

Making the Best of a Bad Situation: Prioritized Storage Management in GEMS

J. M. Wozniak^a P. Brenner^a D. Thain^a A. Striegel^a
J. A. Izaguirre^a

^a *Dept. of Computer Science & Engineering*
University of Notre Dame
Notre Dame, IN 46556 USA

{ jwozniak, pbrenne1, dthain, striegel, izaguirr } @ nd.edu

Abstract

As distributed storage systems grow, the response time between detection and repair of the error becomes significant. Systems built on shared servers have additional complexity because of the high rate of service outages and revocation. Managing high replica counts in this environment becomes very costly in terms of the storage required and bandwidth consumption for file copies. The storage challenge for this situation can thus be phrased as an attempt to function inexpensively with respect to cost constraints such as: disk utilization, network bandwidth consumption, and server CPU time. The GEMS (Grid Enabled Molecular Simulation) storage system provides a replicated and shared workspace for large scale molecular dynamics simulations, and exemplifies the above issues. The GEMS framework offers a solution to this problem by accessing metadata, prioritizing observed faults, and repairing them in an intelligent manner. In this paper, we provide observations from the operation of GEMS and compare its error handling to like storage systems.

Key words: replicated, shared, simulation, storage

1 Introduction

An important observation about large systems is that their response to user and internal system operations is delayed, scaling with the size of the system. The exact nature of the scaling, of course, depends on the algorithms used, but some increase in latency is to be expected. In an active replication system, system operations to maintain integrity become more frequent, and often more expensive individually. This is true even if we assume that the underlying

components do not become less reliable, which is a dubious assumption in Grid computing.

A typical solution to increase system robustness is to decrease the reliance on a central server, either by deploying redundant central servers or by shifting the system architecture to an independent, peer-to-peer model. Yet, reducing the reliance on a centralized server tends to increase the number of computers involved in a given transaction, thus again increasing latency. This architectural technique also increases the algorithmic complexity of system operations.

As systems get decentralized, fault-tolerance must be built in. An additional complexity is often added to Grid enabled storage systems: the distributed ownership of the resource fabric. A shared system allows providers to pull the plug on their volunteered system at any moment. Thus, such a system incorporates as little confidence in a storage device as possible, and fault tolerance is a foundational design feature.

What do we mean by low confidence? Such a shared resource, in the worst case, a commodity workstation, may be rebooted several times a day. However, this does not directly indicate permanent data loss, as the vast majority of machines come back online after a reboot. It does mean that its data is temporarily unavailable, and that a replica must be handy. Even users of mature storage systems, as discussed in the following sections, must sometimes wait for data. However, a properly replicated system will nearly eliminate such outages.

A consequence of the above observations is that repairing storage failures is slow and expensive. Losing, say, a 250GB disk full of simulation data will cause the controlling servers to search for and replicate that much data over the network. This assigns a real cost to the repair of this fault, which would be vain if the system suddenly came back up. If it does not come back, since the lost replicas may take a day or more to copy and restructure, other permanent fault combinations in the meantime could result in permanent data loss for one or more files. This puts a real cost on not just the storage space used by extra replicas, but also on the actions taken to respond to changes in the system, and stresses the need to order repairs effectively.

Earlier this year, a new project was officially started at Notre Dame to store large quantities of biomolecular simulation data. This work stems from a recent NSF grant which intends to combine experimental and simulated scientific data on a shared storage Grid. A variety of requirements for the storage system were quickly identified from the research domain, including data storage requirements relating to the computational needs of researchers in a high performance setting and scientific requirements that were to improve the utility of the storage system for exploring the valuable content of the data. A prototype

