# A Coherent Distributed File Cache with Directory Write-Behind

TIMOTHY MANN, ANDREW BIRRELL, ANDY HISGEN, CHARLES JERIAN, and GARRET SWART Digital Equipment Corporation

Extensive caching is a key feature of the Echo distributed file system. Echo client machines maintain coherent caches of file and directory data and properties, with write-behind (delayed write-back) of *all* cached information. Echo specifies ordering constraints on this write-behind, enabling applications to store and maintain consistent data structures in the file system even when crashes or network faults prevent some writes from being completed. In this paper we describe the Echo cache's coherence and ordering semantics, show how they can improve the performance and consistency of applications, and explain how they are implemented. We also discuss the general problem of reliably notifying applications and users when write-behind is lost; we addressed this problem as part of the Echo design, but did not find a fully satisfactory solution.

Categories and Subject Descriptors: D.4.3 [**Operating Systems**]: File Systems Management—distributed file systems

General Terms: Design, Experimentation, Measurement, Performance, Reliability, Security

Additional Key Words and Phrases: Coherence, file caching, write-behind

#### 1. INTRODUCTION

Echo is a distributed file system that incorporates replication, caching, global naming, and distributed security. Figure 1 gives a block diagram of the Echo system.<sup>1</sup>

-*Replication*. Echo replicates servers for availability and disks for data integrity, allowing a wide range of configurations and tolerating both server crashes and network faults. The interconnections between disks, servers, and client machines can be replicated as well.

© 1994 ACM 0734-2071/94/0500-012300.00

ACM Transactions on Computer Systems, Vol. 12, No. 2, May 1994, Pages 123-164.

<sup>&</sup>lt;sup>1</sup>Echo is no longer under development or in use, but to avoid awkwardness we speak of it in the present tense throughout this paper.

A version of this paper was issued as Research Report 103, Systems Research Center, Digital Equipment Corporation, Palo Alto, California, 1993.

Authors' address: Systems Research Center, Digital Equipment Corporation, 130 Lytton Avenue, Palo Alto, CA 94301.

Permission to copy without fee all or part of this material is granted provided that the copies are not made or distributed for direct commercial advantage, the ACM copyright notice and the title of the publication and its date appear, and notice is given that copying is by permission of the Association for Computing Machinery. To copy otherwise, or to republish, requires a fee and/or specific permission.



Fig. 1. Block diagram of the Echo system.

- --Caching. Echo client machines keep coherent write-back caches of data and properties for both files and directories in volatile memory, thereby reducing the latency seen by applications on file-system operations and reducing the read load and peak write load on servers. The caches use ordered write-behind; that is, updates buffered in the cache are automatically flushed after a time delay and are written to server disks in a well-defined partial order, thus limiting the damage that can occur when a client machine crashes and its volatile memory is erased. We use the term clerk for the module in the client operating system that performs these functions.
- -Global naming. Echo supports a globally scalable, hierarchical name space. The upper levels of the hierarchy are implemented using a replicated name service that achieves very high availability with loose consistency semantics, while the lower levels are implemented by the file system with tight consistency, but somewhat lower availability.
- -Distributed security. Echo has adopted a security model in which servers do not trust client machines; instead, client machines use a cryptographic protocol to authenticate themselves as acting on behalf of the users who have logged in to them [Lampson et al. 1991].

This paper concentrates on Echo's client cache. Separate papers give a complete overview of Echo [Birrell et al. 1993] and discuss various other

ACM Transactions on Computer Systems, Vol. 12, No. 2, May 1994.

aspects in detail [Hisgen et al. 1989; 1990; 1992; 1993; Mann et al. 1989; Swart et al. 1993].

In the next section, we discuss the motivation for Echo's cache design; then in Section 3 we present the cache's coherence and ordering semantics and the facilities for reporting lost write-behind. We show how these semantics are useful to some applications in Section 4, and how directory write-behind can improve the performance of some applications in Section 5. In Section 6 we give details of the cache implementation. Section 7 compares the Echo cache with related work, while Section 8 summarizes our conclusions.

# 2. DESIGN MOTIVATION

In this section we discuss four goals that motivated the Echo cache design and explain how they shaped the resulting system. We wanted the Echo distributed file system to meet or exceed the standards set by single-machine file systems in performance, semantics, fault tolerance, and security.

Echo's *performance* goal was to provide service comparable to or better than a single-machine file system on similar hardware. We strove not only for high throughput, but also for low latency on individual operations. We also wanted to maximize the number of client machines that could be handled per server.

This ambitious goal led us to an aggressive cache design. Most of our workstations dedicate more than 16 Mbytes of main memory to the cache. Echo clerks cache both files and directories, so that users see good performance not only when reading file data, but also when opening files and listing directories. Clerks do write-behind on both files and directories, so that users see low latency not only when writing to files, but also when creating, deleting, renaming, or changing access permissions on files and directories. Write-behind also opens up the opportunity to reduce the load on servers (thus increasing their capacity) by optimizing out sequences of operations that cancel before they reach the server—for example, writing the same file page repeatedly, or creating and deleting a temporary file—though we have not gone far in implementing such optimizations.

Echo's *semantic* goal was to present a single-system view to applications. That is, so far as possible, the distributed file system should present the same interface and semantics to an application program as does a single-machine file system. Even an application distributed across many machines should see one consistent file system, not many slightly inconsistent ones. Existing applications should not require changes to work with the new file system, and programmers should not need to learn a new set of skills to write new applications.

This goal led us to implement a file cache with strict coherence and a replication scheme with strict consistency between copies. We chose not to explore the optimistic approach pioneered by the LOCUS [Popek et al. 1981; Walker et al. 1983] file system and more recently used in Coda [Kistler and Satyanaraynan] and Ficus [Guy and Popek 1991; Page et al. 1991], where replicas or cached copies of files are allowed to diverge during periods when

ACM Transactions on Computer Systems, Vol 12, No 2, May 1994.

the network is not fully connected. Instead, we forbid both write/write conflicts (where two clients make different changes to copies of the same file) and read/write conflicts (in which one client changes a copy of a file while another client continues to read the old version).

Our desire to to present a single-system image also led us to emulate a single-machine UNIX<sup>2</sup> file system more carefully than some other distributed file systems have done. For example, unlike NFS [Sandberg et al. 1985; Sun Microsystems 1989], Echo keeps a file that has been unlinked from the name space in existence as long as any application program has it open. Also unlike NFS, an Echo clerk allows applications to write data into its cache only if it knows there will be disk space available on the server to hold the data when the cache is flushed.

Echo's *fault-tolerance* goal was to mask faults in all system components wherever feasible, and to fail as cleanly as possible when faults occur that cannot be masked. Fault tolerance is becoming increasingly important in distributed systems as they are built from more and more individual machines, each of which can fail independently.

The most visible fault-tolerance feature in Echo is the replication of servers and disks, which we do not discuss in this paper, but fault-tolerance considerations affected our cache design as well. The token directory that is needed to maintain cache coherence is replicated in the main memory of two different servers, so that the failure of one server will not invalidate client caches and force their write-behind to be discarded. If a server loses touch with a clerk, the server is allowed to revoke the clerk's tokens, but not before a mutually agreed-upon time-out (or lease [Gray and Cheriton 1989]) has expired. When a lease expires, the clerk holding it is required to make cached data that the lease covers inaccessible to clients; thus, two applications that read from different caches at the same time can never see inconsistent data, even if one is running on a machine that has been partitioned away from the rest of the network. To provide clean failure semantics, we guarantee that write-behind will reach disk in a partial order known to applications, and we allow applications to add their own constraints to this partial order. When we must discard write-behind, we give any application that may have read or written the discarded data an error return on any further reads or writes it attempts, thus halting its progress before it can observe the anomaly.

Echo's *security* goal was to protect the privacy and integrity of stored information without requiring all machines in the distributed system to trust all others and without compromising the system's performance. This goal led us to place a trust boundary between servers and clerks. Servers do not trust the operating systems on client machines to provide proper security. Instead, clerks must authenticate themselves to servers as acting on behalf of particular users, using a cryptographic protocol. An authenticated clerk is given only the privileges of the user it is acting for; it cannot touch data that the user is not permitted to access, and it cannot corrupt data structures in the file-sys-

<sup>&</sup>lt;sup>2</sup>UNIX is a registered trademark.

ACM Transactions on Computer Systems, Vol 12, No 2, May 1994.

tem implementation.<sup>3</sup> Because access checks can be expensive, we also developed a form of access check caching: An Echo server needs to do access checking only when a clerk requests a cache coherence token, not on every read or write. This security machinery was easy to add to the system during the design phase, but would have been hard to retrofit had we chosen to omit security from the initial design.

# 3. COHERENCE AND ORDERING SEMANTICS

Echo is *single-copy equivalent*; that is, at any moment, the file-system data is in a single, well-defined state. An operation that changes this state is called a *write*. An operation that returns information about the state without changing it is called a *read*. Each operation is *logically performed* (carried out) at a distinct point in real time, reading or changing the state as it exists at that moment. (We take the view that every operation is logically performed at a different time because we want these times to define a total order on the operations.) If an application issues a system call requesting an operation at time  $t_1$  and the call returns at time  $t_2$ , the operation will have been logically performed at some time t with  $t_1 \leq t \leq t_2$ . Moreover, in the absence of faults, Echo caching is *transparent*; that is, if no network faults or machine crashes occur, nothing other than a write operation changes the file system state.

When faults do occur, Echo remains single-copy equivalent, but the caching is no longer fully transparent. A fault can change the state of the file-system data by causing updates that were written behind to be discarded before they reach disk. Such writes are logically *undone* at a unique point in real time, moving the file-system data into a new state that is reflected in all reads and writes performed after that point.

These semantics are the same as those of a single-machine file system with write-behind, but differ from those of NFS, in which caches are incoherent. With NFS, a process on one machine can see old data when reading a file if a process on another machine has the file open for writing, and can see old data in a directory if a process on another machine has modified the directory recently.

The Echo cache remains single-copy equivalent even if network faults cause one or more machines to be partitioned away from the rest of the system. If a partition (or server crash) prevents a clerk from accessing the file system's current state to perform an application read or write request, the clerk either blocks waiting for the state to become accessible or returns an error indication. (These two options are similar, respectively, to the *hard* and *soft* options in mounting an NFS volume.) In our environment, we found it most practical to have the clerks block such operations for a limited time (up to two minutes), and then give up and return an error indication if the problem remains. This policy works well in Echo because Echo's fault-tolerance features correct many problems automatically in well under two

<sup>&</sup>lt;sup>3</sup>The interface between clerks and servers is in terms of logical operations, such as creating and deleting files, not in terms of direct modifications to bytes in directory pages or disk blocks.

ACM Transactions on Computer Systems, Vol. 12, No. 2, May 1994.

minutes. If a problem is still present after two minutes, it is likely to be the kind that needs to be fixed manually, which may take a long time. So it seems good to unblock applications and give them error returns at that point.

# 3.1 Ordering Constraints

Because crashes and network faults can cause write-behind to be discarded, the order in which writes reach disk is important. If Echo allowed writes to reach disk in an arbitrary order, applications that store mutable data in files could find their data inconsistent after a crash in which write-behind is lost. Therefore, Echo specifies a partial order on writes and guarantees that the actual order in which writes reach disk will be consistent with this partial order. Echo also provides a primitive that lets applications augment the partial order with additional constraints. With a knowledge of Echo's ordering guarantees, a carefully coded application can assure itself that the data structures it stores in the file system will remain consistent even when some write-behind is lost in a crash. Section 4 gives some examples; the remainder of this section describes the constraints themselves. Section 7 includes a brief comparison of this approach to fault tolerance with the alternative of providing atomic transactions in the file-system interface. In this section we begin by defining some necessary terms, then state Echo's ordering guarantees—first informally, then formally—and finally present two examples.

We say that a write is *stable* when it has reached disk; we say that a write is *discarded* when it is logically undone and will never reach disk. All writes are eventually either stable or discarded (but not both). We call a write *unstable* when it has been logically performed, but is not yet either stable or discarded.

Stated informally, Echo's guarantees are as follows:

- (A) If a write is requested by one client machine and the results are observed by another, the write is stable.
- (B) Writes to a given object become stable in the same order they are logically performed, except that an unbroken sequence of overwrites requested by a single client may be reordered. (An *overwrite* is a write operation to a file that does not change its length.)<sup>4</sup>
- (C) Writes are stable when they are forced to disk by *fsync*. (When the UNIX system call *fsync(f)* returns, all writes to file *f* that returned before *fsync* was called are guaranteed to be on disk. Echo allows *fsync* on directories as well as on ordinary files.)

