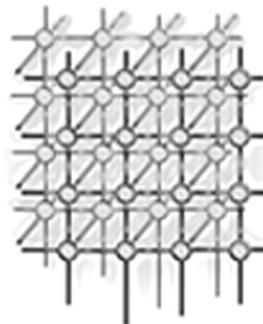


# Global-scale Distributed I/O with ParaMEDIC



P. Balaji<sup>\*†,‡</sup>, W. Feng<sup>§</sup>, H. Lin<sup>§</sup>, J. Archuleta<sup>§</sup>,  
S. Matsuoka<sup>||</sup>, A. Warren<sup>¶</sup>, J. Setubal<sup>¶</sup>, E. Lusk<sup>‡</sup>,  
R. Thakur<sup>‡</sup>, I. Foster<sup>‡,\*\*</sup>, D. S. Katz<sup>\*\*††</sup>, S. Jha<sup>††</sup>,  
K. Shinpaugh<sup>§</sup>, S. Coghlan<sup>‡</sup> and D. Reed<sup>‡‡</sup>

<sup>‡</sup>*Argonne National Laboratory*

<sup>§</sup>*Department of Computer Science, Virginia Tech*

<sup>¶</sup>*Virginia Bioinformatics Institute, Virginia Tech*

<sup>||</sup>*Department of Mathematical and Computing Science, Tokyo Institute of Technology*

<sup>\*\*</sup>*Computation Institute, University of Chicago & Argonne National Laboratory*

<sup>††</sup>*Center for Computation & Technology, Louisiana State University*

<sup>‡‡</sup>*Extreme Computing Group, Microsoft Research*

---

## SUMMARY

Achieving high performance for distributed I/O on a wide-area network continues to be an elusive holy grail. Despite enhancements in network hardware as well as software stacks, achieving high-performance remains a challenge. In this paper, our worldwide team took a completely new and non-traditional approach to distributed I/O, called *ParaMEDIC: Parallel Metadata Environment for Distributed I/O and Computing*, by utilizing application-specific *transformation* of data to orders of magnitude smaller metadata before performing the actual I/O. Specifically, this paper details our experiences in deploying a large-scale system to facilitate the discovery of missing genes and constructing a genome similarity tree by encapsulating the mpiBLAST sequence-search algorithm into ParaMEDIC. The overall project involved nine computational sites spread across the U.S. and generated more than a petabyte of data that was “teleported” to a large-scale facility in Tokyo for storage.

KEY WORDS: Distributed I/O, Bioinformatics, BLAST, Grid computing, Cluster computing

---

\*Correspondence to: Argonne National Laboratory, 9700 S. Cass Avenue, Argonne, IL 60439

†E-mail: balaji@mcs.anl.gov



## 1. Introduction

With the rapid growth in the scale and complexity of scientific applications over the past few decades, the requirements for compute, memory, and storage resources are now greater than ever before. With the onset of petascale and exascale computing, issues related to managing such grand-scale resources, particularly related to data I/O, need to be carefully studied. For example, applications including genomic sequence search and the emergent field of metagenomics, large-scale data mining, data visual analytics, and communication profiling on ultrascale parallel computing platforms generate massive amounts of data that needs to be managed for later processing or archival.

Adding to the complexity of this problem is the issue of resource locality. While system sizes have certainly grown over the past few years, most researchers do not have local access to systems of the scale required by their applications. Therefore, researchers access such large systems remotely to perform the required computations and move the generated data to their local systems after the computation is complete. Similarly, many applications tend to require multiple resources simultaneously for efficient execution. For example, applications that perform large computations and generate massive amounts of output data are becoming increasingly common. While several large-scale supercomputers provide either the required compute power or the storage resources, very few provide both. Thus, data generated at one site often has to be moved to a different site for storage and/or analysis.

In order to alleviate issues related to moving such massive data across sites, considerable monetary and intellectual investment has been put into high-speed distributed network connectivity [35, 5, 2]. However, the utility of these investments is limited in the light of three primary observations: (1) such infrastructure is scarce and does not provide end-to-end connectivity to a very high percentage of the scientific community, (2) the amount of data generated by many applications is so large that even at 100% network efficiency, the I/O time for these applications can significantly dominate their overall execution time, and (3) based on recent trends and published results, existing distributed I/O mechanisms have not been able to achieve a very high network utilization for “real data” on high-speed distributed networks, particularly for single-stream data transfers [3, 38].

To resolve such issues on a *global* scale, we proposed a new, non-traditional approach for distributed I/O known as *ParaMEDIC* (Parallel Metadata Environment for Distributed I/O and Computing) [11, 12, 13]. ParaMEDIC uses application-specific semantic information to process the data generated by treating it as a collection of high-level abstract objects, rather



than as a generic byte-stream. It uses such information to transform the data into orders of magnitude smaller metadata before transporting it over the distributed environment and regenerating it at the target site. All data transformation, movement, and regeneration are done while the application is executing, giving the illusion of an ultrafast teleportation device for large-scale data over distributed environments.

At a high level, ParaMEDIC is similar to standard compression algorithms. However, the term “compression” typically has the connotation that the data is processed as a generic byte-stream. Since ParaMEDIC uses a more abstract application-specific representation of the data to achieve a **much** larger reduction in the data size, we use the terminology of “metadata transformation” in this case.

Because ParaMEDIC utilizes application semantics to generate metadata, it loses some portability compared to traditional byte-stream-based distributed I/O. For example, an instance of ParaMEDIC’s metadata transformation in the context of the mpiBLAST sequence search application is described in Section 3.2. By giving up some portability, however, ParaMEDIC can potentially attain tremendous savings in the amount of actual distributed I/O performed, consequently resulting in substantial performance gains. Further, through the use of a generic framework with an application plug-in model, different applications can use the overall framework in an easy and flexible manner.

In this paper, we demonstrate how we used ParaMEDIC to tackle two large-scale computational biology problems—discovering missing genes and adding structure to genetic sequence databases—on a worldwide supercomputer [14]. The overall worldwide supercomputer comprised nine different supercomputers distributed at seven sites across the U.S. and one large-scale storage facility located in Japan. The overall experiment consisted of sequence searching the entire microbial genome database against itself, generating approximately a petabyte of data that was transported to Tokyo for storage. We present several insights gained from this large-scale run, which will be valuable to other researchers performing such large, global-scale distributed computation and I/O.

## ***2. Large-Scale Computational Biology: A Peek at Compute and Storage Requirements***

In this section we discuss different aspects of computational biology with a focus on the compute and storage requirements of large-scale applications in this domain.

