

The Zebra Striped Network File System

JOHN H. HARTMAN

University of Arizona

and

JOHN K. OUSTERHOUT

Sun Microsystems Laboratories, Inc.

Zebra is a network file system that increases throughput by striping the file data across multiple servers. Rather than striping each file separately, Zebra forms all the new data from each client into a single stream, which it then stripes using an approach similar to a log-structured file system. This provides high performance for writes of small files as well as for reads and writes of large files. Zebra also writes parity information in each stripe in the style of RAID disk arrays; this increases storage costs slightly, but allows the system to continue operation while a single storage server is unavailable. A prototype implementation of Zebra, built in the Sprite operating system, provides 4–5 times the throughput of the standard Sprite file system or NFS for large files and a 15–300% improvement for writing small files.

Categories and Subject Descriptors: D.4.2 [**Operating Systems**]: Storage Management—*allocation/deallocation strategies; secondary storage*; D.4.3 [**Operating Systems**]: File Systems Management—*access methods; distributed file systems; file organization*; D.4.5 [**Operating Systems**]: Reliability—*fault-tolerance*; D.4.7 [**Operating Systems**]: Organization and Design—*distributed systems*; D.4.8 [**Operating Systems**]: Performance—*measurements*; E.5 [**Data**]: Files—*organization/structure*

General Terms: Design, Measurement, Performance, Reliability

Additional Key Words and Phrases: Log-based striping, log-structured file system, parity computation, RAID

1. INTRODUCTION

Zebra is a network file system that uses multiple file servers to provide greater throughput and availability than can be achieved with a single server. Clients stripe file data across servers so that different pieces of data

This work was supported in part by NSF grant CCR-89-00029, NASA/ARPA grant NAG2-591, NSF grant MIP-87-15235, ARPA contract N00600-93-C-2481, and the California MICRO Program.

Authors' addresses: J. H. Hartman, Department of Computer Science, Gould-Simpson Building, The University of Arizona, Tucson, AZ 85721; email: jhh@cs.arizona.edu; J. K. Ousterhout, Sun Microsystems Laboratories, Inc., 2550 Garcia Avenue, MS UMTV29-232, Mountain View, CA 94043-1100; email: john.ousterhout@eng.sun.com.

Permission to make digital/hard copy of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage, the copyright notice, the title of the publication, and its date appear, and notice is given that copying is by permission of ACM, Inc. To copy otherwise, to republish, to post on servers, or to redistribute to lists, requires prior specific permission and/or a fee.

© 1995 ACM 0734-2071/95/0800-0274 \$03.50

ACM Transactions on Computer Systems, Vol. 13, No. 3, August 1995, Pages 274–310.

are stored on different servers. Striping makes it possible for a single client to keep several servers busy, and it distributes the load among the servers to reduce the likelihood of hot spots. Zebra also stores parity information in each stripe, allowing it to continue operation while any one server is unavailable.

In current network file systems the read and write bandwidth for a single file is limited by the performance of a single server, including its memory bandwidth and the speed of its processor, network interface, I/O buses, and disks. It is possible to split a file system among multiple servers; but each file must reside on a single server, and it is difficult to balance the loads of the different servers. For example, system directories often lie on a single server, making that server a hot spot.

In the future, new styles of computing such as multimedia and parallel computation are likely to demand much greater throughput than today's applications, making the limitations of a single server even more severe. For example, a single video playback can consume a substantial fraction of a file server's bandwidth even when the video is compressed. A cluster of workstations can easily exceed the bandwidth of a file server if they all run video applications simultaneously, and the problems will become much worse when video resolution increases with the arrival of HDTV. Another example is parallel applications. Several research groups are exploring the possibility of using collections of workstations connected by high-speed low-latency networks to run massively parallel applications [Anderson et al. 1995; Freeh et al. 1994]. These "distributed supercomputers" are likely to present I/O loads equivalent to traditional supercomputers, which cannot be handled by today's network file servers.

A striping file system offers the potential to achieve very high performance using collections of inexpensive computers and disks. Several striping file systems have already been built, such as Swift [Cabrera and Long 1991] and Bridge [Dibble et al. 1988]. These systems are similar in that they stripe data within individual files, so that only large files benefit from the striping. Zebra uses a different approach borrowed from log-structured file systems (LFS) [Rosenblum and Ousterhout 1991]. Each client forms its new data for all files into a sequential log that it stripes across the storage servers. This not only improves large-file performance through striping, but it also improves small-file writes by batching them together and writing them to the servers in large, efficient transfers. It also reduces network overhead, simplifies the storage servers, and spreads write traffic uniformly across the servers.

Zebra's style of striping also makes it easy to use redundancy techniques from RAID disk arrays to improve availability and data integrity [Patterson et al. 1988]. One of the fragments of each stripe stores parity for the rest of the stripe, allowing the stripe's data to be reconstructed in the event of a disk or server failure. Zebra can continue operation while a server is unavailable. Even if a server is totally destroyed, Zebra can reconstruct the lost data.

We have constructed a prototype implementation of Zebra as part of the Sprite operating system [Ousterhout et al. 1988]. Although it does not incorporate all the reliability and recovery aspects of the Zebra architecture, it does demonstrate the performance benefits. For reads and writes of large

files the prototype achieves up to 2.6MB/second for a single client with five servers, which is 4–5 times the throughput of either NFS or the standard Sprite file system on the same hardware. For small files the Zebra prototype improves performance by more than a factor of 3 over NFS. The improvement over Sprite is only about 15%, however. This is because both Zebra and Sprite require the client to notify the file server of the file opens and closes, and when writing small files these notifications dominate the running time. With the addition of file name caching to both systems, Zebra is expected to have more of an advantage over Sprite.

The rest of the article is organized as follows. Section 2 describes the computing environment for which Zebra is intended, and the types of failures it is designed to withstand. Section 3 describes the RAID and log-structured file system technologies used in Zebra and introduces Zebra's logging approach. Section 4 describes the structure of Zebra, which consists of clients, storage servers, a file manager, and a stripe cleaner. Section 5 shows how the components of the system work together in normal operation; communication between the components is based on *deltas*, which describe file block creations, updates, and deletions. Section 6 describes how Zebra restores consistency to its data structures after crashes. Section 7 shows how the system provides service while components are down. Section 8 gives the status of the Zebra prototype and presents some performance measurements. Section 9 discusses related work. Section 10 concludes.

2. ZEBRA APPLICABILITY

Zebra makes several assumptions concerning its computing environment and the types of failures that it will withstand. Zebra is designed to support UNIX workloads as found in office and engineering environments. These workloads are characterized by short file lifetimes, sequential file accesses, infrequent write-sharing of files by different clients, and many small files [Baker et al. 1991]. This environment is also notable because of the behavior it does not exhibit: namely, random accesses to files. Zebra is, therefore, designed to handle sequential file accesses well, perhaps at the expense of random file accesses. In particular, this means that Zebra may not be suitable for running database applications, which tend to update and read large files randomly. This is not to say that the Zebra design precludes good performance on such a workload, but that the current design has not been tuned to improve random-access performance.

Zebra is also targeted at high-speed local area networks; it assumes that in a data transfer between a client and server the point-to-point bandwidth of the network is not a bottleneck. Zebra is also not designed to handle network partitions. New point-to-point network architectures, such as ATM, typically include redundant links that reduce the probability of a network partition and make partitions less of a concern in the design of a network file system for use on a local area network.

Zebra also assumes that clients and servers will have large main-memory caches to store file data. These caches serve two purposes: to allow frequently

used data to be buffered and accessed in memory, without requiring an access to the server or the disk, and to buffer newly written file data prior to writing it to the server or the disk. The former filters out accesses to data that are frequently read, whereas the latter filters out short-lived data and allows Zebra to batch together many small writes by application programs into large writes to the servers.

Zebra is designed to provide file service despite the loss of any single machine in the system. Multiple server failures are not handled; the loss of a second server causes the system to cease functioning, and data may be lost if disks fail catastrophically on two servers at the same time. Any number of clients may fail, however, without affecting the availability of file data. A client crash may lose newly written data cached on that client; but it cannot lose data older than a time limit, nor can it lose data written by another client. This is analogous to losing the data stored in a UNIX file system cache when the machine crashes.

3. STRIPING IN ZEBRA

Zebra distributes file data over several file servers while ensuring that the loss of a single server does not affect the availability of the data. To do this, Zebra borrows from two recent innovations in the management of disk storage systems: RAID technology (Redundant Arrays of Inexpensive Disks) [Patterson et al. 1988] and log-structured file systems (LFS) [Rosenblum and Ousterhout 1991]. RAID technology allows Zebra to provide scalable file access performance while using parity instead of redundant copies to guard against server failures. The log-structured approach simplifies the parity implementation, reduces the impact of managing and storing parity, and allows clients to batch together small writes to improve the efficiency of writing to the servers.

3.1 RAID

RAID is a storage system architecture in which many small disks work together to provide increased performance and data availability. A RAID appears to higher-level software as a single very large and fast disk. Transfers to or from the disk array are divided into blocks called *striping units*. Consecutive striping units are assigned to different disks in the array as shown in Figure 1 and can be transferred in parallel. A group of consecutive striping units that spans the array is called a *stripe*. Large transfers can proceed at the aggregate bandwidth of all the disks in the array, or multiple small transfers can be serviced concurrently by different disks.

Because a RAID has more disks than a traditional disk storage system, disk failures will occur more often. Furthermore, a disk failure anywhere in a RAID can potentially make the entire disk array unusable. To improve data integrity, a RAID reserves one of the striping units within each stripe for parity instead of data (see Figure 1); each bit of the parity striping unit contains the exclusive OR of the corresponding bits of the other striping units in the stripe. If a disk fails, each of its striping units can be recovered using

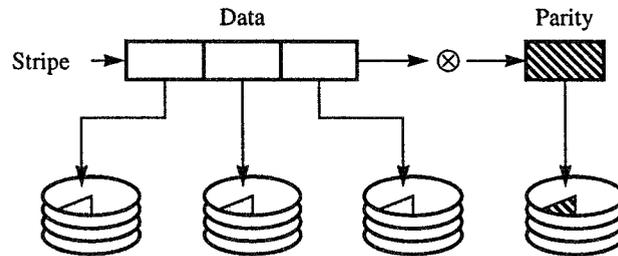


Fig. 1. Striping with parity. The storage space of a RAID disk array is divided into stripes, each stripe containing a striping unit on each disk of the array. All but one of the striping units hold data; the remaining striping unit holds parity information that can be used to recover after a disk failure.

the data and parity from the other striping units of the stripe. The file system can continue operation during recovery by reconstructing data on-the-fly.

A RAID offers large improvements in throughput, data integrity, and availability, but it presents potential problems. First, the parity mechanism makes small writes expensive. If all write operations are in units of whole stripes, then it is easy to compute the new parity for each stripe and write it along with the data. This increases the cost of writes by only $1/(N - 1)$ relative to a system without parity, where N is the number of disks in the array. However, the overhead of small writes is much higher. In order to keep the stripe's parity consistent with its data, it is necessary to read the current value of the data block that is being updated, read the current value of the corresponding parity block, use this information to compute a new parity block, then rewrite both parity and data. This makes small writes in a RAID about four times as expensive as they would be in a disk array without parity, since they require two reads and two writes to complete. Unfortunately, the best size for a striping unit appears to be tens of kilobytes or more [Chen and Patterson 1990], which is larger than the average file size in many environments [Baker et al. 1991; Hartman and Ousterhout 1993]; therefore, writes will often be smaller than a full stripe.

The second problem with a disk array is that all the disks are attached to a single machine, so its memory and I/O system are likely to be a performance bottleneck. For example, it is possible to attach multiple disks, each with a bandwidth of 1–2MB/second, to a single SCSI I/O bus, but the SCSI bus has a total bandwidth of only 2–10MB/second. Additional SCSI buses can be added, but data must be copied from the SCSI channel into memory and from there to a network interface. On the DECstation 5000/200 machines used for the Zebra prototype, these copies to and from the SCSI and network controllers can only proceed at about 6–8MB/second. The Berkeley RAID project has built a special-purpose memory system with a dedicated high-bandwidth path between the network and the disks [Drapeau et al. 1994], but even this system can support only a few dozen disks at full speed.

