

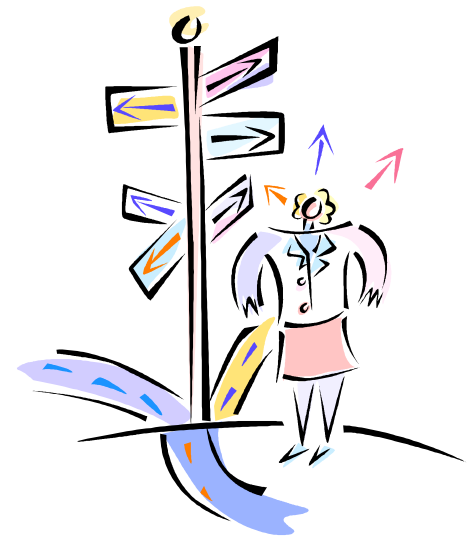
# Fault Tolerance and Recovery of Scientific Workflows on Computational Grids

Gopi Kandaswamy  
Anirban Mandal  
Daniel A. Reed

Resilience'08, May 22, 2008

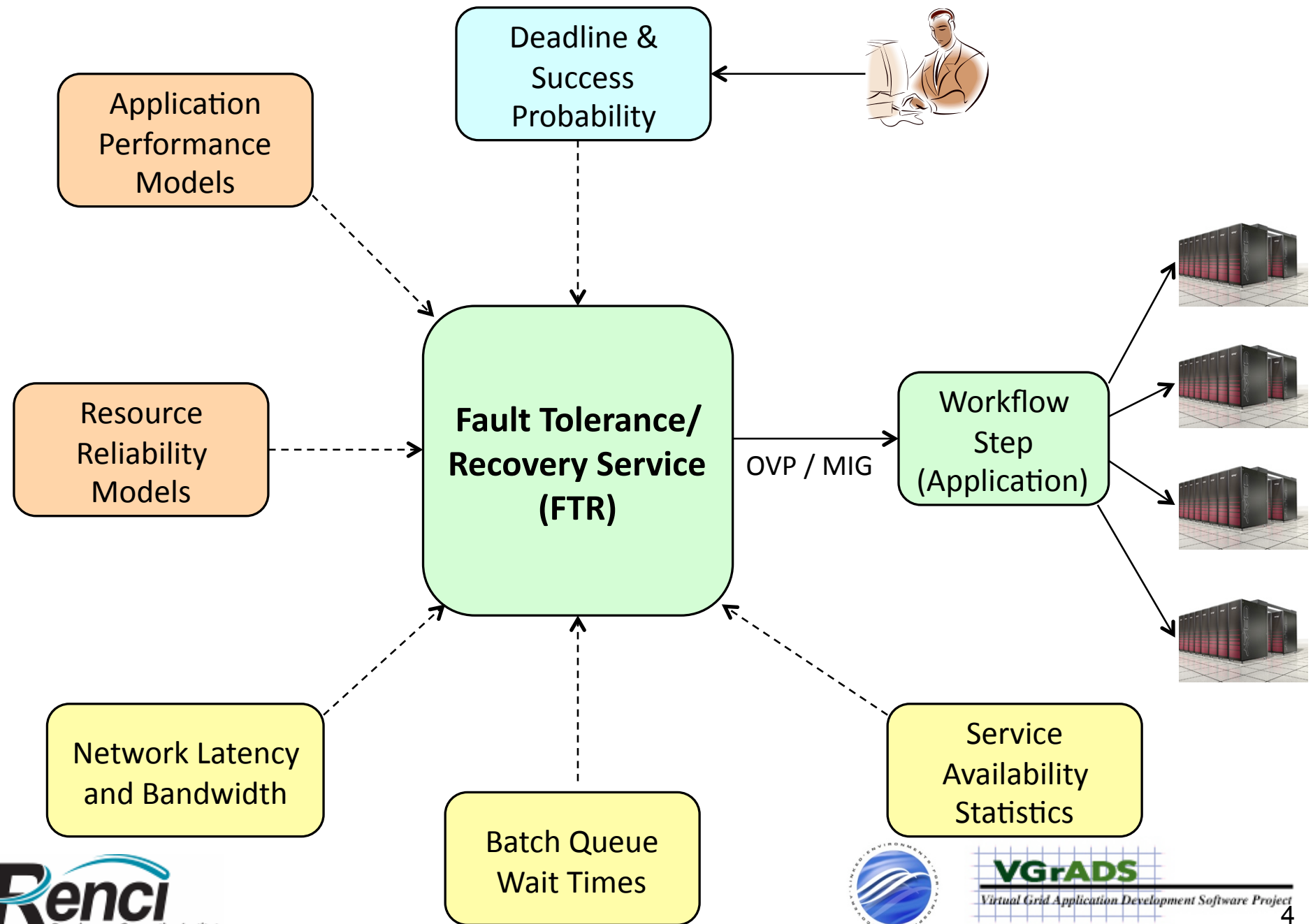
# Presentation Outline

- Motivation and rationale
- Fault tolerance and recovery service
- Algorithms
  - migration, over-provisioning
- Evaluation with LEAD
- Conclusion



# Motivation

- **Reliability and performance are related**
  - failure is the limiting case of poor performance
  - both involve measures of behavior over time
- **Large, complex workflows are sensitive to failures**
  - faults are the norm, rather than exceptions
    - distributed systems, services and resources
  - completion “guarantees” are problematic
    - workflow completion is probabilistic in the presence of faults
- **Many time-critical workflows are deadline driven**
  - severe weather events, disaster response, ...



