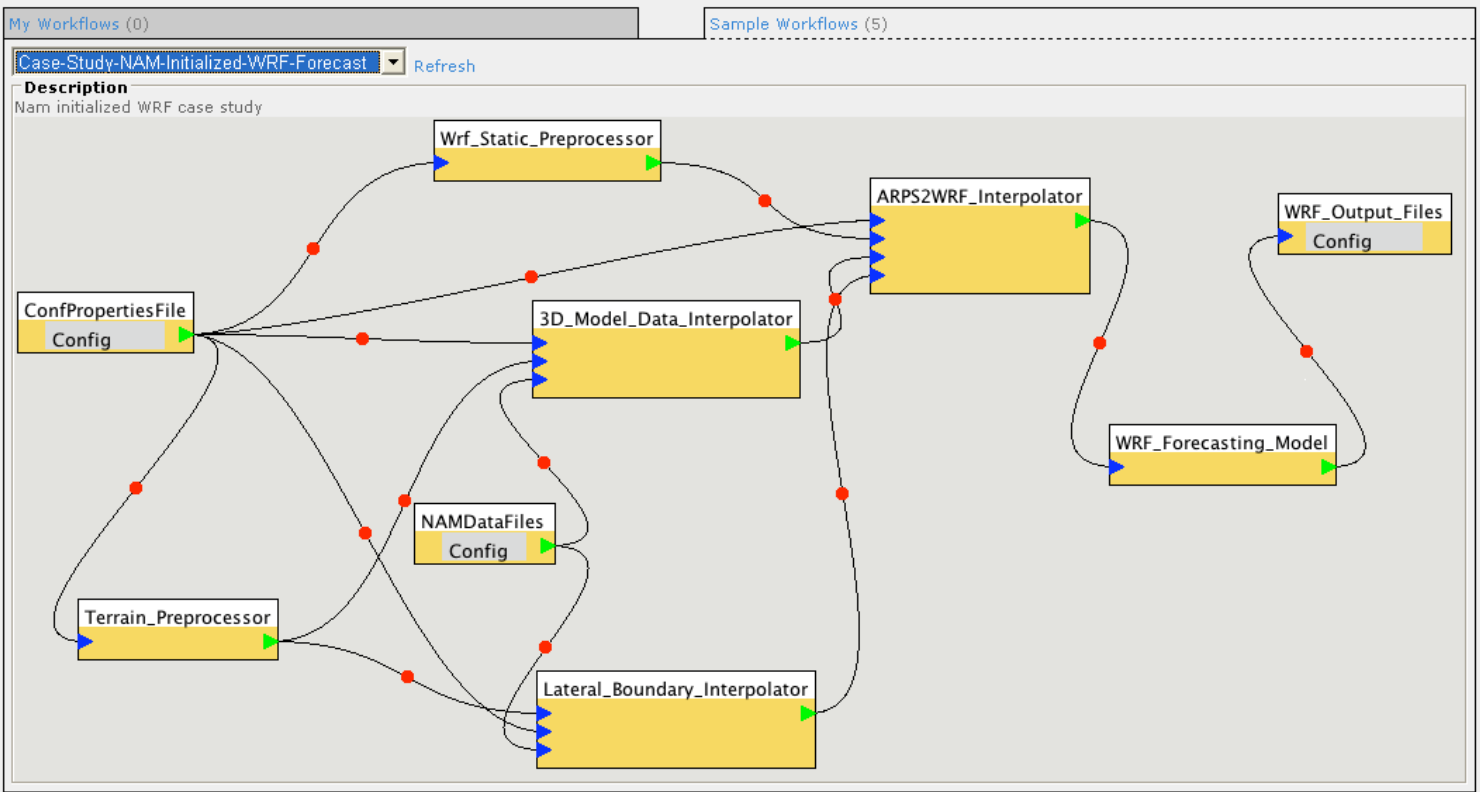
THE LEAD TEAM

Specify a name, description, and select workflow

Name: VGrADSWorkflow

Description:

Workflow



< Back Next > Cancel Launch



Experiment Builder Portlet

Experiment Wizard

User: VGrADS Demo Project: SC-Testing
 Name: VGrADSWorkflow
 Description:
 Workflow: Case-Study-NAM-Initialized-WRF-Forecast

Select options for NAMDataFiles required for this experiment

Description: Choose a set of NAM data files for generating interpolated boundary condition files. For a 6 hour forecast, please choose 3 NAM forecast time steps ending with "f00", "f03", "f06" and for a 12 hour forecast please choose 5 NAM forecast time steps ending in "f00", "f03", "f06", "f09" and "f12". Note that all time step files should be from the same model run.