The fundamental problem with using a disk array to improve server bandwidth is that the server itself becomes a performance bottleneck. In

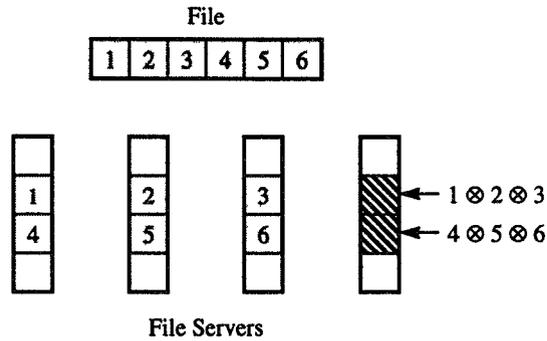


Fig. 2. File-based striping for a large file. The file is divided into stripe fragments that are distributed among the servers. One fragment of each stripe contains the parity of the stripe's contents.

order to eliminate the bottlenecks presented by centralized resources, multiple paths must exist between the source or sink of data and the disks so that different paths can be used to reach different disks. For example, this might be done by spreading the disks among different machines on a single very high speed network, or even by using different networks to reach different disks. Unfortunately, this turns the disk array into a distributed system and introduces issues such as who should allocate disk space or compute parity. Nonetheless, this distribution is necessary to avoid the bottleneck presented by having multiple data paths share the same resource. One of our goals for Zebra was to solve these distributed-system problems in a simple and efficient way.

3.2 File-Based Striping in a Network File System

A striped network file system is one that distributes file data over multiple file servers in the same way that a RAID distributes data over multiple disks. This allows several servers to participate in the transfer of a single file. The terminology we use to describe a striped network file system is similar to RAID's: a collection of file data that spans the servers is called a *stripe*, and the portion of a stripe stored on a single server is called a *stripe fragment*.

The most-obvious way to organize a striped network file system is to stripe each file separately, as shown in Figure 2. We refer to this method as *file-based striping*. Each file is stored in its own set of stripes. As a result, parity is computed for each file because each stripe contains data from only one file. Although conceptually simple, file-based striping has drawbacks. First, small files are difficult to handle efficiently. If a small file is striped across all the servers, as in Figure 3(a), then each server will only store a very small piece of the file. This provides little performance benefit, since most of the access cost is due to network and disk latency; yet it incurs overhead on every server for every file access. Thus it seems better to handle small files differently than large files and to store each small file on a single

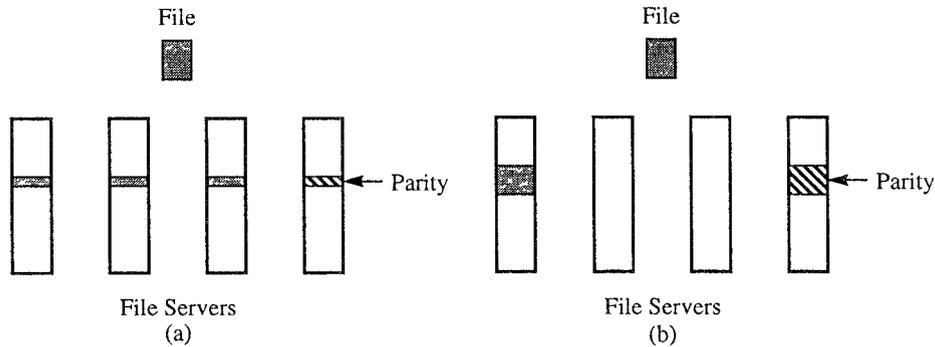


Fig. 3 File-based striping for a small file. In (a) the file is striped evenly across the servers, resulting in small fragments on each server. In (b) the entire file is placed on one server, but the parity requires as much space as the file data

server, as in Figure 3(b). This leads to problems in parity management, however. If a small file is stored on a single server, then its parity will consume as much space as the file itself, resulting in high storage overhead. In addition, the approach in Figure 3(b) can result in unbalanced disk utilization and server loading.

The second problem with file-based striping is that it requires a parity fragment to be updated each time an existing file block is modified. As with RAIDs, small updates such as this require two reads (the old data and the old parity) followed by two writes (the new data and the new parity). Furthermore, the two writes must be carried out atomically. If one write should complete but not the other (e.g., because a client or server crashed), then the parity will be inconsistent with the data; if this parity is used later for reconstructing lost data, incorrect results will be produced. There exist protocols for ensuring that two writes to two different file servers are carried out atomically [Bernstein and Goodman 1981], but they are complex and expensive.

3.3 Log-Structured File Systems and Log-Based Striping

Zebra uses techniques from log-structured file systems (LFS) [Rosenblum and Ousterhout 1991] to avoid the problems of file-based striping. LFS is a disk management technique that treats the disk like an append-only log. When new file blocks are created or existing file blocks modified, the new data are batched together and written to the end of the log in large sequential transfers. The metadata for the affected files are also updated to reflect the new locations of the file blocks. LFS is particularly effective for writing small files since it can write many files in a single transfer. In contrast, traditional file systems require at least two independent disk transfers for each file. Rosenblum and Ousterhout reported a tenfold speedup over traditional file systems for writing small files. LFS is also well suited for RAIDs because it

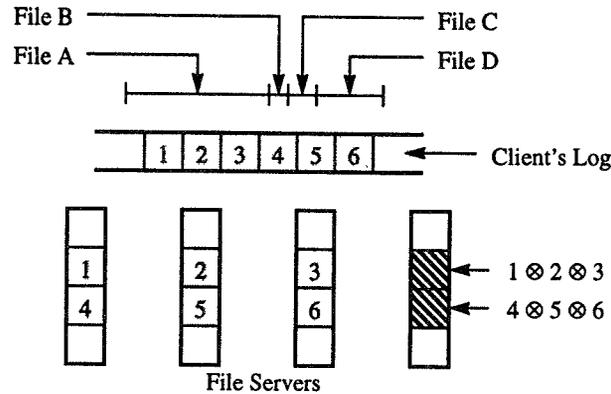


Fig. 4. Log-based striping in Zebra. Each client forms its new file data into a single append-only log and stripes this log across the servers. In this example, File A spans several servers whereas File B is stored entirely on a single server. Parity is computed for the log, not for individual files.

batches small writes together into large sequential transfers and thus avoids the expensive parity updates associated with small random writes.

Zebra can be thought of as a log-structured network file system: whereas LFS uses the logging approach at the interface between a file server and its disks, Zebra uses the logging approach at the interface between a client and its servers. Figure 4 illustrates this approach, which we call *log-based striping*. Each Zebra client organizes its new file data into an append-only log, which it then stripes across the servers. The client computes parity for the log, not for individual files. Each client creates its own log, so that each stripe in the file system contains data written by a single client.

Log-based striping has several advantages over file-based striping. The first is that the servers are used efficiently regardless of file sizes: large writes are striped, allowing them to be completed in parallel, and small writes are batched together and written to the servers in large transfers; no special handling is needed for either case. Second, the parity mechanism is simplified. Each client computes parity for its own log without fear of interactions with other clients. Small files do not have excessive parity overhead because parity is computed for the logs, not individual files. Furthermore, once a stripe is complete, its parity is never updated because file data are not overwritten in place.

The preceding description of log-based striping leaves several questions unanswered. For example, how can files be shared among client workstations if each client is writing its own log? Zebra solves this problem by introducing a central *file manager*, separate from the storage servers, that manages metadata such as directories and file attributes and supervises interactions between clients. Also, how is free space reclaimed from the logs? Zebra solves this problem with a *stripe cleaner*, which is analogous to the cleaner in a log-structured file system. Section 4 provides a more-detailed discussion of these issues and several others.

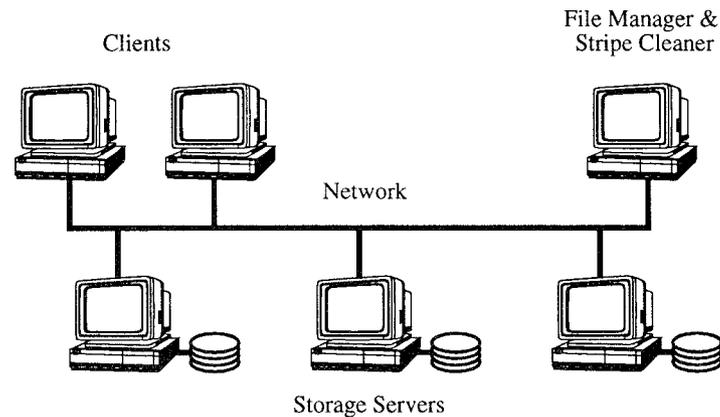


Fig. 5. Zebra schematic. Clients run applications; storage servers store data. The file manager and the stripe cleaner can run on any machine in the system, although it is likely that one machine will run both.

4. ZEBRA COMPONENTS

The Zebra file system contains four main components as shown in Figure 5: *clients*, which are the machines that run application programs; *storage servers*, which store file data; a *file manager*, which manages the file and directory structure of the file system; and a *stripe cleaner*, which reclaims unused space on the storage servers. There may be any number of clients and storage servers but only a single file manager and stripe cleaner. More than one of these components may share a single physical machine; for example, it is possible for one machine to be both a storage server and a client. If care is not taken, the single file manager and stripe cleaner may both be potential single points of failure and performance bottlenecks. Section 7 describes how the system is able to continue operation even if the file manager or stripe cleaner crashes, and the bottleneck issue is addressed in Section 8, which provides performance measurements of the prototype. The remainder of this section describes each of the components in isolation. Section 5 shows how the components work together to implement operations such as reading and writing files, and Sections 6 and 7 describe how Zebra deals with crashes.

We will describe Zebra under the assumption that there are several storage servers, each with a single disk. However, this need not be the case. For example, storage servers could each contain several disks managed as a RAID, thereby giving the appearance to clients of a single disk with higher capacity and throughput. Doing so would also provide additional redundancy: the parity maintained in the RAID would protect against disk failures, whereas the parity maintained by Zebra would protect against server failures as well. It is also possible to put all the disks on a single server; clients would treat it as several logical servers, all implemented by the same physical machine. This approach would still provide many of Zebra's benefits: clients would still batch small files for transfer over the network, and it would still

be possible to reconstruct data after a disk failure. However, a single-server Zebra system would limit system throughput to that of the one server, and the system would not be able to operate when the server is unavailable.

4.1 Clients

Clients are machines where application programs execute. When an application reads a file, the client must determine which stripe fragments store the desired data, retrieve the data from the storage servers, and return them to the application. As will be seen in the following, the file manager keeps track of where file data are stored and provides this information to clients when needed. When an application writes a file, the client appends the new data to its log by creating new stripes to hold the data, computing the parity of the stripes, and writing the stripes to the storage servers. Clients' logs do not contain file attributes, directories, or other metadata. This information is managed separately by the file manager, as described in Section 4.3.

4.2 Storage Servers

The storage servers are the simplest part of Zebra. They are just repositories for stripe fragments. As far as a storage server is concerned, a stripe fragment is a large block of bytes with a unique identifier. The identifier for a fragment consists of an identifier for the client that wrote the fragment, a sequence number that identifies the stripe uniquely among all those written by the client, and an offset for the fragment within its stripe. All fragments in Zebra are the same size, which should be chosen large enough to minimize the network and disk overheads of transferring data among the clients and the storage servers. The Zebra prototype uses 512KB fragments.