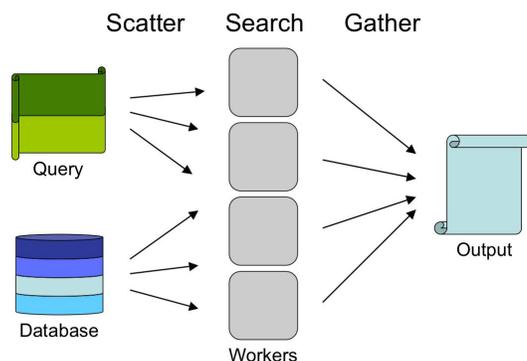


Figure 1: High-level Algorithm of mpiBLAST

### 2.1. Sequence Searching

With the advent of rapid DNA sequencing, the amount of genetic sequence data available to researchers has increased exponentially [9]. The GenBank database, a comprehensive database that contains genetic sequence data for more than 260,000 named organisms, has exhibited exponential growth since its inception over 25 years ago [15]. This information is available for researchers to search new sequences against and infer homologous relationships between sequences or organisms. This helps a wide range of projects, from assembling the Tree of Life [20] to pathogen detection [22] and metagenomics [23].

Unfortunately, the exponential growth of sequence databases necessitates faster search algorithms to sustain reasonable search times. The Basic Local Alignment Search Tool (BLAST), the de facto standard for sequence searching, uses heuristics to prune the search space and decrease search time with an accepted loss in accuracy [7, 8]. mpiBLAST parallelizes BLAST using several techniques including database fragmentation, query segmentation [18], parallel input-output [29], and advanced scheduling [40]. As shown in Figure 1, mpiBLAST uses a *master-worker* model and performs a *scatter-search-gather-output* execution flow. During the scatter, the master splits the database and query into multiple pieces and distributes them among worker nodes. Each worker then searches the query segment against the database fragment that it was assigned. The results are gathered by the master, formatted, and output to the user. Depending on the size of the query and that of the database, such output generated can be large. For example, as shown in Table I, an all-to-all search of the nucleotide database can generate as much as 30 TB of data. Thus, for environments with limited I/O capabilities, such as distributed systems, the output step can cause significant overheads.



## 2.2. *Discovering Missing Genes*

Genome annotation is the process of associating information with a genomic sequence. Part of this process includes determining (by computational analysis) the location and structure of protein-encoding and RNA-encoding genes, also known as making gene calls. It is important to be as accurate and as sensitive as possible in making gene calls: avoiding false positives and missing real genes. Gene prediction in prokaryotes (bacteria and archaea) typically involves

evaluating the coding potential of genomic segments that are delimited by conserved nucleotide motifs. The most widely used gene-finding programs [19, 30] build a sequence model based on statistical properties of genes known to be (very likely) real genes in a given genome. This model is then used to evaluate the likelihood that an individual segment codes for a gene; using this method, some genes with anomalous composition are almost always missed. Another popular method for locating genes is to compare genomic segments with a database of gene sequences found in similar organisms. If the sequence is conserved, the segment being evaluated is likely to be a coding gene (this is the “similarity method”). Genes that do not fit a given genomic pattern and do not have similar sequences in current annotation databases may be systemically missed.

One way to detect missed genes is to use the similarity method and compare raw genomes against each other, as opposed to comparing a raw genome to a database of known genes [41]. If gene  $a$  in genome  $A$  and gene  $b$  in genome  $B$  have been missed and  $a$  is similar to  $b$ , then this method will find both. However, this involves performing an all-to-all comparison of the entire database against itself (in our case study, the entire microbial genome database against itself). This task is heavily compute and data intensive, requiring thousands of compute processors and generating on the order of a petabyte of output data that needs to be stored for processing.

## 2.3. *Adding Structure to Genetic Sequence Databases*

One of the major issues with sequence search is the structure of the sequence database itself. Currently, these databases are unstructured and stored as a flat file, and each new sequence that is discovered is simply appended to the end of the file. Without more intelligent structuring,

Table I: Estimated Output of an All-to-All NT Search

Query Size (KB)	Number of Queries	Estimated Output (GB)
0-5	3,305,170	1,139
5-50	87,506	593
50-150	25,960	23,555
150-200	26,524	3,995
>200	9,840	<i>Not run</i>
Total	3,455,000	>29,282



a query sequence has to be compared to every sequence in the database, forcing the best case to take just as long as the worst case. By organizing and providing structure to the database, searches can be performed more efficiently by discarding irrelevant portions entirely.

One way to structure the sequence database is to create a sequence similarity tree where “similar” sequences are closer together in the tree than dissimilar sequences. The connections in the tree are created by determining how “similar” the sequences are to each other through sequence searches. To create *every* connection, however, the entire database has to be searched against itself, resulting in an output size of  $N^2$  values (where  $N$  is the number of sequences in the database).

### 3. Overview of ParaMEDIC-enhanced mpiBLAST

In our previous work [12, 13], we provided a detailed description of ParaMEDIC. Here we present a brief summary of that work.

#### 3.1. The ParaMEDIC Framework

ParaMEDIC provides a framework for *decoupling* computation and I/O in applications relying on large quantities of both. Specifically, it does not hinder application computation. As the output data is generated, however, the framework differs from traditional distributed I/O in that it uses application-semantic information to process the data generated by treating it as a collection of high-level application-specific objects rather than as a generic byte-stream. It uses such information to transform the data into orders of magnitude smaller metadata before transporting it over the distributed environment and regenerating the original data at the target site.

As shown in Figure 2, ParaMEDIC provides several capabilities, including support for data encryption and integrity as well as data transfer in distributed environments (either directly via TCP/IP communication or through global file-systems). However, the primary semantics-based metadata creation is done by the *application plug-ins*. Most application plug-ins are specific to each application and thus rely on knowledge of application semantics. These plug-ins provide two functionalities: (1) processing output data generated by the application to create metadata and (2) converting metadata back to the final output. Together with application-specific plug-ins, ParaMEDIC also provides application-independent components such as data compression, data integrity, and data encryption. These can be used in conjunction with the application-specific plug-ins or independently.