# FTR Service

- **Resource models**
  - MDS for static resource characteristics
  - NWS for network latency and bandwidth between resources
  - QBETS for batch queue wait time prediction on resources
  - simple reliability models
- **Application models**
  - based on simple parametric historical performance models
- **Deadline of workflow steps**
- **Success probability**
  - expected probabilistic completion guarantee
- **Grid services' availability**
  - core middleware services like WS-GRAM, GridFTP

# FTR Algorithms

- **Notation**

$p_i$  : failure probability of resource  $i$  (eg. 1 hr. failure probability)

$h_i$  : expected execution cost of application on resource  $i$

- expected queue wait time
- expected computation time
- expected communication time

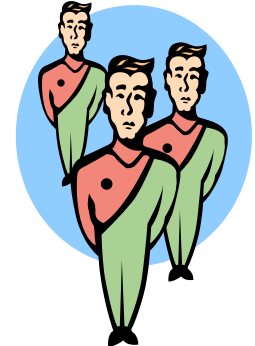
$x$  : required success probability

$d$  : required deadline

$m_i$  : failure probability of application (based on reliability model)

- Resource  $i$  represents (queue, #nodes) combinations
- Use a simple reliability model
  - assumption : resource failures are independent over time
  - resource failures follow a binomial distribution

# Over Provisioning



- Find
  - degree and resources for over-provisioning
- Number of application copies
  - meet a deadline  $d$  with a success probability  $x$
- Solve the following optimization problem

For given  $[1..M]$  resources, find a partition  $P = \{s_1, s_2 \dots s_n\}$  of  $[1..M]$  such that

$$1 - m_{s_1} * m_{s_2} * \dots * m_{s_n} \geq x \wedge |P| \text{ is minimum} \wedge \max \{h_{s_1} \dots h_{s_n}\} \leq d$$

Probability of failure

Minimum number of resources meeting deadline

# Migration



- Find the best migration path
  - Best resource chain
- Solve the following optimization problem

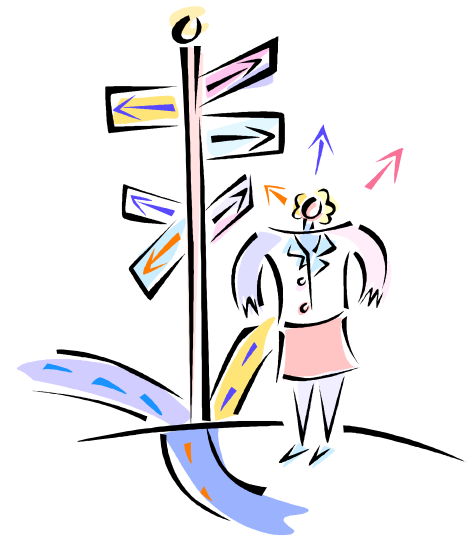
For given  $[1..M]$  resources, find a partition  $P = \{s_1, s_2 \dots s_n\}$  such that

$$1 - m_{s_1} * m_{s_2} * \dots * m_{s_n} \geq x \wedge |P| \text{ is minimum } \wedge \text{sum}(h_{s_1} + \dots h_{s_n}) \leq d,$$



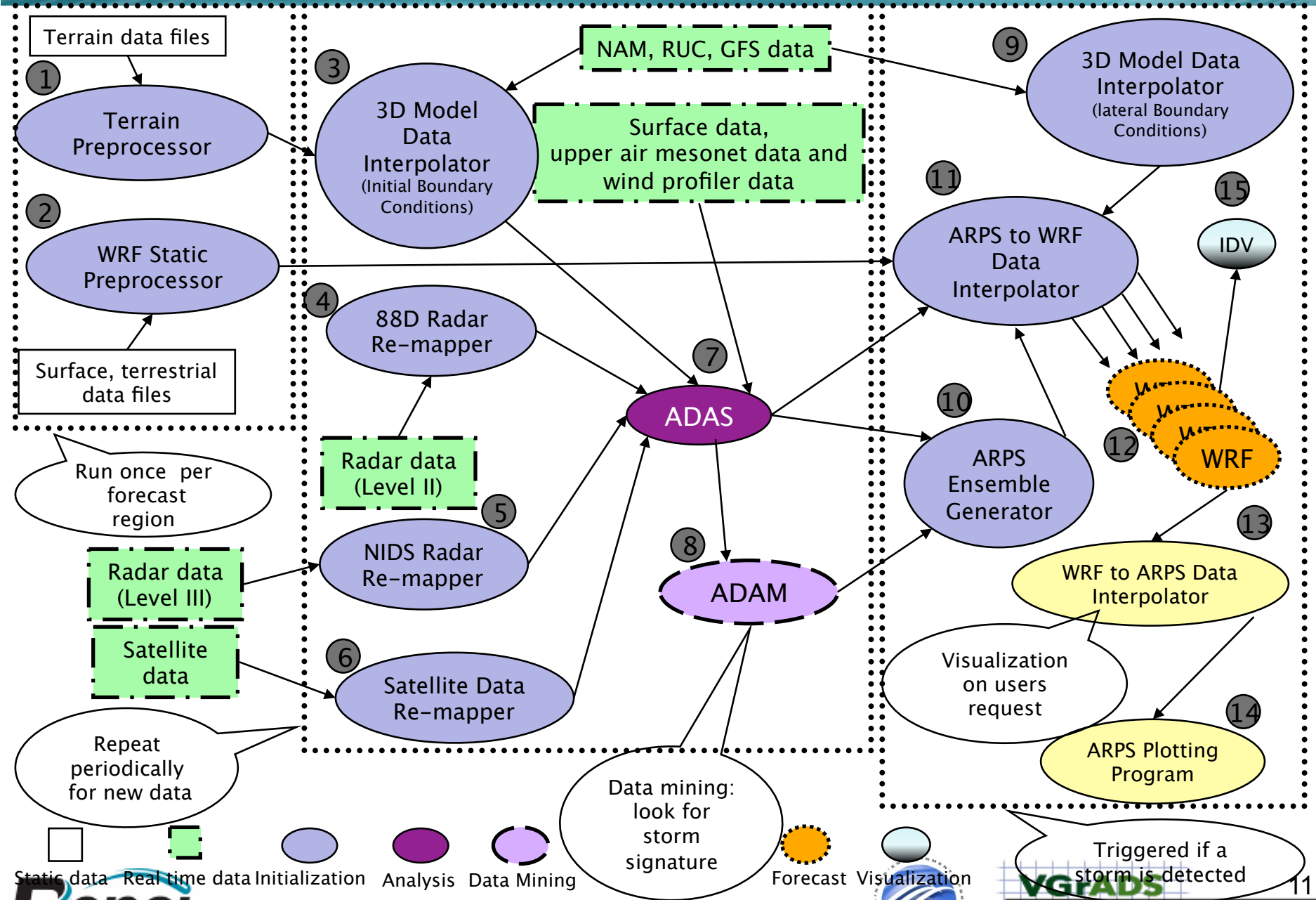
# Presentation Outline

- Motivation and rationale
- Fault tolerance and recovery service
- Algorithms
  - migration, over-provisioning
- **Evaluation with LEAD**
- **Conclusion**