Storage servers provide five operations:

- Store a fragment*: This operation allocates space for the fragment, writes the fragment to disk, and records on disk the fragment identifier and disk location for use in subsequent accesses. The operation is synchronous: it does not complete until the fragment has been safely stored. The fragment must not already exist unless it is a parity fragment, in which case the new copy of the fragment replaces the old. This is done in a nonoverwrite manner to avoid corruption of a parity fragment in the event of a crash.
- Append to an existing fragment*: This operation is similar to storing a fragment except that it allows a client to write out a fragment in pieces if it does not have enough data to fill the entire fragment at once (this can happen, for example, if an application invokes the `fsync` system call to force data to disk). Appends are implemented atomically so that a crash during an append cannot cause the previous contents of the fragment to be lost.
- Retrieve a fragment*: This operation returns part or all of the data from a fragment. It is not necessary to read the entire fragment; a fragment identifier, offset, and length specify the desired range of bytes.
- Delete a fragment*: This operation is invoked by the stripe cleaner when the fragment no longer contains any useful data. It makes the fragment's disk space available for new fragments.

—*Identify fragments*: This operation provides information about the fragments stored by the server, such as the most-recent fragment written by a client. It is used to find the ends of the clients' logs after a crash.

Stripes are immutable once they are complete. A stripe may be created with a sequence of append operations, but nonparity fragments are never overwritten, and once the stripe is complete it is never modified except to delete the entire stripe. A parity fragment, however, can be overwritten if data are appended to a partial stripe (see Section 5.2).

4.3 File Manager

The file manager is responsible for all the information in the file system except for file data. We refer to this information as *metadata*: it includes file attributes such as protection information, block pointers that tell where file data are stored, directories, symbolic links, and special files for I/O devices. The file manager performs all the usual functions of a file server in a network file system, such as name lookup and maintaining the consistency of client file caches. However, the Zebra file manager does not store any file data; where a traditional file server would manipulate data, the Zebra file manager manipulates block pointers. For example, consider a read operation. In a traditional file system, the client requests the data from the file server; in Zebra, the client requests block pointers from the file manager, then it reads the data from the storage servers.

In the Zebra prototype, we implemented the file manager using a Sprite file server with a log-structured file system. For each Zebra file there is one file in the file manager's file system, and the "data" in this file are an array of block pointers that indicate where the blocks of data for the Zebra file are stored. This allows Zebra to use almost all the existing Sprite network file protocols without modification. Clients open, read, and cache Zebra metadata in the same manner that they cache "regular" Sprite files. There is nothing in the Zebra architecture that requires Sprite to be used as the network file system, however: any existing network file server could be used in the same way by storing block pointers in files instead of data.

The performance of the file manager is a concern because it is a centralized resource. In our implementation, clients must contact the file manager on each open and close, so communication with the file manager is a performance bottleneck when clients are accessing small files. We believe that this problem can be solved by caching naming information on clients so that the file manager need not be contacted for most opens and closes. Name caching has been used successfully in several network file systems, including AFS [Howard et al. 1988], LOCUS [Walker et al. 1983], and Echo [Hisgen et al. 1989]. There have been several published studies of the effectiveness of name caching, and they all indicate that a relatively small directory cache can absorb a large fraction of directory accesses. A study of directory reference patterns in a time-shared UNIX system [Floyd and Ellis 1989] found that a cache of 10 directories, occupying 14KB of space, would have a hit ratio of 85%. A hit ratio of 95% was attainable with a cache of only 30 directories

requiring 41KB of memory. Sheltzer et al. [1986] found that in the LOCUS network file system a 40-directory cache produced a hit ratio of 87–96%. A more-recent study of directory reference patterns in a network file system by Shirriff and Ousterhout [1992] found that a 10-directory cache had a hit ratio of 91%, whereas a 20-directory cache had a hit ratio of 97%. Despite this evidence of the benefits of name caching, we decided not to implement name caching in the Zebra prototype because it would have required major modifications to the Sprite file system. Nonetheless, we would expect any production version of Zebra to incorporate name caching, due to the large benefits that can be attained from relatively small caches.

The centralized nature of the file manager also makes its reliability a concern; this issue is addressed in Section 7.

4.4 Stripe Cleaner

When a client writes a new stripe, it is initially full of live data. Over time, though, blocks in the stripe become free, either because their files are deleted or because the file blocks are overwritten. If an application overwrites an existing block of a file, Zebra does not modify the stripe containing the block; instead it writes a new copy of the block to a new stripe. The only way to reuse free space in a stripe is to *clean* the stripe, so that it contains no live data whatsoever, and then delete the entire stripe. The storage servers can then reuse the stripe's disk space for new stripes.

The Zebra stripe cleaner runs as a user-level process and is very similar to the segment cleaner in a log-structured file system. It first identifies stripes with large amounts of free space; then it reads the remaining live blocks out of the stripes and appends them to the end of the log of the client on which the cleaner is running, thus copying the blocks to a new stripe. Once this has been done, the stripe cleaner deletes the stripe's fragments from the storage servers. Section 5.5 describes the cleaning algorithm in more detail.

5. SYSTEM OPERATION

This section describes several of the key algorithms in Zebra to show how the pieces of the system work together in operation. Most of these algorithms are similar to the approaches used in log-structured file systems, RAIDs, or other network file systems.

5.1 Communication via Deltas

A client's log contains two kinds of information: *blocks* and *deltas*. A block is simply a collection of data from a file (i.e., the information that is read and written by applications). Deltas identify changes to the blocks in a file, and are used to communicate these changes among the clients, the file manager, and the stripe cleaner. For example, a client puts a delta into its log when it writes a file block, and the file manager subsequently reads the delta to update the metadata for that block. Deltas contain the following information:

—*File identifier*: a unique identifier for a file, analogous to an i-number in a UNIX file system.

- File version*: identifies the time when the change described by the delta occurred. A file's version number increments whenever a block in the file is written or deleted. The version numbers allow deltas in different logs to be ordered during crash recovery.
- Block number*: identifies the block to which the delta applies.
- Old block pointer*: gives the fragment identifier and offset of the block's old storage location. If the delta is for a new block, then the old block pointer has a special null value. The old block pointer is used by the stripe cleaner to keep track of the live data within stripes, and by the file manager to detect races caused by the simultaneous cleaning and modification of a file, as described in Section 5.6.
- New block pointer*: gives the fragment identifier and offset for the block's new storage location. If the delta is for a block deletion, then the new block pointer has a special null value.

Deltas are created whenever blocks are added to a file, deleted from a file, or overwritten. Deltas for these events are called *update deltas*. Deltas are also created by the stripe cleaner when it copies live blocks out of stripes; this type of delta is called a *cleaner delta*. In addition, *reject deltas* are created by the file manager to resolve races between stripe cleaning and file updates.

Deltas provide a simple and reliable way for the various system components to communicate changes to files. Because deltas are stored in the client logs, and the logs are reliable, each component is ensured that any delta it writes will not be lost. When a client modifies a block of a file it only needs to write the block and the update delta to the log to ensure that both the file manager and the stripe cleaner learn of the modification. After crashes, the file manager and stripe cleaner replay deltas from the client logs to recover their state.

5.2 Writing Files

For Zebra to run efficiently, clients must collect large amounts of new file data and write them to the storage servers in large batches (ideally, whole stripes). The existing structure of the Sprite file caches made batching relatively easy. When an application writes new data they are placed in the client's file cache. The dirty data are not written to a server until either (a) they reach a threshold age (30 seconds in Sprite), (b) the cache fills with dirty data, (c) an application issues an `fsync` system call to request that data be written to disk, or (d) the file manager requests that data be written in order to maintain consistency among client caches. In many cases, files are created and deleted before the threshold age is reached, so their data never need to be written at all [Baker et al. 1991; Hartman and Ousterhout 1993].

When information does need to be written to disk, the client forms the new data into one or more stripe fragments and writes them to the storage servers. For each file block written, the client also puts an update delta into its log and increments the file's version number.

To benefit from multiple storage servers it is important for a client to transfer fragments to all of the storage servers concurrently. We added support for asynchronous remote procedure calls to Sprite to allow clients to do this. A client can also transfer the next stripe fragment to a storage server while the server is writing the current stripe fragment to disk, so that both the network and the disk are kept busy. The client computes the parity as it writes the fragments, and at the end of each stripe, the client writes the parity to complete the stripe. In the Zebra prototype the client also sends the stripe's deltas to the file manager and stripe cleaner. This improves performance by avoiding the disk accesses that would occur if the file manager and stripe cleaner were to read the deltas from the log. This optimization does not reduce the reliability of the system, however, because if the client crashes before sending the deltas, then the file manager and stripe cleaner will read the deltas from the log on their own.

If a client is forced to write data in small pieces (e.g., because an application invokes `fsync` frequently), then it fills the stripe a piece at a time, appending to the first stripe fragment until it is full, then filling the second fragment, and so on until the entire stripe is full. When writing partial stripes, the client has two choices for dealing with parity. It can delay writing the parity until the stripe is complete. This is the most-efficient alternative and is relatively safe (the client has a copy of the unwritten parity, so information will be lost only if both a disk is destroyed and the client crashes). For even greater protection, the client can update the stripe's parity fragment each time it appends to the stripe. Parity fragments written in this way include a count of the number of bytes of data in the stripe at the time the fragment was written, which is used to determine the relationship between the parity and the data after crashes. Parity updates are implemented by storage servers in a nonoverwrite fashion, so either the old parity or the new parity is always available after a crash. This is done by writing the new parity fragment to an unused location on disk, then updating the storage server's on-disk data structures to record the new location of the fragment.

The rate at which applications invoke `fsync` will have a large effect on Zebra's performance (or any other file system's) because `fsync` requires synchronous disk operations. Baker et al. [1992] found that, under a transaction-processing workload, up to 90% of the segments written on an LFS file system were partial segments caused by an `fsync`. Such a workload would have poor performance on Zebra as well. Fortunately, they found that on nontransaction-processing workloads, `fsync` accounted for less than 20% of the segments written.

The ability of Zebra clients to write directly to the storage servers opens a potential security hole. The storage servers do not implement a file abstraction; therefore, it is impossible for the servers to prevent a client from modifying a file for which it does not have permission, or from filling the storage servers with garbage. Zebra is able to prevent either of these occurrences, however, because only the file manager can modify the file system's metadata. A client modifies a file block by writing a new copy of the block to

its log, along with an update delta that describes the change. The file manager uses the information in the delta to update the file's metadata, and can easily ignore the delta if the client does not have permission to modify the file. Similarly, if a client tries to fill the storage servers with garbage, the file manager will not update the file system metadata. In both cases, the file manager issues a reject delta to indicate that the update delta was ignored, allowing the stripe cleaner to reclaim the new block written by the client. The mechanism for rejecting update deltas is described in greater detail in Section 5.6. The net result is that a malicious client cannot jeopardize the integrity of the file system, and, at worst, forces the stripe cleaner to run more often.

5.3 Reading Files

File reads in Zebra are carried out in almost the same fashion as in a nonstriped network file system. The client opens and closes the file in the same way as for a non-Zebra file; in Sprite this means a remote procedure call to the file manager for each open or close. Reading data is a two-step operation in the Zebra prototype. A client must first fetch a file's block pointers from the file manager before it can read the file blocks from the storage servers. This results in at least one extra RPC relative to a nonstriped file system. The effect of these extra RPCs is negligible for large files because as many as 2048 block pointers can be returned in an RPC, allowing the block pointers up to 8MB of data to be returned in a single RPC. For small files, however, the effect is more pronounced, since the RPC to fetch the block pointers takes 2ms even if the file manager has the block pointers cached. A better solution to reduce this additional latency is to return the block pointers for small files in the reply to the RPC to open the file. The current prototype does not implement this optimization, but it does allow clients to cache block pointers, avoiding the need to fetch them from the file manager each time a file is read.

For large files being accessed sequentially, Zebra prefetches data far enough ahead to keep all the storage servers busy. As with writing, asynchronous RPCs are used to transfer data from all the storage servers concurrently and to read the next stripe fragment on a given server from disk while transferring the previous one over the network to the client.

