# Topic 18
# Parallel I/O and Storage Technology

Peter Brezany, Marianne Winslett, Denis A. Nicole, and Toni Cortes

Topic Chairpersons

## Introduction

Input and output (I/O) is a major performance bottleneck for large-scale scientific applications running on parallel platforms. For example, it is not uncommon that performance of carefully tuned parallel programs can slow dramatically when they read or write files. This is because many parallel applications need to access large amounts of data, and although great advances have been made in the CPU and communication performance of parallel machines, similar advances have not been made in their I/O performance. The densities and capacities of disks have increased significantly, but improvement in performance of individual disks has not followed the same pace. For parallel computers to be truly usable for solving real, large-scale problems, the I/O performance must be scalable and balanced with respect to the CPU and communication performance of the system. Parallel I/O techniques can help to solve this problem by creating multiple data paths between memory and disks. However, simply adding disk drives to an I/O system without considering the overall software design will improve performance only marginally.

The parallel I/O and storage research community is pursuing solutions in several different areas in order to solve the problem. Active areas of research include disk arrays, network-attached storage, parallel and distributed file systems, theory and algorithms, compiler and language support for I/O, runtime libraries, reliability and fault tolerance, large-scale scientific data management, database and multimedia I/O, real-time I/O, and tertiary storage. The MPI-IO interface, defined by the MPI Forum as part of the MPI-2 standard, aims to provide a standard, portable API that enables implementations to deliver high I/O performance to parallel applications.

The Parallel I/O Archive at Dartmouth, http://www.cs.dartmouth.edu/pario is an excellent resource for further information on the subject. It has a comprehensive bibliography and links to various I/O projects.

## Papers in this Track

This year nine papers, from three continents and five countries, were submitted. All papers were reviewed by three or more referees. Using the referees' reports guidelines, the program committee picked three papers for publication and presentation at the conference. These were presented in one session.

The first paper, by Hakan Ferhatosmanoglu, Divyakant Agrawal and Amr El Abbadi, explores optimal partitioning techniques for data stored in large spatial databases for different types of queries, and develops multi-disk allocation techniques that maximize the degree of I/O parallelism obtained during the retrieval. The authors show that hexagonal partitioning has optimal I/O cost for circular queries compared to all possible non-overlapping partitioning techniques that use convex regions. The second paper, by Xavier Molero, Federico Silla, Vicente Santonja and José Duato, proposes several strategies for dealing with short control messages and analyzes their impact on the performance of storage area networks. This analysis is carried out for a fully adaptive routing algorithm in the context of different network topology environments. The third paper, by Jonathan Ilroy, Cyrille Randriamaro and Gil Utard, deals with improving the performance of the MPI-IO implementation running on top of the Parallel Virtual File System (PVFS). Their optimization mainly focuses on the collective I/O functionality of MPI-IO.