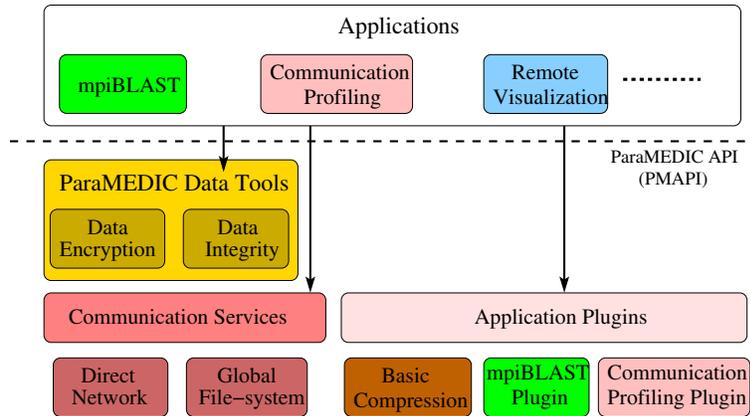


Figure 2: ParaMEDIC Architecture

**Trading Computation with I/O:** The amount of computation required in ParaMEDIC is higher than what is required by the original application. After the output is generated by the application processes, it has to be further processed to generate the metadata, sent to the storage site, and processed yet again to regenerate the final output. However, the I/O cost achieved can potentially be significantly reduced by using this framework. In other words, ParaMEDIC trades (a small amount of) additional computation for (potentially large) reduction in I/O cost. With respect to the additional computational cost incurred, ParaMEDIC is quite generic with respect to the metadata processing required by the different processes. For many applications, it is possible to tune the amount of post-processing performed on the output data, with the trend being, the more the post-processing computation, the better the reduction in the metadata size. That is, an application plug-in can perform more processing of the output data to reduce the I/O cost.

### 3.2. Integration with *mpiBLAST*

In a *cluster* environment, most of the *mpiBLAST* execution time is spent on the search itself, since the BLAST string-matching algorithm is computationally intensive. In comparison, the cost of formatting and writing the results is minimal, especially when many advanced clusters are configured with high-performance parallel file-systems. In a *distributed* environment, however, the output typically needs to be written over a wide-area network to a remote file-system. Hence, the cost of writing the results can easily dominate the execution profile of *mpiBLAST* and become a severe performance bottleneck. By replacing the traditional

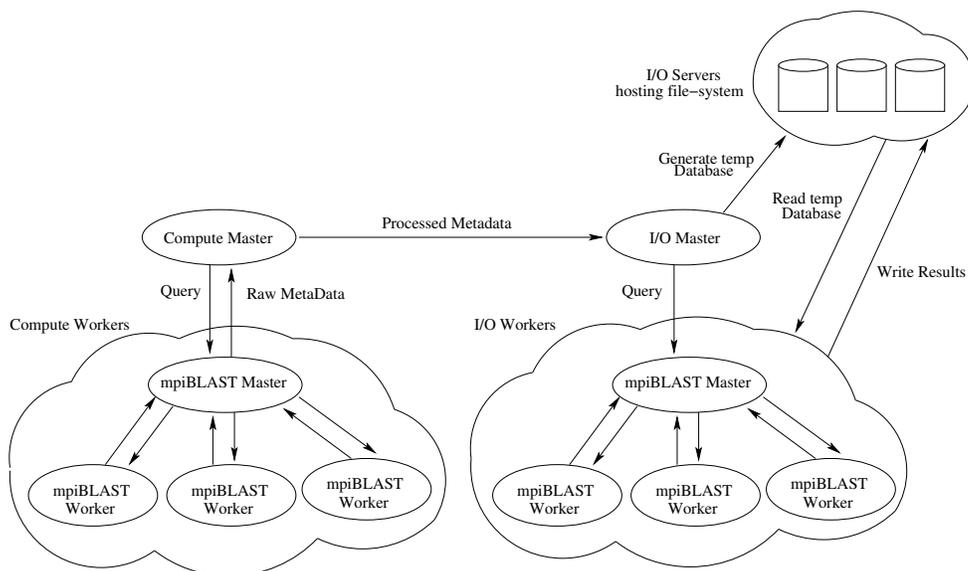


Figure 3: ParaMEDIC and mpiBLAST Integration

distributed I/O framework with ParaMEDIC (as shown at the top of Figure 3), we can provide large reduction in the amount of data communication performed. For example, as we will see in Section 5, a mpiBLAST-specific instance of ParaMEDIC reduces the volume of data written across a wide-area network by *more than two orders of magnitude*.

Figure 3 depicts how mpiBLAST is integrated with ParaMEDIC. First, on the compute site (the left cloud in Figure 3), once the output is generated by mpiBLAST, the mpiBLAST application plug-in for ParaMEDIC processes this output to generate orders of magnitude lesser metadata. Specifically, the output of mpiBLAST consists of alignment information and scores corresponding to the top matches it found for each sequence in the entire database. Thus, while the search time largely depends on the size of the database, once the search is complete, the output only depends on how closely the input query sequence matches the top matching sequences in the database. Based on this observation, the metadata basically contains information identifying the top matching sequences in the database, and not the other alignment or score information.

ParaMEDIC transfers this metadata to the storage site. At the storage site, a temporary (and much smaller) database that contains only the top matching sequences is created by extracting the corresponding sequence data from a local database replica. ParaMEDIC then



reruns mpiBLAST at the storage site by taking as input the same query sequence and the temporary database to generate and write the complete output (including the alignments and scores) to the storage system. The overhead in rerunning mpiBLAST at the storage site is small, since the temporary database that is searched is substantially smaller, with only about 250 sequences in it, compared to the several millions of sequences in large DNA databases.

#### 4. *ParaMEDIC on a Worldwide Supercomputer*

To accommodate the compute and storage requirements of the computational biology applications discussed in Section 2, we utilize a worldwide supercomputer that, in aggregate, provides the required compute power and storage resources. The worldwide supercomputer comprises nine high-performance computing systems at seven different sites across the U.S. and a large-scale storage facility in Japan, to create a single high-performance *distributed* computing system. The specifics of each individual system are in Table II. In the next subsections, we address the issues with working on such a large-scale distributed system that are not immediately apparent on smaller-scale systems.

##### 4.1. *Dealing with Dynamic Availability of Compute Clients and Other Faults*

Several systems in our worldwide supercomputer operate in batch-mode. Users submit jobs to system queues and are scheduled to execute on the available resources. That is, compute resources are not available in a dedicated manner but become available when our job is scheduled for execution and become unavailable when our job terminates.

To handle this issue, we segment the overall work to be performed into small tasks that can be performed independently (i.e., sequentially, concurrently, or out-of-order). The management of tasks is done by a centralized server running on a dedicated resource. As each job is executed, it contacts this server for the next task, computes the task, transforms the output to metadata, and transmits the metadata to the storage site. This approach has two benefits. First, the clients are completely stateless. That is, if a client job terminates before it has finished its computation or metadata transmission to the storage site, the servers handle this failure by reassigning the task to a different compute client. The second advantage is if the metadata corresponding to a task that is either not received completely or is corrupted, the server just discards the data and reassigns the task to another compute node. Thus, I/O failures are never catastrophic.