The Zebra prototype does not attempt to optimize reads of small files: each file is read from its storage server in a separate operation, just as for a nonstriped file system. However, it is possible to prefetch small files by reading entire stripes at a time, even if they cross file boundaries. If there is locality of file access so that groups of files are written together and then later read together, this approach might improve read performance. We speculate that such locality exists, but we have not attempted to verify its existence or capitalize on it in Zebra.

The separation of metadata management and data storage in Zebra introduces a potential security problem because the storage servers do not offer any protection for the data they store. A client can read any block of data on

the servers simply by constructing the proper block pointer. Although the current Zebra design assumes that clients are trusted, this assumption would probably not be valid for a production version of the system. One possible solution is to extend the storage server interface and functionality to allow clients to associate a “security identifier” with each file block they write. The storage servers would maintain an access control list for each identifier, specifying which clients are allowed to read blocks with that identifier. This would allow Zebra to ensure that clients can only read blocks if they are authorized to do so, while requiring only minimal modifications to the storage servers.

5.4 Client Cache Consistency

If a network file system allows clients to cache file data and allows files to be shared between clients, then cache consistency is a potential problem. For example, a client could write a file that is cached on another client; if the second client subsequently reads the file, it must discard its stale cached data and fetch the new data. We chose to use the Sprite approach to consistency, which involves flushing or disabling caches when files are opened [Nelson et al. 1988], because it was readily available, but any other approach could have been used as well. The only changes for Zebra occur when a client flushes a file from its cache. Instead of just returning dirty data to a file server, the Zebra client must write the dirty blocks to a storage server, and then the file manager must process all the deltas for the blocks so that it can provide up-to-date block pointers to other clients.

5.5 Stripe Cleaning

The first step in cleaning is to select one or more stripes to clean. To do this intelligently, the stripe cleaner needs to know how much live data are left in each stripe. Deltas are used to compute this information. The stripe cleaner processes the deltas from the client logs and uses them to keep a running count of space utilization in each existing stripe. For each delta, the cleaner increments the utilization of the stripe containing the new block (if any) and decrements the utilization of the stripe that contained the old block (if any). In addition, the cleaner appends all the deltas that refer to a given stripe to a special file for that stripe, called the *stripe status file*, whose use will be described later. The stripe status files are stored as ordinary Zebra files. Note that a single delta can affect two different stripes; a copy of the delta is appended to the status files for both stripes.

During cleaning the stripe cleaner first looks for stripes with no live data. If any are found, then the cleaner deletes the stripes’ fragments from the storage servers and deletes the corresponding stripe status files. If there are no empty stripes, and more free space is needed, then the cleaner chooses one or more stripes to clean. The policy it uses for this is identical to the one described by Rosenblum and Ousterhout [1991]; that is, a cost-benefit analysis is done for each stripe, which considers both the amount of live data in the stripe and the age of the data.

There are two issues in cleaning a stripe: identifying the live blocks and copying them to a new stripe. The stripe status files make the first step easy: the cleaner reads the deltas in the stripe's status file and finds blocks that have not yet been deleted. Without the stripe status files, this step would be much more difficult, because the deltas that cause blocks to become free could be spread throughout the stripes in the file system.

Once the live blocks have been identified, the stripe cleaner (which executes as a user-level process) copies them to a new stripe using a special kernel call. The kernel call reads one or more blocks from storage servers, appends them to its client log along with the corresponding cleaner deltas, and writes the new log contents to the storage servers. The kernel call for cleaning blocks has the same effect as reading and rewriting the blocks except that (a) it does not open the file or invoke cache consistency actions, (b) it need not copy data out to the user-level stripe cleaner process and back into the kernel again, (c) it does not update last-modified times or version numbers for files, and (d) it generates cleaner deltas instead of update deltas.

One concern about the stripe cleaner is how much of the system's resources it will consume in copying blocks. We do not have measurements of Zebra under real workloads, but we expect the fraction of the system resources consumed by the stripe cleaner to be comparable to those for other log-structured file systems running the same workloads, since Zebra's file layout and cleaning algorithm are similar. In a transaction-processing benchmark on a nearly full disk, Seltzer et al. [1993] found that cleaning accounted for 60–80% of all write traffic and significantly affected system throughput. Unfortunately, that study was unable to account fully for the surprisingly poor LFS performance [Ousterhout 1995], leading to the publication of a more-extensive study [Seltzer et al. 1995]. The new study found that LFS performance on a transaction-processing benchmark was at most 10% worse than FFS. The reasons for the LFS performance degradation are still not fully explicable, indicating that further study is warranted.

Despite the controversy surrounding LFS performance on transaction-processing workloads, several studies have shown the cleaning cost to be minimal on more-typical workstation workloads. Seltzer et al. [1993] found LFS cleaning costs to be negligible on a software development benchmark. Rosenblum measured production usage of LFS on Sprite for several months and found that only 2–7% of the data in stripes that were cleaned were live and needed to be copied [Rosenblum and Ousterhout 1991]. Based on these measurements, we believe that the cleaning overhead will be low for typical workstation workloads, but more work may be needed to reduce the overheads for transaction-processing workloads.

5.6 Conflicts Between Cleaning and File Access

It is possible for an application to modify or delete a file block at the same time that the stripe cleaner is copying it. Without any synchronization, a client could modify the block after the cleaner reads the old copy but before the cleaner rewrites the block, in which case the new data would be lost in

favor of the rewritten copy of the old data. In the original LFS, this race condition was avoided by having the cleaner lock files to prevent them from being modified until after cleaning was finished. Unfortunately, this produced lock convoys that effectively halted all normal file accesses during cleaning and resulted in significant pauses.

Zebra's stripe cleaner uses an optimistic approach similar to that of Seltzer et al. [1993]. It does not lock any files during cleaning nor invoke any cache consistency actions. Instead, the stripe cleaner copies the block and issues a cleaner delta, assuming optimistically that its information about the block is correct and that the block has not been updated recently. If, in fact, the block is updated while the cleaner is cleaning it, an update delta will be generated by the client that made the change. Regardless of the order in which these deltas arrive at the file manager, the file manager makes sure that the final pointer for the block reflects the update delta, not the cleaner delta. This approach results in wasted work by the cleaner in the unusual case where a conflict occurs, but it avoids synchronization in the common case in which there is no conflict.

The file manager detects conflicts by comparing the old block pointer in each incoming delta with the block pointer stored in the file manager's metadata; if they are different, it means that the block was simultaneously cleaned and updated. Table I shows the various scenarios that can occur. The first two scenarios represent the cases where there is no conflict: the delta's old block pointer matches the file manager's current block pointer, so the file manager updates its block pointer with the new block pointer in the delta. If an update delta arrives with an old block pointer that does not match, it can only mean that the block was cleaned (any other update to the block is prevented by the cache consistency protocol); the file manager updates its block pointer with the new block pointer from the delta. If a cleaner delta arrives with an old block pointer that does not match, it means that the block has already been updated, so the cleaned copy is irrelevant: the cleaner delta is therefore ignored.

In both cases where the file manager detects a conflict, it generates a reject delta, which is placed in the client log of the machine on which the file manager is running. The old block pointer in the reject delta refers to the cleaned copy of the block, and the new pointer is null to indicate that this block is now free. The reject delta is used by the stripe cleaner to keep track of stripe usage; without it the stripe cleaner would have no way of knowing that the file manager ignored the block generated by the cleaner, leaving the space it occupies unused.

It is also possible for an application to read a block at the same time that it is being cleaned. For example, suppose that a client has retrieved a block pointer from the file manager, but the block is moved by the cleaner before the client retrieves it. If the client then tries to use the out-of-date block pointer, one of two things will happen. If the block's stripe still exists, then the client can use it safely, since the cleaner did not modify the old copy of the block. If the stripe has been deleted, then the client will get an error from the storage server when it tries to read the old copy. This error indicates that the

Table I. File Manager Delta Processing

Type of Delta	Block Pointer Matches?	Update Pointer ?	Issue Reject Delta?
Update	Yes	Yes	No
Cleaner	Yes	Yes	No
Update	No	Yes	Yes
Cleaner	No	No	Yes

When a delta arrives at the file manager, the old block pointer in the delta is compared to the current block pointer. If they do not match (the bottom two scenarios), then a conflict has occurred.

block pointer is out of date: the client simply discards the pointer and fetches an up-to-date version from the file manager.

5.7 Adding a Storage Server

Zebra's architecture makes it easy to add a new storage server to an existing system. All that needs to be done is to initialize the new server's disk(s) to an empty state and notify the clients, file manager, and stripe cleaner that each stripe now has one more fragment. From this point on clients will stripe their logs across the new server. The existing stripes can be used as is even though they do not cover all the servers; in the few places where the system needs to know how many fragments there are in a stripe (such as reconstruction after a server failure), it can detect the absence of a fragment for a stripe on the new server and adjust itself accordingly. Over time the old stripes will gradually be cleaned, at which point their disk space will be used for longer stripes that span all the servers. Old stripes are likely to be cleaned before new ones since they contain less live data. If it should become desirable for a particular file to be reallocated immediately to use the additional bandwidth of the new server, this can be done by copying the file and replacing the original with the copy.

5.8 Removing a Storage Server

Removing a storage server is a three-step process. First, the system administrator must verify that there is enough free space in the system to accommodate the loss of a server. If not, files must be deleted until the total amount of free space exceeds the storage capacity of the server. Second, the clients, file manager, and stripe cleaner are notified that stripes now have one less fragment. Once this is done, any new stripes created will not use the server that is being decommissioned. Third, the stripe cleaner is instructed to clean all the old stripes. This has the effect of moving live data from the unwanted server to the remaining servers. When the stripe cleaner has finished, the

unwanted server will not contain any live data and can be safely removed from the system.

6. RESTORING CONSISTENCY AFTER CRASHES

There are two general issues that Zebra must address when a client or server machine crashes: consistency and availability. If a crash occurs in the middle of an operation, then data structures may be left in a partially modified state after the crash. For example, the file manager might crash before processing all the deltas written by clients; when it reboots, its metadata will not be up-to-date with respect to information in the clients' logs. This section describes how Zebra restores internal consistency to its data structures after crashes. The other issue is availability, which refers to the system's ability to continue operation even while a component is down. Zebra's approach to availability is described in Section 7.

In many respects the consistency issues in Zebra are the same as in other network file systems. For example, the file manager will have to restore consistency to all its on-disk data structures. Since the file manager uses the same disk structures as a nonstriped file system, it can also use the same recovery mechanism. In the Zebra prototype, the metadata is stored in a log-structured file system, so we use the LFS recovery mechanism described by Rosenblum and Ousterhout [1991]. The file manager must also recover the information that it uses to ensure client cache consistency; for this Zebra uses the same approach as in Sprite, which is to let clients reopen their files to rebuild the client cache consistency state [Nelson et al. 1988]. If a client crashes, then the file manager cleans up its data structures by closing all the client's open files, also in the same manner as Sprite.

However, Zebra introduces certain consistency problems that are not present in other file systems. These problems arise from the distribution of system state among the storage servers, file manager, and stripe manager; each problem is a potential inconsistency between system components. The first problem is that stripes may become internally inconsistent (e.g., some of the data or parity may be written but not all of it); the second problem is that information written to stripes may become inconsistent with metadata stored on the file manager; the third problem is that the stripe cleaner's state may become inconsistent with the stripes on the storage servers. These problems are discussed separately in the remainder of this section.

The solutions to all the consistency issues are based on logging and checkpoints. Logging means that operations are ordered so that it is possible to tell what happened after a particular time and to revisit those operations in order. Logging also implies that information is never modified in place, so if a new copy of information is incompletely written, the old copy will still be available. A checkpoint defines a system state that is internally consistent. To recover from a crash, the system initializes its state to that of the most-recent checkpoint, then reprocesses the portion of the log that is newer than the checkpoint.