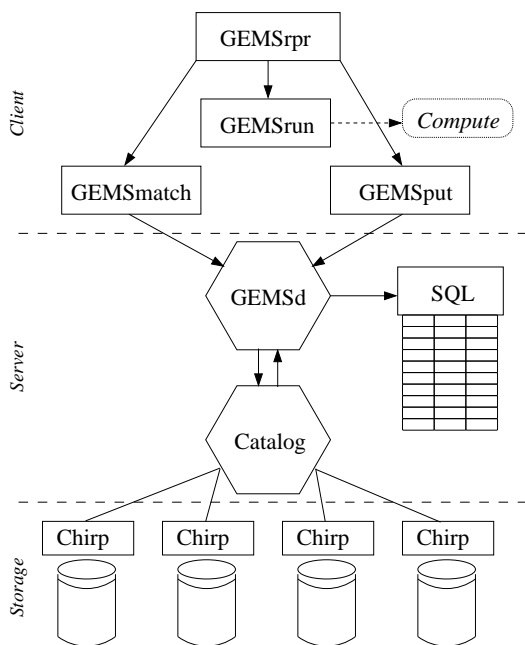


Fig. 1. GEMS Toolset Framework.

was designed and quickly brought online to begin experimentation.

The system is designed to maximize the opportunities for scientific research by cataloging complete metadata for data in the system, allowing for automatically spawned simulations and automatic data reuse, and exploiting relationships among data sets for research purposes. More abstractly, GEMS provides a framework for *analyzing*, *sharing*, and *archiving* important biomolecular data. The storage fabric for this system is a widely distributed uncontrolled network of servers that, individually, offer only best effort reliability. This provides an interesting observatory for the investigation of large scale fault tolerant systems built upon partially reliable servers.

The resulting system, named GEMS, for Grid Enabled Molecular Simulation, consists of client tools and a network of services that combine to provide the requested functionality to researchers in this actively evolving field. Figure 1 provides the skeleton of the system. This figure breaks GEMS into the client tools used to query the system, and the server's relationship to the file storage servers, which run Chirp [Thain, 2004].

The three client tools are GEMSmatch, GEMRun, and GEMSpur which can be accessed from the shell but are ultimately intended to be accessed through GEMSRpr. GEMSmatch is used to query the metadatabase with a variety of simulation parameters for existing files and replica locations. Original files and their simulation metadata are added to the system by GEMSpur, at which point the maximum number of replicas is set for each file. Data is retrieved from the system or rebuilt by GEMRun, which implements a virtualized data

model specific to our target application, biomolecular simulation (cf. [Foster et al., 2002]). New simulations are currently sent to the Grid via APST [Casanova et al., 2000]. To provide a research-oriented workspace, all three operations may be accessed through GEMSRpr, which stands for Result Production Request. This high-level tool is currently accessed through the shell but intended to be controlled by a web portal.

We have previously discussed [Wozniak et al., 2005] how GEMS manages high levels of disk loss, including loss of half of the storage servers simultaneously. A method of reallocating the replication count of files was presented to reduce the probability of permanent data loss, resulting in a high degree of disk utilization control. In this paper, we intend to focus on improvements to the failure handling model to handle *subsequent massive failures*.

From the ground up, GEMS has been intended to respond to various faults that occur in large distributed systems. In this paper, we specifically focus on the effect of failures in the underlying Chirp servers to deliver requested data, or, even an acknowledgement. GEMS interprets these events and builds up an error handling system in response. So GEMS is fault-tolerant in the sense that it can continue to respond to queries and file requests despite the problematic storage layer upon which it operates. As discussed below, GEMS has knowledge of the metadata and data stored on its servers, and can *prioritize* response operations to reduce the probability of permanent loss of *important* data, and attempt to minimize the cost of replicating files unnecessarily.

In this paper, we compare the fault responses of GEMS to that of other systems, and demonstrate how GEMS meets the requirements of its research-based user pool. In the next section, we describe the GEMS model as it compares to other commonly used storage architectures and existing storage software systems. We then move into GEMS' internal management techniques, as Section 3 explains what constitutes a fault in the GEMS model, and Section 4 explains the GEMS fault response system and its priority assignment model. Experimental results are provided in 5, and we offer some concluding remarks in Section 6.

