



iRODS and the RENCI Data Working Group

*Howard Lander
Michael Shoffner*

renci

RESEARCH \ ENGAGEMENT \ INNOVATION

The Renaissance Computing Institute

- Formed in 2004 as a collaborative institute involving the University of North Carolina at Chapel Hill, Duke University and North Carolina State University.
- RENCI develops and deploys advanced technologies to enable research discoveries and practical innovations.
- This ***science of cyberinfrastructure*** is essential to continuing scientific discovery and innovation.

RENCI Resources

- A diverse group of people including domain scientists in oceanography, meteorology, chemistry, informatics and computer science.
- A diverse set of projects and collaborators spanning the domains listed above and more.
- Several compute clusters with an aggregate peak computing power of approximately 30 Teraflops.
- More than one Pb of spinning disk.
- An ideal laboratory to develop the science of cyberinfrastructure

The Data Working Group

- Chartered in May 2010, as an outgrowth of discussions that started in late 2009.
- Motivated by the realization that RENCI had a number of ongoing projects with significant data challenges.
- Existing projects and knowledge were confined to project specific stove pipes. No way to run an Institute!

RENCI Data Working Group

- Is responsible for providing leadership and strategic guidance for RENCi in the data technology area.
- Includes data architecture, technology research, development and operations, and dissemination and education.
- RDWG focuses on large scale research-based data challenges such as very large scale data sets, distributed data sets, multi-institutional data collections and novel analysis and visualization approaches.

Procedures and Practices

- Meetings every two weeks.
- Provide consulting services and discussion forum for new projects and proposals.
- Catalog data needs, architectures, successes and failures of existing projects. Goal is to establish a set of design patterns for management of large amounts of scientific data.
- Maintain an archive of NSF style data management plans to assist proposal writers.

The Data Working Group and iRODS

- A close collaborative relationship between RENCI and the DICE Center.
- Arcot Rajasekar and Reagan Moore are RDWG members and regular contributors.
- We have several projects with iRODS involved:
- National Climatic Data Center: Next Few Slides.
- RENCI Sequencing Initiative: Charles Schmitt.

National Climatic Data Center Project

- NCDC is in Asheville, NC. Worlds largest archive of weather data. Some data is over 150 years old and there is data collected by Thomas Jefferson and Benjamin Franklin.
- One of the data sets is an archive of radar precipitation estimates.
- RENCi and NCDC are collaborating on a pilot program produce a repeatable scalable workflow with this data set.
- Project has a computational component and a data management component.

National Climatic Data Center Project

- Computation occurs at RENCI on our Blue Ridge cluster.
- Combines 9 overlapping precipitation estimates to produce a single mosaic estimate. Period of the study is 10 years.
- Radar mosaic is augmented with “truth on the ground” to produce a high resolution gridded data set. Result set is known as “Q2”. Must be returned to NCDC, but is small compared to the input data.
- So what’s the problem?

National Climatic Data Center Project

- RENCI wants to save copy of Q2 and share it with other collaborators.
- Input data for calculation is low 10's of Tb's.
- Input data is not at RENCI: it's behind a firewall at NCDC.
- The computation is not one calculation: it's hundreds to thousands of “embarrassingly parallel” tasks. Easily separated without much interdependency.
- Too many jobs to launch at once and too much data to move at once.
- Can iRODS help?

National Climatic Data Center Project

- Saving Q2 and sharing is easy. Replication and federation.
- First usage so far is data transfer. iRODS data transfer using iput is much faster than scp. NCDC uses iRODS client to the iren data grid at RENCI.
- scp: 2.8 MB/s
- iput: 32.8 MB/s
- Big improvement! Fast enough?

National Climatic Data Center Project

- Naïve case: transfer all the data, then run all the jobs. Answer: Nope, still not fast enough. 32.8 MB/s is less than 3 Tb per day. Tie up the network completely for 10 days for 30Tb.
- Still have the problem of overrunning our shared computational queue. There must be a better idea. If only ...

National Climatic Data Center Project

- Tie file transfer and job submission together in iRODS.
- iRODS would estimate download time for input data and remaining run time for job. When these 2 times are equal, iRODS would begin downloading the needed input data. When the data has arrived, iRODS would start the job.
- iRODS could maintain a job queue, to handle this process for multiple concurrent jobs.
- May require iRODS/Globus integration.
- Similar to double/multiple buffering in graphics.

RENCI Sequencing Initiative

- Consists of several RENCi collaborations.
- Deep Sequencing Studies for Stimulant Dependence with Kirk Wilhelmsen (UNC School of Medicine).
- National Institutes of Health Exome Project with Kari North (UNC Epidemiology) and Ethan Lange (UNC Genetics).

Contact information

Howard Lander <howard@renci.org>
Michael Shoffner <shoffner@renci.org>