The combination of these two techniques allows Zebra to recover quickly after crashes. It need not consider any information on disk that is older than the most-recent checkpoint. Zebra is similar to other logging file systems such as LFS, Episode [Chutani et al. 1992], and the Cedar File System [Hagmann 1987] in this respect. In contrast, file systems without logs, such as the BSD Fast File System [McKusick et al. 1984], cannot tell which portions of the disk were being modified at the time of a crash, so they must rescan all the metadata in the entire file system during recovery.

6.1 Internal Stripe Consistency

When a client crashes, it is possible for fragments to be missing from stripes that were in the process of being written. The file manager detects client crashes and recovers on behalf of the client: it queries the storage servers to identify the end of the client's log and verifies that any stripes that could have been affected by the crash are complete. If a stripe is missing a single fragment, then the missing data can be reconstructed using the other stripes in the fragment. If a stripe is missing more than one fragment, then it is discarded along with any subsequent stripes in the same client's log. This means that data being written at the time of a crash can be lost or partially written, just as in other file systems that maintain UNIX semantics.

When a storage server crashes and recovers, two forms of stripe inconsistency are possible. If a stripe fragment was being written at the time of the crash, then it might not have been completely written. To detect incomplete stripe fragments, Zebra stores a simple checksum for each fragment. After a storage server reboots, it verifies the checksums for fragments written around the time of the crash and discards any that are incomplete.

The other inconsistency after a storage server crash is that it will not contain fragments for new stripes written while it was down. After the storage server reboots, it queries other storage servers to find out what new stripes were written. Then it reconstructs the missing fragments as described in Section 7.2 and writes them to disk. The storage servers in the prototype do not perform this reconstruction after a crash.

6.2 Stripes Versus Metadata

The file manager must maintain consistency between the client logs and its metadata. To do this it must ensure that it has processed all the deltas written by clients and updated its metadata accordingly. During normal operation, the file manager keeps track of its current position in each client's log, and at periodic intervals it forces the metadata to disk and writes a checkpoint file containing the current log positions. If a client crashes, the file manager checks with the storage servers to find the end of the client's log and make sure it has processed all the deltas in the log. If the file manager crashes, then when it reboots it processes all the deltas that appear in the client logs after the positions stored in the last checkpoint, thereby bringing the metadata up-to-date. A checkpoint is relatively small (a few hundred bytes) as it only contains current log positions for each client, but it does have

a performance impact because the metadata is flushed before the checkpoint is written. Decreasing the checkpoint interval improves the file manager's recovery time at the expense of normal operation; we anticipate that a checkpoint interval on the order of several minutes will provide acceptable recovery time without significantly affecting the system performance.

There are two complications in replaying deltas, both of which are solved with version numbers. Some of the deltas may have already been processed and applied to the metadata. This will happen if the file manager crashes in the interval between writing the metadata out to disk and writing a checkpoint. If an update delta is encountered that has already been applied, then its version number will be less than that of the file, and it is therefore ignored. As in normal operation, a cleaner delta is applied only if its old block pointer matches the file manager's current block pointer.

The other complication is that a file could have been modified by several different clients, resulting in deltas for the file in several client logs. The file manager must replay the deltas for each file in the same order that they were originally generated. If the file manager encounters a delta during replay whose version number is greater than the file's version number, it means that there are deltas in some other client log that must be replayed first. In this case the file manager must delay the processing of the delta, and the other unprocessed deltas in that client's log, until all the intervening deltas have been processed from the other client logs.

6.3 Stripes Versus Cleaner State

In order for the stripe cleaner to recover from a crash without completely reprocessing all the stripes in the file system, it checkpoints its state to disk at regular intervals. The state includes the current utilizations for all the stripes plus a position in each client log, which identifies the last delta processed by the stripe cleaner. Any buffered data for the stripe files are flushed before writing the checkpoint.

When the stripe cleaner restarts after a crash, it reads in the utilizations and log positions, then starts processing deltas again at the saved log positions. If a crash occurs after appending deltas to a stripe status file but before writing the next checkpoint, then the status file could end up with duplicate copies of some deltas, since the stripe cleaner will process those deltas both before and after the crash. These duplicates are easily weeded out when the cleaner processes the status files.

7. AVAILABILITY

Our goal for Zebra is for the system to continue to provide service even if some of its machines have crashed. A single failure of a storage server, the file manager, or the stripe cleaner should not prevent clients from accessing files, nor should any number of client failures affect the remaining clients. Each of the system components is discussed separately in the following sections. The prototype does not implement all these features, as noted.

7.1 Client Crashes

The only way that one client can prevent other clients from accessing files is through the cache consistency protocol: if a client has a file open and cached, then other clients' access to the file is restricted to prevent inconsistencies. After a client crash, the file manager closes all the open files on the client, thus allowing those files to be cached by other clients.

7.2 Storage Server Crashes

Zebra's parity mechanism allows it to tolerate the failure of a single storage server using algorithms similar to those described for RAIDs [Patterson et al. 1988]. To read a file while a storage server is down, a client must reconstruct any stripe fragment that was stored on the down server. This is done by computing the parity of all the other fragments in the same stripe; the result is the missing fragment. Writes intended for the down server are simply discarded; the storage server will reconstruct them when it reboots, as described in Section 6.1. In the prototype, clients are capable of reconstruction, but only under manual control. Clients do not yet automatically reconstruct fragments when a server crashes.

Reconstruction is relatively inexpensive during large sequential reads: all the fragments of the stripe are needed anyway, so the only additional cost is the parity calculation. For small reads reconstruction is expensive because it requires reading all the other fragments in the stripe. If small reads are distributed uniformly across the storage servers, then reconstruction doubles the average cost of a read.

7.3 File Manager Crashes

The file manager is a critical resource for the entire system because it manages all the file system metadata. If the metadata are stored nonredundantly on the file manager, then the file system will be unusable whenever the file manager is down, and the loss of the file manager's disk will destroy the file system. These problems are avoided by using the Zebra storage servers to store the file manager's metadata. Instead of using a local disk, the file manager writes the metadata to a virtual disk implemented as a Zebra file. The metadata is stored in the virtual disk file which, in turn, is stored in the file manager's client log and striped across the storage servers with parity, just like any other Zebra file. This provides higher performance for the metadata than storing it on a local disk, and improves its availability and integrity. This approach also allows the file manager to run on any machine in the network, because it does not depend on having local access to a disk. If the file manager's machine should break, then the file manager can be restarted on another machine. Of course, if the file manager crashes, Zebra will be unavailable until the file manager restarts, but it should be possible to restart the file manager quickly [Baker and Sullivan 1992]. A similar approach has been proposed by Cabrera and Long [1991] for the Swift file system for making its storage mediator highly available.

7.4 Stripe Cleaner Crashes

The technique used to make the stripe cleaner highly available is similar to that used for the file manager. The key is that access to the stripe cleaner's state must not be confined to the machine on which the stripe cleaner runs. If this is not the case (as it would be if the stripe cleaner stored its state on a local disk), the stripe cleaner would be vulnerable to a failure of the machine on which it is running. For this reason, the stripe cleaner stores its state in a collection of Zebra files, so that the files are stored in the stripe cleaner's client log and striped across the servers. If the machine on which the stripe cleaner is running fails, the stripe cleaner is simply restarted on a different machine.

8. PROTOTYPE STATUS AND PERFORMANCE

The implementation of the Zebra prototype began in April, 1992. As of March, 1995, Zebra supports all the usual UNIX file operations; the cleaner is functional; and clients can write parity and reconstruct fragments. The file manager and cleaner both checkpoint their states and are able to recover after a failure. The prototype does not implement all the crash recovery and availability features of Zebra, however; clients do not automatically reconstruct stripe fragments when a storage server crashes; storage servers do not reconstruct missing fragments after a crash; and the file manager and stripe cleaner are not automatically restarted. We have simplified the prototype by choosing not to implement name caching or support for concurrent write-sharing.

The rest of this section contains preliminary performance measurements made with the prototype. The measurements show that Zebra provides an improvement factor of 4–5 in throughput for large reads and writes relative to either NFS or the Sprite file system, but its lack of name caching prevents it from providing much of a performance advantage for small files. We estimate that a Zebra system with name caching would also provide substantial performance improvements for small writes.

8.1 Experimental Setup

For our measurements we used a cluster of DECstation-5000 Model 200 workstations connected by an FDDI ring (maximum bandwidth 100Mb/second). The workstations are rated at about 20 integer SPECmarks, and each contains 32MB of memory. In our benchmarks, the memory bandwidth is at least as important as CPU speed; these workstations can copy large blocks of data from memory to memory at about 12MB/second, but copies to or from disk controllers and FDDI interfaces run at only about 8MB/second. This limits the bandwidth of Sprite RPCs over the FDDI to 3.1MB/second, despite the capacity of the network itself to support higher bandwidths. Each storage server is equipped with a single RZ57 disk with a capacity of about one gigabyte and an average seek time of 15ms. Data can be read from the disk to the host at about 1.6MB/second, and written at about 1.1MB/second.

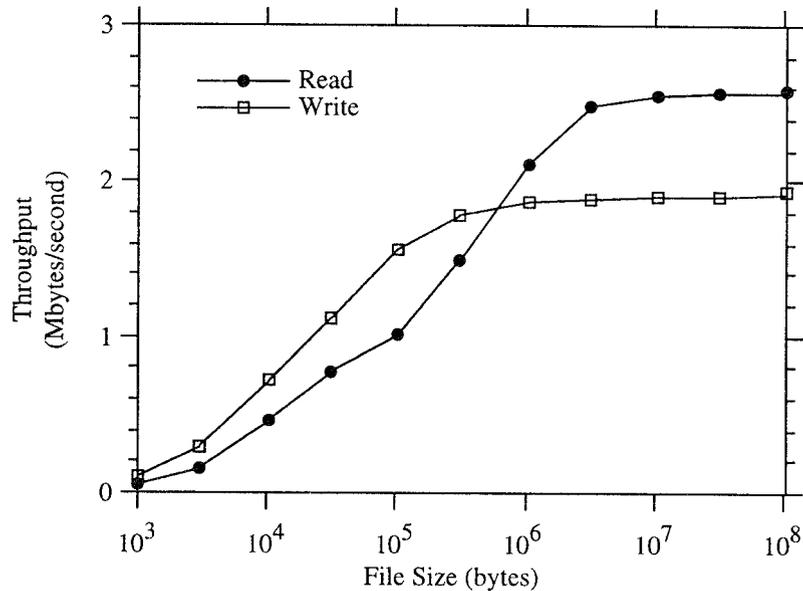


Fig. 6. Throughput versus file size. A single client reads or writes files of varying size to five storage servers. Writing small files is faster than reading due to Zebra's ability to batch small writes; writing of large files is slower than reading due to the parity computation.

There were a total of nine workstations available for running these experiments. The minimum configuration tested consisted of one client, one storage server, and one file manager. In the maximum configuration there were three clients, five storage servers, and one file manager. The relatively small number of storage servers available eliminates the possibility of the FDDI ring being a performance bottleneck; the five servers are capable of transferring data at a maximum rate of 8MB/second, which is well below the FDDI's maximum bandwidth.

During the measurements the file manager did not generate checkpoints, nor was the stripe cleaner running. Each data point was collected by running the benchmark ten times and averaging the results. Standard deviations are reported, but not shown in the graphs because most are too small to be discernible.

For comparison we also measured a standard Sprite configuration and an Ultrix/NFS configuration. The Sprite system used the same collection of workstations as the Zebra experiments, except that the standard Sprite network file system is used instead of Zebra, and the Sprite log-structured file system was used as the disk storage manager on the file server. The NFS configuration used the same client configuration as Zebra, but the file server had a slightly faster CPU and slightly faster disks. The NFS server also included a 1MB PrestoServe nonvolatile RAM card for buffering disk writes.