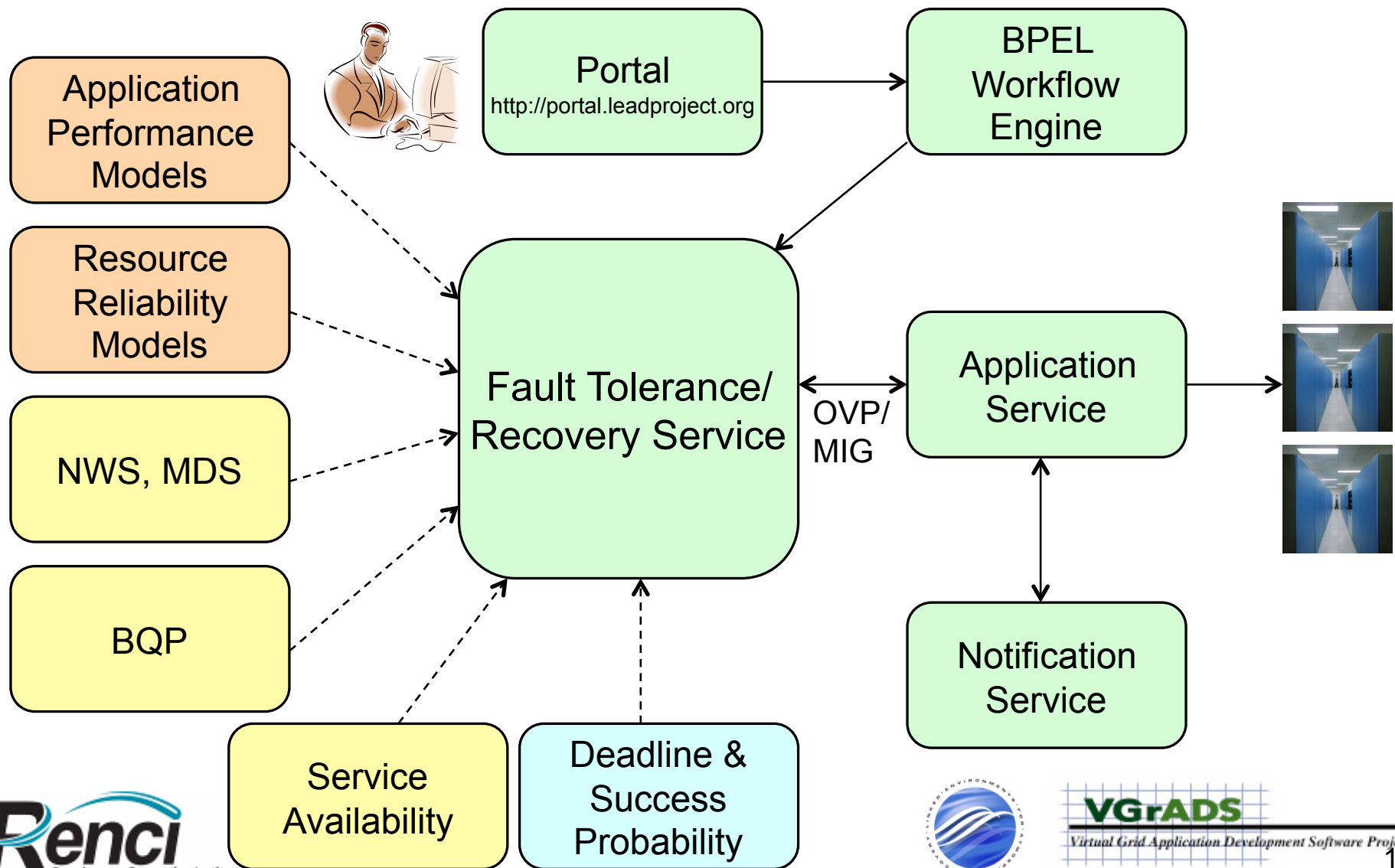


# LEAD – Linked Environment and Atmospheric Discovery

- **Integrated scalable framework for dynamic and adaptive meso-scale weather prediction**
  - computations continually steered by new weather data
  - responds to decision-driven inputs from users
  - steers remote observing technologies to optimize data collection for problem at hand
  - consists of analysis, visualization and data-mining tools
- **Framework consists of**
  - Teragrid resources at NCSA, UC and IU
  - weather data repositories (static and dynamic)
  - web portal for user interaction
  - <http://portal.leadproject.org>



# Simplified Architecture: LEAD



LEAD Portal - Mozilla Firefox

File Edit View History Bookmarks Tools Help

https://portal-dev.leadproject.org/gridsphere/gridsphere?cid=experiment

MyUNC UNC-CH Weather Google Scholar CNET Google News Slashdot Google Maps WWW Board Gmail

LEADPORTAL  
LINKED ENVIRONMENTS FOR ATMOSPHERIC DISCOVERY

SPONSORED BY THE NATIONAL SCIENCE FOUNDATION

HOME MY WORKSPACE ABOUT LEAD DATA SEARCH EXPERIMENT VISUALIZE EDUCATION RESOURCES HELP

Introduction Experiment Builder

Experiment Builder Portlet

User: Gop... swamy Project: TestProject Add Project ...

Click pencil icon

Experiments

Experiment Name	Description	Created On	Last Updated On	Status
test2	test2	Tue Apr 17 08:53:49 EDT 2007	Tue Apr 17 08:53:49 EDT 2007	FINISHED
test3	test3	Tue Apr 17 10:06:12 EDT 2007	Tue Apr 17 10:06:11 EDT 2007	FINISHED
test5	test5	Wed Apr 18 09:12:56 EDT 2007	Wed Apr 18 09:12:56 EDT 2007	FINISHED
test1	test1	Mon Apr 23 09:55:03 EDT 2007	Mon Apr 23 09:55:03 EDT 2007	FINISHED

Copyright © 2006 Linked Environments for Atmospheric Discovery. All rights reserved. CONTACT US | LOGOUT

POWERED BY TeraGrid

Done portal-dev.leadproject.org

LEAD Portal - Mozilla Firefox

File Edit View History Bookmarks Tools Help

https://portal-dev.leadproject.org/gridsphere/gridsphere?gs\_mode=EDIT&cid= Google

MyUNC UNC-CH Weather Google Scholar CNET Google News Slashdot Google Maps WWW Board Gmail

