# Building the Bioscience  Gateway

**Lavanya Ramakrishnan**
lavanya@renci.org

**Renaissance Computing Institute**
**Duke University**
**North Carolina State University**
**University of North Carolina -  Chapel Hill**

NORTH CAROLINA
**BIOPORTAL**

**www.ncbioportal.org**

# Building Science Communities

- **National Evolutionary Synthesis Center**
  - $15M, five year project, Duke, NCSU and UNC-CH
  - national center and resource
    - research, data federation and outreach
    - sabbaticals and teaching release
  - RENCI support
    - *data models, portals and Grid infrastructure*
- **The Carolina Center for Exploratory Genetic Analysis**
  - develop collaborative experiences and plans
    - preliminary data to apply for a P50 grant
  - develop a prototype informatics infrastructure
    - data models, methods, tools and portals
  - facilitate use of best practices for existing projects
- **North Carolina Bioportal**
  - leverage state-wide investment in bioinformatics and grid
  - undergraduate education, graduate education, faculty research

# North Carolina Bioportal
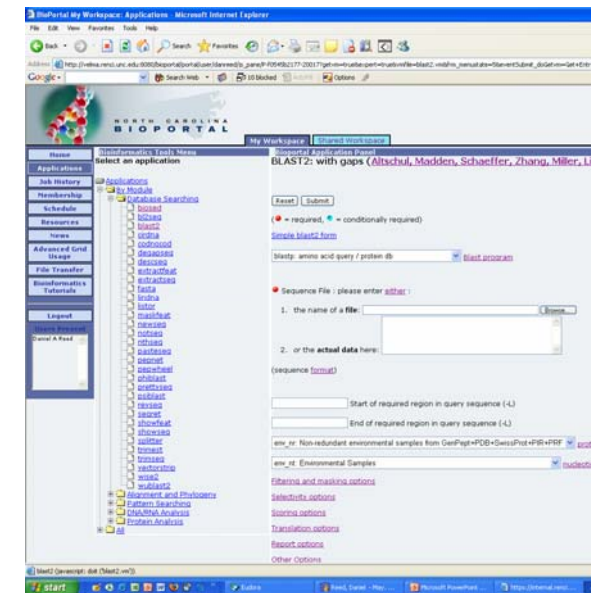
- **Features**
  - access to common bioinformatics tools
  - extensible toolkit and infrastructure
    - OGCE and National Middleware Initiative (NMI)
    - leverages emerging international standards
  - remotely accessible or locally deployable
  - packaged and distributed with documentation
- **National reach and community**
  - TeraGrid deployment
    - scheduled for summer 2005
- **Education and training**
  - hands-on workshops across North Carolina
    - clusters, Grids, portals and bioinformatics

# Bioportal Computing Infrastructure

- **34 node Linux cluster**
  - one head node
  - 32 compute nodes
  - one storage node
- **Configuration**
  - 3.06 GHz dual Xeon processors
  - 4 GB memory/node
  - 8 GB memory on storage node
- **1.73 TB storage array**
  - 14 x 146GB U320 SCSI Drives
  - RAID 5 partitioned
- **Software stack**
  - ROCKS cluster software
  - Globus toolkit + Torque/Maui + MyProxy
  - OGCE portal software

# Portlets



OGCE - GridFTP, Globus
CHEF/Sakai – Resources,
Schedule, News

Bioinformatics Applications

Shared Workspace
(Discussion, Schedule, Chat)

Job History

# Current Bioportal Applications

- **Applications**
  - ~140 distinct codes
- **Application Suites**
  - EMBOSS
    - European Molecular Biology Open Software Suite
  - GLIMMER
    - gene identification in microbial DNA
  - HMMER
    - Hidden Markov Model program for profile-based sequence analysis
  - NCBI
    - diverse set of tools
  - PHYLIP
    - PHYLogeny Inference Package for inferring phylogenies
- **Others (incomplete list)**
  - ClustalW, FASTA

