# Distributed Data Sharing with PetaShare for Collaborative Research in CyberTools

Ismail Akturk[1,2], Mehmet Balman[1,3], Xinqi Wang[1,3], Tevfik Kosar[1,3], Tyler Barker[1,3,] Erik Schnetter[1,5], Raju Gottumukkala[6], Ramesh Kolluru[6], Somnath Roy[1,4], Sumanta Acharya[1,4], Werner Benger[1]

[1]Center for Computation and Technology, Louisiana State University
[2]Department of Electrical and Computer Engineering, Louisiana State University
[3]Department of Computer Science, Louisiana State University
[4]Department of Mechanical Engineering, Louisiana State University
[5]Department of Physics & Astronomy, Louisiana State University
[6]NIMSAT Institute, University of Louisiana at Lafayette

**Abstract**. The unbounded increase in the size of data generated by scientific applications necessitates collaboration and sharing data among the nation's education and research institutions. Simply purchasing high-capacity, high-performance storage systems and adding them to the existing infrastructure of the collaborating institutions does not solve the underlying and highly challenging data handling problems. Scientists are compelled to spend a great amount of time and energy on solving basic data-handling issues, such as how to find out the physical location of data, how to access it, and/or how to move it to visualization and/or compute resources for further analysis. The NSF funded PetaShare project aims to enable transparent handling of underlying data sharing, archival, and retrieval mechanisms for a wide range of applications including several CyberTools science drivers. The ultimate goal is to enable scientists to focus on their primary research problems and collaborate with their colleagues, assured that the underlying infrastructure will manage the low-level data handling issues.

**Key words**: storage management, data management, collaborative research, CyberTools, PetaShare.

## 1 Introduction

The NSF funded PetaShare[1] project aims to enable transparent handling of underlying data sharing, archival and retrieval mechanisms, and make data available to scientists for analysis and visualization on demand. The goal is to enable scientists to focus on their primary research problems, assured that the underlying infrastructure will manage the low-level data handling issues such as data migration, replication, synchronization, data-coherence, and metadata management.

_____
[1] http://www.petashare.org

Unlike existing approaches, PetaShare treats data storage resources and the tasks related to data access as first class entities just like computational resources and compute tasks, and not simply the side effect of computation. Besides of providing globally unified name space, PetaShare also aims to provide sufficient metadata service to the collaborative research community. The key technologies which are developed as part of PetaShare include data-aware storage systems and data-aware schedulers [1].

PetaShare has been deployed at seven Louisiana research institutions connected with high speed LONI[2] network. These institutions include Louisiana State University, Tulane University, LSU Health Sciences Center in New Orleans, LSU Shreveport, University of New Orleans, University of Louisiana at Lafayette, and Louisiana Tech University. PetaShare manages approximately 300 TB of disk storage distributed across these institutions as well as 400 TB tape storage centrally located in downtown Baton Rouge.

## 2 Research Areas and Applications Supported by PetaShare

For collaborative scientific applications, the archival and sharing of data are basic components of data management. While in the past data had been shared through floppy discs, CD-ROM's or DVD, such approaches are no longer feasible with data of state-of-the art simulations, since those easily spread several dozens of gigabytes at least, reaching into the terabyte-range and beyond with current demands on simulation accuracy. PetaShare provides a convenient way of archiving and sharing of data in distributed storage resources by providing global name space that enables transparent access to data such that users are not required to know or remember the physical location of data.

PetaShare allows collaboration of scientists who are performing multidisciplinary research that spans various fields of engineering and basic sciences. Application areas supported by PetaShare include numerical relativity, incident management and interdependency analysis, fluid mechanics, coastal and environmental modeling, geospatial analysis, bioinformatics, medical imaging, petroleum engineering, and high energy physics.

### 2.1 Experimental and Computational Fluid Mechanics

Having efficient mixers and better design of stirred tank reactors (STRs) requires a detailed understanding of the associated flow behavior which involves identification of large-scale mixing structures and the dynamics of their growth and dispersion with the inherent instabilities of the STR flows. Mixing inside a stirred tank is a dynamic process which involves convection, stretching and twisting of fluid elements over the entire operational time. There a number of periodic and random fluctuations through-out the flow domain contributing to the whole mixing process.