2 GEMS Distributed Storage Model

GEMS has commonality with distributed file systems, databases, and peer-to-peer sharing systems. Users of these different systems have different expectations when it comes to error states, and internally, faults are treated differently.

Users of file systems and databases generally expect all-or-nothing responses.

File system users expect that if one file on their machine is present and correct, then they will all be there. This locality assumption comes from the standard assumption that when a physical storage device fails, all of its data is permanently lost, and that partial failures are handled by the operating system and not exposed to the user. This carries over into expectations of network file systems, which often attempt to emulate the behavior of a local file system.

While GEMS does store whole files and their directory information, it does not attempt to provide file system behavior. Files that are expected to be found in the same directory may be found on different hosts, destroying any assumptions about locality in the delivery of a file. This affects the design of clients to GEMS, which have to be able to utilize lists of replica locations. While GEMS does not necessarily expose internal errors to the user, the user should be aware of the actions that GEMS may take in the case of data loss, which include contacting various remote hosts, or submitting jobs to a compute Grid to rebuild lost data.

NFS [Sandberg et al., 1985] is a the standard point of comparison for file retrieval latency and API semantics for remote storage. GEMS is not designed to compete with a finely tuned cluster storage system, and since it is not a filesystem, it does not attempt to mimic any system's API semantics. GEMS does function in an RPC fashion from the client's perspective, and is stateless. On failure, GEMS does not block or disconnect, but continues to attempt to satisfy a user request, recreating data if possible. If permanent data loss has occurred, this will be reported and must be handled at a higher level.

Similarly, AFS [Howard et al., 1988] shares NFS semantics, but provides some additional error handling functionality. Read-only replicas of stored data may be easily configured, and roll back to previous data is built in. However, AFS servers must maintain cache consistency for their clients. GEMS does not need to cache data because of its write-once characteristic, greatly reducing the server's responsibilities.

Database systems similarly are expected to provide all-or-nothing responses. GEMS may only be able to partially fulfill a query and will resort to recomputation or partial delivery. GEMS does provide a metadatabase which may be queried, and normal database user expectations apply; when querying for output data files, the returned result is a fault-tolerant plan on how to retrieve the requested information from remote storage devices, and feed it into an external program.

GEMS can be compared to a peer-to-peer system because it allows users to query each others' data sets, and it allows users to pool storage resources. Peer-to-peer users have very low or non-existent expectations about reliability, for example, a file that is seen once could quickly vanish forever. However, GEMS

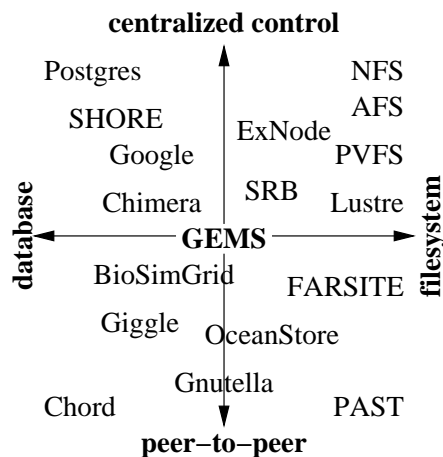


Fig. 2. Related Work

is a write-once system, so that data that is published to the system cannot be modified, which is quite different from the expectations of most storage users, but not so surprising in a setting in which data is “published” to other users. GEMS data is much more reliable than an *ad hoc* peer-to-peer network, because it abstracts the storage layer from the client.

The most obvious point of comparison for a redundant data storage system is a hardware RAID [Patterson et al., 1988], or a network-based RAID, such as Zebra [Hartman and Ousterhout, 1995]. In a typical hardware RAID setting, faults may be detected as blocks are read from the device, where in GEMS, errors are actively probed by a server component. Zebra detects faults over the network, which is a technically difficult observation to make. GEMS and RAID both benefit from hot-pluggable hardware. As devices are added, both systems can discover and begin utilizing new storage. When a new server is added to GEMS, it may immediately begin receiving data from a client, just as RAID may immediately use a new disk. However, intentionally removing a disk from a RAID system is a very expensive operation, which is not compatible with a Grid shared storage system. GEMS offers some additional benefits in that it has knowledge of the importance and replication status of the stored data, so upon replica loss, replication does not necessarily immediately consume network bandwidth for all lost data.