- **Standard bioinformatics databases**
  - NCBI Aggregate (95 GB)
    - three formats: native, BLAST and WUBLAST
  - GenBank (206 GB)
  - GenPept (3 GB)
  - PDB (6.3 GB)
  - Prints (72 MB)
  - RepBase (8.6 MB)
  - UniProt (12 GB)
  - PFam (8.7 GB)
  - ProSite (16 MB)
  - TransFac (36 MB)
- **Database update mechanism**
  - follows the schedule of the distribution source
  - currently NCBI Aggregate is the only one updated nightly

# PISE

- **Pasteur Institute Software Environment (PISE)**
  - generates web interfaces for molecular biology tools
    - XML specification for command line interfaces
  - see www.pasteur.fr/recherche/unites/sis/Pise

- **Rationale and objectives**
  - simplify specification of program interfaces
    - homogeneous specification mechanisms
  - reuse of existing software interfaces
    - independent development and integration
  - extension for integration with graphical interfaces
    - complexity hiding and commonality

- **Bioportal program described in PISE**
  - semi-automated GUI synthesis from XML via Perl
- **Output is a generated command line, for example**
  - blastall -p blastp -d env_nr -i query.dat.blast2.1116248106513 -a 2

# An Example PISE XML



```
<parameter ismandatory="1" iscommand="1"
    issimple="1" type="Excl">
  <name>blast2</name>
  <attributes>
    <prompt>Blast program</prompt>
    <format>
      <language>perl</language>
      <code>"blastall -p $value"</code>
    </format>
    <vdef><value>blastp</value></vdef>
    <group>1</group>
    <vlist>
      <value>blastn</value>
      <label>blastn: nucleotide query / nucleotide
db</label>
      <value>blastp</value>
      <label>blastp: amino acid query / protein
db</label>
      <value>blastx</value>
      <label>blastx: nucleotide query translated / protein
db</label>
      <value>tblastn</value>
      <label>tblastn: protein query / translated
nucleotide db</label>
      <value>tblastx</value>
      <label>tblastx: nucleotide query transl. / transl.
nucleotide db</label>
      <value>psitblastn</value>
      <label>psitblastn: protein query / transl. nucleotide
db</label>
    </vlist>
  </parameter>
```

```
<parameter type="Integer">
    <name>start_region</name>
    <attributes>
      <prompt>
      Start of required region in query sequence (-L)
      </prompt>
    </attributes>
  </parameter>
```
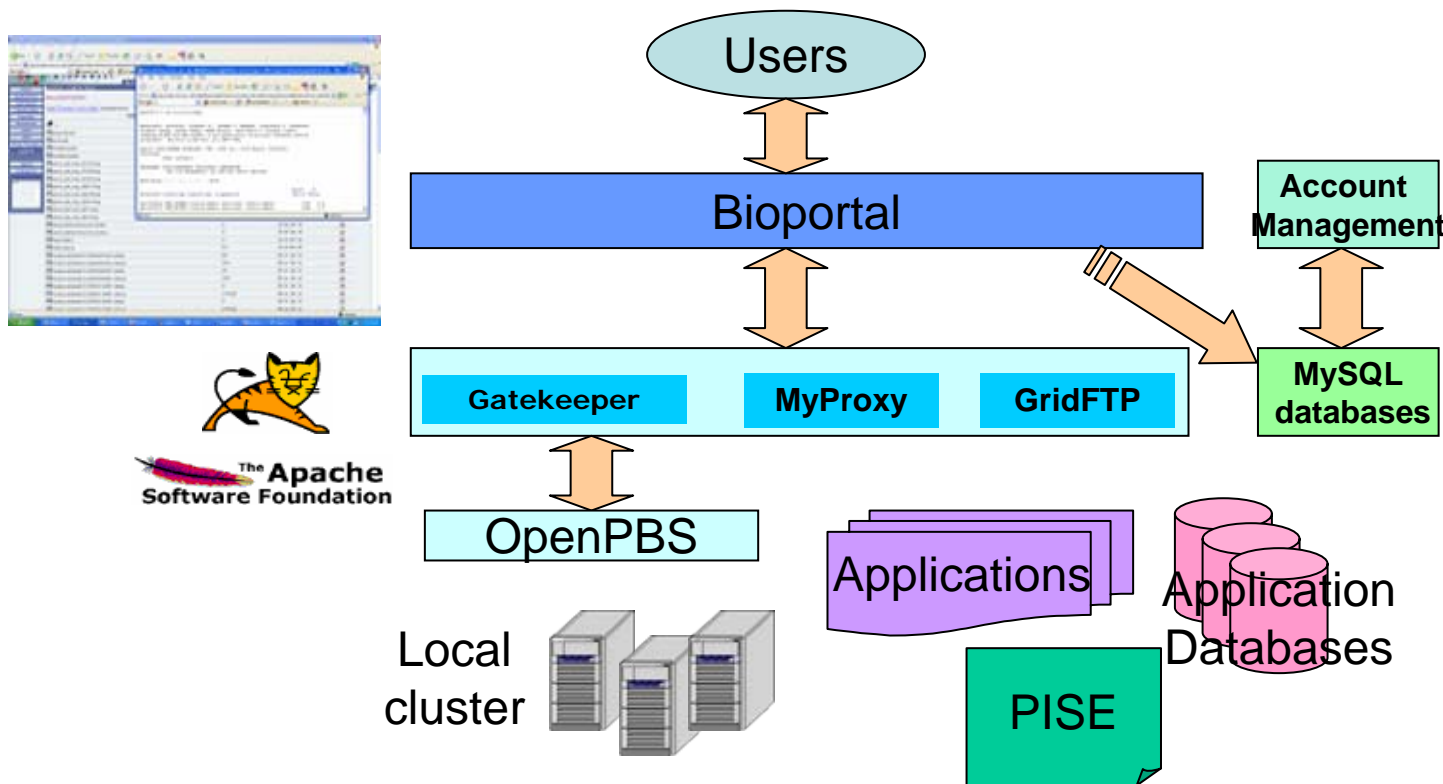
# North Carolina Bioportal