Select	Name	Description	Timestamp
<input checked="" type="checkbox"/>	Default option 0	The value of the default option is [gsiftp://grid-hg.ncsa.teragrid.org/gpfs_scratch1/drlead/data/eta40grb.2006051700f00]	
<input checked="" type="checkbox"/>	Default option 1	The value of the default option is [gsiftp://grid-hg.ncsa.teragrid.org/gpfs_scratch1/drlead/data/eta40grb.2006051700f03]	
<input checked="" type="checkbox"/>	Default option 2	The value of the default option is [gsiftp://grid-hg.ncsa.teragrid.org/gpfs_scratch1/drlead/data/eta40grb.2006051700f06]	
<input checked="" type="checkbox"/>	Default option 3	The value of the default option is [gsiftp://grid-hg.ncsa.teragrid.org/gpfs_scratch1/drlead/data/eta40grb.2006051700f09]	
<input checked="" type="checkbox"/>	Default option 4	The value of the default option is [gsiftp://grid-hg.ncsa.teragrid.org/gpfs_scratch1/drlead/data/eta40grb.2006051700f12]	
<input type="checkbox"/>	North American Model/CONUS 40 km in Form for ADAS (conduit) (eta40grb.2006110706f00)	NCEP North American Model 40 km with forecast hours in separate files : AWIPS 212 (R) Regional - CONUS - Double Resolution. Model runs are made at 00Z, 06Z, 12Z, and 18Z and have analysis and forecasts every 3 hours out to 84 hours. Horizontal = 185 by 129 points, resolution 40.63 km, LambertConformal projection. Vertical = surface, 1000 to 50 hPa pressure levels, layers, and depth.	2006-11-07 11:00Z
<input type="checkbox"/>	North American Model/CONUS 40 km in Form for ADAS (conduit) (eta40grb.2006110706f03)	NCEP North American Model 40 km with forecast hours in separate files : AWIPS 212 (R) Regional - CONUS - Double Resolution. Model runs are made at 00Z, 06Z, 12Z, and 18Z and have analysis and forecasts every 3 hours out to 84 hours. Horizontal = 185 by 129 points, resolution 40.63 km, LambertConformal projection. Vertical = surface, 1000 to 50 hPa pressure levels, layers, and depth.	2006-11-07 11:00Z
<input type="checkbox"/>	North American Model/CONUS 40 km in Form for ADAS (conduit) (eta40grb.2006110706f06)	NCEP North American Model 40 km with forecast hours in separate files : AWIPS 212 (R) Regional - CONUS - Double Resolution. Model runs are made at 00Z, 06Z, 12Z, and 18Z and have analysis and forecasts every 3 hours out to 84 hours. Horizontal = 185 by 129 points, resolution 40.63 km, LambertConformal projection. Vertical = surface, 1000 to 50 hPa pressure levels, layers, and depth.	2006-11-07 11:00Z
<input type="checkbox"/>	North American Model/CONUS 40 km in Form for ADAS (conduit) (eta40grb.2006110706f09)	NCEP North American Model 40 km with forecast hours in separate files : AWIPS 212 (R) Regional - CONUS - Double Resolution. Model runs are made at 00Z, 06Z, 12Z, and 18Z and have analysis and forecasts every 3 hours out to 84 hours. Horizontal = 185 by 129 points, resolution 40.63 km, LambertConformal projection. Vertical = surface, 1000 to 50 hPa pressure levels, layers, and depth.	2006-11-07 11:00Z
<input type="checkbox"/>	North American Model/CONUS 40 km in Form for ADAS (conduit) (eta40grb.2006110706f12)	NCEP North American Model 40 km with forecast hours in separate files : AWIPS 212 (R) Regional - CONUS - Double Resolution. Model runs are made at 00Z, 06Z, 12Z, and 18Z and have analysis and forecasts every 3 hours out to 84 hours. Horizontal = 185 by 129 points, resolution 40.63 km, LambertConformal projection. Vertical = surface, 1000 to 50 hPa pressure levels, layers, and depth.	2006-11-07 11:00Z
<input type="checkbox"/>	North American Model/CONUS 40 km in Form for ADAS (conduit) (eta40grb.2006110712f00)	NCEP North American Model 40 km with forecast hours in separate files : AWIPS 212 (R) Regional - CONUS - Double Resolution. Model runs are made at 00Z, 06Z, 12Z, and 18Z and have analysis and forecasts every 3 hours out to 84 hours. Horizontal = 185 by 129 points, resolution 40.63 km, LambertConformal projection. Vertical = surface, 1000 to 50 hPa pressure levels, layers, and depth.	2006-11-07 17:00Z



Introduction Experiment Builder

Experiment Builder Portlet

SUCCESSFULLY CREATED NEW EXPERIMENT. YOU CAN MONITOR YOUR EXPERIMENT USING THE WORKFLOW COMPOSER.

User: VGrADS Demo Project: SC-Testing Add Project ...

Experiments

New Experiment

Experiment Name	Description	Created On	Last Updated On	Status
VGrADSWorkflow	No description	Wed Nov 08 03:15:29 EST 2006	Wed Nov 08 02:15:32 EST 2006	STARTED
MyFirstOne	No description	Wed Nov 08 01:10:04 EST 2006	Wed Nov 08 00:10:06 EST 2006	RUNNING
NAM Test	testing nam workflow	Wed Nov 08 01:01:32 EST 2006	Wed Nov 08 00:01:35 EST 2006	RUNNING

Project

Scheduling with Batch Queues

- Last Year: VGrADS supported scheduling using estimated batch queue waiting times
 - Batch queue estimates are factored into communication time
 - E.g., the delay in moving from one resource to another is data movement time + estimated batch queue waiting time
 - Unfortunately, estimates can have large standard deviations
- This Year: limiting variability through two strategies:
 - Resource reservations: partially supported on the TeraGrid and other schedulers
 - In advance queue insertion: submit jobs before data arrives based on estimates
 - Can be used to simulate advance reservations
- Exploiting this requires a preliminary schedule indicating when the resources are needed
 - Problem: how to build an accurate schedule when exact resource types are unknown

Preliminary Scheduling Solution

- Use performance models to specify alternative resources
 - For step B, I need the equivalent of 200 Oopteron, where 1 Oopteron = 3 Itanium = 1.3 Power 5
 - Equivalence from performance model
- This permits an accurate preliminary schedule because the performance model standardizes the time for each step
 - Scheduling can then proceed with accurate estimates of when each resource collection will be needed
 - Makes advance reservations more accurate
 - Data will arrive neither too early or too late
- It may provide a mixture to meet the computational requirements, if the specification permits
 - Give me a loose bag of tight bags containing the equivalent of 200 Oopteron, minimize the number of tight bags and the overall cost
 - Solution might be 150 Oopteron in one cluster and 150 Itaniums in another