With these summarized comparisons, we can describe GEMS in the context of modern data-Grids [Finkelstein et al., 2004]. The GEMS architecture combines a few commonly observed architectures. The storage sites may be thought of as peer-to-peer clusters that combine to contribute to a complete storage system, each offering file space and receiving replicas through a metadatabase. Since the system is actively monitored for data loss by a variety of server components, it can accurately be described as an agent controlled system. The active components can reside on various servers, which lends greatly to fault tolerance and scalability, but creates the increased possibility of partial

failure in the system and requires a shared critical resource for inter-component communication.

This hybrid architecture allows our fault handling policy to utilize the metadata archive to make intelligent choices concerning replication and repair. It also meets other requirements of the system, such as the ability for contributing storage servers to revoke storage with limited impact, and simplifies the security aspects of such a shared system.

Another modern model for distributed storage is the replica locator model, as exemplified in the Giggle project [Chervenak et al., 2002]. Giggle provides a highly distributed service to redirect filename requests in a highly efficient way. However, the requirements of the molecular simulation environment demand responses to metadata requests. Since the server components also have access to the metadata, these are used to control storage policy.

3 GEMS Failure to Maintenance Continuum

The GEMS framework is subject to many of the same error modes that plague its cousin Grid distributed systems. Sample sources of faults include: an errant user, the desktop hardware, the OS, the TCP connection, the local switch, numerous routers, the storage server software and its hardware, etc. Even though the apparent severity of one of the mentioned failures may differ by orders of magnitude (the loss of one node due to hard drive failure versus the loss of an entire DNS domain due to switch failure), the GEMS design embraces the spectrum of failure with one unified resolution policy. *Observed faults are viewed not as pre-emptive exceptions, but events on par with typical maintenance operations.*

To give this concept substance we introduce the concept of a Problem. A Problem is an object in GEMS which can be assigned a priority based on its severity and queued for resolution. Problems are generated through Auditors which continuously audit the system state with respect to the metadata, physical reality, and user requests as reflected in the metadata. On the failure side of the Problem continuum, one example would be the loss of a local switch providing access to ten nodes, say, *Cluster A*. The Auditor would recognize that the metadata pointing to *Cluster A* is now inconsistent and for each unavailable file a Problem will be instantiated with all of the pertinent information necessary for calculation of a problem priority and resolution of the problem. On the maintenance side of the Problem continuum, consider a client who submits a record of ten files into the database with a requested replication of 5 each. The GEMSD daemon will make sure the first copy is submitted to stable storage by the client tool but the remaining requested copies will

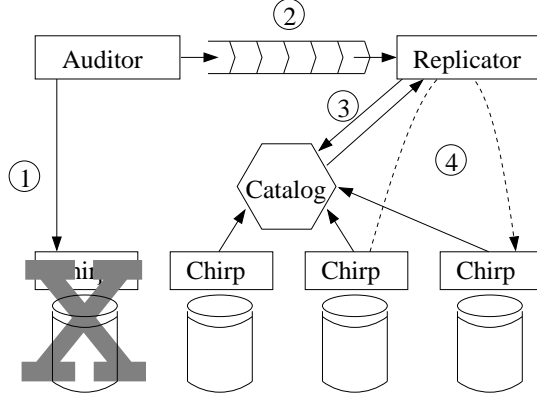


Fig. 3. GEMS Prioritized Repair

be handled through the instantiation of Problems as the metadata shows an insufficient number of replications with respect to the user's request. Hence, Problems can be subtyped to provide a useful abstraction for error handling.

With this unified response to both failures and maintenance tasks we can exploit varied risk/cost targets through specification of those metrics which comprise the prioritization function.