To state Echo's stability ordering constraints formally, we define two relations,  $\rightarrow$  and  $\Rightarrow$ . Intuitively, the relation  $\rightarrow$  expresses data dependency; two operations are related by  $\rightarrow$  if and only if the first could have

<sup>&</sup>lt;sup>4</sup>We chose to allow reordering of overwrites simply to reduce the bookkeeping burden in the clerk implementation; an application that needs to order its overwrites can do so using the *forder* primitive described later in this section.

ACM Transactions on Computer Systems, Vol 12, No. 2, May 1994

affected the result of the second. The relation  $\Rightarrow$  is the partial order in which writes are guaranteed to reach disk. Viewed as a set of ordered pairs, the relation  $\Rightarrow$  is a subset of  $\rightarrow$ . The formal definitions of  $\rightarrow$  and  $\Rightarrow$  are given next; some readers may prefer to skip ahead to the examples in Figures 2 and 3 before reading these definitions.

- (1) If  $o_1$  and  $o_2$  are two operations on the file system, we say that  $o_1 \rightarrow o_2$  if  $-o_1$  is a write operation,
  - $-o_1$  and  $o_2$  have an operand in common,
  - $-o_1$  and  $o_2$  both return successfully,<sup>5</sup>
  - $-o_1$  was logically performed before  $o_2$ , and
  - $-o_1$  was not already discarded when  $o_2$  was logically performed.
- (2) The  $\rightarrow$  relation is transitive; if  $o_1 \rightarrow o_2$  and  $o_2 \rightarrow o_3$ , then  $o_1 \rightarrow o_3$ .
- (3) We say o<sub>1</sub> ⇒ o<sub>2</sub> if
   -o<sub>1</sub> → o<sub>2</sub>, and
   -o<sub>1</sub> and o<sub>2</sub> are both write operations, but they are not both overwrites (data writes to a file that do not change its length).
- (4) The  $\Rightarrow$  relation is transitive; if  $o_1 \Rightarrow o_2$  and  $o_2 \Rightarrow o_3$ , then  $o_1 \Rightarrow o_3$ .
- (5) If  $o_1 \Rightarrow o_2$ , and  $o_1$  is now discarded, then  $o_2$  is now discarded.
- (6) If  $o_1 \rightarrow o_2$ , and  $o_1, o_2$  were requested by applications running on different client machines, then when  $o_2$  is logically performed,  $o_1$  is stable. (This corresponds to informal guarantee (A) above.)
- (7) If  $o_1 \Rightarrow o_2$ , and  $o_2$  is now stable, then  $o_1$  is now stable (informal guarantee (B)).
- (8) If the write operation *fsync* returns successfully, then it is stable (informal guarantee (C)).

Echo's read and write operations include most of the usual UNIX file-system operations, plus a new operation for adding constraints (described below). Read operations include getting the properties of a file or directory (the UNIX *stat* system call), opening a file, reading file data, listing a directory, looking up a pathname component in a directory, and so forth. System calls that involve looking up a pathname are viewed formally as several operations that are logically performed in sequence, consisting of a component lookup in each directory along the path, followed by an operation on the final object found. All read operations have just one operand.

Write operations include writing file data, creating a file or directory, renaming, *fsync*, and so forth. We view *fsync* as a write operation because of the role it plays in constraining the order in which other operations are made stable, even though it does no writing of its own. A file-data write that both overwrites existing data and appends new data is viewed formally as two operations, a pure overwrite logically followed by an append. Also, a file-data overwrite is guaranteed to be failure-atomic only if it modifies bytes in at

<sup>&</sup>lt;sup>5</sup>We say an operation *returns successfully* if it returns control to its caller without indicating an error. Write-behind or other work queued by the operation may still fail later.

ACM Transactions on Computer Systems, Vol. 12, No. 2, May 1994.

most one block of the file, where the block size is a parameter set by the Echo implementation, currently 1024 bytes. An overwrite that does not meet this criterion may be implemented as an arbitrary sequence of shorter overwrite operations, where each of the shorter overwrites is failure-atomic, but the sequence as a whole is not failure-atomic or even ordered. All other write operations are failure-atomic. Write operations may have multiple operands; for example, if we rename a file from /a/f to /b/f when some other file named /b/f already exists, this operation modifies the old parent /a, the new parent /b, the file /a/f being renamed, and the old file /b/f being displaced, for a total of four operands.

Echo adds one more write operation to the usual set, called *forder*. Like *fsync*, the *forder* operation does not modify any of its operands, but does follow the stability rules for write operations. Unlike *fsync*, however, *forder* returns immediately, not waiting for the operations ordered before it to become stable. The *forder* operation is useful for adding constraints to the order in which other writes are made stable, without delaying its caller as *fsync* does; we give some examples in Section 4.

Figures 2 and 3 illustrate some of the ordering constraints. In Figure 2 a series of write operations is applied to a single file **f**. All of the operations are ordered by  $\rightarrow$ , but only those joined by arrows in the figure are ordered by  $\Rightarrow$ . First, two different records within the file are overwritten; these writes are not ordered by  $\Rightarrow$ . Next, a new record is appended; this write is ordered after both of the overwrites. Three more overwrites occur next; these are ordered after the append, but not among themselves. Finally, there is an *forder* call followed by another overwrite. The *forder* itself does not change the file in any way, but causes the final overwrite to be ordered after the earlier ones.

In Figure 3 several files are renamed. The first two renames are not ordered because they affect different files in different directories; they have no operands in common. The third rename follows the first because both affect directory /d. The fourth rename follows all of the others because it affects both directories /d and /e.

Echo limits the scope of its ordering constraints by requiring all of the operands of a given operation to be in the same *volume*. An Echo volume is roughly similar to a UNIX "file system": It is a subtree of the global hierarchical name space, stored on a single (possibly replicated) set of disks, and managed by a single (possibly replicated) server. Echo volumes are tied together by *junctions*, which differ from UNIX "mount points" in that they are stored stably in the joined volumes and are global to all clients of the file system, rather than being established individually (and perhaps differently) by each client machine at run time. Like UNIX, Echo does not permit files to be renamed or hard-linked from one volume to another. It also does not permit *forder* operations whose operands are not all in the same volume. This restriction greatly simplifies the Echo implementation and Echo's interfaces for reporting discarded write-behind (discussed in Section 3.3), but it requires applications to keep their file-based data structures within a single volume or to take extra care when cross-volume dependencies arise.

ACM Transactions on Computer Systems, Vol 12, No. 2, May 1994.



# 3.2 Limits on Buffered Write-Behind

Echo's cache uses a *write-behind* policy, not a simple *write-back* policy. That is, the clerk writes back dirty data after a fixed maximum time delay, even if the stability guarantees of Section 3.1 do not require it to. We chose to do this because dirty data is vulnerable to loss as long as it is buffered on a client machine. We cannot eliminate this vulnerability short of going to a writethrough policy, but we can reduce the probability that data will be lost by limiting the length of time dirty data is buffered. With simple write-back, a dirty block can stay buffered on a client workstation indefinitely, if other clerks do not request the block and if no programs call *fsync* or *sync* on it. With write-behind, a workstation that has not been used for a while will typically have no dirty data buffered.

The time limits on write-behind in Echo are characterized by two parameters. The Echo clerk calls the server to flush each write to disk within  $t_{flush}$ seconds from the time the write was logically performed in the cache, and the clerk blocks all further client write calls on a volume if any write call to the server made more than  $t_{block}$  seconds ago is still waiting for a response. The clerk also logs a diagnostic message if  $t_{block}$  is exceeded. These time limits are part of the defined interface advertised to Echo users. The second limit is needed to deal with the possibility that an Echo server might be unable to

ACM Transactions on Computer Systems, Vol. 12, No 2, May 1994.

keep up with the rate of requests submitted by clerks: If the clerk guarantees to call the server within  $t_{flush}$  seconds, but the server may take arbitrarily long to respond, then the guarantee is of no value to users. In our installation, we set  $t_{flush} = 15$  and  $t_{block} = 300$ . We saw the  $t_{block}$  limit hit only in the early days of the system, when the servers had performance problems and had bugs that occasionally caused them to deadlock.

Besides these time limits, the amount of dirty data that an Echo clerk buffers is, of course, also limited by the amount of memory the clerk has available. If all available cache space fills up with dirty blocks or pending directory modifications, some of this material must be written back before any more work can be done. The Echo clerk uses a variety of heuristics to prevent its cache from being clogged with too much dirty data and to determine what to write back first when the cache begins to fill. But none of these heuristics provide any useful semantic guarantees to applications, so we do not discuss them further in this paper.

# 3.3. Reporting Lost Write-Behind

In Section 3.1 we said that a carefully coded application can take advantage of Echo's ordering guarantees to ensure that its file-based data structures will remain consistent even when write-behind is lost "in a crash." thus suggesting that applications can deal with lost write-behind by restarting and perhaps invoking crash recovery. But not all lost write-behind in Echo is caused by crashes. If a network fault cuts off communication between a clerk and a server, the server revokes the clerk's cache coherence tokens in order to allow other clerks to proceed. This forces the clerk to discard its write-behind in order to maintain cache coherence, but does not halt programs running on the clerk's machine. (Write-behind is also discarded if a double server crash causes the entire token direction to be lost.) By contrast, NFS and singlemachine UNIX discard write-behind only if the machine holding the writebehind crashes, which of course halts all applications running on that machine. Thus, Echo applications have an additional problem to deal with: finding out when writes they depend on have been discarded and taking appropriate action. Echo provides some facilities to help; we describe them next.

To make this discussion more precise, we say that an application process P depends on a write w if P issued the write, or if P has done a read or write operation o such that  $w \rightarrow o$ . In these cases (and in no others), P has directly observed the effect of w on the file system's state.<sup>6</sup> Thus, if P does not expect other processes to be changing the part of the file system that P is using, discarding w will make P's internal state inconsistent with that of the file system. If P continues running after w is discarded and reads from the file

<sup>&</sup>lt;sup>6</sup>Other processes may have observed the effect of w indirectly—for example, they may have been told about it by P—but the system cannot detect this. In general, it is not safe for a process to communicate information about the file system's state outside itself until that information is stable

ACM Transactions on Computer Systems, Vol 12, No 2, May 1994.

system, it may observe the inconsistency and be confused; if it writes, it may write data that is inconsistent with the new file-system state.

Echo provides several different ways of dealing with this problem, as different *recovery modes* that can be selected for each process. If a process that is in *standard recovery mode* depends on discarded write-behind on a given volume, Echo gives an error return on any further operations the process tries to invoke on that volume. A process that receives such an error return should immediately abort with an error message, effectively converting the error to a "crash." We considered sending the affected processes an asynchronous UNIX signal too, so that they would halt quickly even if they did not access the file system frequently; this seems like a good idea, but we did not find time to try it.

If a process that is in *self-recovery mode* depends on discarded write-behind on a given volume, Echo marks all of its open files in the affected volume so that new operations attempted on them will give an error return. The process's working directory is also marked in this way if it is in the affected volume, so that new operations that specify relative pathnames will give an error return. But operations that specify absolute pathnames (including changing the working directory) are allowed without restriction. Our vision was that a process could use this mode by making all of its file references (after initialization) through open files and relative pathnames, until it received its first error return reporting lost write-behind. At that point the process would run its own recovery code, which would reopen its files and working directory using absolute pathnames.

Self-recovery mode turned out not to be very useful in practice. Existing programs make no distinction in their usage of absolute and relative pathnames, so they generally do not behave reasonably in this mode: Either they use absolute pathnames in unfortunate places and thus fail to notice lost write-behind, or they use relative pathnames everywhere and thus fail to recover. Moreover, an application that maintains state in the file system and wants to recover from lost write-behind errors is easily structured as two processes: a child process that does all of the real work and that halts when lost write-behind is detected, plus a parent process that restarts the child after each such halt. The child runs the same recovery code regardless of whether it was restarted due to lost write-behind or due to a machine crash and reboot. Therefore, self-recovery mode was not useful to newly written programs either.

A third recovery mode, which we might call *null recovery mode*, would have been useful for interactive shells. In this mode, if write-behind that a process depends on is lost, all of the process's open files in the affected volume are marked so that new operations attempted on them give an error return. But new files can be opened by name without restriction, and the working directory remains valid. This mode would be useful for interactive UNIX shells because they do not keep files open for long and do not remember file-system state. An unmodified UNIX shell, which simply prints a message and continues when it receives an error return reporting lost write-behind, would work nicely in null recovery mode. Unfortunately, we did not imple-

ACM Transactions on Computer Systems, Vol. 12, No. 2, May 1994

ment null recovery mode in Echo, so we ran shells in self-recovery mode instead.<sup>7</sup> Users did not like the results: they were confused, not enlightened, when relative pathnames stopped working in a shell because lost writebehind had caused its working directory to be marked invalid. But running shells in standard recovery mode would have been even worse: Users would have been quite unhappy if their shells had crashed or lost all access to the file system whenever write-behind was lost.