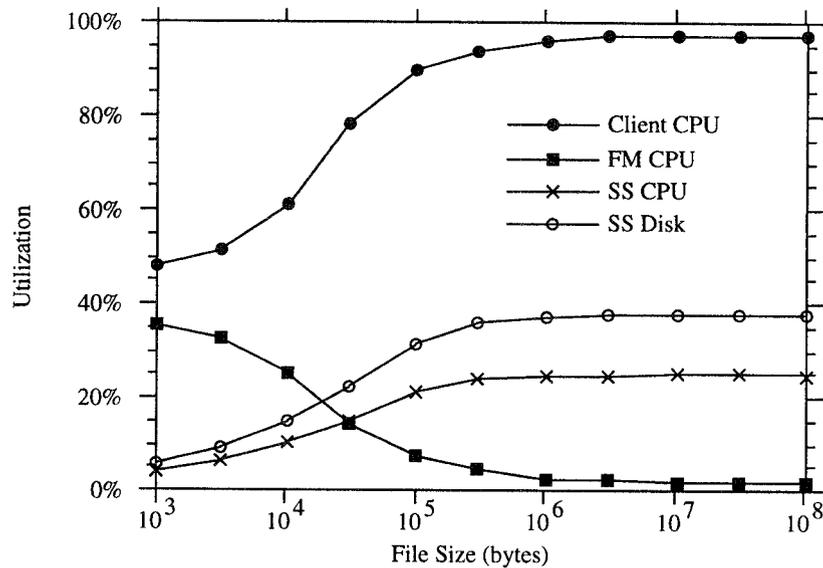


Fig. 7. Write resource utilizations. For small files the time required to open and close the files causes low client utilization and high file manager utilization (see Section 8.4); for large files the client CPU saturates. The storage server utilizations were measured on one of the five servers in the system. The maximum standard deviation for the measurements is 3%.

8.2 Performance Versus File Size

The first experiments varied file size and measured file system throughput and resource utilizations while reading or writing files. In each experiment there was one client, one file manager, four data servers (servers that store data fragments as opposed to parity fragments), and one parity server. An application that wrote or read files ran on the client, and the elapsed time and resource utilizations were measured. In order to measure the steady-state performance of the system, startup and end effects for files smaller than 300KB in size were avoided by having the application read or write 1000 files in each test. For files of size 300KB or greater, each test read or wrote enough files to transfer at least 100MB of data. Figure 6 shows the results. The standard deviations for the read and write measurements are less than 78KB/second and 50KB/second, respectively. As can be seen, throughput increases dramatically as file size increases. For large files, reading is faster than writing; this is because the client CPU is saturated when accessing large files, and writing has the additional overhead of computing parity. For small files, writing is faster than reading; this is because Zebra's log-based striping writes many small files to the servers at a time.

Although Zebra batches small file writes, Figure 6 shows that write performance decreases as file size decreases, indicating that there are still significant per-file overheads associated with writing files that batching does not eliminate. The write bandwidth for 100MB files is more than ten times the bandwidth for 1KB files. Further evidence of these per-file overheads can be found in the resource utilizations during the write tests, shown in Figure 7.

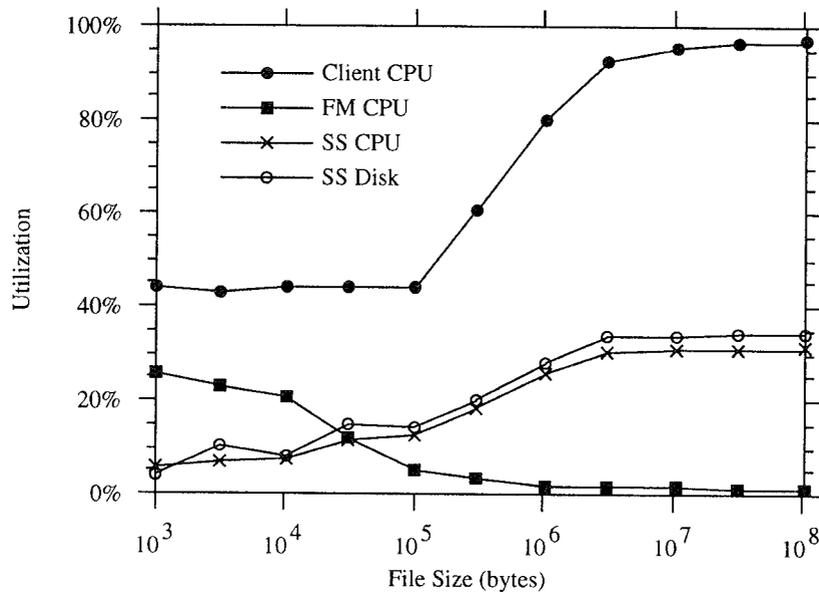


Fig. 8. Read resource utilizations. The curves are similar in shape to those for writing, except that the knees occur at larger file sizes. The storage server utilizations were measured on one of the five servers in the system. For small files the loads were not equal on all the servers, causing the fluctuations in the curves. The standard deviations for all measurements are less than 2%.

The utilizations show that the bottleneck when writing large files is the client CPU, which is over 97% utilized, whereas the file manager CPU is less than 2% utilized. The client is spending all its time copying data between the application program and the kernel, between the kernel and the network interface, and in performing the parity computation. This is a favorable result, because it indicates that Zebra's write performance for large files will track client performance improvements, and that the file manager should be able to support at least 50 clients running this workload.

For small-file writes, however, the bottleneck is no longer the client CPU. As can be seen, when writing 1KB files the client CPU is less than 50% utilized, and the file manager CPU is more than 35% utilized. The source of the high overhead on the file manager is the processing of file open and close requests from the client, and is described in more detail in Section 8.4. In short, each open or close of a file by a client results in a request/response message exchange with the file manager. This not only increases the overhead on the file manager, but reduces the overall performance of the benchmark because the client is idle while the file manager processes the open and close requests.

A similar situation occurs when reading files, as shown in Figure 8. The bottleneck when reading large files is the client CPU. The cost of opening and closing files when reading small files decreases client CPU utilization and increases file manager CPU utilization. The basic shapes of the curves are the same as in Figure 7, but the knees of the curves occur at larger file sizes.

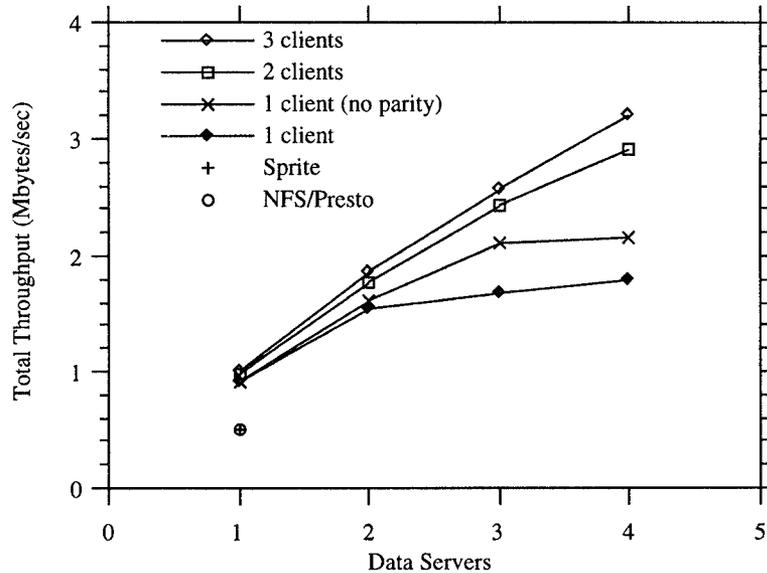


Fig. 9. Total system throughput for large-file writes. Each client ran a single application that wrote a 12MB file and then flushed the file to disk. In multiserver configurations data were striped across all the servers with a fragment size of 512KB. Each Zebra configuration also included a parity server in addition to the data servers. The maximum standard deviations are 0.06MB/second for Zebra, 0.01MB/second for Sprite, and 0.16MB/second for NFS.

This is because Zebra can batch many small writes together, but it cannot do the same for reads. Thus, larger files are required to use the servers efficiently.

8.3 Performance Versus Number of Servers

For the next set of experiments the file size was fixed at 12MB, and the number of servers and clients was varied. The first benchmark consists of an application that writes a single very large file (12MB) and then invokes `fsync` to force the file to disk. We ran one or more instances of this application on different clients (each writing a different file) with varying numbers of servers, and computed the total throughput of the system (total number of bytes written by all clients divided by elapsed time). Figure 9 graphs the results.

Even with a single client and server, Zebra runs at about twice the speed of either NFS or Sprite. This is because Zebra uses large blocks, and its asynchronous RPC allows it to overlap disk operations with network transfers. The limiting factor in this case is the server's disk system, which can write data at about 1.1MB/second. As servers are added in the single-client case, Zebra's performance increases by a factor of 2 to 1.9MB/second with four servers. The nonlinear speedup in Figure 9 occurs because of startup effects caused by Sprite's write-back cache. The client does not begin to write its cache to the servers until it is full, causing the benchmark to run in two

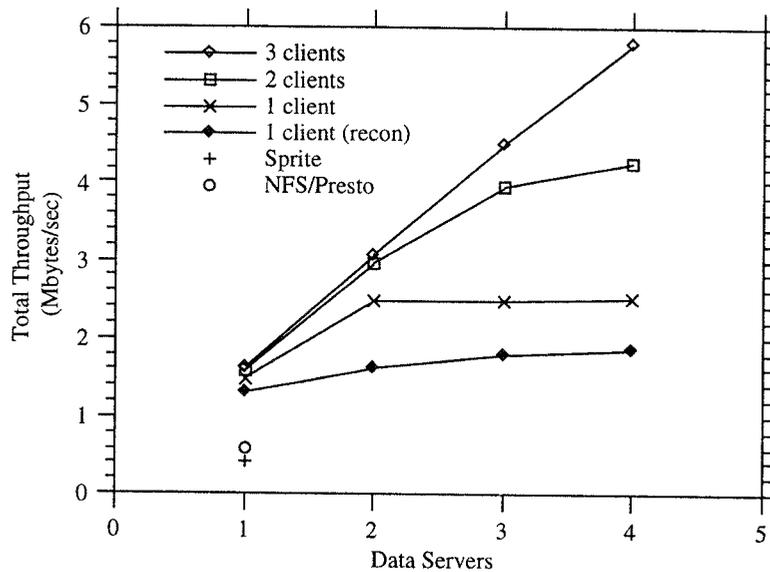


Fig. 10. Throughput for large-file reads. Each client ran a single application that read a 12MB file. In multiserver configurations data were striped across all the servers with a fragment size of 512KB. The line labeled “1 client (recon)” shows reconstruction performance: one server was unavailable, and the client had to reconstruct the missing stripe fragments. In addition to the servers storing file data, each Zebra configuration had a server storing parity. The maximum standard deviations are 0.17MB/second for Zebra, 0.01MB/second for Sprite, and 0.01MB/second for NFS.

phases. In the first phase, the application fills the kernel’s file cache by writing the file; in the second phase, the client’s kernel flushes its cache by transferring stripes to the servers. These phases are not overlapped, and only the second phase benefits from additional storage servers. Performance with two or more clients is limited entirely by the servers, so it scales linearly with the number of servers.

Figure 9 also shows the throughput for a single client when it does not generate parity. Zebra incurs almost no overhead for parity aside from the obvious overhead of writing more data to more servers. If there is only one data server in the system, then the server is the bottleneck, and the client has plenty of time to compute and write the parity. Once there are more than two data servers in the system, the client becomes the bottleneck, and the cost of writing the parity begins to have an effect.