4 Prioritization of Failure Handling Events

The prioritization of Problems is central to the priority queue fault tolerance model. In this work we do not argue for a specific priority function nor discuss the suitability of those parameters which are included in the function. We do however present our prototype's implementation of the priority calculation and mention select parameters from our priority function in order to substantiate a functional implementation. The flexible prioritization function allows the specific user community to tune the system to their risk/cost demands.

As shown in Figure 3, the Auditor iteratively probes Chirp servers to verify their contents against the GEMS metadatabase. When any of a variety of faults are detected (1), the relevant error information is encapsulated into a Problem object and enqueued (2). When the Replicator pops the queue, it determines a response, usually by creating an additional replica. Existing copies of the file in question are located by consulting the metadatabase and a suitable host for the new replica is found by reading the Chirp catalog (3). The file is then copied over (4).

Each Problem object contains a set of members which are necessary for the Problem's resolution and which also serve as the parameters for the priority calculations. From the set of priority members a priority is calculated internally by the Problem on instantiation. Of key interest is the ability for the

Problem to recalculate its own priority. For example, the priority queue could trigger all Problems in the queue to recalculate their priority (stale problem member data is addressed in the following section).

A good demonstration of the priority recalculation utility leads us into the discussion of priority function parameters. Classical job starvation concerns would prompt us to include “time in queue” as a parameter with varies in proportion to an increased priority. This way, Problems involving a host failure which has only been recognized for a short period are not handled immediately, which gives the host system a chance to recover, eliminating an unnecessary file transfer. A second parameter of importance for our scientific user base is file type. Whereas an output file may be automatically regenerated because the metadata contains enough information to derive the output files, the input configuration files are irreplaceable without human intervention. To again reduce the probability of permanent data loss we increase priority in with respect to the number of remaining replications. Additional parameters such as file size allow us the ability to fine tune the prioritization for improved response to massive failure scenarios. An important observation in biomolecular simulation data sets is that input files tend to be small, and output files are large. This allows us to aggressively replicate the input files while being more cavalier about output files, which saves storage space and bandwidth.

In general, our priority function for a Problem P takes the form:

$$\begin{aligned}
 P.\text{priority}() &= k_1 \times (\text{replicas missing}) \\
 &+ k_2 \times (\text{max replicas}) \\
 &+ k_3 \times (\text{age}) \\
 &+ k_4 \times (\sqrt{\text{size}})^{-1}.
 \end{aligned}$$

This turns a system design, or policy design, problem into an optimization problem. We can now flexibly pick parameters to meet the needs of various systems, if risks and costs can be approximated, or desired priorities are given:

Missing	Max	Age	Size	Priority
2	4	1 hr.	50KB (input)	100
1	4	1 hr.	50KB (input)	30
1	3	2 hr.	10MB (output)	20
2	3	6 hr.	10MB (output)	150

Table 1
Sample Priorities

Solving for k provides us with a useful policy, which can be tested by applying it to other hypothetical Problems or by experiment. An interesting case

occurs when the priority is evaluated and negative, or lower than the Replicator threshold. Assuming $k_3 > 0$, this means that the Problem will not be immediately handled for a few hours, giving the host machine a chance to recover.

4.1 Stale Problems

System state data is inherently stale in a distributed framework, and is expensive to obtain, as discussed in the introduction. Therefore we expect that the state which was recognized when the problem was instantiated may not reflect the system state when the problem is handled. In response our failure handler must validate that the Problem member data is still current. If the Problem data remained in sync with the current metadata the Problem is resolved as normal, however, if the data is no longer current the Problem priority status is re-evaluated. If the Problem priority is reduced, as in the case of a returning storage server, a new Problem will be created and queued, and if the Problem is eliminated, the Problem is discarded. If the priority of the new Problem exceeds or matches that of the stale problem it is, of course, resolved immediately.