Self-recovery mode is ugly largely because it is an atternpt to fit a new concept into an existing UNIX-like file-system interface. Had we been free to change this interface, we might have adopted a cleaner approach to self-recovery using explicit *failure handles*. In this approach, each read or write operation accepts a failure handle as an additional parameter. The semantics of lost write-behind reporting are changed to replace *process* with *failure handle* in the concept of dependency on writes: If a write w is issued using failure handle h, or if an operation o is issued using h and  $w \rightarrow o$ , then h depends on w. If a write is discarded, subsequent operations issued using any failure handle that depends on that write will give an error return. Standard recovery mode then becomes a special case in which a process chooses a new failure handle on its first access to each volume and uses it for all subsequent accesses.

The Echo file-system interface allows applications to obtain locks on files. Echo provides Berkeley-style advisory locks, and also uses a form of internal lock to implement the UNIX feature that keeps a file from being deleted as long as at least one process has it open, even if it is removed from the name space. In the implementation, such locks are obtained and cached in much the same way as file data. The locks are discarded along with write-behind when a client machine crashes or has its cache coherence tokens revoked by the server due to a network fault. Thus, when a client machine crashes, its locks are released, allowing other machines to obtain them and to make progress. Locking operations do not participate in the  $\rightarrow$  or  $\Rightarrow$  relations, but lock acquisitions do count as file accesses for the purpose of lost writebehind reporting. That is, if a process depends on lost write-behind on a volume, it receives error returns on attempts to acquire new locks in the volume. Conversely, if a lock on a given open file is discarded, all further attempts to read or write it receive error returns; there is no timing window during which the lock is released, but processes can still read or write the file.

When write-behind is discarded, besides notifying the affected processes, it is also a good idea to notify the user directly. Programs that have not been specially coded to deal with the possibility of lost write-behind may not handle it cleanly. For example, a process that writes to the file system will usually exit without calling *fsync* to make its changes stable, so if those

 $<sup>^{7}\,\</sup>mathrm{In}$  fact, for historical reasons, we ran most programs in self-recovery mode, which was clearly a mistake.

ACM Transactions on Computer Systems, Vol. 12, No. 2, May 1994

changes are discarded later, there is no process left to notify of the error.<sup>8</sup> In other cases, a process that ought to abort when write-behind is lost may instead print a message and continue, or may simply ignore the error returns. Unfortunately, it is hard to notify users about discarded write-behind in a way that makes sense to them. Echo prints a console message when this happens, but users find these messages cryptic, or miss them entirely because the console window is iconized. Worse, we have no way at all of notifying a user when write-behind is lost because his or her machine crashed. (Single-machine UNIX and NFS also have this problem.) We have no really good ideas for improving this situation; inherently, when we accept write-behind and store it unstably in client machine memory, we violate the simple abstraction of a stable, on-disk file system that programs and users expect to see. We can ameliorate the problems this causes, by providing ordering guarantees and notification of discarded write-behind, but we cannot eliminate them.

# 4. USING THE SEMANTICS

In this section we describe some ways in which applications can make use of Echo's caching and ordering semantics.

First and most important, a file system with coherent caching is much easier to use when writing a distributed application than one with incoherent caching. For example, suppose you would like to build a distributed, parallel make, that is, a tool that speeds up the recompilation of large programs by farming out the compilation of individual source files to different machines and then linking the results together on one machine. With a coherent distributed file system, such a tool can work much like an ordinary, nondistributed make: The compilation steps can simply write their results into the file system, and the link step can simply read them, knowing it will get the correct data. With a system like NFS that has incoherent caches, things are not so simple. Although recent implementations of NFS give close-to-open coherence, meaning that, once a file is closed, readers that open it later will get the correct data, NFS does not provide coherence on directories. So in this example, when the linker tries to open the files that the compilers have written, it may not find them all in its cached copy of the directory. Perhaps this example sounds contrived, but we are aware of real-world instances of the same phenomenon. Our colleagues in the Hector lexicography project [Glassman et al. 1992] tried to use NFS to maintain a shared file of dictionary definitions being read and updated by multiple lexicographers. After

ACM Transactions on Computer Systems, Vol. 12, No. 2, May 1994.

<sup>&</sup>lt;sup>8</sup>Because of this problem, we considered changing the semantics of process exit to include calling *fsync* on all of the files and directories that the process modified; however, we judged that the additional safety would not be worth the performance penalty. In particular, this change would slow down the operation of many UNIX shell scripts and Makefiles, which run a sequence of processes that communicate through the file system to perform a single task. Furthermore, existing programs are not prepared to handle the error returns that could be generated by adding *fsync* calls to the process exit routine; the UNIX *exit* primitive is declared as a void function and is specified not to return under any circumstances.

running into intractable bugs caused by NFS's incoherent caching, they were forced to rewrite their code to get the data from a centralized server, instead of going directly to the file system. Their original design would have worked fine with Echo.

Echo guarantees that if the data is written on one machine and read on another, that data is stable. (This is a consequence of the rules pertaining to the  $\rightarrow$  relation given above.) Thus, a distributed application whose component processes communicate only through the file system can be sure that, when one Echo clerk's write-behind is lost, only the processes running on that clerk's machine will be affected. The lost data will not have propagated to processes on other machines. The effect will be much the same as if the affected machine crashed.<sup>9</sup>

Making use of Echo's ordering semantics is more difficult than making use of its cache coherence, but the effort is worthwhile for applications that store data structures in the file system and want to make sure these structures do not become corrupt if write-behind is lost. One way to use the semantics is as follows: First, carefully order the application's file-system write calls so that, if a crash halts the application at any point, the data structures it has written will be consistent (or, at worst, automatically repairable) when it is restarted. This first step produces an application that would be robust against crashes if run on a file system with no write-behind. Second, check whether Echo's ordering constraints forbid all reorderings of writes that would leave data structures in an inconsistent (or unrepairable) state. If not, add *forder* calls to the application to forbid the unwanted reorderings. In the worst case, this can always be done by adding enough *forder* calls to forbid any reordering at all. A few examples follow.

Figure 4 gives a simple example of writing a file and then replacing it atomically with a new version. The arrows in the figure show which operations are related by  $\Rightarrow$  (omitting arrows that can be inferred from the transitivity of  $\Rightarrow$ ). In this example, Echo's built-in ordering constraints provide exactly what is needed, with no calls to *forder*. First, we create a file /d / f and write some data to it. Next, we create a file /d / f.new, which is intended to be a new version of /d / f, and write some new data to it. Finally, we rename /d / f.new to replace /d / f. Because the rename is a write operation on both the new file (changing its name) and the old file (deleting it), it is ordered by  $\Rightarrow$  after all of the other operations. Therefore, we can be assured that, even if some write-behind is lost, the new file will replace the old one only if its intended contents have reached disk.

An ordinary UNIX file system does not provide this guarantee. In most UNIX systems, directory operations are write-through while file data is write-behind, so chances are excellent that /d/f.new will be renamed before its contents have reached disk. If there is a crash before the contents

<sup>&</sup>lt;sup>9</sup>Providing this guarantee exacts a certain cost: Applications that communicate across machines through the file system can do so only at disk speed, not at network speed. Note that NFS effectively provides this guarantee as well, at the same cost.

ACM Transactions on Computer Systems, Vol 12, No 2, May 1994



reach disk, upon reboot the name /d / f will point to a garbage file, and the name's old referent will have been destroyed.

Figure 5 gives an example of replacing a directory /d with a new version. The left-hand column shows an initial attempt at coding the procedure. First, we create the new directory /d.new and write two files into it, /d.new / f1 and /d.new / f2. Then we rename the old version of /d to /d.old and rename the new version to /d. Finally, we can remove /d.old and its contents (not shown). If this procedure is interrupted by a crash, we do a small amount of recovery upon restart—renaming /d.old back to /d if /d does not exist, and then removing /d.new or /d.old if they exist.

This initial attempt would work if write-behind could not be reordered—after recovery,  $/\mathbf{d}$  would be a complete copy of either the old version or the new version—but reordering introduces a problem. Because the final rename operation does *not* have any of the new files in  $/\mathbf{d}$ .new as operands, the  $\Rightarrow$  relation does not order the data writes to those files before the rename. (The file creations are ordered because they have  $/\mathbf{d}$ .new itself as an operand, but the data writes do not.) So when the rename reaches disk, there is no guarantee that the file data is on disk, and an inopportune crash or network fault could leave  $/\mathbf{d}$  as a directory full of garbage files.

This problem is easily fixed by inserting an *forder* call, as shown in gray in the right-hand column of Figure 5. The *forder* call establishes a synchronization barrier, such that all operations on /d.new, /d.new / f1, or /d.new / f2 that were logically performed before the *forder* are ordered by  $\Rightarrow$  before any operations on these operands that are logically performed later. In particular, the three ordering arrows shown in gray ensure that the contents of the files reach disk before the final rename operation.

Applications that make complex changes to data structures stored in the file system may use write-ahead logging. The application first appends an intentions record to a log file and then proceeds to make its changes in any order. When the application is restarted after a crash, it reads the log file and redoes the changes recorded in the log. From time to time, the application reclaims log space (and speeds up the next crash recovery) by trimming off a prefix of the log, after making sure that the changes up to that point have reached disk.

Figure 6 shows how Echo's ordering constraints can be used to improve the performance of write-ahead logging. With a conventional UNIX file system,

ACM Transactions on Computer Systems, Vol. 12, No. 2, May 1994.



Fig 5. Replacing a directory with a new version.

the only way to be sure that a log record reaches disk before the changes it describes is to call *fsync* on the log file after writing the record and before making the changes. Doing this slows down the update by introducing an additional wait for the disk. Although an application using write-ahead logging generally does need to force its log to disk at certain *commit points*, that is, points where the application reports to the user that the updates requested are stable, there are typically many writes to the log between commit points, so one does not want to force the log after every log write. With Echo, one can use *forder* in place of *fsync*, eliminating the need to wait. In Figure 6, after writing the log record, we issue an *forder* whose operands are the log file and each of the files or directories that are to be updated. This *forder* ensures that none of the updates will reach disk before the log record does. When the application does need to force the log at a commit point, it simply calls *fsync* on the log file.

To make a complete, usable write-ahead logging system from this example, we would need to add a means for trimming the log. It is quite straightforward to do this in a simple way that requires updates to be forced to disk with *fsync* before the log is trimmed. Alternatively, it is possible to reduce latency by using *forder* instead of *fsync*. We leave the details of log trimming with *forder* as an exercise for the interested reader.

Echo's ordering constraints make fsync useful not only to ensure that previously written data is stable, but also that previously *read* data is stable. Figure 7 illustrates this. First, one application writes some data to file **f**; later, a second application does a read that returns this data. Suppose the second application wants to make sure the data is stable before printing results, aborting instead if the data is discarded. It can do so by calling *fsync* on **f** at some point after the read and aborting if the call gives an error return. Why does this work? Because the read returned the written data, the read is

ACM Transactions on Computer Systems, Vol. 12, No. 2, May 1994



ordered after the write by  $\rightarrow$ , as shown by the broken arrow in Figure 7. Therefore, the application process *depends on* the write in the sense defined in Section 3.3, so Echo's lost write-behind reporting rules guarantee that, if the write was already discarded when the *fsync* is logically performed, the *fsync* call will give an error return. Hence, if the *fsync* call returns without error, the write must not have been discarded at the time the *fsync* was performed: so by Echo's ordering rules, the *fsync* is ordered after the write by  $\Rightarrow$  (as shown by the solid arrow in the figure), and the write is stable.

Echo's advisory locks provide a convenient building block for implementing other replicated services on top of the Echo replicated file system. One technique is to store the service's data in an Echo volume and to provide access to the data through a pair of servers. In normal operation, one server acts as the primary and holds an exclusive advisory lock on some agreed-upon file in \_\_e volume. The second server acts as a backup and is blocked waiting to acquire its own exclusive lock on the same file. If the primary server crashes or is partitioned away from the file system, Echo revokes its lock. As a result, the backup server acquires the lock and begins running as the new primary; it can then start a new backup if desired. Because Echo's lost write-behind reporting removes the old primary's ability to access the shared volume at the same time it revokes the lock, the new primary can be certain that the old one is no longer modifying the shared data when it takes over. Conversely, an existing primary can be sure that no new primary has taken over as long as it is still able to access the shared data.