Table II: Specification of the Different Systems that Formed our Worldwide Supercomputer

Name	Location	Cores	Architecture	Memory (GB)	Network	Storage (TB)	Distance from Storage (km)
SystemX	Virginia Tech	2200	IBM PowerPC 970FX	4	InfiniBand (8Gbps)	NFS (30)	11,000
Breadboard Compute Cluster	Argonne National Lab	128	AMD Opteron	4	10G Ethernet (10Gbps)	NFS (5)	10,000
BlueGene/L Supercomputer	Argonne National Lab	2048	IBM PowerPC 440	1	Proprietary (1.4Gbps)	PVFS (14)	10,000
SiCortex Supercomputer	Argonne National Lab	5832	SiCortex SC5832	3	Proprietary (11.2Gbps)	NFS (4)	10,000
Jazz Compute Cluster	Argonne National Lab	700	Intel Xeon	1-2	1G Ethernet (1Gbps)	GFS (10) PVFS (10)	10,000
NSF TeraGrid Cluster	University of Chicago	320	Intel Itanium2	4	Myrinet 2000 (2Gbps)	NFS (4)	10,000
NSF TeraGrid Cluster	San Diego Super-computer Center	60	Intel Itanium2	4	Myrinet 2000 (2Gbps)	GPFS (50)	9,000
Oliver Compute Cluster	Louisiana State University, Lafayette	512	Intel Xeon	4	InfiniBand (8Gbps)	Lustre (12)	11,000
Open Science Grid	United States	200	AMD Opteron Intel Xeon	1-2	1G Ethernet (1Gbps)	-	11,000
TSUBAME Storage Cluster	Tokyo Institute of Technology	72	AMD Opteron	16	1G Ethernet (1Gbps)	Lustre (350)	0



---

#### 4.2. *Architectural Heterogeneity*

One of the key impediments to large-scale distributed systems is system heterogeneity. Many distributed systems, such as the one used in this paper, cannot obtain a homogeneous environment in either hardware or software, and efficient use of the system requires overcoming this obstacle. The worldwide supercomputer used in this paper contains six different processor architectures (IBM PowerPC 970FX, IBM PowerPC 440, AMD Opteron, SiCortex MIPS64, Intel Xeon, and Intel Itanium2), five different network interconnects (Gigabit Ethernet, 10-Gigabit Ethernet, InfiniBand, IBM proprietary 3D toroidal network, and SiCortex Kautz graph), and eight variations of the Linux operating system.

In order to deal with this issue, all data being transferred over the network has to be converted to an architecture-independent format. Since the total amount of data that is generated and must be moved to the storage site is enormous, this can have a significant impact on traditional distributed I/O. However, with ParaMEDIC, only metadata generated by processing the actual output is transferred across the wire. Since this metadata is orders of magnitude smaller as compared to the actual output, such byte manipulation to deal with heterogeneity has minimal impact on overall performance.

#### 4.3. *Utilizing the Parallelism in Compute Nodes*

In traditional file I/O, there are two levels of parallelism. First, multiple I/O servers are aggregated into a parallel file-system to take advantage of the aggregate storage bandwidth of these servers. Second, multiple compute clients, that process different tasks, write data to such file-systems in parallel as well. Most parallel file-systems are optimized for such access to give the best performance.

With ParaMEDIC, there are three I/O components: (1) compute clients that perform I/O, (2) post-processing servers that process the metadata to regenerate the original output, and (3) I/O servers that host the file-system. Similar to the traditional I/O model, the first and third components are already parallelized. That is, multiple streams of data being written in parallel by different compute clients and the I/O servers parallelize each stream of data that is being written to them. However, in order to achieve the best performance, it is important that the second component, post-processing servers, be parallelized as well.

Parallelizing the post-processing servers adds its own overhead and complexity mainly with respect to synchronization between the different parallel processes. To avoid this, we use an embarrassingly parallel approach for these servers. Each incoming stream of data is allocated to a separate process till a maximum number of processes is reached, after which the incoming data requests are queued till a process becomes available again. Thus, different processes do not

---



have to share any information and can proceed independently. The advantage of this approach is its simplicity and lack of synchronization required between different parallel post-processing servers. The disadvantage, however, is that the number of data streams generated from the post-processing servers is equal to the number of incoming data streams. That is, if only two tasks are active at one time, only two streams of data are written to the actual storage system. Thus, the performance might not be optimal. However, in most cases, we expect the number of incoming streams to be sufficiently large to not face such performance issues.

#### 4.4. *Handling Large Communication Latencies with Disconnected I/O*

As seen in Table II, the computational sites are between 9,000 and 11,000 kilometers away from the storage site. At these distance, communication latencies are in tens of milliseconds. Such large latencies can severely limit the effectiveness of a synchronous remote file-systems that can be used for distributed I/O, since each synchronization operation has to make round-trip hops on the network. To overcome the bottleneck incurred by such high-latency, our worldwide supercomputer utilizes a lazy asynchronous I/O approach. By caching the output data locally before performing the actual output, clients can perform their computations while disconnected from the remote file-system. After a substantial amount of data is generated a bulk transfer of the metadata occurs, thereby maximally utilizing the bandwidth available between the sites and mitigating the effect of high-latency.

An issue with this approach of disconnected computation is fault tolerance. Once a task is assigned to a compute client, the server is completely disconnected from this client. After the computation is complete, the client reconnects and sends the generated metadata. Although, this two-phase synchronization model is more error prone and requires additional state information in the server, it makes the compute clients *truly stateless* even when the actual computation is going on.

### 5. *Experiments, Measurements and Analysis*

In this paper, we use ParaMEDIC to search the entire microbial genome database against itself. Several supercomputing centers within the U.S. perform the computation, while the data generated is stored at a large storage resource in Tokyo. This section compares the performance of ParaMEDIC with vanilla mpiBLAST with respect to the amount of data communicated and the I/O time taken, as well as the storage bandwidth utilization.



### 5.1. Data I/O Overheads

Figure 4(left) illustrates the amount of data transmitted between the compute and the storage sites for different number of post-processing threads, and Figure 4(right) shows the factor of reduction in the amount of data. Each post-processing thread processes one segment of data that has the output for 10,000 query sequences. Most segments have approximately similar output sizes, so the amount of data communicated over the distributed network increases linearly with the number of segments, and hence the number of post-processing threads. ParaMEDIC, on the other hand, processes the generated data and converts it into metadata before performing the actual transfer. Thus, the actual data that is transferred over the network is much smaller. For example, with 288 threads, mpiBLAST communicates about 100 GB of data, while ParaMEDIC only communicates about 108 MB—a *900-fold* reduction.