**LEADPORTAL**  
LINKED ENVIRONMENTS FOR ATMOSPHERIC DISCOVERY

SPONSORED BY THE NATIONAL SCIENCE FOUNDATION

HOME MY WORKSPACE ABOUT LEAD DATA SEARCH **EXPERIMENT** VISUALIZE EDUCATION RESOURCES HELP

Introduction Experiment Builder

Experiment Builder Portlet

Customize

- ☐ I have SPRUCE tokens and I would like to have the option of running SPRUCE workflows.
- ☐ Use the VGrADS Scheduler when running my workflows.
- ☐ Submit my workflows to the Workflow Configuration Service (WCS)
- ☒ Use the Fault Tolerant Recovery (FTR) service when submitting workflow

Submit

Back

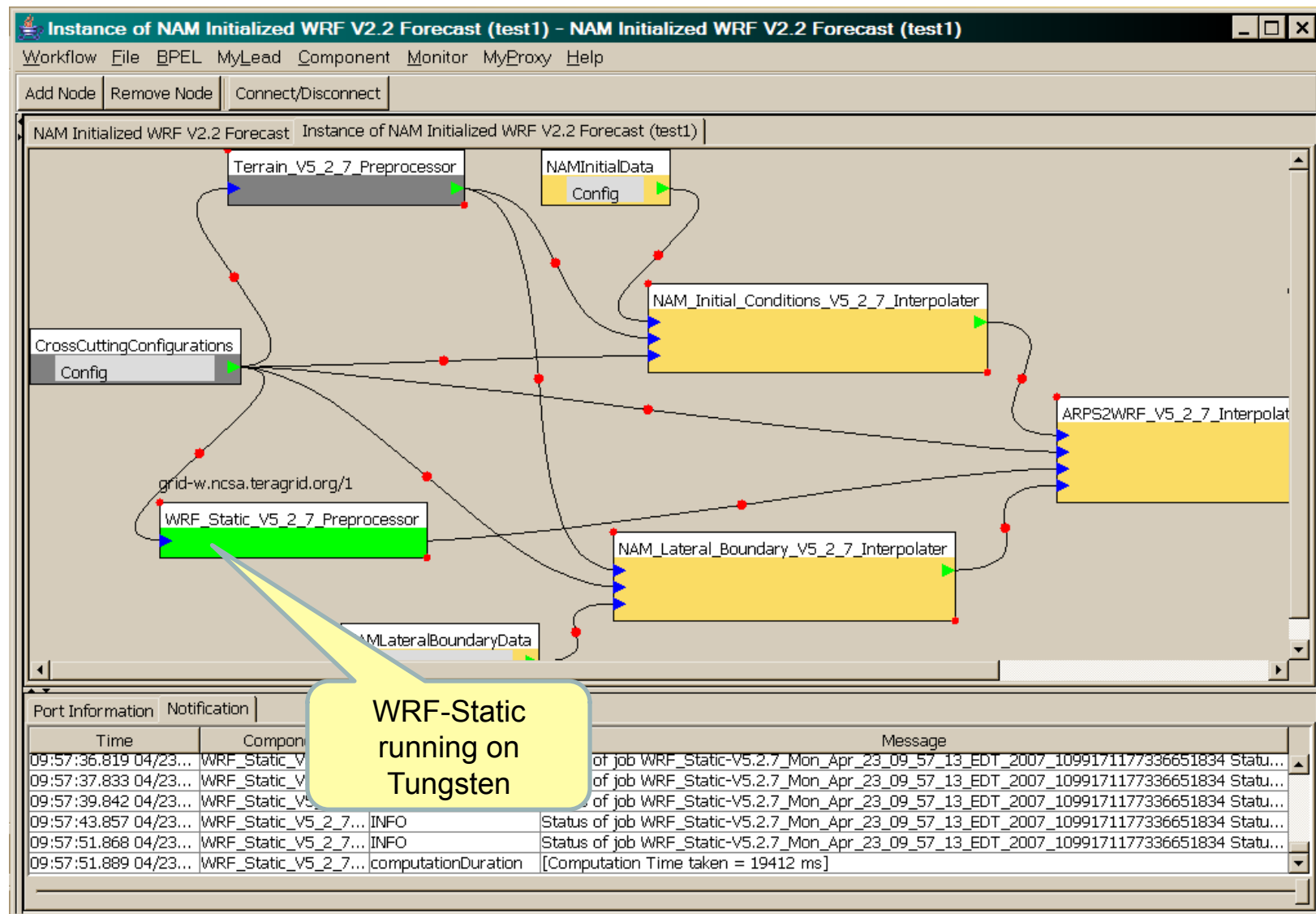
Enable FTR check box

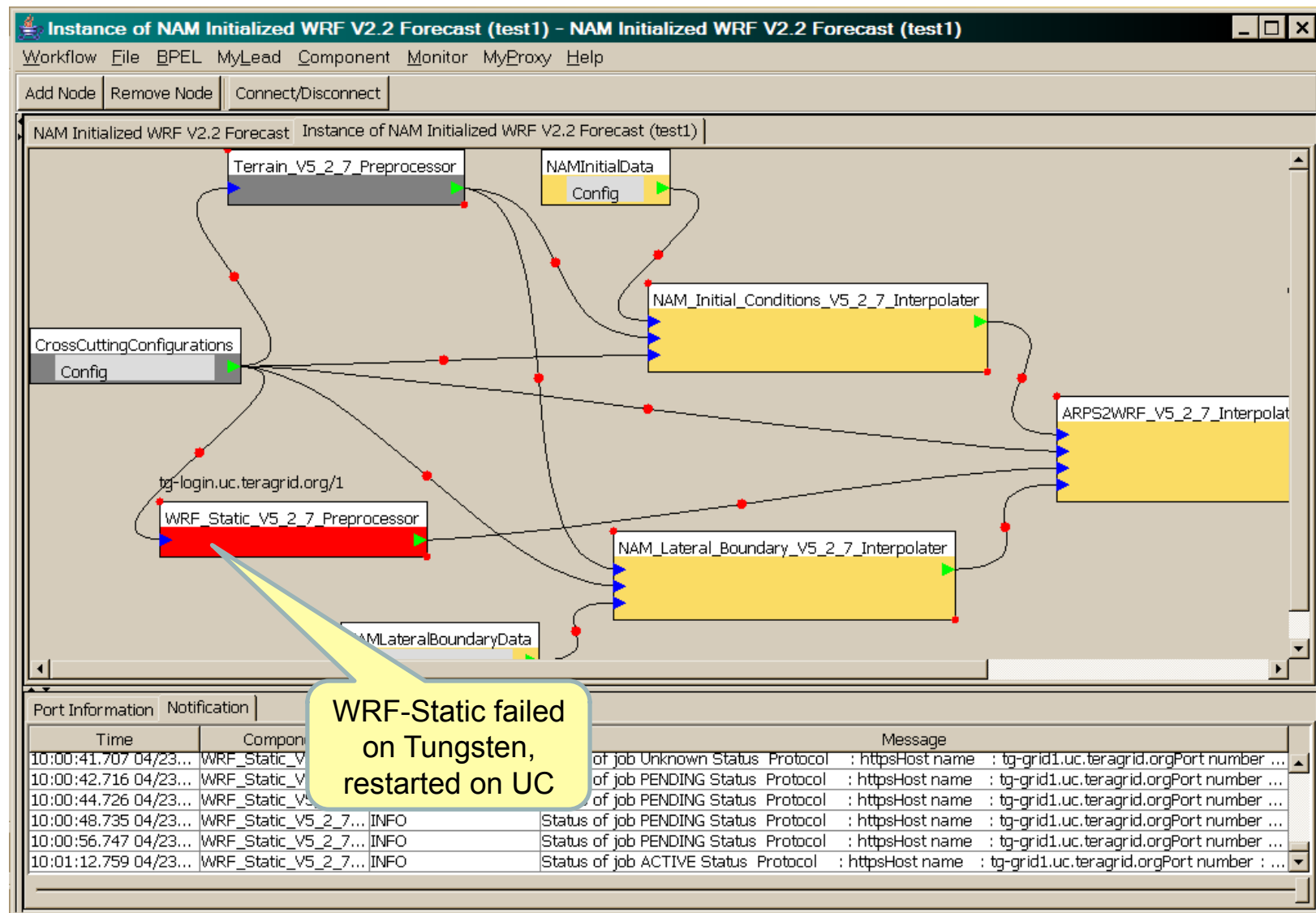
COPYRIGHT © 2006 LINKED ENVIRONMENTS FOR ATMOSPHERIC DISCOVERY. ALL RIGHTS RESERVED. CONTACT US | LOGOUT