A number of the techniques for using Echo's ordering constraints that we have just described have been used in real applications. In particular, the Vesta software configuration management system [Chiu and Levin 1993; Levin and McJones 1993], developed by colleagues at our research center, uses both the file and directory replacement techniques. (Vesta uses slight variants of the techniques presented, as it needs to create new files and directories full of files atomically and not to replace existing ones.) Vesta could also have used the write-ahead logging technique, but unfortunately, the Vesta procedures that need to do logging involve updates to more than one Echo volume.

# 5. PERFORMANCE

In this section we give benchmark measurements that show how Echo's directory write-behind can improve the performance of some applications.

ACM Transactions on Computer Systems, Vol. 12, No. 2, May 1994.

Fig. 7. Checking that data read is stable

write(f, data)



The benchmarks include compiling 100 one-line C programs, running the *make* phase of the Andrew benchmark [Howard et al. 1988], creating a new version of a software package in the Vesta system, and creating and deleting files and directories with some simple loops. We ran each benchmark both with the normal Echo clerk and with a modified version that does synchronous write-through to disk on directory changes instead of doing writebehind; all of the benchmarks run considerably faster with writebehind providing better performance on directory operations than a production-quality UNIX system, even though Echo without write-behind is much slower on directory operations than the UNIX system.

These benchmarks measure the kind of file-system performance that an ordinary application observes. Most applications that write data to files on UNIX-like systems do not concern themselves with when the data actually reaches disk; they are quite satisfied to report success and exit as soon as the system has their data buffered in memory for write-behind. Hence, except where specifically stated, the running times measured in these benchmarks do not include the time to flush writes to disk with *sync* or *fsync*.

In all of the benchmarks discussed below, the Echo clerk and the benchmark programs themselves were running on a four-processor Firefly [Thacker et al. 1987], with each processor using a 3-MIPS MicroVAX chip. The Echo servers were running on six-processor Fireflies and used DEC RA-90 disks. Client and server machines were connected by the 100-Mbit/s Autonet network [Schroeder et al. 1990]. The one UNIX benchmark was run on a one-processor VAX, again using the 3-MIPS MicroVAX chip, but with DEC RA-82 disks.

Table I shows the results of a simple compilation benchmark. One hundred C source files, each containing a single one-line function, were compiled under the control of the UNIX *make* program. Thus, during each run, *make* tested for the existence of 100 object files and created 100 processes running the C compiler. Each C compiler instance read one source file and wrote one object file. The table gives the average elapsed time for three different versions of the benchmark. The first line of the table shows the time to run the benchmark in Echo's normal configuration, with write-behind for both files and directories. The second line shows the running time when the Echo clerk was modified to write through changes to directories, but to continue to write behind changes to file data. (In this benchmark, file creation was the only operation involving a directory change.) The benchmark took 42 percent longer to run in this case. The third line shows the running time when the clerk was forced both to write through directory changes and to write back

ACM Transactions on Computer Systems, Vol. 12, No 2, May 1994

	Mean time (20 runs)	Range	Relative time
Full write-behind	240 s	+12/-4	1.00
Directory write-through	$341 \mathrm{~s}$	+18/-6	1.42
Write-back on close	$366 \mathrm{s}$	+10/-7	1.52

Table I. Elapsed Time to Compile 100 One-Line C Programs

changes to file data whenever a file was closed. In this case, the benchmark took 52 percent longer to run.

Table II shows the results of a similar benchmark, but with a more realistic work load. The *make* phase (Phase V) of the well-known Andrew file-system benchmark from CMU was run in the same three Echo configurations. This phase involves compiling 17 C source files, creating two libraries, and linking one application program.<sup>10</sup> As the table indicates, the benchmark runs took 4 percent longer with directory write-through, and 8 percent longer with both directory write-through and file write-back on close.

What do these results mean? In both benchmarks, directory write-behind gives a clearly measurable improvement in running time. Also, file write-back on close gives a clearly measurable penalty<sup>11</sup>. It is clear that the benefits of directory write-behind depend on how much time an application spends modifying directories. In the first benchmark, the files being compiled are very short, so little time was spent running the compiler. Thus, speeding up file creation gives a large improvement. The second benchmark compiled larger files and also spent time running the linker and library builder, so there was much less benefit in speeding up file creation. It seems likely that real applications will see benefits closer to the 4 percent of the Andrew benchmark than the 42 percent of the simple compilation benchmark. But it is important to note that these benchmarks were run on relatively slow machines by today's standards. In the future, as processors get faster and disk speeds lag further and further behind, applications are likely to see more and more benefit from directory write-behind.

Table III shows the results of a benchmark run on the Vesta software configuration management system. This benchmark was developed by the Vesta research group to evaluate the performance of Vesta's code repository. It simulates storing a new version of a software package in the repository, in the case where all of the source files in the package have actually changed since the previous version, and where the user has already compiled and

<sup>&</sup>lt;sup>10</sup> We report only on Phase V of the Andrew benchmark because the other phases either are very short (Phase I), do no writing at all (Phases III and IV), or focus on large-volume data copying with relatively few directory modifications (Phase II). Running these phases would provide general performance information about Echo, but our goal in this section is restricted to evaluating the effectiveness of directory write-behind.

<sup>&</sup>lt;sup>11</sup>We think of directory write-through with file-data write-behind as the base case in this comparison, because it is what single-machine UNIX file system provide, as do the newer distributed file access protocols emerging as successors to NFS (see Section 7).

ACM Transactions on Computer Systems, Vol 12, No. 2, May 1994.

	Mean time (20 runs)	Range	Relative time
Full write-behind	$546 \mathrm{~s}$	+6/-4	1.00
Directory write-through	$566 \mathrm{~s}$	+5/-5	1.04
Write-back on close	$592 \mathrm{~s}$	+32/-8	1.08

Table II. Elapsed Time for the Make Phase of the Andrew Benchmark

linked the sources to produce a new set of derived object files. Thus, a whole new set of files must be copied into the repository, but there is little else to be done. In this test, there were 8 source files comprising 16 Kbytes of data and 10 derived object files comprising 1478 Kbytes of data.

As Table III shows, the benchmark took 2.43 times as long with directory write-behind turned off and 2.54 times as long with file write-back on close. This result may seem strange. Why should directory write-through slow down the benchmark so much when only 18 files are being written? We can explain this by recalling that the Vesta implementation makes heavy use of Echo's write-behind ordering guarantees to ensure that its data structures remain consistent despite crashes, as discussed in Section 4. In particular, Vesta makes sure that a file has been completely written before it appears in the repository under its permanent name and that all of the files in a directory tree are completely written before the tree appears under its permanent name. So, in this benchmark, there was at least one rename operation ordered after each file was written. When these renames were written through, the files had to be written first, making the benchmark run almost as slowly as if the files themselves had been written through (or written back on close).

Thus, this benchmark makes an interesting point about file-system design choices. When we were designing Echo, we first decided to do directory write-behind and then later decided to give ordering guarantees on the write-behind. But, as we have said, operation reordering and lost writebehind are problems even on conventional UNIX systems where only file data is written behind. So one might consider adding ordering guarantees even to systems without directory write-behind. Unfortunately, as this benchmark shows, directory write-through and ordering guarantees are a bad combination; together they can reduce the performance of file writes to the point where they might as well be write-through too.

As a final benchmark, we did some absolute speed tests against a version of UNIX running with local disks on hardware similar to Echo's. Table IV displays some cases where Echo showed an advantage: creating 100 empty files, deleting 100 files, creating 100 directories, and deleting 100 directories. In each case the 100 operations were performed by a single program executing 100 system calls. Comparing the first two columns of the table, the UNIX system was fast at deleting files, but in every other case, Echo was considerably faster than UNIX. In fact, in some cases Echo was faster even when the timing included the time to flush all of the write-behind using the *sync* 

ACM Transactions on Computer Systems, Vol. 12, No. 2, May 1994

	Mean time (18 runs)	Range	Relative time
Full write-behind	9.9 s	+1.3/-0.4	1.00
Directory write-through	$24.1~{ m s}$	+0.4/-0.3	2.43
Write-back on close	$25.1~\mathrm{s}$	+2.4/-1.0	2.54

Table III. Elapsed Time for the Vesta Benchmark

system call, as shown in the table's third column. In all cases, however, Echo was much slower than UNIX when write-behind was turned off (fourth column).

Some additional information about the internals of UNIX and Echo helps to explain the numbers in Table IV. All of the measured operations are writethrough to disk in the UNIX implementation, so the numbers in the "UNIX with local disk" column each include the time for at least one disk write and. in most cases, more than one. It appears that file deletion requires only one synchronous disk write on the UNIX system we measured, since 1.72 s is roughly the time required for an RA-82 disk drive to rotate 100 times. The numbers in the "Echo with write-behind" column primarily reflect the time needed to update the clerk's in-memory data structures; no disk writes were required. The numbers in the last column, "Echo with write-through," are the sum of the time required for the clerk to update its in-memory structures, the time for 100 remote procedure calls, the time for the server to process these calls, and the time for 100 disk writes. The tests reported in the third column. "Echo with write-behind plus sync," involved the same amount of work as those in the fourth column, but the measured times are much smaller because of pipelining-the client, server, and disk drive were all working in parallel. In addition, group commit to the log could have come into play on the server [Hisgen et al. 1993], reducing the total number of disk writes.

As is evident from these measurements, both the Echo clerk and server required considerably more CPU time to do comparable operations than their UNIX counterparts. The server in particular was heavily CPU bound. We do not believe these performance problems were caused by any fundamental flaw in the Echo approach or algorithms; they were merely an artifact of our prototype implementation.

What overall conclusions can we draw from the benchmarks in this section? We have seen that the speedup an application gets from directory writebehind depends strongly on the amount of computation or other work it does between directory writes. From the final benchmark, we see that the Echo prototype benefits more from directory write-behind than most file systems would, because Echo's servers are relatively slow. These two factors make it difficult to tell just how worthwhile it would be to add directory write-behind to a production file system. There would undoubtedly be some benefits, however. Even though the implementations of the Echo clerk and server are both relatively slow, directory write-behind made some operations faster in Echo than in UNIX. Therefore, it is reasonable to expect that adding direc-

ACM Transactions on Computer Systems, Vol. 12, No. 2, May 1994.

	UNIX with local dısk	Echo with write-behind	Echo with write- behind plus <i>sync</i>	Echo with write-through
Create files	4.95	2.94	7.89	13.1
Delete files	1.72	3.28	4.18	13.8
Create directories	18.6	2.17	9.33	18.7
Delete directories	7.21	1.53	5.45	15 0

Table IV. Elapsed Times for 100 Operations, in Seconds (mean over 20 runs)

tory write-behind to a fast file-system implementation would make it even faster. Also, even a fast file server can become slow when it is loaded down with requests from many clients. Though we did not measure this effect, it seems clear that write-behind can give clients faster responses to their requests when the underlying server is slowed by heavy load.

#### 6. IMPLEMENTATION

In this section we discuss some significant features of the Echo caching implementation.

#### 6.1 Coherence

Tokens are the mechanism Echo uses to keep its caches coherent. The basic token algorithm is simple. Before an Echo clerk can hold data in its cache, it must obtain an appropriate token from the Echo server that stores the data.<sup>12</sup> Holding a *read token* on a data item permits the clerk to read the data from the server and cache the data, but not to modify it. Holding a *write token* on a data item permits the clerk to read to write the changes back to the server. Reading or writing without a token is not permitted. The server keeps track of which clerks have which tokens and prevents conflicts that could cause cache incoherence. In particular, if any clerk holds a write token on the same data. Thus, the caches are single-copy equivalent: Either there are no write tokens outstanding on a data item, in which case the server's copy and all of the cached copies of the item are identical; or there is exactly one clerk with a write token, in which case that couly one accessible.

Whenever a clerk needs a token that it does not have, it makes a remote procedure call (RPC) to the server. If there are no conflicting tokens, the server grants the new token immediately. If there are any conflicting tokens, the server calls each clerk that holds one, asking it to give the token back. When a clerk is asked to give up a read token, it does so immediately. When asked to give up a write token, the clerk immediately writes back any changes it has made to the data that the token covers, and then gives back the token. In either case, the clerk discards its cached copy of the data.

 $<sup>^{12}</sup>$  For the moment, we ignore the fact that Echo servers may be replicated; the impact of replication on the token mechanism is discussed in Section 6.2

ACM Transactions on Computer Systems, Vol. 12, No. 2, May 1994.