_____
[2] http://www.loni.org

The dynamic description of the flow field can be obtained through visualization of the unsteady flow-field and quantification of the mixedness during the entire operational time. The visualization challenges need accessing the entire set of time-dependent data [2] and computing pathlines of tracer particles. The current dataset contains flow information over 3.6 million grid points distributed in 2088 blocks for a time-span of 100 revolutions of the impeller blade. Therefore, PetaShare is being used as a storage resource of this dataset generated on LONI Linux clusters which can be accessed by several researchers inside and outside of the campus for analysis using new visualization tools. LONI allows fast communication between the visualization servers and PetaShare storage that is essential for efficient sharing and processing of dataset of interest.

## 2.2 Numerical Relativity

The numerical relativity group at Louisiana State University is using PetaShare to store data of simulations and associated metadata. PetaShare is being used for two reasons: PetaShare provides a large amount of storage that is readily available to the researchers, and PetaShare is, in principle, available from all LONI machines as well as workstations and laptops.

In addition to provided storage space, associating metadata with the data stored in PetaShare is another important concern. To exploit the metadata management infrastructure of PetaShare, a scheme is being developed, called PetarXiv, for storing simulation data and metadata in PetaShare, creating a repository shared between all numerical relativity group members.

PetarXiv allows scientists to archive data from simulations, making use of the extensive storage space in PetaShare. After a simulation is finished, output data can be uploaded to PetaShare. Along with output data, metadata is stored, which contains more information to uniquely identify the output. To retrieve previously archived data, the metadata attributes can be used as search queries to locate and identify the data. These files can then be retrieved from PetaShare to the users' machines or other computing resources for further analysis.

## 2.3 Incident Management and Interdependency Analysis

The National Incident Management and Advanced Technologies Institute[3] conducts research in development to improve disaster response primarily in three areas: Critical infrastructure interdependency analysis, government-industry partnerships, decision support tools for emergency response.

The decision support tools rely extensively on urgent, reliable and secure access to potentially terabytes of heterogeneous data (in the form of raster and vector geospatial data, and multimedia) ranging from geospatial imagery, LIDAR, DEMs, databases of critical infrastructures, public and private infrastructure, demographics, recent and historical hazard data.

[3] http://www.nimsat.org/

The heterogeneous data has to be accessed by applications running on LONI systems and multiple users on their desktops for geospatial analysis and visualization. PetaShare will help the group in various aspects of data management including providing high availability storage, urgent/prioritized access to data during emergencies and the ability to allocate data for computing, visualization and storage.

**2.4 Scientific Visualization**

PetaShare provides a convenient solution as it does not only provide storage space to keep such large data, but also provides light-weighted client tools that have zero learning curve that allows sharing data among all involved cooperation partners. Instead of passing around CD-ROMs, one can therefore easily upload data into PetaShare, a huge improvement as all data will be immediately available at all involved institutions and new data can be analyzed as soon as the application scientist makes them available.

Current simulation data under investigation stem from diverse application domains, in particular computational fluid dynamics from mechanical engineering, astrophysical data from cosmological evolutions of colliding galaxies, simulations of merging black holes and neutron stars from numerical relativity as well as observational and computational data sets describing hurricanes and cyclones in the gulf of Mexico and Southern Pacific.

## 3 Conclusions

PetaShare provides an infrastructure for transparent handling of underlying data sharing, archival, and retrieval mechanisms for a wide range of applications and researchers including several CyberTools science drivers. It enables researchers and science drivers to focus on their primary research problems, assured that the underlying infrastructure will manage the low-level data handling issues. PetaShare also acts as a bridge between different science drivers as well between individual work packages in CyberTools by allowing them easily share data and by acting as a backbone in the dataflow of the end-to-end application scenarios.

## References

w
1. Kosar, T., and Balman, M. A new paradigm: Data-aware scheduling in grid computing. *Future Generation Computer Systems In Press, DOI: 10.1016/j.future.2008.09.006* .
2. Harhad, F., Roy, S., Acharya, S., and Benger, W. Visualization Challenges in Stirred Tank Fluid Simulations. *5-th High-end Visualization Workshop, March2009, Louisiana State University*