POWERED BY TeraGrid

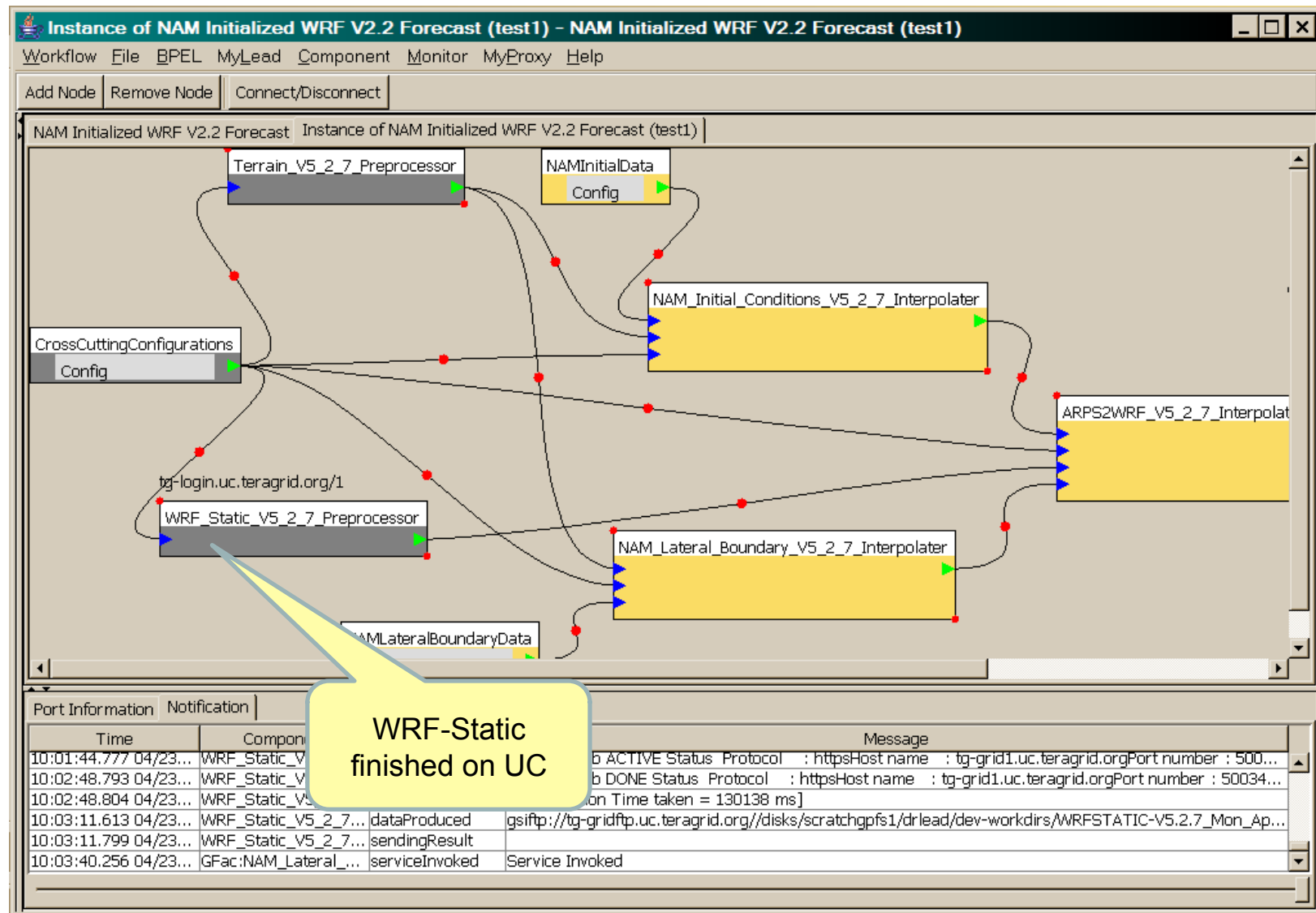
Done portal-dev.leadproject.org

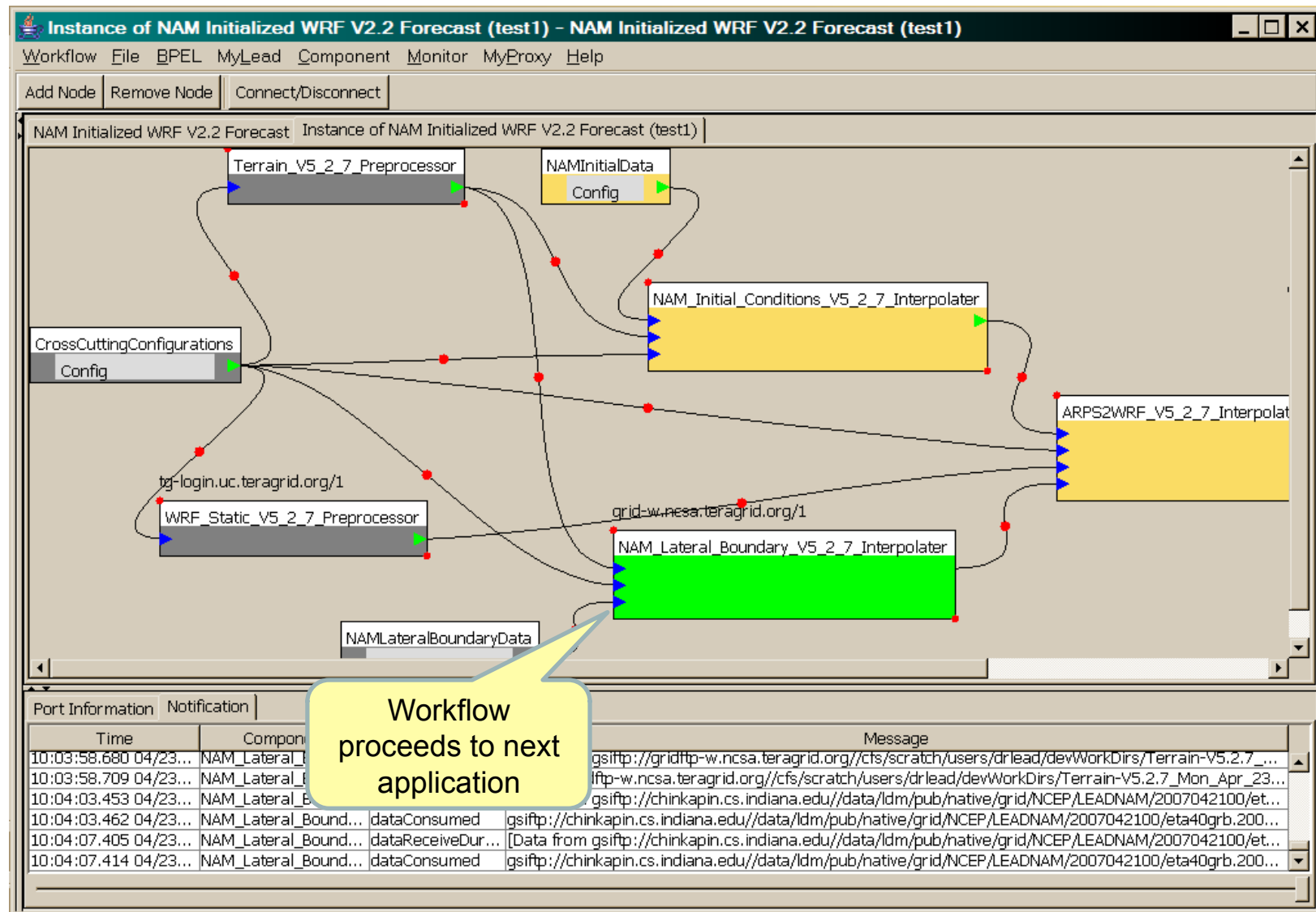


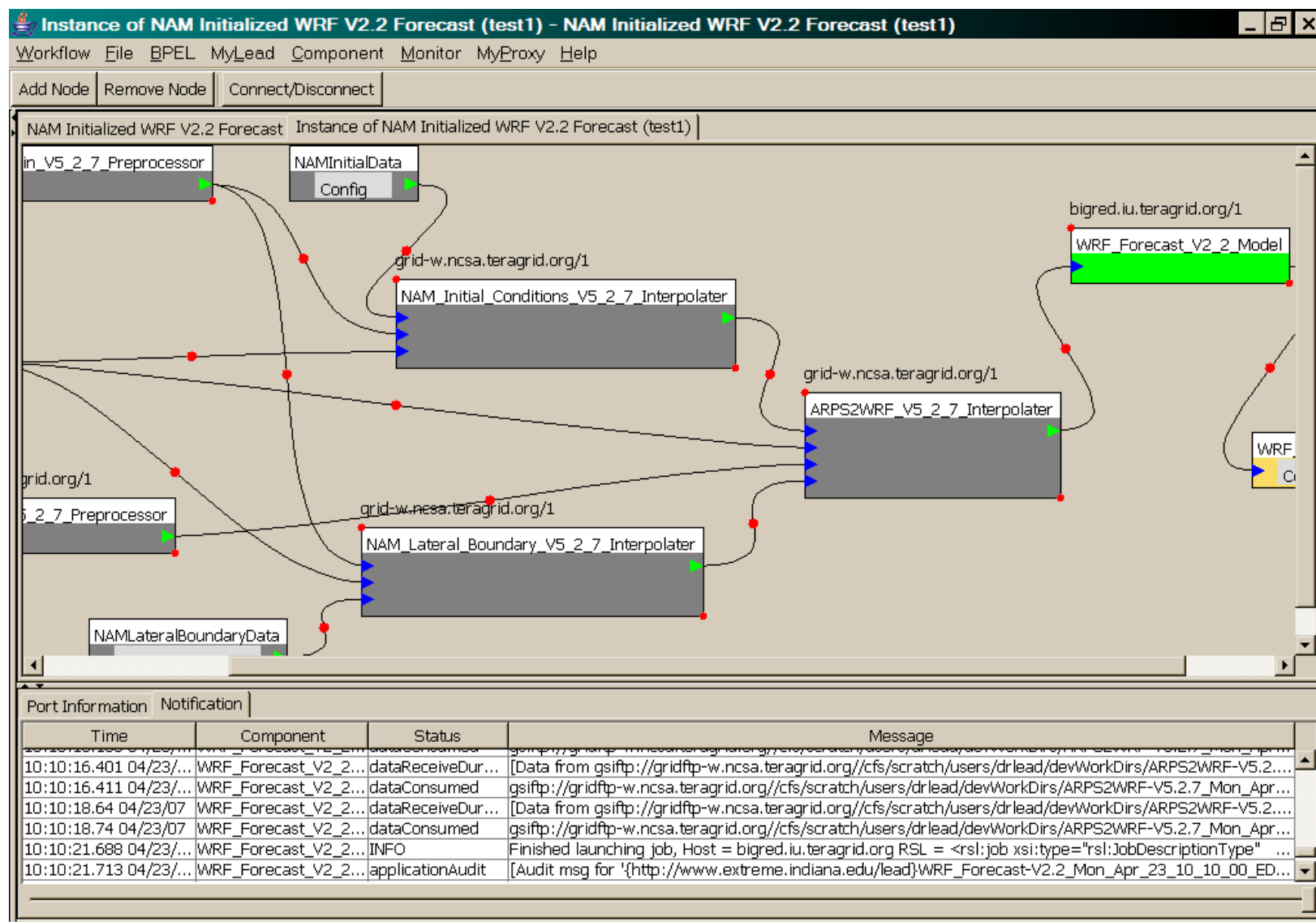




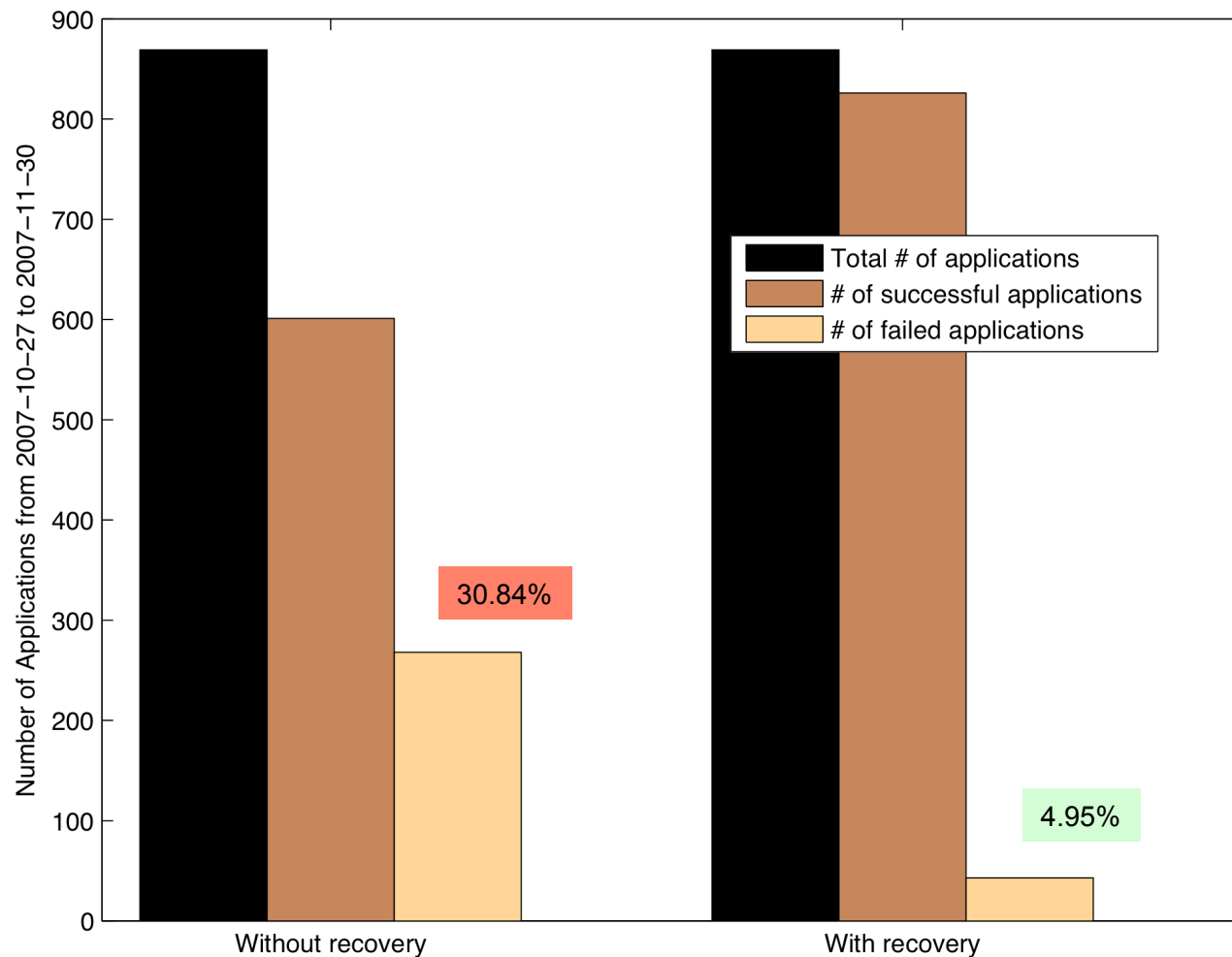




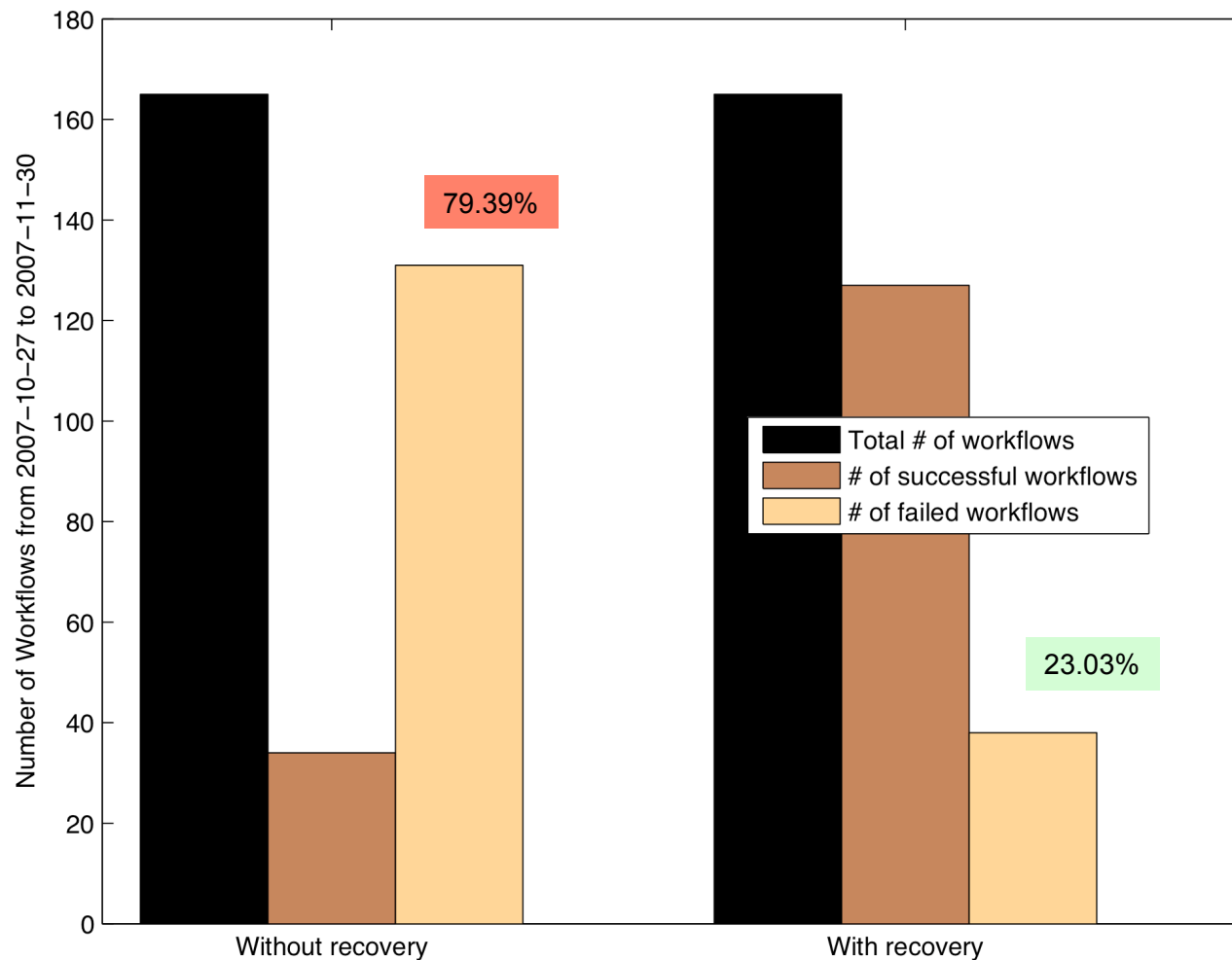