In our implementation, the data item covered by a token is a whole file or whole directory. However, clerks can read, write, and cache individual file blocks or directory entries. This token granularity is a convenient choice, but makes for artificially poor performance if two or more applications on different machines are concurrently accessing different parts of the same file or directory, and at least one of them is writing. In such a case, the object's token is continually moving back and forth between the two machines, and their writes are continually being flushed to disk, even if there is no real communication going on between the applications.

This behavior is not a problem on files in our environment, because our users present Echo with a UNIX-like work load in which shared, mutable random-access files are rare. In environments where files are often accessed in this way—VMS, for example—one could change Echo's token mechanism to work on byte ranges instead of on whole files, using the technique Burrows describes in his thesis [Burrows 1988].

We did have a few problems with widely shared, mutable directories, but we were able to work around them. For example, much of the shared software at our installation is kept in *packages*, which are stored as subdirectories of /proj / packages. The directory /proj / topaz / bin contains symbolic links to the executable programs that are stored in these packages. Every user has **/proj / topaz / bin** on his or her shell's search path, and some of the programs in it are frequently used. As a result, every machine needs a read token for /proj / packages most of the time. But the tool that ships new versions of packages to **/proj / packages** works by first creating a new directory with the new contents and then using two rename operations to replace the old version with the new, as in the example of Figure 5, Section 4. These rename operations, of course, require the write token on /proj/ packages. When Echo was first put into use, users often noticed a delay in service when a package was shipped, as their clerks waited to reacquire their read tokens on **/proj / packages**. At first, we were puzzled by this delay, but we soon discovered the cause. The final rename calls in a package installation forced the write token on /proj / packages to be acquired, but since all of the file creations and writes were ordered before the renames, the renames could not be performed, and the token could not be released, until all of this work was completed. This could take a long time-and during this time, any client that needed the read token on /proj / packages would have to wait. Fortunately, modifying the package tool to call *fsync* just before the final rename operations proved to be a satisfactory work-around. With this change, the write token on /proj / packages needs to be held only long enough to do the renames themselves. The package tool runs a bit more slowly, but all other users see better performance, so the trade-off is well worthwhile. If we had encountered more instances of this problem, we could have fixed it in a more general way, perhaps by extending the token mechanism to make directory tokens cover individual entries or alphabetical ranges.

Some write operations in Echo modify more than one file or directory, and, therefore, require more than one token, so care is needed to avoid the possibility of deadlock or livelock. For example, a deadlock could occur if two

ACM Transactions on Computer Systems, Vol. 12, No. 2, May 1994

different clerks both needed the same two write tokens, each clerk had one, and neither clerk would release its token before acquiring the other. On the other hand, a livelock could occur if each clerk was willing to release its first token before acquiring the other; the tokens might bounce back and forth indefinitely with neither clerk able to get both at once.

We solved these problems by defining a fixed order in which tokens must be acquired and held. Unfortunately, the tree structure of files and directories does not define a natural total order that matches the order in which operations find their operands—some operations work down the tree, some work up, and some operate on two unrelated directories—so we had to choose an arbitrary order. We chose to sort the operands by their unique numeric identifiers.

A write operation with several operands (such as rename) therefore proceeds in two phases. In the first phase, the clerk finds all of the operands; this may involve reading a number of directories and acquiring a number of tokens. During this phase, the clerk is willing to release any of the tokens it has acquired immediately. In the second phase, the clerk rechecks that it is holding each of the tokens it needs and reacquires any it has lost since the first phase, proceeding in the defined order for deadlock avoidance. When a token is found or reacquired in the second phase, the clerk marks the object it covers as *dirty* by incrementing a counter associated with the object. After the operation is completed and written back to the server, the dirty counter is decremented. Thus, whenever the clerk has write-behind for an object, the object's dirty counter is greater than zero. If the server asks the clerk to give back a token for an object whose dirty counter is greater than zero, the clerk does so only after driving the counter to zero by sending all of the object's buffered write operations to the server. (If the last of these operations is still in phase two, the clerk, of course, finishes processing it locally before sending it to the server.) If the clerk finds it must reacquire a token during the second phase, it checks to see whether the object that the token covers has been modified since the first phase; if so, the clerk aborts what it is doing, decrements all of the dirty counters it has incremented, and retries the operation from the beginning. This abort and retry are necessary to ensure that the current operation and the operation that modified the object are serialized.

Token-based schemes are flexible. Once we had the basic token scheme designed, it was easy to extend it to handle more details of the UNIX file-system interface.

Our first addition was the *open token*. In a UNIX file system, a file can be deleted from the name space while one or more processes still have it open. When this occurs, the file continues to exist (even though it is nameless) until it is no longer open. To implement this behavior in Echo, a machine's clerk acquires and holds an open token on every file that is open on that machine. Unlike read and write tokens, open tokens are never called back by the server; they do not conflict with any other type of token.

The Echo clerk uses two heuristics to reduce the overhead of acquiring and releasing open tokens. First, the clerk does not release an open token just

ACM Transactions on Computer Systems, Vol. 12, No. 2, May 1994

because no process on its machine still has the file open. Thus, if the same file is repeatedly opened and closed, the token is acquired only once. Second, whenever an application opens a file for reading, the clerk requests both a read token and an open token in one call to the server; this reduces overhead when the file is actually read, which is likely to happen soon. The clerk takes similar action when a file is opened for writing.

Two further heuristics are needed for releasing open tokens, so that deleted files are not kept in existence when no one has them open. First, when an application deletes a file from the name space, if no process on the same machine has the file open, the clerk releases its open token in the same call to the server that deletes the file. Finally, if the server calls back the token that permits deleting a file (a special kind of write token; see Section 6.3) and no application on the machine has the file open, the clerk gives up its open token in the response to the call.

More additions to the set of tokens are discussed in Section 6.3.

#### 6.2 Fault Tolerance

Fault tolerance in the Echo clerk is implemented using *leases, sessions*, and *token replication*.

Echo uses *leases* to allow servers to reclaim tokens from clerks that have crashed, while at the same time ensuring that caches remain single-copy equivalent even when network faults cut off communication between servers and clerks.<sup>13</sup> A lease is an agreement between server and clerk that a clerk's token will remain valid for a given period of time. If the clerk does not renew its lease on a token before the lease expires, the server is free to revoke the token, even if network faults prevent the server from communicating with the clerk. But, until the lease expires, the server is not allowed to revoke a token unilaterally; it must ask the clerk to release the token and wait for a response.

Without leases, there is no safe way for a server to reclaim tokens held by a clerk it cannot communicate with. The clerk machine may have crashed, or it may remain out of communication with the server for a long time, so we certainly want to have the server take the tokens back when other clerks ask to access the files they cover. Yet, if the server unilaterally revokes the tokens while the clerk machine that holds them is still running, then single-copy equivalence may be violated: Applications running on the machine that is out of touch with the server may read (or write) values in the local cache that disagree with the values seen on machines that are still in touch.

But with leases, a clerk can maintain single-copy equivalence simply by checking, each time a client process asks it to read or write cached data, that the lease on the token that covers the data is still in force; if the lease has expired, the clerk does not return the cached data. The clerk does not immediately discard its cached data and write-behind in this case; if the clerk

 $<sup>^{13}</sup>$  The term *lease* was first used by Gray and Cheriton [1989]; we compare their system with ours in Section 7.

is able to contact the server again later and the server has not needed to revoke the clerk's tokens, the server will allow the clerk to reestablish its lease. In the meantime, the clerk can either block client requests or give error returns; as mentioned in Section 3, we chose to have the clerk block requests if the server has been out of touch for less than two minutes and return errors otherwise.

In effect, leases use real time as a communication channel, so that a clerk knows when its tokens are definitely valid and when the server may have revoked them, even if the network connecting the clerk and server is broken. (For this application, the clerk and server clocks do not actually have to be synchronized; it is sufficient for them to run at the same rate within some known error bound, which is fortunate since there is no way to synchronize clocks on two machines that cannot communicate!)

Echo uses *sessions* to reduce the amount of network traffic needed to keep leases up-to-date. Whenever a clerk wants to begin caching data from a new server, it calls the server to establish a session. Individual tokens do not have leases; instead, each token is associated with a particular session, and there is a lease on the session as a whole. This technique dramatically reduces the number of lease renewal messages that must be sent; it also reduces the bookkeeping burden on both clerks and servers. If a session's lease expires and the server needs to revoke one of the tokens associated with the session because another clerk needs it, the server revokes the entire session, invalidating all of the tokens associated with it. Revoking the whole session may seem rather draconian, but it is the most efficient thing to do if the clerk that owns the session has actually crashed; also, it somewhat simplifies the implementation of Echo's semantic guarantees, by making it easy to detect when a process might depend on lost write-behind.

We chose a lease period of 30 s for Echo sessions. A longer lease period would further reduce the number of lease renewal messages sent; a shorter period would speed the recovery of tokens from failed clerks. We found 30 s a satisfactory compromise for our purposes.

Echo provides a fast way for the system to recover when a server crashes by *replicating the token directory* of each server on a backup server. This approach is natural for Echo because (as mentioned in Section 1) Echo already uses replicated servers for high availability. Server replication in Echo uses a primary/backup scheme; at any given moment, one server is the primary for a given set of disks, and all requests from clerks are directed to it. A backup server takes over if the primary crashes. The backup has a path to at least one replica of the data that can function even when the primary is down. Token replication is easy to add on top of this base. Whenever a clerk asks the primary for a token, the primary makes a nested RPC to the backup, which records the token in its copy of the directory; only when this call returns does the primary return the token to the clerk. Both the primary and backup keep their token directories in main memory, not on disk, so that they are fast to access.

One alternative to this scheme would be for the primary server to record the token directory on its (possibly replicated) disk. The main advantage of

ACM Transactions on Computer Systems, Vol. 12, No 2, May 1994

this alternative is that the token database is not lost even if both servers crash. We did not find this advantage compelling, since it is much slower to write the tokens to disk than to record them in memory on two servers. Another apparent advantage of writing tokens to disk is that it does not require two servers, so it works even when servers are not replicated. But token replication does not require the backup server to be a dedicated machine or to have access to a disk, so there is little cost in configuring a backup server for use only by the token machinery.

Another alternative would be for a server that crashes to recover tokens from clerks when it reboots, as is done in the Sprite file system [Nelson et al. 1988]. This scheme could also be used when a backup server takes over from the primary. With this scheme, normal operation is slightly faster, because the primary does not have to call the backup on each token acquisition. But failure recovery is much slower, because many clerks must be contacted and a substantial number of tokens must be recovered from each—two minutes is a typical time for this process in Sprite [Baker and Sullivan 1992]. Also, due to Echo's model and implementation of security (discussed in the next section), the recovering server would have to check the clerks' token lists for consistency with each other and for legality according to file access control lists, slowing down recovery even more. Moreover, if the lists are all legal but are not consistent, there is no way to determine which clerk has made a false claim, so in the worst case, a faulty clerk could cause a nonfaulty one to lose its cache tokens and associated write-behind.

Weighing the advantages and disadvantages, we prefer token replication as the first line of defense against server crashes. However, we encountered enough double server crashes in Echo to convince us that token recovery would have been useful as a second line of defense. It would have considerably reduced the disruption to users in these cases.<sup>14</sup>

#### 6.3 Security

As mentioned in Section 1, Echo servers do not trust Echo clerks. For each operation a clerk requests from a server, the clerk must authenticate itself as acting on behalf of some user who is authorized to perform the operation. For example, to read part of a file into its cache, the clerk has to authenticate itself as acting on behalf of some user, and the clerk has to check the file's ACL (access control list) to see whether that user has read access to the file. A clerk is able to authenticate itself as a particular user if (and only if) that user has logged into the clerk machine.<sup>15</sup> The authentication and login

<sup>&</sup>lt;sup>14</sup>Perhaps the Sprite research group has reached a similar conclusion; Sprite has recently been modified to replicate tokens in a segment of the server's own memory that is usually preserved across reboots, recovering tokens from clients only when this memory is lost [Baker and Sullivan 1992].

<sup>&</sup>lt;sup>15</sup>Thus, a user who logs into a machine obviously must trust its clerk, since the clerk is able to request any operation it pleases on behalf of the user. Because the clerk is part of the machine's operating system, it seems reasonable to require this kind of trust.

ACM Transactions on Computer Systems, Vol. 12, No 2, May 1994.

protocols Echo uses were developed as part of a separate security architecture project and are described in detail by Lampson et al. [1991].