- **Open Grid Computing Environment (OGCE)**
  - shared development based on NMI toolkit
  - standard web services
  - adopting portal standards (JSR168)
  - used by cyberinfrastructure projects
    - LEAD, NEES, PACI, DOE, TeraGrid …

# Bioportal User Interactions

Users

User requests an account

Bioportal

Account Management System

**Behind the scenes**
1. Create Unix account
2. Create a certificate request
3. Sign the certificate request
4. Update MyProxy
5. Add entry to gridmap file
6. Create a portal account

**Grid Gatekeeper**

**MyProxy**

**GridFTP**

**MySQL databases**

OpenPBS

Applications

Application Databases

PISE

Linux cluster
34 nodes
1.73 TB storage

# Bioportal User Interactions
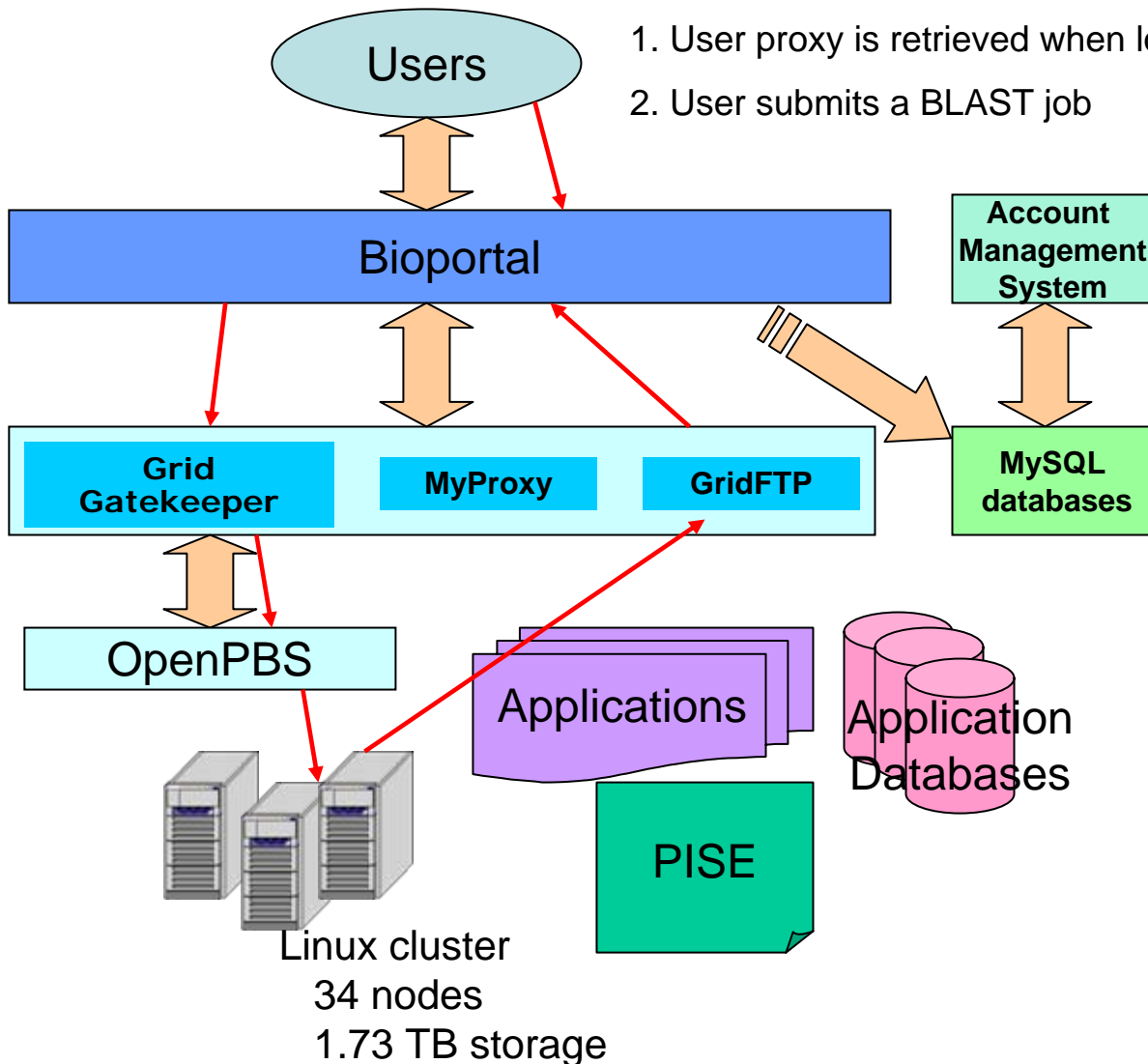
Users

Bioportal

**Account Management System**

Grid Gatekeeper

MyProxy

GridFTP

**MySQL databases**

OpenPBS

Applications

Application Databases

PISE

Linux cluster
34 nodes
1.73 TB storage

1. User proxy is retrieved when logged in

2. User submits a BLAST job

3. User can view job history

4. User can view job output

**Behind the scenes**
1. Job history database is updated
2. Job submitted to the gatekeeper
3. Enqueued in OpenPBS queue
4. Job executes
5. Output files viewed via GridFTP

# Bioportal Experiences

- **Security**
  - account creation and management
  - Grid Security Infrastructure (GSI) ,SSL
- **Job Management**
  - unique job directory
  - manage job  files -14 day policy
- **Application Domain Issues**
  - conflicts with Globus RSL
  - size and policy of database updates

# Bioportal: What's Next

- **Community Engagement**
  - workshops, experiences and deployments
  - Software and documentation
- **Infrastructure**
  - dynamic job scheduling across multiple sites
    - load driven based on community use
  - fully automated database updates, possible distributed replication
    - driven by user needs and available disk space
- **Portal tool suite**
  - expand application and databases based on user feedback
    - phylogeny, morphology, microarray analysis, …
  - different file format support
- **Leverage national presence**
  - TeraGrid/NCSA bioinformatics portal
  - NESCent evolutionary biology portal