Figure 10 shows Zebra’s throughput for reading large files. Zebra’s performance for reading is better than for writing because the servers can read data from their disks at the full SCSI bandwidth of 1.6MB/second. Thus a single client can read a file at 1.5MB/second from a single server, and three clients can achieve a total bandwidth of 5.8MB/second with four data servers. Two servers can saturate a single client, however, causing the single-client curve in Figure 10 to level off at 2.5MB/second. At that speed the client is spending most of its time copying data from a network buffer into the file cache and then from the file cache to the application. This overhead could be

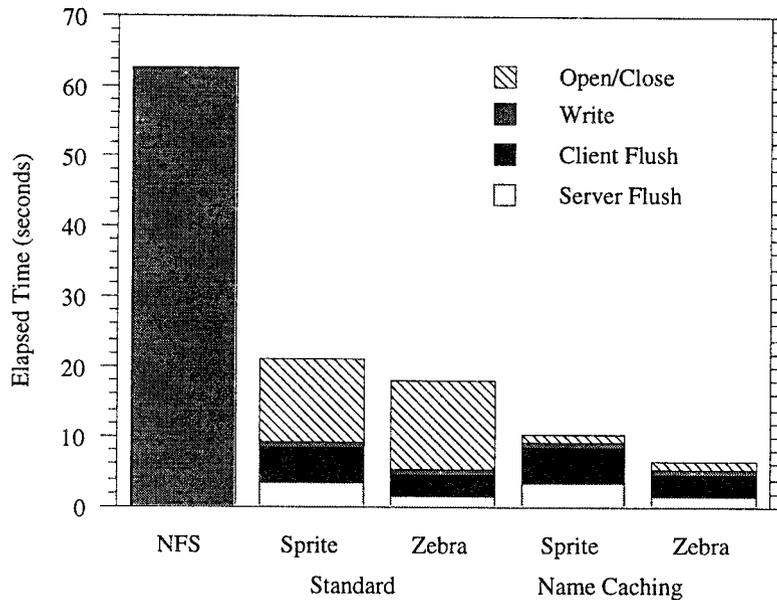


Fig. 11. Performance for small writes. A single client created 2048 files, each 1KB in length, then flushed all the files to a single server. The elapsed time is divided into four components: the time to open and close the files, the time for the application to write the data, the time for the client to flush its cache to the server's cache, and the time for the server to flush its cache to disk. For NFS, each file was flushed as it was closed. The two rightmost bars are estimates for Sprite and Zebra if name caching was implemented. The maximum standard deviations for the components are 0.98 seconds for NFS, 0.24 seconds for Sprite, and 0.54 seconds for Zebra.

reduced significantly by modifying the Sprite kernel to use the FDDI interface's DMA capability to transfer incoming network packets directly into the file cache, thus eliminating one of the data copies.

The performance of reads that require reconstruction is shown in the line labeled "1 client (recon)" in Figure 10. In this test one of the storage servers was unavailable, and the client had to reconstruct any stripe fragments stored on that server by reading all the other fragments in each stripe and computing their parity. With only one data server, the throughput during reconstruction is only slightly less than in normal operation; this is because each parity block in a system with only one data server is a mirror image of its data block, and therefore, reconstruction does not require any additional computation by the client. The throughput does not increase much with additional servers because the client CPU has saturated due to additional copying and exclusive-or operations to reconstruct the missing data.

8.4 Small-File Performance

Figure 11 shows the elapsed time for a single client to write small files. In the NFS and Sprite tests the client was writing to a single file server, whereas the Zebra test used two storage servers (one stored parity) and one file manager. Although Zebra is substantially faster than NFS for this bench-

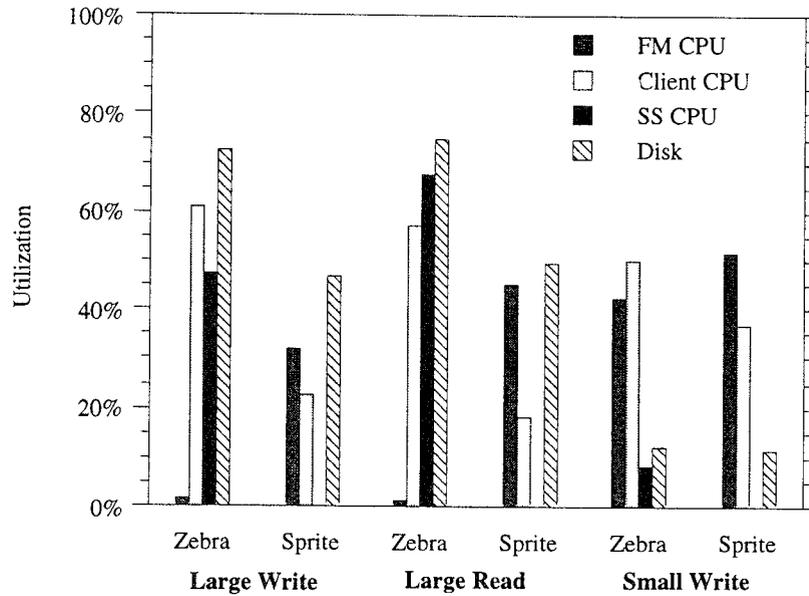


Fig. 12. Resource utilizations. Utilizations of the file manager (FM) CPU and disk, client CPU, storage server (SS) CPU, and the disk during the previous three benchmarks. The Zebra system consisted of a single client, a single file manager, and two storage servers, one of which stored parity; the Sprite system consisted of a single client and a single file server, which served as both file manager and storage server. The standard deviations for all measurements are less than 1%.

mark, it is only about 15% faster than Sprite. The main reason for this is that neither Zebra nor Sprite caches naming information; each open and close requires a separate RPC to either the file server or file manager, and the figure shows that most of the time is spent in these RPCs. The rightmost bars in the figure estimate the times for Sprite and Zebra if name caching were implemented; the estimates were made by running the same benchmark directly on a Sprite file server. These estimates show that the addition of name caching would reduce the time required to open and close the files by almost 90%; this result agrees with the published studies of directory reference patterns. Zebra is significantly faster than Sprite during the cache-flush portion of the benchmark. Both systems merge the small files into large blocks for writing, but Sprite does not do it until the data have reached the server: each file is transferred over the network in a separate message exchange. Zebra batches the file together before transferring over the network, which is more efficient because it amortizes the overhead associated with a network transfer over more bytes of data.

8.5 Resource Utilization

Figure 12 shows the utilization of various system components during the previous three benchmarks, both for Zebra and for Sprite. For large reads

and writes the Zebra file manager's CPU is almost idle; the system could scale to dozens of storage servers before the file manager becomes a performance bottleneck. When compared to Sprite, Zebra has higher utilizations of the client CPU, server CPU, and server disk; this is because Zebra is running the benchmark more quickly.

For small writes both Zebra and Sprite spend most of their time in synchronous RPCs to open and close files. In both systems the sum of client CPU utilization and file manager CPU utilization is nearly 100%; it cannot exceed 100% because the RPCs do not allow much overlap in processing between the two CPUs. In both Zebra and Sprite it appears that the server CPU will saturate with the addition of a second client; without name caching the server CPU will be a performance bottleneck.

In all the benchmarks the Zebra client has higher CPU utilization than the file manager; the opposite is true for the Sprite system. This indicates that Zebra is better able to take advantage of client performance improvements, because the overall performance of the benchmark is more heavily dependent on the client performance than in Sprite. For the large read and write benchmarks the Zebra file manager is less than 5% utilized, whereas the Sprite file server is more than 30% utilized.

9. RELATED WORK

Most of the key ideas in Zebra were derived from prior work in disk arrays and log-structured file systems. However, there are many other related projects in the areas of striping and availability.

RAID-II [Drapeau et al. 1994], DataMesh [Wilkes 1992], and TickerTAIP [Cao et al. 1993] all use RAID technology to build high-performance file servers. RAID-II uses a dedicated high-bandwidth data path between the network and the disk array to bypass the slow memory system of the server host. DataMesh is an array of processor/disk nodes connected by a high-performance interconnect, much like a parallel machine with a disk on each node. TickerTAIP is a refinement of DataMesh that focuses on distributing the functions of the traditionally centralized RAID controller across multiple processors, thus removing the controller as a single point of failure. In all these systems the striping is internal to the server, whereas in Zebra the clients participate in striping files.

RADD (Redundant Array of Distributed Disks) [Schloss and Stonebraker 1990] is similar to RAID in that it uses parity to withstand the loss of a disk, but it differs by separating the disks geographically to decrease the likelihood of losing multiple disks. Furthermore, RADD does not stripe data; the data stored on each disk are logically independent; thus RADD does not improve the performance of individual data accesses.

Several other striping file systems have been built. Most, such as sfs [LoVerso et al. 1993], Bridge [Dibble et al. 1988], and CFS [Pierce 1989], stripe across I/O nodes in a parallel computer; to our knowledge only one, Swift [Cabrera and Long 1991], stripes across servers in a network file

system. All these systems use file-based striping, so they work best with large files. Swift's performance while reading and writing large files improves nearly linearly as the number of servers increases to three; but the CPUs and disks for Swift are much slower than those for Zebra, so its absolute performance is lower than Zebra's. The Swift prototype has recently been reimplemented to incorporate the reliability mechanisms described in the Swift architecture [Long et al. 1994]. The prototype can now support a variety of parity organizations. Measurements show that the parity computation incurs a significant overhead, so that the performance of a five-server system with parity enabled is only 53% of the original Swift prototype with the same number of servers.

There have also been several recent research efforts to improve the availability of network file systems, such as Locus [Walker et al. 1983], Coda [Satyanarayanan et al. 1990], Deceit [Siegel et al. 1990], Ficus [Guy et al. 1990], and Harp [Liskov et al. 1991]. All these systems replicate data by storing complete copies, which results in higher storage and update costs than Zebra's parity scheme. Harp uses write-behind logs with uninterruptible power supplies to avoid synchronous disk operations and thereby reduce the update overhead. In addition, some of the systems, such as Locus and Coda, use the replicas to improve performance by allowing a client to access the nearest replica; Zebra's parity approach does not permit this optimization.

Another approach to highly available file service is to design file servers that can quickly reboot after a software failure [Baker and Sullivan 1992]. The idea is to reboot the file server so quickly that file service is not interrupted. This alternative does not require redundant copies or parity, but neither does it allow the system to continue operation in the event of a hardware failure.

Zebra borrows its log structure from LFS [Rosenblum and Ousterhout 1991], a high-performance write-optimized file system. A recent paper by Seltzer et al. [1993] has shown that adding extents to FFS [McKusick et al. 1984] results in a file system (called EFS [McVoy and Kleiman 1991]) that has comparable performance to LFS on large reads and writes. However, EFS does not improve performance for small files as does LFS and therefore Zebra, nor does it address the parity and striping issues presented by a striped network file system.

The create and delete deltas used by Zebra are similar to the active and deleted sublists used in the Grapevine mail system to manage entries in a registration database [Birrell et al. 1982]. Grapevine used timestamps whereas Zebra uses version numbers, but they each allow the system to establish an order between different sources of information and to recover from crashes.

10. CONCLUSIONS

Zebra takes two ideas that were originally developed for managing disk subsystems—striping with parity and log-structured file systems—and ap-

plies them to network file systems. The result is a network file system with several attractive properties:

- Performance*: Large files are read or written 4–5 times as fast as other network file systems, and small files are written 15–300% faster.
- Scalability*: New disks or servers can be added incrementally to increase the system's bandwidth and capacity. Zebra's stripe cleaner reorganizes data automatically over time to take advantage of the additional bandwidth.
- Cost-effective servers*: Storage servers do not need to be high-performance machines or have special-purpose hardware, since the performance of the system can be increased by adding more servers. Zebra transfers information to storage servers in large stripe fragments, and the servers do not interpret the contents of stripes; therefore, the server implementation is simple and efficient.
- Availability*: By combining ideas from RAID and LFS, Zebra can use simple mechanisms to manage parity for each stripe. The system can continue operation while one of the storage servers is unavailable and can reconstruct lost data in the event of a total failure of a server or disk.
- Simplicity*: Zebra adds very little complexity over the mechanisms already present in a network file system that uses logging for its disk structures. Deltas provide a simple way to maintain consistency among the components of the system.