We use two kinds of caching to make access control decisions fast. First, there is caching within the authentication protocol implementation, so that most RPCs are authenticated with no extra packets or cryptographic overhead. Authentication caching is beyond the scope of this paper. Second, we extended the Echo token mechanism to provide caching for ACL checks. We extended the set of tokens so that there is a separate token for each kind of access permission a clerk may have on a file or directory. Thus, the server needs to do ACL checking only on token acquisition requests. On actual read or write requests, the server checks only that the clerk has the appropriate tokens, which is considerably faster (and which the server would have to do anyway, since it does not trust the clerk).

In the extended set of tokens, the read token described in Section 6.1 is replaced by three tokens: *InfoToken*, *SearchToken*, and *ReadToken*. Holding an InfoToken on an object allows a clerk to read and cache a record of status information about the object, corresponding to what is returned by the *stat* system call in UNIX. There are no access control restrictions on obtaining an InfoToken. Holding a SearchToken on a directory allows looking up individual names in the directory and caching the results, but not reading the entire directory. Holding a ReadToken allows reading file data or reading a directory in its entirety. A clerk can obtain a SearchToken or ReadToken only if acting on behalf of a user with the corresponding access permission.

There is no token corresponding to execute-only file access, because there is no way for the server to enforce the distinction between reading and executing. For a user to execute a program, the clerk on his machine must be able to read it; but if the clerk is able to read the program, there is no way for the server to keep the clerk from letting the user read it. Therefore, a clerk is allowed to obtain a ReadToken when acting on behalf of a user with executeonly access.

The write token is replaced by three tokens: *WriteToken, ChangeAccessToken*, and *ChangeParentToken*. Holding a WriteToken on a file allows writing data or changing the file's length; on a directory, it allows modifications such as creating, deleting, or renaming objects named in the directory. Holding a ChangeAccessToken on an object allows changing the object's ACL. To obtain either of these tokens, a clerk must be acting on behalf of a user with the corresponding access permission: write access for WriteToken, and ownership for ChangeAccessToken. Finally, a clerk must hold a ChangeParentToken on an object to rename it or to delete it. There are no access control restrictions on obtaining a ChangeParentToken.

The ChangeAccessToken has a special property: A clerk that holds this token on an object has blanket permission to acquire any other tokens on the same object that it asks for, with no ACL check. This property is needed because the clerk may have written behind a change to the ACL that makes it legal to obtain the tokens it is asking for, even if the server's copy of the ACL says it is not. There is no security hole here. If the clerk is acting on behalf of a user who can change the ACL, then the server must allow the

ACM Transactions on Computer Systems, Vol. 12, No. 2, May 1994.

clerk to change the ACL whenever it asks, so there is nothing to be gained by forcing the clerk to write back such a change before it can take advantage of the changed permissions.

Which tokens does a clerk need to perform which file-system operations? Echo provides a complete set of UNIX-like file-system operations, and it would be tedious to list the tokens required for all of them here, so we explain the most complicated case as an example in the next paragraph. The reader can work out what tokens are needed for the other, simpler operations from the token definitions given above.

Suppose that a clerk is asked to rename a directory from /a/b to /x/y/c/d and that there is already an empty directory named /x/y/c/d, which will be effectively deleted by the rename operation. To look up the pathnames involved, the clerk will need a ReadToken or SearchToken on each directory along each of the two paths. To check that directory d is empty, the clerk will need a ReadToken or InfoToken on it. The clerk can (if asked to) release these tokens before actually calling the server to carry out the rename; in fact, it can release the token on each directory before obtaining a token for the next one on the path—pathname lookup is not atomic. To actually carry out the rename, the clerk must hold the following tokens:

- -the ChangeParentToken on b;
- -the ChangeParentToken on d;
- -the WriteToken on a;
- -the WriteToken on **c**; and
- —the InfoTokens on x and y, assuming that / is the root directory of an Echo volume.

The only surprise here is the need for the InfoTokens. When a UNIX-like file system renames a directory, it must ensure that this does not create a cycle; that is, a directory must not be renamed so as to become a child of one of its own descendants. So, when a directory rename is requested, the Echo clerk scans the directory tree upward from the new parent directory to the root of the volume and refuses to do the rename if it finds the directory to be moved on the path. The InfoToken is needed on each of these directories to determine its identity and find its parent, and these InfoTokens must be held throughout the rename operation to ensure that a concurrent rename requested by some other clerk cannot change one of the parents after it has been checked.

Table V summarizes the compatibility rules among all types of cache coherence tokens. Where "No" appears at the intersection of a row and column in the table, it is not permitted for two different clerks to hold the two tokens named at the head of the row and column on the same object at the same time; where "Yes" appears, it is permitted. OpenToken applies only to files, and SearchToken only to directories, so there is no table entry at their intersection.

How did we choose these compatibility rules? It is clear that ReadToken and WriteToken cannot be compatible, because WriteToken grants the right

ACM Transactions on Computer Systems, Vol. 12, No 2, May 1994.

	Open	Info	Search	Read	Write	Change Access	Change Parent
Open	Yes	Yes		Yes	Yes	Yes	Yes
Info	Yes	Yes	Yes	Yes	No	No	No
Search		Yes	Yes	Yes	No	No	No
Read	Yes	Yes	Yes	Yes	No	No	No
Write	Yes	No	No	No	No	No	No
Change Access	Yes	No	No	No	No	No	No
ChangeParent	Yes	No	No	No	No	No	No

Table V Compatibility Matrix for Cache Coherence Tokens

to change information that ReadToken grants the right to cache. Similar reasoning explains many other incompatible token pairs; in particular, Info-Token grants the right to cache information (such as length, ACL, and link count) that can be changed by clerks holding WriteToken, ChangeAccessToken, or ChangeParentToken, so these tokens cannot be compatible with InfoToken. It would have been possible for us to make a few more pairs of tokens compatible than we did, since they protect disjoint pieces of information: for example, we could have allowed one clerk to hold the WriteToken on a file while another holds its ChangeParentToken. But we did not find a use for such additional compatibility in our implementation: To simplify our handling of cached object status information, we adopted the convention that, whenever a clerk holds any token on an object (other than the OpenToken), it also holds the InfoToken. Thus, continuing the previous example, if a clerk holds the WriteToken on an object, it also holds the InfoToken, and so no other clerk could hold the object's ChangeParentToken even if we made WriteToken and ChangeParentToken compatible.

Our security implementation departs from single-machine UNIX semantics in one detail. In UNIX, a file's ACL is checked only when the file is opened. So, if a process has a file open for writing, it can continue to write the file indefinitely, even if the file's ACL is changed to remove the process's write access; and similarly for reading. We could have emulated this feature by splitting the OpenToken into OpenReadToken and OpenWrite Token, with the property that holding such a token allows a clerk to obtain the corresponding ReadToken or WriteToken without an ACL check. We chose not to do so because it seems too complicated and because we feel that this UNIX feature is not a good idea in a fault-tolerant distributed system. On a single-machine UNIX system, one can always force all processes to close their open files by rebooting. But there is nothing corresponding to rebooting in Echo, so there is no way to ensure that no process is holding a file open when its ACL is changed to remove access. Unfortunately, it turns out that a few applications we wanted to run depend on the UNIX semantics, so we adopted a compromise: Echo servers allow the owner of a file to get a read or write token even if the file's ACL does not give the owner the corresponding access. This feature allows the owner (and only the owner) to continue accessing an open file after its permissions are taken away. (It could also allow a file's owner to open the file without having access permission, but the Echo clerk

ACM Transactions on Computer Systems, Vol 12, No 2, May 1994.

does not allow this.) As with the special rule for ChangeAccessToken, there is no security hole here, because a file's owner has permission to change the ACL and give himself access. Restricting this feature to the owner has not been a problem for any UNIX application we have tried. NFS has essentially the same feature for the same reason and also restricts it to a file's owner.

# 6.4 Resource Reservations

Creating or enlarging a file or directory requires resources on the file server —chiefly disk space. Echo allows clerks to reserve resources in advance, so that, when a clerk buffers an operation for write-behind, the user can be assured that the necessary resources will be available when the write is performed.

Disk space reservations are fairly simple. Each write operation that the clerk sends to the server uses up some of the clerk's reserved disk space. The clerk is linked with a library routine provided by the server that tells it how much space each operation requires. (For complex operations that require varying amounts of space, the routine computes a simple, conservative estimate.) Whenever an application asks the clerk to do a write, the clerk first checks whether it has the necessary space reserved; if not, it asks the server for more. The clerk specifies two numbers in its request: a minimum amount and a desired amount. The minimum amount is what is needed to complete the current application request. If the server grants less than the minimum, the clerk returns a "Disk full" error to the application. The desired amount is the amount the clerk would like to reserve ahead to reduce the number of future reservation requests it has to make.

A clerk's disk-space reservation is associated with its session, so if the session's lease runs out, the server can reclaim the reserved space. Unlike tokens, however, we do not replicate space reservations; instead, we use a lazy form of recovery for them. (We chose to do this because write requests are more frequent than token acquisitions, and we did not want to slow them down by adding an extra RPC to the backup.) After a primary server crashes and the backup takes over, each clerk's reservation is recovered the next time the clerk sends a write request to the server. The server responds to the write request by asking the clerk to send its current reservation value and to retry the request. There is a minor security hole in this mechanism. A malicious clerk could lie about its reservation, claiming to have reserved much more space than it really did. This could cause the server's disk space to be overcommitted, so that other clerks would not be able to get back their legitimate reservations. We decided not to worry about this problem-it is unlikely to arise in practice, and we have not found a solution less expensive than replicating the reservations or recording them on disk, both of which seem impractical.

Unique identifiers for files and directories are another resource that clerks reserve. Echo servers do not trust clerks to generate these identifiers correctly, because a malicious clerk might give the same identifier to two different files, thereby corrupting the file system's tree structure. Therefore,

ACM Transactions on Computer Systems, Vol. 12, No 2, May 1994

so that file and directory creation can be write-behind, clerks are allowed to reserve a stock of these identifiers ahead of time. A clerk's reserved identifiers are associated with its session, so the server can reclaim the space needed to keep track of them if the session's lease expires. A clerk is automatically granted all of the coherence tokens on each object that it creates using a reserved identifier; this avoids an extra call to the server upon object creation.

# 6.5 Ordering Constraints

The Echo clerk's implementation of ordering constraints falls out naturally from the way it does write-behind. For each cached object that has been locally modified, the clerk keeps a representation of the object's current state (or at least the modified portion), plus a write-behind queue-a list of operations that have not yet been written back to the server, in the order they were logically performed. An operation with more than one operand is on the write-behind queue of each operand. A sequence of file overwrites unbroken by any other write operation appears as a single element in the file's write-behind queue. An forder operation appears in the write-behind queue of each of its operands, but is treated as a no-op when it reaches the head of the queue. This queue data structure corresponds precisely to the  $\Rightarrow$  relation for write-behind ordering that we defined in Section 3.1. Before the clerk sends any write operation to the server, it checks that the operation is at the head of the write-behind queue for each of its operands. If not, the clerk first sends out the operation's predecessors, after recursively checking that each of them is at the head of the write-behind queue for each of its operands.

To improve performance, we use a pipeline to send write requests from clerk to server. Each write request is an RPC, but several calls may be issued in parallel by separate threads. Each call carries a sequence number, however, so that the server knows the order in which they must be physically carried out. Thus, under load, the clerk is able to send a continuous stream of requests to the server; it does not have to wait for one request to complete and the reply to arrive over the network before it sends out another.

# 6.6 Reporting Lost Write-Behind

To identify the processes that depend on discarded write-behind, the Echo clerk uses a simple, conservative technique. When a process accesses a volume for the first time, the clerk's current session on that volume is recorded as part of the process state. When a process accesses a volume it has accessed before, the saved session identifier is compared with the current session identifier for the volume; if the identifiers differ, the process may depend on write-behind that was discarded, so it receives an error return.

This technique is more conservative than necessary, however. First, there may not actually have been any write-behind lost in the old session. Second, even if write-behind was lost, some processes that accessed the volume may not actually have been affected by the loss. It would be easy to fix the first

ACM Transactions on Computer Systems, Vol. 12, No 2, May 1994

problem with a small amount of added bookkeeping in the clerk, but fixing the second would require substantial extra bookkeeping.

Whether fixing these problems is worthwhile depends on how often they occur, which, in turn, depends on how often network faults occur and how often there is actually write-behind buffered when one does occur. When we were designing Echo, we believed that token revocation due to network faults would be very rare, so the extra bookkeeping would not be worthwhile; unfortunately, this turned out not to be the case. Fixing the first problem would have had a substantial payoff in Echo, because the workstations in our installation are lightly loaded, so at any given moment a workstation usually has no write-behind buffered (Section 3.2). But, clearly, this property could easily change if Echo users started doing more computing or started running a different mix of applications.