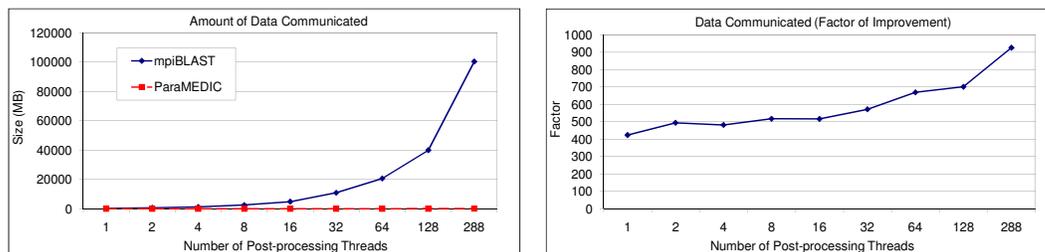


Figure 4: Data I/O Overheads: (left) Total Amount of Data Communicated and (right) Factor of Improvement

We also illustrate the I/O time in Figure 5. As shown, ParaMEDIC outperforms mpiBLAST by more than two orders of magnitude. This result is attributed to multiple aspects. First, given the shared network connection between the two sites, the achievable network performance is usually much lower than within the cluster. Thus, with mpiBLAST transferring the entire output over this network, its performance would be heavily impacted by the network performance. Second, the distance between the two sites causes the communication latency to be high. Thus, file-system control messages tend to take a long time to be exchanged, resulting in further loss in performance. Third, for mpiBLAST, since the wide-area network is a bottleneck, the number of simultaneously transmitted data streams does not matter;

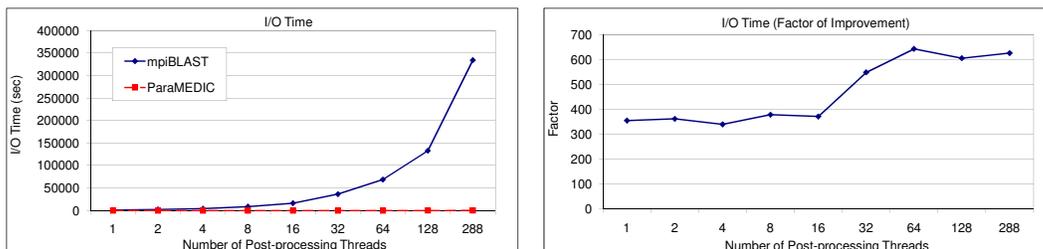


Figure 5: I/O Time Measurements: (left) Total I/O Time and (right) Factor of Improvement

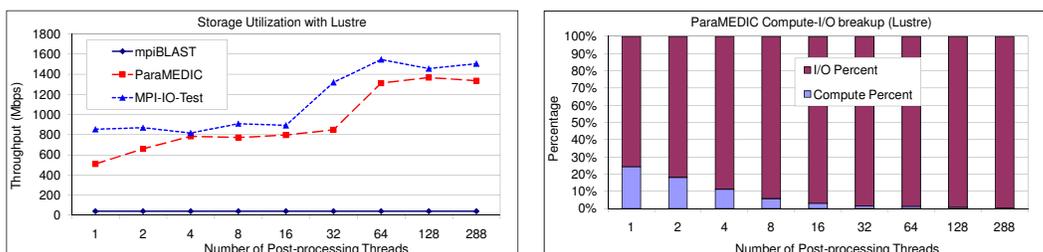


Figure 6: Storage Bandwidth Utilization Using Lustre: (left) Storage Utilization Improvement and (right) Computation and I/O Time

communication is serialized in the network. However, with ParaMEDIC, since the wide-area network is no longer a bottleneck, it can more effectively utilize the parallelism in the data streams to better take advantage of the storage bandwidth available at the storage site, as described in more detail in Section 5.2.

### 5.2. Storage Bandwidth Utilization

Figure 6(left) illustrates the storage bandwidth utilization achieved by mpiBLAST, ParaMEDIC, and the MPI-IO-Test benchmark, which is used as an indication of the peak

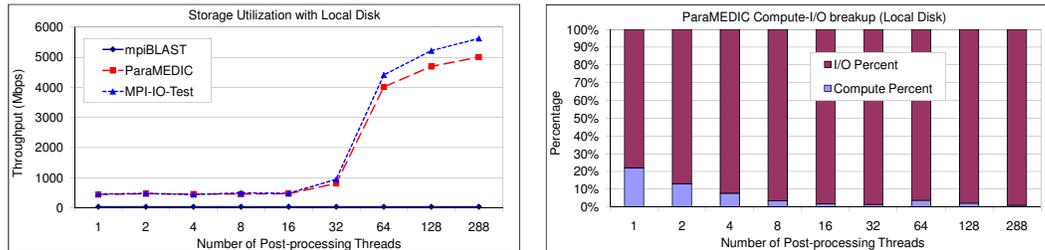


Figure 7: Storage Utilization Using Local File-System: (left) Storage Utilization Improvement and (right) Computation and I/O Time

performance capability of the I/O subsystem. We notice that the storage utilization of mpiBLAST is very poor compared to ParaMEDIC. The reason is that, for mpiBLAST, the I/O is limited by the wide-area network bandwidth. Thus, though more than 10,000 processors are performing the compute part of the task, the network connecting the compute servers in the U.S. and the storage system in Tokyo becomes the bottleneck.

On the other hand, ParaMEDIC uses more than 90% of the storage system capability (shown by MPI-IO-test). When the number of processing threads is low (x-axis in the figure), ParaMEDIC uses about half the storage capability. However, as the number of processing threads increases, the I/O utilization of ParaMEDIC increases as well.

Figure 6(right) illustrates the percentage breakup of the time spent in ParaMEDIC's post-processing phase. A significant portion of the time spent is in the I/O part. This shows that in spite of using a fast parallel file system such as Lustre, ParaMEDIC is still bottlenecked by the I/O subsystem. In fact, our analysis has shown that in this case the bottleneck lies in the 1-gigabit Ethernet network subsystem connecting the storage nodes. Thus, we expect that even for systems with faster I/O subsystems, ParaMEDIC will further scale up and continue to use a significant portion of the I/O bandwidth provided.