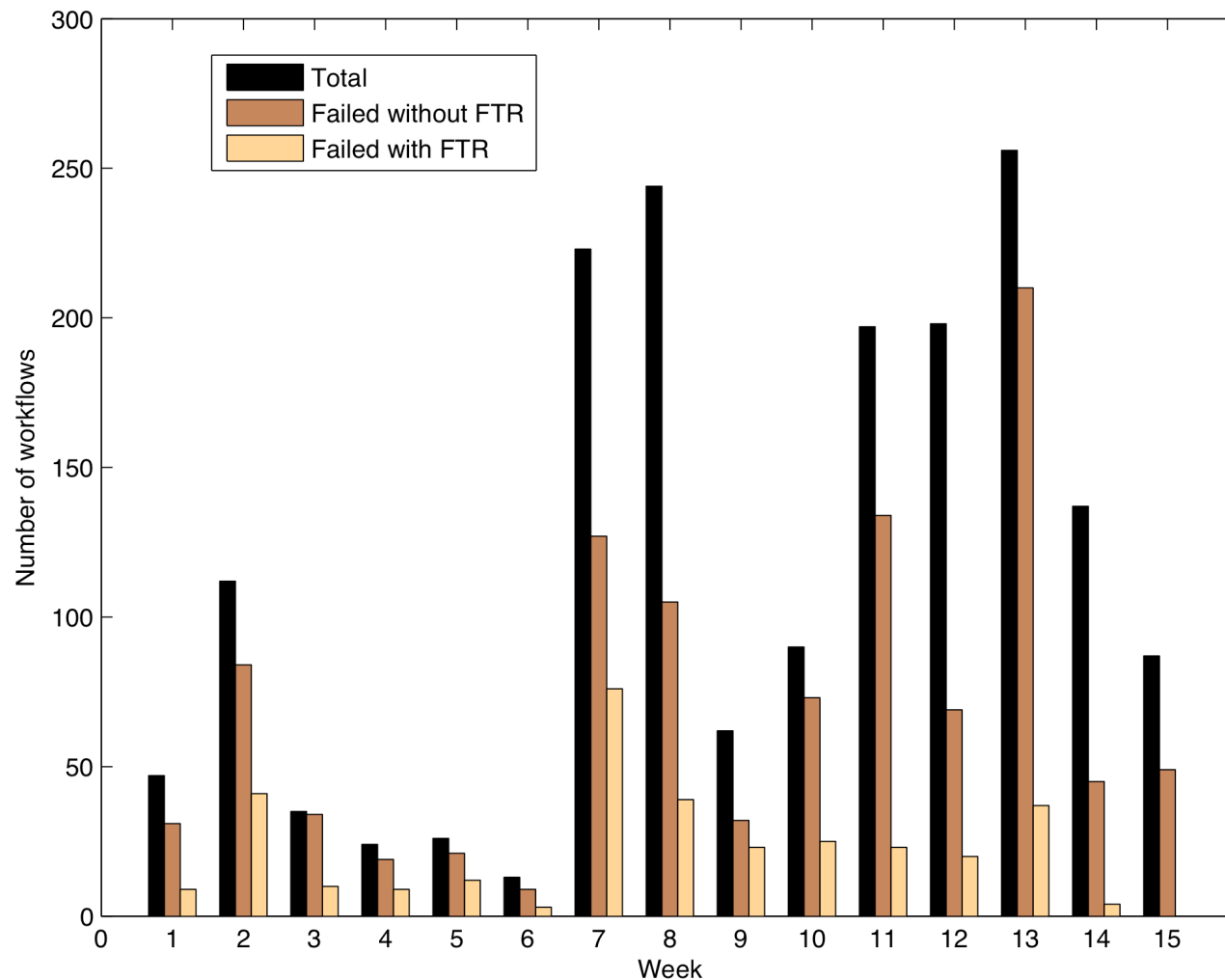
# Results: Workflow Step Failure Rate



# Results: Workflow Failure Rate



# Weekly Statistics (past 5 months)



# Conclusions

- **Developed a fault tolerance and recovery (FTR) service**
  - delivers reliable execution of workflows on grids
    - under deadline and success probability constraints
  - uses migration and over-provisioning techniques
- **Deployed FTR with LEAD production infrastructure**
  - transparent to users
- **Results from LEAD workflows show**
  - Reduction of application failure rate from 31% to 5%
  - Reduction of workflow failure rate from 80% to 23%
- **Future work**
  - accurate reliability estimates of resources
  - other fault-tolerance techniques for different workflow types

# Questions ?

Thank you..



<http://portal.leadproject.org>