# 6.7 Advisory Lock Tokens

The Echo clerk implements advisory locks using another set of tokens for each file and directory. If any process on a machine holds a shared advisory lock on an object, that machine's clerk must have a SharedToken for the object. If any process holds an exclusive lock, the clerk must have an ExclusiveToken. Table VI gives the compatibility matrix for advisory lock tokens.

As with cache coherence tokens, a clerk may hold an advisory lock token even when no process on its machine needs it. Thus, if advisory locks on an object are repeatedly acquired and released by applications running on one machine, no communication with the server is necessary.

When a clerk requests an advisory lock token that conflicts with one held by a second clerk, the server calls back the clerk holding the conflicting token. If the conflicting token is not in use (i.e., no process on the second clerk's machine is holding an advisory lock that needs it), the second clerk returns it immediately. If the token is in use, however, the callback blocks until the token is no longer in use. The blockage propagates back through the server to the first clerk, and ultimately causes the application that called the first clerk requesting the conflicting lock to block until the lock is released. Other threads within the server and clerks continue to run.

A problem with this implementation is that conflicting locks result in RPCs that remain outstanding for a long time. In the RPC implementation we are using, such calls consume excessive resources in the RPC run-time library. We know of two ways to fix this problem, but have not actually implemented either of them. The most obvious fix is to change the RPC system (or to add a layer on top of it) so that long-term outstanding calls are not expensive. Another possibility is to take the states that are currently represented by long-term outstanding RPCs and instead represent them with additional token types. We designed a protocol that does this using two additional token types for each kind of advisory lock token: WantSharedToken, WantExclusiveToken, KeepSharedToken, and KeepExclusiveToken. If a clerk wants an advisory lock token, but cannot be given it immediately, the server gives the

ACM Transactions on Computer Systems, Vol. 12, No 2, May 1994.

Table VI. Compatibility Matrix for Advisory Lock Tokens

	Shared	Exclusive	
Shared	Yes	No	
Exclusive	No	No	

clerk the corresponding WantToken instead of blocking its call. If a clerk asks for an advisory lock token and there are no other clerks contending for it, the clerk is also granted the corresponding KeepToken. When the server wants to take back the lock token, if the clerk cannot give it back immediately, it gives back the KeepToken to the server instead of blocking the server's call. When a clerk no longer needs a lock token, it does not keep it unless it also holds the corresponding KeepToken. When a server obtains a lock token back from a clerk, if another clerk has the corresponding WantToken, the server takes back the WantToken and grants the token the clerk wanted in its place (along with the KeepToken, if appropriate).

A clerk's advisory lock tokens are associated with its session. Thus, if the session's lease expires, the server can reclaim the tokens and grant them to another clerk. Error returns due to discarded locks are generated using the same mechanism as those due to lost write-behind. This implementation gives the desired semantics for advisory locks, as described in Section 3.3.

# 7. RELATED WORK

Throughout this section (and in other parts of this paper) we use Echo terminology to discuss related systems, even when the papers describing those systems use different terminology. In particular, other systems often use different terms for what we call *volumes*, *clerks*, and *tokens*.

Like Echo, the Sprite file system uses tokens to maintain coherence in a distributed file cache [Nelson et al. 1988]. Sprite's caching differs in several respects from Echo's, however. Sprite clerks cache only files, not directories. When a Sprite application process asks to open a file, the local clerk sends the request on to the server machine that stores the file, and the pathname lookup is done there. Sprite clerks do *prefix caching* so that most requests can be sent directly to the correct server, without the need to broadcast or consult a name server first [Welch and Ousterhout 1986]. A prefix cache is not a full-fledged directory cache; it simply maps prefixes of absolute pathnames to the servers that store the files whose names begin with those prefixes. For example, if the directory subtree rooted at **/usr** is stored as a single volume on file server **fred**, a clerk's prefix cache might contain the mapping from **/usr** to **fred**, but the clerk would have no other information about **/usr** and no information at all about longer pathnames such as **/usr / include**.

A Sprite clerk is required to hold a read token on a file to have it open for reading or a write token to have it open for writing. This differs from the Echo scheme, where read and write tokens are needed only when a clerk is actually caching file data. If two Sprite clerks want to hold the same file open

ACM Transactions on Computer Systems, Vol. 12, No. 2, May 1994

in conflicting modes, the server detects the conflict and turns off caching for the file, requiring all reads and writes to go to the server. Doing this seems like a good idea if two machines are really communicating data through the file: If one machine is only writing and the other is only reading, caching is pointless, and the overhead of moving tokens back and forth between the machines is wasted. On the other hand, if two machines are accessing different parts of a random-access file, turning off caching is not the best solution; a better idea would be to give each machine a token on just the range of bytes within the file it is using.

When we were designing Echo, we looked for a way to avoid useless caching and token-passing overhead when two applications are actively communicating through the file system, but we did not find a solution we were satisfied with. We could certainly have adopted a scheme like Sprite's, so that caching on a file would be turned off when applications on two different machines have it open in conflicting modes. However, this scheme may turn off caching in some cases where it would be beneficial, and it does not catch cases where files are shared sequentially: One process opens a file, writes it, and closes it, then another process opens it, reads it, and closes it; and this pattern repeats. It is also unclear how to extend the scheme to directories. As with files, one would like to turn off caching for a directory when processes on different machines are actively communicating by modifying it. But processes do not announce the start and end of their access to a directory by opening and closing it, so one would have to use a heuristic to decide when to turn directory caching on and off. Perhaps such a heuristic could also detect sequentially shared files. Further research seems to be needed in this area; however, when Baker et al. [1991] compared the performance of the Sprite and Echo schemes for handling write-shared files, they found that the choice made little difference on the work-load traces they were able to gather.

Sprite does not replicate its token directory. When a Sprite file server reboots after a crash, it reconstructs its token directory by contacting all of its clients and asking them what tokens they hold. The pros and cons of this approach were discussed in Section 6.2.

Sprite does not use leases. If a server is unable to contact a clerk (i.e., if its RPC to the clerk times out), the server invalidates the clerk's tokens immediately. Therefore, for the reasons discussed in Section 6.2, Sprite does not provide strict single-copy equivalence. However, because a Sprite clerk communicates with its server on every file open and because Sprite clerks do not cache directories, a Sprite client can access stale data during a partition only in files that the client has open at the time their tokens are revoked, and only as long as those files remain open. Thus, there seems to be little or no practical disadvantage in this departure from single-copy equivalence. In fact, one can argue that there is a practical advantage in Sprite's approach, because, after a clerk fails or is partitioned, Sprite can reclaim that clerk's tokens more quickly than a lease-based system can: The delay is at most one RPC time-out, rather than the full lease time-out.

Sprite blocks client file operations interruptibly when communication between the Sprite clerk and server is broken. A process requesting such an

ACM Transactions on Computer Systems, Vol 12, No 2, May 1994

operation blocks until either communication is restored or a user interrupts it from the keyboard. This behavior is similar to the *hard*, *interruptible* mount option on NFS volumes. It differs from the Echo behavior (described in Section 3), where such operations are blocked noninterruptibly for up to two minutes and then unblocked with an error return. Sprite may have made the better choice here; we considered switching to the Sprite behavior after gaining some experience with our initial design, but we did not have time to implement the change.

The Andrew File System (AFS) [Howard et al. 1988; Kazar 1988] caches both files and directories on client machines. On directory updates, AFS does write-through, not write-behind. AFS caches files in their entirety, not block by block; when a client process tries to open for reading a file that is not already in the cache, the AFS clerk reads the complete file from the server. When a client process writes a file, the AFS clerk buffers the changes locally until the file is closed and then writes them through to the server. Thus, AFS does not provide single-copy equivalence when multiple client machines hold the same file open at the same time and at least one is writing it.

The AFS file cache is kept coherent using a scheme similar to Echo's; what AFS terms a *callback* is similar to an Echo token. AFS's tokens time out—they must be periodically renewed to remain valid—but AFS servers do not give the same lease guarantee to their clerks that Echo servers do. If an AFS server tries to ask a clerk to give a token back, but cannot contact the clerk over the network, it immediately revokes the clerk's tokens without waiting for the token time-out to expire. If this happens and the clerk in question is still running, single-copy equivalence is violated; all processes running on the clerk's machine can see stale data until the time-out expires. As with Sprite, this violation of strict semantics is defensible as a practical choice.

Burrows has implemented a file caching service, called MFS [Burrows 1988], that is similar to Echo's in many respects. MFS caches both files and directories, with write-behind for both, and with coherence implemented using a token scheme similar to Echo's. MFS does not use leases—if an MFS token server is unable to contact a clerk, it immediately revokes all of the clerk's tokens—and thus, MFS does not provide strict single-copy equivalence in the face of network partitions. MFS tokens do not time out at all, so there is no bound on how long a clerk that is partitioned away from its server may operate with stale data before it detects that the server has revoked its tokens; however, the clerk will certainly find out on its next cache miss.

MFS does not provide any guarantees about the order in which updates are written back to the file server; in particular, updates to a directory can be reordered arbitrarily, as long as the reordered update sequence would (in the absence of faults) leave the directory in the same state as did the original logical sequence. For example, if a user creates file /d/a and then file /d/b, but a crash causes some write-behind to be lost, file /d/b may exist while /d/a is never created. Or, if a user writes some data into file /d/c, then renames it to /d/d, and a crash causes some write-behind to be lost, the file may have been renamed without the data having been written to it.

ACM Transactions on Computer Systems, Vol 12, No 2, May 1994

MFS does an excellent job of saving work for the server by canceling sequences of operations that have no net result before they leave the clerk's write-behind buffer. For example, if an application creates a temporary file, writes into it, reads it, and deletes it within a short time, MFS avoids sending any of the operations to the server. (Actually, this sequence changes the last-modified time of the directory in which the file was created, so MFS should tell the server about that change, but apparently this was not done.) Our original plans for Echo called for us to do this kind of optimization, but we did not find time to implement it.<sup>16</sup> MFS's lack of ordering guarantees makes it easier to implement these optimizations than it would have been in Echo.

Burrows introduced the concept of *byte-range tokens* in MFS. Logically, each byte of an MFS file has its own token, and tokens on different bytes within a file do not conflict. Changing the length of a file requires holding write tokens on the bytes being added or deleted. A file's other properties (such as owner and last-modified time) are covered by another token. In data structures and interfaces, sets of tokens on a file are represented as ranges; this representation is very compact except in pathological cases. As mentioned above, we believe that byte-range tokens would be a useful addition to the Echo token scheme under work loads where shared, mutable randomaccess files are common.

Gray and Cheriton [1989] coined the term *lease* in a paper about their file caching server. We developed the lease concept simultaneously and independently of Gray and Cheriton, but chose to adopt their terminology when we learned of their work. Gray and Cheriton's system does not use write-behind at all, so it needs only one kind of token. If a clerk has a token on a file, it may cache a clean copy as long as its lease remains valid. If a clerk wants to write a file, it sends the write request directly to the file server, which recalls all tokens on the file before performing the write. (Alternatively, the server can delay the write until all of the leases have expired, refusing to honor any lease renewal requests that come in during the waiting period.) Gray and Cheriton's system does not have sessions; each cached file has an independent lease.

QuickSilver [Schmuck and Wyllie 1991] provides a transactional interface to its file system. Operations on multiple files and directories can be grouped into a transaction that is committed or aborted as a unit. Although we chose not to explore transactional file semantics in Echo, we view this area with interest and would be pleased to see research in it succeed in producing a practical system. The transactional programming model is cleaner and seemingly easier to use than Echo's model of ordered write-behind, while still allowing write-behind for transactions that have not yet been committed. Moreover, QuickSilver transactions can include operations on both files and nonfile objects, and can cover files in different volumes managed by different

<sup>&</sup>lt;sup>16</sup>As implemented, Echo collapses multiple overwrites to the same file bytes if they are not ordered by  $\Rightarrow$ , but it does no other work-canceling optimizations.

ACM Transactions on Computer Systems, Vol. 12, No. 2, May 1994.

servers. A drawback of the QuickSilver approach is that it requires considerably more machinery to implement, giving rise to concern about how well it may perform compared to more conventional approaches. Also, many applications that use the file system need to be modified to work properly with QuickSilver. By default, every file-system action taken in a QuickSilver program is part of a single atomic transaction that commits when the program exits, but this behavior is not appropriate for many programs, for example, long-running text editors.<sup>17</sup>

For a more extensive bibliography on file systems that do caching on client machines, see Burrows' thesis [Burrows 1988].