In Figure 7(left), we remove the file-system network bottleneck and directly perform I/O on the local nodes. Storage utilization achieved in this case is significantly higher than going over the network. Even in this case, ParaMEDIC completely uses the storage capability with more than 90% efficiency. Figure 7(right) shows the percentage breakup of the time spent. Similar



to the case with the Lustre file-system, a significant portion of the time is still spent on I/O. Thus, again, ParaMEDIC can be expected to scale and fully use even faster storage resources.

## 6. Discussion

Although this paper deals only with enhancing the mpiBLAST application through ParaMEDIC, the idea is relevant for many other applications as well. For example, applications that have natively been built for distributed environments such as SETI@home [37] and other BOINC applications [1] can easily use similar ideas and can benefit aspects in which such techniques are possible. In the field of communication profiling with MPE [4], we have also done some preliminary work that uses metadata transformation of profiled data through ParaMEDIC. Specifically, based on the observation that most scientific applications have a very uniform and periodic communication pattern, we perform a Fourier transform on the data to identify this periodicity and use this as an abstract block. The metadata comprises one complete abstract block and just the differences for all other blocks. Our preliminary numbers in this field have demonstrated between two and five-fold reduction in the I/O time using ParaMEDIC. Work on other application fields including earthquake modeling and remote visualization is ongoing as well, with promising preliminary results.

## 7. Related Work

Efficient I/O access for scientific applications in distributed environments has been an ongoing subject of research for various parallel and distributed file-systems [24, 34, 16, 36]. There has also been work on explicit data transfer protocols such as GridFTP [6]. Other efforts include providing remote data access through MPI-IO [31]. RIO [21] introduced a proof-of-concept implementation that allows applications to access remote files through ROMIO [39]. RFS [27] enhanced the remote write performance with *active buffering*, by optimizing overlap between applications computation and I/O. Studies have also been done in translating MPI-IO calls into operations of lower level data protocols such as Logistic Network [28]. However, all these approaches deal with data as a byte-stream. Conversely, our approach focuses on aggressively reducing the amount of I/O data to be communicated by taking advantage of application semantics and dealing with data as high-level abstract units, rather than a stream of bytes.

Semantic-based data transformation is not new. Several semantic compression algorithms have been investigated in compressing relational databases [25, 10, 26]. Leveraging the table semantics, these algorithms first build descriptive models of the database using data mining techniques such as clustering and then strip out data that can be regenerated. In the



multimedia field, context-based coding techniques (similar to semantics-based approaches) have been widely used in various video compression standards [32, 33, 17]. With aid of context modeling, these techniques efficiently identify redundant information in the media. Although sharing the same goal of reducing data to store or transfer with ParaMEDIC, these data compression studies do not address the remote I/O issue.

Thus, ParaMEDIC utilizes ideas from different fields to provide a novel approach for distributed I/O.

## 8. Conclusion

Rapid growth of computational power is enabling computational biology to tackle increasingly large problems such as discovering missing genes and providing structure to genetic sequence databases. As the problems grow larger, however, so does the data consumed and produced by the applications. For many applications, the required compute power and storage resources cannot be found at a single location, precipitating the transfer of large amounts of data across the wide-area network. ParaMEDIC mitigates this issue by pursuing a non-traditional approach to distributed I/O. By trading computation for I/O, ParaMEDIC utilizes application semantics information to transform the output to orders of magnitude smaller metadata. In this paper, we presented our experiences in solving large-scale computational biology problems by utilizing nine different high-performance compute sites within the U.S. to generate a petabyte of data that then was transferred to a large-scale storage facility in Tokyo using ParaMEDIC's distributed I/O capability. We demonstrated that ParaMEDIC can achieve a performance improvement of several orders of magnitude compared to traditional I/O. In future, we plan to evaluate semantic-based compression for other applications.

## 9. Acknowledgments

We thank the following people for their technical help in managing the large-scale run and other experiments associated with this paper: R. Kettimuthu, M. Papka and J. Insley from the University of Chicago; N. Desai and R. Bradshaw from Argonne National Laboratory; G. Zelenka, J. Lockhart, N. Ramakrishnan, L. Zhang, L. Heath, and C. Ribbens from Virginia Tech; M. Rynge and J. McGee from Renaissance Computing Institute; R. Fukushima, T. Nishikawa, T. Kujiraoka and S. Ihara from Tokyo Institute of Technology; S. Vail, S. Cochrane, C. Kingwood, B. Cauthen, S. See, J. Fragalla, J. Bates, R. Cagle, R. Gaines and C. Bohm from Sun Microsystems; and H. Liu from Louisiana State University/LONI.



---

## 10. *Author Bios*

**Pavan Balaji** holds a joint appointment as an Assistant Computer Scientist at the Argonne National Laboratory and as a research fellow of the Computation Institute at the University of Chicago. He had received his Ph.D. from the Computer Science and Engineering department at the Ohio State University. His research interests include high-speed interconnects, efficient protocol stacks, parallel programming models and middleware for communication and I/O, and job scheduling and resource management.

**Wu-chun Feng** (a.k.a. “Wu”) holds a joint appointment as an Associate Professor in the Department of Computer Science and Department of Electrical & Computer Engineering and as Director of the Synergy Laboratory at Virginia Tech. He also founded and leads “The Green500 List” effort. Previous professional stints include IBM T.J. Watson Research Center, NASA Ames Research Center, Orion Multisystems, Vosaic, University of Illinois at Urbana-Champaign, Purdue University, The Ohio State University, and most recently, Los Alamos National Laboratory. He received his Ph.D. in Computer Science from the University of Illinois at Urbana-Champaign. His research interests include high-performance computing from systems software to middleware to applications software, accelerator-based parallel computing, K-12 pedagogy in computer science, green supercomputing, and bioinformatics.

**Heshan Lin** is a Senior Research Associate at the Department of Computer Science, Virginia Tech. He had received his Ph.D. from the Department of Computer Science at the North Carolina State University. His research interests include data-intensive parallel and distributed computing, parallel I/O, high-performance bioinformatics and cloud computing.

**Jeremy Archuleta** is a graduate student in the Department of Computer Science at Virginia Polytechnic Institute and State University. He received his undergraduate degree in Electrical Engineering and Computer Science at the University of California, Berkeley and a masters degree from the School of Computing at the University of Utah. His research interests encompass a broad range of topics in computer science but focus predominantly in applying high-performance computing and software engineering methodologies towards bioinformatics applications.