4.2 Failed Failure Responses

One challenge in the field of fault tolerance has been the determination of when to abort attempts at repair if the repair operation itself continues to fail. Although the GEMS response to the failure of one system is chiefly the replication of a remaining copy of the file on a second machine to a new destination on a third machine, one can imagine that by the time the Problem is being handled either the second machine with the remaining copy has also gone down. The typical answer would be to retry at regular intervals with an eventual timeout. The GEMS model improves upon this by allowing a Problem to count its history of failed repairs, and adjust its priority accordingly. Problems are re-queued when repair fails, so that other repairs may be attempted. Since the age of a Problem is also recorded, extremely old Problems can signal for human intervention.

4.3 Starvation and Priority Inversion

Queue structures prompt the user to investigate the fate of those jobs relegated to the back. The GEMS prioritization function as we have implemented it includes consideration of the time spent in the queue to prevent starvation. The

priority queue enables preferential response unavailable to sequential queues but does not guarantee against priority inversion even for a large number of problem handlers running in parallel. Our model allows for Problems to be re-prioritized *while the repair is in progress*, and the Replicator to be interrupted and forced to pop the queue again, by a third agent that probes for priority inversion. Thrashing is avoided by prioritizing a Problem currently being repaired with respect to the size of the data yet to be replicated, and the fact that it is in progress. Since our priority function includes a term which is inversely proportional to the square root of the size, copies in progress are unlikely to be interrupted except in the case of real emergencies, as defined by the tuned coefficients.

4.4 Scaled Response on par with Problem Queue Size

For a sufficient number of Auditor components detecting problems, it is reasonable to propose that the size of the problem queue is indicative of the number of failure handlers required. Furthermore the average priority of problems in the queue provides insight into the severity of the problems and should influence the number of failure handlers. Although the experimental results presented in the work reflect a single failure handler, the authors plan to implement scaled instantiation of failure handling components on other GEMS servers with respect to priority queue metrics.

5 Experimental Results

The authors have developed a functional GEMS prototype as reported in [Wozniak et al., 2005]. Specific to this work we implemented the Problem Queue and report experimental data in demonstration of the error recovery capability.

We intend to demonstrate that the new functionality improves upon the previous functionality, which is similar to functionality obtained in a system with limited access to simulation metadata. Such a data-agnostic system could only repair errors in the order observed, one at a time, and could not respond dynamically to a rapidly changing storage fabric. This storage model is not an artificial straw man, indeed, it is the model produced by combining existing virtualized data techniques with a replica locator service. Our prototype can emulate this behavior if we make all the observed priorities equivalent, and perform in-order response.

Our first example simply shows that the prototype can respond to problems in

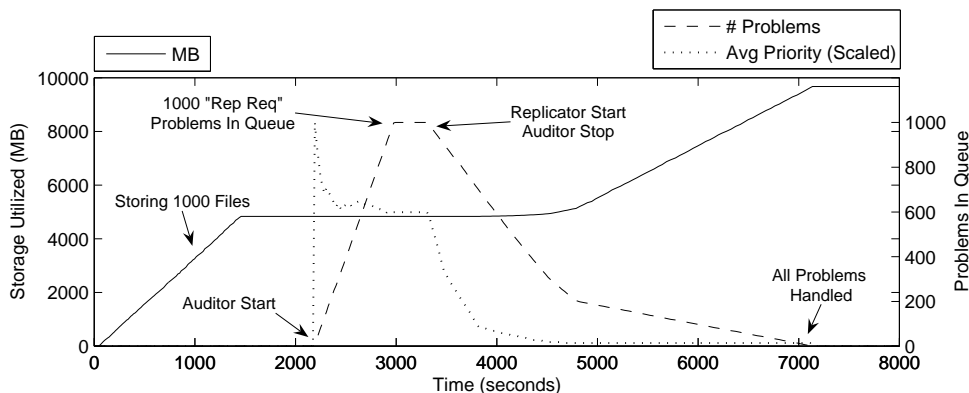


Fig. 4. Queue Performance: *Priority Dominated by File Size*

the order specified by the priority function. In this case, small files have priority over large files, which is typical in a biomolecular simulation environment where the input files and configuration files are irreplaceable, and the output files represent derived data and are only required to speed the retrieval time.