There are at least five areas where Zebra could benefit from additional work:

- Name caching*: Without name caching, Zebra provides only about a 15% speedup for small writes in comparison to a nonstriped Sprite file system. A system with name caching would provide a much greater speedup.
- Transaction processing*: Zebra is expected to work well on the same workloads as LFS, which includes most workstation applications. However, there is a significant amount of controversy surrounding the performance of LFS under a transaction-processing workload. More work is needed to understand this area.
- Metadata*: It was convenient in the Zebra prototype to use a file in an existing file system to store the block pointers for each Zebra file, but this approach suffers from a number of inefficiencies. The system could be improved if the metadata structures were redesigned from scratch with Zebra in mind.
- Small reads*: It would be interesting to verify whether there is enough locality in small-file reads for prefetching of whole stripes to provide a substantial performance improvement.
- Security and protection*: The current design does little to provide security and protection for the files it stores. Malicious clients cannot overwrite existing files, but they can read files for which they should not have permission. The addition of security identifiers to file blocks and access

control lists to the storage servers appears to be a simple solution that would greatly improve Zebra's security.

Overall Zebra offers higher throughput, availability, and scalability than today's network file systems at the cost of only a small increase in system complexity.

ACKNOWLEDGMENTS

Paul Leach, Felipe Cabrera, Ann Drapeau, Ken Shirriff, Bruce Montague, and Mary Baker provided useful comments on various drafts of the article. Ken Lutz, Peter Chen, Peter Belleau, and Ares Ho built the timer boards that proved invaluable in debugging the system and running the experiments. Our thanks to the anonymous referees whose comments greatly improved this article.

REFERENCES

- ANDERSON, T. E., CULLER, D. E., AND PATTERSON, D. A. 1995. A case for NOW (Networks of Workstations). *IEEE Micro*, 15, 1 (Feb.), 54-64.
- BAKER, M. AND SULLIVAN, M. 1992. The Recovery Box: Using fast recovery to provide high availability in the UNIX environment. In *Proceedings of the Summer 1992 USENIX Conference* (June). USENIX Assoc., Berkeley, Calif., 31-43.
- BAKER, M., ASAMI, S., DEPRIT, E., AND OUSTERHOUT, J. 1992. Non-volatile memory for fast, reliable file systems. In *Proceedings of the 5th International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS)* (Boston, Mass., Oct.). ACM, New York, 10-22.
- BAKER, M. G., HARTMAN, J. H., KUPFER, M. D., SHIRRIFF, K. W., AND OUSTERHOUT, J. K. 1991. Measurements of a distributed file system. In *Proceedings of the 13th Symposium on Operating Systems Principles (SOSP)* (Asilomar, Calif., Oct.). *ACM SIGOPS Oper. Syst. Rev.* 25, 5, 198-212.
- BERNSTEIN, P. A. AND GOODMAN, N. 1981. Concurrency control in distributed database systems. *ACM Comput Surv.* 13, 2 (June), 185-222.
- BIRRELL, A. D., LEVIN, R., NEEDHAM, R. M., AND SCHROEDER, M. D. 1982. Grapevine. An exercise in distributed computing. *Commun. ACM* 25, 4 (Apr.), 260-274.
- CABRERA, L.-F. AND LONG, D. D. E. 1991. Swift: Using distributed disk striping to provide high I/O data rates. *Comput. Syst.* 4, 4 (Fall), 405-436.
- CAO, P., LIM, S. B., VENKATARAMAN, S., AND WILKES, J. 1993. The TickerTAIP parallel RAID architecture. In *Proceedings of the 20th Annual International Symposium of Computer Architecture* (May). ACM/IEEE, New York, 52-63.
- CHEN, P. M. AND PATTERSON, D. A. 1990. Maximizing performance in a striped disk array. In *Proceedings of the 17th Annual International Symposium of Computer Architecture* (May). ACM/IEEE, New York, 322-331.
- CHUTANI, S., ANDERSON, O. T., KAZAR, M. L., LEVERETT, B. W., MASON, W. A., AND SIDEBOTHAM, R. N. 1992. The Episode File System. In *Proceedings of the Winter 1992 USENIX Conference* (Jan.). USENIX Assoc., Berkeley, Calif., 43-60.
- DIBBLE, P. C., SCOTT, M. L., AND ELLIS, C. S. 1988. Bridge: A high-performance file system for parallel processors. In *Proceedings of the 8th International Conference on Distributed Computing Systems (ICDCS)*. IEEE, New York, 154-161.
- DRAPEAU, A. L., SHIRRIFF, K., HARTMAN, J. H., MILLER, E. L., SESHAN, S., KATZ, R. H., LUTZ, K., PATTERSON, D. A., LEE, E. K., CHEN, P. M., AND GIBSON, G. A. 1994. RAID-II: A high-bandwidth network file server. In *Proceedings of the 21st Annual International Symposium of Computer Architecture* (Apr.). ACM/IEEE, New York, 234-244.

- FLOYD, R. A. AND ELLIS, C. S. 1989. Directory reference patterns in hierarchical file systems. *IEEE Trans. Knowl. Data Eng.* 1, 2 (June), 238–247.
- FREEH, V. W., LOWENTHAL, D. K., AND ANDREWS, G. R. 1994. Distributed filaments: Efficient fine-grain parallelism on a cluster of workstations. In *Proceedings of the 1st USENIX Symposium on Operating Systems Design and Implementation (OSDI)* (Nov.). USENIX Assoc., Berkeley, Calif., 201–213.
- GUY, R. G., HEIDEMANN, J. S., MAK, W., PAGE, T. W., JR., POPEK, G. J., AND ROTHMEIER, D. 1990. Implementation of the Ficus replicated file system. In *Proceedings of the Summer 1990 USENIX Conference* (Anaheim, Calif., June). USENIX Assoc., Berkeley, Calif., 63–71.
- HAGMANN, R. 1987. Reimplementing the Cedar file system using logging and group commit. In *Proceedings of the 13th Symposium on Operating Systems Principles (SOSP)* (Nov.). *ACM SIGOPS Oper. Syst. Rev.* 21, 5, 155–162.
- HARTMAN, J. H. AND OUSTERHOUT, J. K. 1993. Letter to the editor. *ACM SIGOPS Oper. Syst. Rev.* 27, 1 (Jan.), 7–10.
- HISGEN, A., BIRRELL, A., MANN, T., SCHROEDER, M., AND SWART, G. 1989. Availability and consistency tradeoffs in the Echo distributed file system. In *Proceedings of the 2nd Workshop on Workstation Operating Systems* (Sept.). IEEE, New York, 49–54.
- HOWARD, J. H., KAZAR, M. L., MENEES, S. G., NICHOLS, D. A., SATYANARAYANAN, M., SIDEBOTHAM, R. N., AND WEST, M. J. 1988. Scale and performance in a distributed file system. *ACM Trans. Comput. Syst.* 6, 1 (Feb.), 51–81.
- LISKOV, B., GHEMAWAT, S., GRUBER, R., JOHNSON, P., SHRIRA, L., AND WILLIAMS, M. 1991. Replication in the Harp file system. In *Proceedings of the 13th Symposium on Operating Systems Principles (SOSP)* (Asilomar, Calif., Oct.). *ACM SIGOPS Oper. Syst. Rev.* 25, 5, 226–238.
- LONG, D. D. E., MONTAGUE, B. R., AND CABRERA, L.-F. 1994. Swift/RAID: A distributed RAID system. *Comput. Syst.* 7, 3 (Summer), 333–359.
- LO VERSO, S. J., ISMAN, M., NANOPOULOS, A., NESHEIM, W., MILNE, E. D., AND WHEELER, R. 1993. sfs: A parallel file system for the CM-5. In *Proceedings of the Summer 1993 USENIX Conference* (Cincinnati, Ohio, June). USENIX Assoc., Berkeley, Calif., 291–305.
- MCKUSICK, M. K., JOY, W. N., LEFFLER, S. J., AND FABRY, R. S. 1984. A fast file system for UNIX. *ACM Trans. Comput. Syst.* 2, 3 (Aug.), 181–197.
- MCVOY, L. W. AND KLEIMAN, S. R. 1991. Extent-like performance from a UNIX file system. In *Proceedings of the Winter 1991 USENIX Conference* (Dallas, Tex., Jan.). USENIX Assoc., Berkeley, Calif., 33–43.
- NELSON, M. N., WELCH, B. B., AND OUSTERHOUT, J. K. 1988. Caching in the Sprite network file system. *ACM Trans. Comput. Syst.* 6, 1 (Feb.), 134–154.
- OUSTERHOUT, J. 1995. A critique of Seltzer's 1993 USENIX paper. Available as <http://www.smli.com/~ouster/seltzer93.html>.
- OUSTERHOUT, J., CHERENSON, A., DOUGLIS, F., NELSON, M., AND WELCH, B. 1988. The Sprite network operating system. *IEEE Comput.* 21, 2 (Feb.), 23–36.
- PATTERSON, D. A., GIBSON, G., AND KATZ, R. H. 1988. A case for redundant arrays of inexpensive disks (RAID). In *Proceedings of the 1988 ACM Conference on Management of Data (SIGMOD)* (Chicago, Ill., June). ACM, New York, 109–116.
- PIERCE, P. 1989. A concurrent file system for a highly parallel mass storage subsystem. In *Proceedings of the 4th Conference on Hypercubes* (Monterey, Calif., Mar.). ACM, New York, 155–160.
- ROSENBLUM, M. AND OUSTERHOUT, J. K. 1991. The design and implementation of a log-structured file system. In *Proceedings of the 13th Symposium on Operating Systems Principles (SOSP)* (Asilomar, Calif., Oct.). *ACM SIGOPS Oper. Syst. Rev.* 25, 5, 1–15.
- SATYANARAYANAN, M., KISTLER, J. J., KUMAR, P., OKASAKI, M. E., SIEGEL, E. H., AND STEERE, D. C. 1990. Coda: A highly available file system for a distributed workstation environment. *IEEE Trans. Comput.* 39, 4 (Apr.), 447–459.
- SCHLOSS, G. A. AND STONEBRAKER, M. 1990. Highly redundant management of distributed data. In *Proceedings of the IEEE Workshop on the Management of Replicated Data* (Nov.). IEEE, New York, 91–95.

- SELTZER, M., BOSTIC, K., MCKUSICK, M. K., AND STAELIN, C. 1993. An implementation of a log-structured file system for UNIX. In *Proceedings of the Winter 1993 USENIX Conference* (San Diego, Calif., Jan.). USENIX Assoc., Berkeley, Calif., 307–326.
- SELTZER, M., SMITH, K. A., BALAKRISHNAN, H., CHANG, J., MCMAINS, S., AND PADMANABHAN, V. 1995. File system logging versus clustering: A performance comparison. In *Proceedings of the Winter 1995 USENIX Conference* (Jan.). USENIX Assoc., Berkeley, Calif., 249–264.
- SHELTZER, A. B., LINDELL, R., AND POPEK, G. J. 1986. Name service locality and cache design in a distributed operating system. In *Proceedings of the 6th International Conference on Distributed Computing Systems (ICDCS)* (May). IEEE, New York, 515–522.
- SHIRRIFF, K. AND OUSTERHOUT, J. 1992. A trace-driven analysis of name and attribute caching in a distributed file system. In *Proceedings of the Winter 1992 USENIX Conference* (Jan.). USENIX Assoc., Berkeley, Calif., 315–331.
- SIEGEL, A., BIRMAN, K., AND MARZULLO, K. 1990. Deceit: A flexible distributed file system. In *Proceedings of the Summer 1990 USENIX Conference* (Anaheim, Calif., June). USENIX Assoc., Berkeley, Calif., 51–61.
- WALKER, B., POPEK, G., ENGLISH, R., KLINE, C., AND THIEL, G. 1983. The LOCUS distributed operating system. In *Proceedings of the 9th Symposium on Operating Systems Principles (SOSP)* (Nov.) *ACM SIGOPS Oper. Syst. Rev.* 17, 5, 49–70.
- WILKES, J. 1992. DataMesh research project, phase 1. In *Proceedings of the USENIX File Systems Workshop* (May). USENIX Assoc., Berkeley, Calif., 63–69.

Received April 1994; revised August 1994; accepted March 1995