**Satoshi Matsuoka** is a full Professor at the Global Scientific Information and Computing Center of Tokyo Institute of Technology. He is the leader of TSUBAME which became the fastest supercomputer in Asia-Pacific in June, 2006. He has also co-lead the Japanese national grid project NAREGI during 2003-2007. Currently he is designing TSUBAME21 likely to become the first Petaflops supercomputer in Japan. He has chaired many ACM/IEEE international conferences, including the technical papers chair role for SC09. He has won many awards including the JSPS Prize from the Japan Society for Promotion of Science in 2006, from his Majesty Prince Akishinomiya.

---



---

**Andrew Warren** is a Ph.D. candidate in the Computer Science department at Virginia Tech and graduate researcher at Virginia Bioinformatics Institute. He received his B.S. in biochemistry and M.S. in Computer Science from Virginia Tech. His research interests include knowledge integration in bioinformatics, genomics and evolutionary analysis, ontology analysis, and data visualization.

**Joao Setubal** is Associate Professor at the Virginia Bioinformatics Institute and at the Department of Computer Science, Virginia Tech. He has a Ph.D. from the Department of Computer Science at the University of Washington. His research interests are in the area of bioinformatics and computational biology, including algorithms and software for genome sequence analysis.

**Ewing “Rusty” Lusk** is director of the Mathematics and Computer Science Division at Argonne National Laboratory and an Argonne Distinguished Fellow. He received his B.A. in mathematics from the University of Notre Dame in 1965 and his Ph.D. in mathematics from the University of Maryland in 1970. He was a professor of computer science at Northern Illinois University before joining Argonne in 1982. His current research interests include programming models for scalable parallel computing, implementation issues for the MPI Message-Passing Interface standard, parallel performance analysis tools, and system software for large-scale machines. He is the author of five books and more than a hundred research articles in mathematics, automated deduction, and parallel computing.

**Rajeev Thakur** is a Computer Scientist in the Mathematics and Computer Science Division at Argonne National Laboratory. He is also a Fellow in the Computation Institute at the University of Chicago and an Adjunct Associate Professor in the Department of Electrical Engineering and Computer Science at Northwestern University. He received a Ph.D. in Computer Engineering from Syracuse University. His research interests are in the area of high-performance computing in general and particularly in parallel programming models and message-passing and I/O libraries.

**Ian Foster** is Director of the Computation Institute, a joint institute of the University of Chicago and Argonne National Laboratory. He is also an Argonne Senior Scientist and Distinguished Fellow, and the Arthur Holly Compton Distinguished Service Professor of Computer Science. Ian received a BSc (Hons I) degree from the University of Canterbury, New Zealand, and a PhD from Imperial College, United Kingdom, both in computer science. His research deals with distributed, parallel, and data-intensive computing technologies, and innovative applications of those technologies to scientific problems in such domains as climate change and biomedicine.

**Daniel S. Katz** is the TeraGrid GIG Director of Science, and holds a joint appointment as a Senior Fellow of the Computation Institute at the University of Chicago and Argonne

---



National Laboratory. He is also an affiliate faculty member of the Center for Computation & Technology and an adjunct associate professor in Electrical and Computer Engineering, at Louisiana State University. He received his Ph.D. from Northwestern University. His research interests include applications and tools for parallel and distributed computing.

**Shantenu Jha** is the Director for Cyberinfrastructure Development at the CCT, and a Research Professor in Computer Science at Louisiana State University (LSU). He is also a theme-leader at the e-Science Institute, Edinburgh and a Visiting Researcher at UC-London. His research interests lie at the triple point of Computational Science, Cyberinfrastructure Development and Computer Science. Shantenu leads the SAGA project and is currently working on co-writing a book on “Abstractions for Distributed Applications and Systems: A Computational Science Perspective” (Wiley).

**Kevin Shinpaugh** is the Director of High Performance Computing at Virginia Tech and is also an Adjunct Professor in Aerospace Engineering. He received his Ph.D. from the Aerospace and Ocean Engineering department at Virginia Tech. His research interests are high performance computing, signal processing optical sensors and spacecraft systems.

**Susan Coghlan** has worked on parallel and distributed computers for over 20 years. She worked at Los Alamos National Laboratory on development of scientific applications and management of ultra-scale supercomputers. In 2000 she co-founded TurboLabs, a research laboratory in Santa Fe, which developed the world’s first dynamic provisioning system for large clusters and data centers. In her current role as Associate Division Director for the Argonne Leadership Computing Facility, she is responsible for the project to procure and deploy the next big system at the ALCF, a 10PF to 20PF Blue Gene/Q.

**Daniel Reed** is the Corporate Vice President for Technology Strategy and Policy and the eXtreme Computing Group (XCG) at Microsoft. Previously, he was the Chancellor’s Eminent Professor at the University of North Carolina, at Chapel Hill, as well as the Director of the Renaissance Computing Institute (RENCI) and the Chancellor’s Senior Advisor for Strategy and Innovation. Before that, he was Gutsell Professor and Head of the Department of Computer Science at the University of Illinois at Urbana-Champaign (UIUC) and Director of the National Center for Supercomputing Applications (NCSA).

## REFERENCES

1. BOINC: Berkeley Open Infrastructure for Network Computing. <http://boinc.berkeley.edu>.
2. CANARIE. <http://www.canarie.ca>.
3. Ethernet Everywhere: Price/Performance as the Key to Cluster & Grid Interconnects. <http://www.cse.scitech.ac.uk/disco/mew14-cd/Talks/Force10-VanCampen.pdf>.
4. MPE: MPI Parallel Environment. <http://www-unix.mcs.anl.gov/perfvis/download/index.htm>.
5. SURFnet. <http://www.surfnet.nl>.