As demonstrated in Figure 4, we start by disabling the GEMS storage management components, the Auditor, Replicator, and GarbageCollector. The simulation data sets, about 3 GB total, are added to the system, so each file has one copy on some Chirp server. The Auditor is then turned on and obvious problems are detected: the files need to be replicated up to the requested level. The graphic shows the average priority of problems, which are scored by, in this case, the size of the file and the replication count. The Replicator is turned on, and the average priority quickly drops as high priority problems are quickly handled by replicating these small files on other Chirp hosts. The average priority slowly levels off as larger, low priority output files are copied over the network, increasing the disk usage level.

6 Conclusion

Researchers that produce and use biomolecular simulation data have needs, both computational and scientific, for a storage system that reliably catalogs and stores simulation data and metadata. Scientifically, the system must be centered around the metadata and provide a virtualized data framework, which makes output data recreateable and new explorations more systematic. Computationally, the metadatabase allows the storage policy to act with substantially more intelligence than existing distributed or replica systems. In this paper, we explained the necessary observations about the data in question and showed that it can be prioritized by value and risk in a variety of tunable

ways. We provided our Problem model, which encapsulates information about any of a variety of observed situations in a sortable context, which facilitates appropriate response with respect to other conditions in the system.

Because the risk/cost parameters are difficult to obtain in advance, we will have to continue experimentation to improve the system performance with respect to its most important metric: permanent loss of input files. In the meantime, we provide some observed behavior that shows that a first-try prioritization scheme is much better than no prioritization, especially important on volatile network of unreliable storage servers.

References

- H. Casanova, G. Obertelli, F. Berman, and R. Wolski. The AppLeS parameter sweep template: user-level middleware for the grid. In *Proceedings of the 2000 ACM/IEEE conference on Supercomputing*. IEEE Computer Society, 2000.
- A. Chervenak, E. Deelman, I. Foster, L. Guy, W. Hoschek, A. Iamnitchi, C. Kesselman, P. Kunszt, M. Ripenu, B. Schwartzkopf, H. Stocking, K. Stockinger, and B. Tierney. Giggle: A framework for constructing scalable replica location services. In *Proc. IEEE Supercomputing*, 2002.
- A. Finkelstein, C. Gryce, and J. Lewis-Bowen. Relating requirements and architectures: A study of data-Grids. *J. Grid Computing*, 2, 2004.
- I. Foster, J. Voeckler, M. Wilde, and Y. Zhou. Chimera: A virtual data system for representing, querying, and automating data derivation. In *Proceedings of the 14th Conference on Scientific and Statistical Database Management*, Edinburgh, Scotland, July 2002.
- J. H. Hartman and J. K. Ousterhout. The Zebra striped network file system. *ACM Transactions on Computer Systems*, 13(3), 1995.
- J. Howard, M. Kazar, S. Menees, D. Nichols, M. Satyanarayanan, R. Sidebotham, and M. West. Scale and performance in a distributed file system. *ACM Transactions on Computer Systems*, 6(1):51–81, February 1988.
- D. A. Patterson, G. Gibson, and R. Katz. A case for redundant arrays of inexpensive disks (RAID). In *Proc. ACM SIGMOD international conference on management of data*, 1988.
- R. Sandberg, D. Goldberg, S. Kleiman, D. Walsh, and B. Lyon. Design and implementation of the Sun network filesystem. In *Proc. USENIX Summer Technical Conference*, 1985.
- D. Thain. *Chirp User's Manual*. University of Notre Dame, 2004. www.cse.nd.edu/~ccl/software/chirp/.
- J. M. Wozniak, P. Brenner, D. Thain, A. Striegel, and J. A. Izaguirre. Generosity and gluttony in GEMS: Grid enabled molecular simulations. In *Proc. IEEE High Performance Distributed Computing*, 2005.