6. B. Allcock, J. Bresnahan, R. Kettimuthu, M. Link, C. Dumitrescu, I. Raicu, and I. Foster. The Globus Striped GridFTP Framework and Server. In *SC*, 2005.
7. S. F. Altschul, W. Gish, W. Miller, E. W. Myers, and D. J. Lipman. Basic local alignment search tool. *J Mol Biol*, 215(3):403–410, October 1990.
8. S. F. Altschul, T. L. Madden, A. A. Schaffer, J. Zhang, Z. Zhang, W. Miller, and D. J. Lipman. Gapped BLAST and PSI-BLAST: a new generation of protein database search programs. *Nucleic Acids Res.*, 25:3389–3402, 1997.
9. S.F. Altschul, M.S. Boguski, W. Gish, and J.C. Wootton. Issues in searching molecular sequence databases. *Nat Genet*, 6(2):119–29, 1994.
10. S. Babu, M. Garofalakis, and R. Rastogi. SPARTAN: a model-based semantic compression system for massive data tables. *SIGMOD Rec.*, 30(2), 2001.
11. P. Balaji, W. Feng, J. Archuleta, and H. Lin. ParaMEDIC: Parallel Metadata Environment for Distributed I/O and Computing. In *IEEE/ACM International Conference for High Performance Computing, Networking, Storage and Analysis (SC Storage Challenge Award)*, Reno, Nevada, Nov. 2007.
12. P. Balaji, W. Feng, J. Archuleta, H. Lin, R. Kettimuthu, R. Thakur, and X. Ma. Semantics-Based Distributed I/O for mpiBLAST. In *Proceedings of the ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming*, Salt Lake City, Utah, February 2008.
13. P. Balaji, W. Feng, and H. Lin. Semantics-Based Distributed I/O with the ParaMEDIC Framework. In *Proceedings of the IEEE/ACM International Conference on High Performance Distributed Computing (HPDC)*, Boston, MA, June 2008.
14. P. Balaji, W. Feng, H. Lin, J. Archuleta, S. Matsuoka, A. Warren, J. Setubal, E. Lusk, R. Thakur, I. Foster, D. S. Katz, S. Jha, K. Shinpaugh, S. Coghlan, and D. Reed. Distributed I/O with ParaMEDIC: Experiences with a Worldwide Supercomputer. In *Proceedings of the International Supercomputing Conference (ISC Best Paper Award)*, Dresden, Germany, 2008.
15. D. A. Benson, I. Karsch-Mizrachi, D. J. Lipman, J. Ostell, and D. L. Wheeler. Genbank. *Nucleic Acids Res*, 36(Database issue), January 2008.
16. P. Carns, W. Ligon III, R. Ross, and R. Thakur. PVFS: A Parallel File System For Linux Clusters. In *LSC*, 2000.
17. T. Wiegand D. Marpe, H. Schwarz. Context-based Adaptive Binary Arithmetic Coding in the H.264/AVC Video Compression Standard. *IEEE Transactions on Circuits and Systems for Video Technology*, 13(7), 2003.
18. A. E. Darling, L. Carey, and W. Feng. The Design, Implementation, and Evaluation of mpiBLAST. In *ClusterWorld Conference & Expo and the 4th International Conference on Linux Cluster: The HPC Revolution*, 2003.
19. A. L. Delcher, K. A. Bratke, E. C. Powers, and S. L. Salzberg. Identifying Bacterial Genes and Endosymbiont DNA with Glimmer. *Bioinformatics*, 23(6):673–679, 2007.
20. A.C. Driskell, C. Ané, J.G. Burleigh, M.M. McMahon, B.C. O’Meara, and M. J. Sanderson. Prospects for building the tree of life from large sequence databases. *Science*, 306(5699):1172–1174, November 2004.
21. I. Foster, D. Kohr, R. Krishnaiyer, and J. Mogill. Remote I/O: Fast access to distant storage. In *Workshop on I/O in Parallel and Distributed Systems*, 1997.
22. J. D. Gans, W. Feng, and M. Wolinsky. Whole Genome, Physics-Based Sequence Alignment for Pathogen Signature Design. In *12th SIAM Conference on Parallel Processing for Scientific Computing*, San Francisco, California, February 2006. (electronic version unavailable).
23. S.L. Havre, B.-J. Webb-Robertson, A. Shah, C. Posse, B. Gopalan, and F.J. Brockma. Bioinformatic Insights from Metagenomics through Visualization. *Proceedings of the IEEE Computational Systems Bioinformatics Conference*, pages 341–350, 8-11 Aug. 2005.
24. Howard. An Overview of the Andrew File System. In *USENIX Winter Technical Conference*, February 1988.



25. H. Jagadish, J. Madar, and R. Ng. Semantic compression and pattern extraction with fascicles. In *VLDB '99: Proceedings of the 25th International Conference on Very Large Data Bases*, San Francisco, California, 1999. Morgan Kaufmann Publishers Inc.
26. H. Jagadish, R. Ng, B. Ooi, and A. Tung. ItCompress: An Iterative Semantic Compression Algorithm. In *Proc. of International Conference on Data Engineering*, 2004.
27. J. Lee, X. Ma, R. Ross, R. Thakur, and M. Winslett. RFS: Efficient and flexible remote file access for MPI-IO. In *Cluster*, 2004.
28. J. Lee, R. Ross, S. Atchley, M. Beck, and R. Thakur. MPI-IO/L: Efficient Remote I/O for MPI-IO via Logistical Networking. In *IPDPS*, 2006.
29. H. Lin, X. Ma, P. Chandramohan, A. Geist, and N. Samatova. Efficient Data Access for Parallel BLAST. In *IPDPS*, 2005.
30. A. Lukashin and M. Borodovsky. GeneMark: New Solutions for Gene Finding. *Nucleic Acids Res.*, 26:1107–1115, 1998.
31. Message Passing Interface Forum. *MPI-2: Extensions to the Message-Passing Standard*, July 1997.
32. Generic Coding of Moving Pictures and Associated Audio Information - Part 2: Video. ITU-T and ISO/IEC JTC1, ITU-T Recommendation H.262 ISO/IEC 13 818-2 (MPEG-2), 1994.
33. Coding of Audio-Visual Objects - Part 2: Visual. ISO/IEC JTC1, ISO/IEC 14 496-2 (MPEG-4 Visual version 1), Apr. 1999; Amendment1 (version 2), Feb. 2000; Amendment 4 (streaming profile), Jan. 2001.
34. B. Nowicki. *NFS: Network File System Protocol Specification*. Network Working Group RFC1094, 1989.
35. D. A. Reed. Grids, the TeraGrid, and Beyond. *Computer*, 36(1):62–68, 2003.
36. F. Schmuck and R. Haskin. GPFS: A shared-disk file system for large computing clusters. In *FAST*, 2002.
37. Seti@home: The search for extraterrestrial intelligence. <http://setiathome.ssl.berkeley.edu/>.
38. S. Simms, G. Pike, and D. Balog. Wide Area Filesystem Performance using Lustre on the TeraGrid. In *TeraGrid Conference*, 2007.
39. R. Thakur, W. Gropp, and E. Lusk. Data Sieving and Collective I/O in ROMIO. In *Frontiers of Massively Parallel Computation*, 1999.
40. O. Thorsen, B. Smith, C. P. Sosa, K. Jiang, H. Lin, A. Peters, and W. Feng. Parallel Genomic Sequence-Search on a Massively Parallel System. In *ACM International Conference on Computing Frontiers*, Ischia, Italy, May 2007.
41. A. Warren, J. Archuleta, W. Feng, and J. Setubal. Missing Genes in the Annotation of Prokaryotic Genomes. *BMC Bioinformatics*, 2010.