# 8. CONCLUSIONS

The Echo distributed file system has studied a collection of useful techniques for improving the performance and semantics of distributed file caches and has demonstrated their feasibility. We believe that future file systems will benefit from adopting many of these techniques. In particular, Echo provides fully coherent file and directory caching on clients, with ordered write-behind on updates to both files and directories.

Coherent caching simplifies the task of writing distributed applications that use the file system, by allowing processes running on different machines to communicate simply and reliably by reading and writing files, just as though they were all on the same machine. With incoherent caching (as in NFS), it is quite tricky to write such applications. Together with Echo's location-transparent global naming, coherent caching makes the distributed nature of the file system invisible to applications, while providing much better performance than an uncached file system could. The bookkeeping cost of maintaining cache coherence seems well worth the benefits.

For coherent caching with write-behind to work well, however, the underlying network must have good availability, and the medium used to store the system's token directory must be reliable.

If the network is often broken, clients will often be unable to use the file system, even if they have all of the files they need cached, because they are unable to communicate with the server and, thus, are unable to be sure their caches are coherent. The Echo system uses the Autonet network [Schroeder et al. 1990], which achieves very high availability through redundancy, though Ethernet and most other nonredundant local-area networks have high enough availability for the purpose. On networks that are often unavailable, however, file systems that allow controlled forms of incoherence, such as Coda [Kistler and Satyanaraynan 1991] and Ficus [Guy and Popek 1991; Page et al. 1991], may be more practical.

<sup>&</sup>lt;sup>17</sup>Of course, for some applications, using Echo's semantics requires modifying applications by adding *forder* calls. But these modifications are optional; if they are omitted, the only problem is that the application does not tolerate lost write-behind well, a problem that is most likely just as bad when the application is run on a single-machine UNIX system with write-behind.

ACM Transactions on Computer Systems, Vol. 12, No 2, May 1994

If the system's token directory is lost, all client machines lose their writebehind, disrupting the work of many users. Therefore, the token directory must be either replicated, recoverable from clients after a server crash, or both. We prefer token replication as the first line of defense because it makes for fast recovery, at the cost of updating two replicas on each token acquisition. It seems worthwhile to implement token recovery as well, as a second line of defense when all replicas of the token directory are lost.

Write-behind is an important building block for file systems, whether distributed or not. Nearly every file system we are aware of does write-behind of some kind, though often only on file data. Write-behind certainly reduces latency, by eliminating the need for applications to wait for the disk on every write. It can also improve throughput, by smoothing out load peaks and by enabling sequences of operations that have no net result to be canceled before they leave the write-behind buffer.

Echo does write-behind on directory modifications, including file and directory creation and deletion, and we have shown that this can improve the performance of some applications. For example, it speeds the compilation of programs made up of many small modules, like typical large C programs and libraries. This speedup is not large, but may grow as the disparity between CPU and disk speeds increases. Directory write-behind adds complexity to the file cache implementation, but we implemented it successfully and had few problems with the code once it was in place.

The major drawback of write-behind is that a write can fail long after the process that requested it has gone on to other things or has even exited. Directory write-behind does not make this problem worse in any fundamental way, but it did make both the Echo designers and our users more conscious of the problem. We tried to provide good facilities for allowing applications to tolerate lost write-behind cleanly, and we did make considerable progress in this area; however, lost write-behind remains an ugly problem. Fundamentally, write-behind cannot provide a correct implementation of the natural semantics programs expect to get from a file-system interface: When you write to a file, you want the bits to be stored stably every time, not most of the time. Perhaps the ultimate solution to this problem is to eliminate write-behind, instead using nonvolatile RAM on each file server to shield clients from the latency of synchronous disk writes [Baker et al. 1992].

Echo took a two-pronged approach to dealing with lost write-behind. First, make it easy for an application writer to ensure that, when a crash halts the application and causes some of its write-behind to be lost, the data structures it stores in the file system remain consistent. Second, make it easy to ensure that whenever write-behind is lost, all affected applications are halted (or otherwise notified), so that they can recover cleanly instead of continuing to run with an incorrect notion of what is on disk.

Echo did well on the first prong. It actually provides better semantics than are commonly provided even on nondistributed file systems (such as conventional UNIX). Echo's ordering constraints do the right thing automatically for many simple applications and provide enough power for more complex ones, while giving better performance than if only *fsync* were available. Ordering

ACM Transactions on Computer Systems, Vol. 12, No. 2, May 1994

constraints were easy to implement as a part of doing directory write-behind. On the down side, ordering constraints do appear to be more difficult for application writers to use than transactions would be.

We were less successful on the second prong. Our current design is good at making directly affected processes halt, but it also often halts processes that did not really care about the lost write-behind, causing confusion for users. And processes that indirectly depend on lost write-behind cannot be detected: for example, in a distributed or multiprocess application where only some processes access the file system, those that do not access it may not learn of lost write-behind.

Moreover, we have no solution for the problem of what to do when an application's write-behind is lost after it has exited. These cases escape our two-pronged approach entirely.

Echo's approach to write-behind was motivated by the semantic goal stated in Section 2: We wanted to present application programs with a file-system interface and semantics that were as close as possible to those of a singlemachine file system, so that existing applications would work without change and programmers would not have to learn a new set of skills. But, in the end, we did change the interface and semantics somewhat, by adding the *forder* primitive and defining ordering semantics for the other primitives. In the resulting system, all applications written for a single-machine system can be run with no changes, but some are not fault-tolerant without changes. Other sets of small interface changes might also be worth exploring in the future, either as alternatives or as additions to Echo's. For example, one might allow users to force write-through on selected files using new file attributes, or on all writes done by selected programs using file attributes or environment variables. Such facilities seem attractive for use with programs that are not available in source form or that are too complex to modify. On the other hand, they put the burden on the user to decide when write-behind must be suppressed and when it can be allowed.

# ACKNOWLEDGMENTS

Hania Gajewska, Jim Gettys, Mark Manasse, and Mike Schroeder contributed ideas early in the Echo clerk's design phase. Mike Burrows and Mike Schroeder helped in the selection of material for this paper and provided useful comments on the presentation. The anonymous referees helped to improve the clarity and completeness of the paper.

#### REFERENCES

- BAKER, M., AND SULLIVAN, M. 1992. The recovery box. Using fast recovery to provide high availability in the UNIX environment. In *Proceedings Summer 1992 USENIX Conference* (June) USENIX Association, Berkeley, Calif., pp. 31–43.
- BAKER, M., ASAMI, S., DEPRIT, E., OUSTERHOUT, J., AND SELTZER. M., 1992. Non-volatile memory for fast, reliable file systems. In *Proceedings 5th International Conference on Architectural* Support for Programming Languages and Operating Systems (Boston, Mass, Oct. 12–15). ACM, New York, pp 10–22

ACM Transactions on Computer Systems, Vol 12, No 2, May 1994

- 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 13th Symposium on Operating Systems Principles* (Pacific Grove, Calif., Oct. 13–16). ACM, New York, pp. 198–212.
- BIRRELL, A. D., HISGEN, A., JERIAN, C., MANN, T., AND SWART, G. 1993. The Echo distributed file system. Res. Rep. Systems Research Center, Digital Equipment Corporation, Palo Alto, Calif, Sept.
- BURROWS, M., 1988. Efficient data sharing. Ph.D. thesis, Computer Laboratory, Univ. of Cambridge, U.K. Sept.
- CHIU S.-Y., AND LEVIN, R. 1993. The Vesta repository: A file system extension for software development. Res. Rep. 106, Systems Research Center, Digital Equipment Corporation, Palo Alto, Calif.
- GLASSMAN, L., GRINBERG, D., HIBBARD, C., REID, L. G., AND VAN LEUNEN, M. C. 1992. Hector: Connecting words with definitions. Res. Rep. 92A, Systems Research Center, Digital Equipment Corporation, Palo Alto, Calif., Oct.
- GRAY, C. G., AND CHERITON, D. R., 1989. Leases: An efficient fault-tolerant mechanism for distributed file cache consistency. In *Proceedings 12th Symposium on Operating Systems Principles* (Litchfield Park, Ariz., Dec. 3-6) ACM, New York, pp. 202-210.
- GUY, R. G., AND POPEK, G. J. 1991. Algorithms for consistency in optimistically replicated file systems. Tech. Rep. CSD-910006, UCLA Computer Science Dept. UCLA, Los Angeles, Calif., Mar.
- HISGEN, A., BIRRELL, A., JERIAN, C., MANN, T., AND SWART, G. 1992. Some consequences of excess load on the Echo replicated file system. In *Proceedings 2nd Workshop on the Management of Replicated Data* (Monterey, Calif., Nov. 12–13). IEEE, New York, pp. 92–95.
- HISGEN, A., BIRRELL, A., JERIAN, C., MANN, T., AND SWART, G., 1993. New-value logging in the Echo replicated file system. Res. Rep. 104, Systems Research Center, Digital Equipment Corporation, Palo Alto, Calif., June.
- HISGEN, A., BIRRELL, A., MANN, T., SCHROEDER, M., AND SWART, G. 1989. Availability and consistency tradeoffs in the Echo distributed file system. In *Proceedings 2nd Workshop on Workstation Operating Systems* (Pacific Grove, Calif., Sept. 27–29). IEEE, New York, pp. 49–54.
- HISGEN, A., BIRRELL, A., JERIAN, C., MANN, T., SCHROEDER, M., AND SWART, G., 1990. Granularity and semantic level of replication in the Echo distributed file system. In *Proceedings Workshop* on the Management of Replicated Data (Houston, Tex., Nov. 8-9). IEEE, New York, pp. 2-4.
- HOWARD, J. H., KAZAR, M. L., MENEES, S. G., NICHOLS, D. A., SATYANARAYNAN, 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.
- KAZAR, M. L., 1988. Synchronization and caching issues in the Andrew file system. In Proceedings Winter 1988 USENIX Conference (Dallas, Tex., Feb. 9–12). USENIX Association, Berkeley, Calif., pp. 27–36.
- KISTLER, J. J., AND SATYANARAYNAN, M. 1991. Disconnected operation in the coda file system. In Proceedings 13th Symposium on Operating Systems Principles (Oct.). ACM, New York, pp. 213-225.
- LAMPSON, B., ABADI, M., BURROWS, AND M., AND WOBBER, E. 1991. Authentication in distributed systems: Theory and practice. In *Proceedings 13th Symposium on Operating Systems Princi*ples (Oct.). ACM, New York, pp. 165–182.
- LEVIN, R., AND MCJONES, P. 1993. The Vesta approach to precise configuration of large software systems. Res. Rep. 105, Systems Research Center, Digital Equipment Corporation, Palo Alto, Calif.
- MANN, T., HISGEN, A., AND SWART, G. 1989. An algorithm for data replication. Res. Rep. 46, Systems Research Center, Digital Equipment Corporation, Palo Alto, Calif. June.
- 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.
- PAGE, T. W. JR., GUY, R. G., POPEK, G. J., AND HEIDEMANN, J. S. 1991. Architecture of the Ficus scalable replicated file system. Tech. Rep. CSD-910005, UCLA Computer Science Dept., UCLA, Los Angeles, Calif., Mar.

POPEK G., WALKER, B., CHOW, J., EDWARDS, D., KLINE, C., RUDISIN, G. AND THIEL, G. 1981.

ACM Transactions on Computer Systems, Vol. 12, No. 2, May 1994.

LOCUS: A network transparent, high reliability distributed system. In *Proceedings 8th Symposium on Operating Systems Principles* (Pacific Grove, Calif., Dec. 14–16). ACM, New York, pp. 169–177.

- SANDBERG, R., GOLDBERG, D., KLEIMAN, S., WALSH, D., AND LYON, B. 1985. Design and implementation of the Sun network filesystem. In *Proceedings Summer 1985 USENIX Conference* (Portland, Ore., June. 11-14). USENIX Association, Berkeley, Calif., pp. 119-130.
- SCHMUCK, F., AND WYLLIE, J. 1991. Experience with transactions in QuickSilver. In Proceedings 13th Symposium on Operating Systems Principles (Pacific Grove, Calif., Oct. 13-16). ACM, New York, pp. 239-253.
- SCHROEDER, M. D., BIRRELL, A. D., BURROWS, M., MURRAY, H., NEEDHAM, R. M., RODEHEFFER, T L., SATTERHWAITE, E. H., AND THACKER, C. P. 1990 Autonet: A high-speed, self-configuring local area network using point-to-point links. Res. Rep. 59, Systems Research Center, Digital Equipment Corporation, Palo Alto, Calif., Apr.
- SUN MICROSYSTEMS, INC. 1989. NFS: Network file system protocol specification. RFC 1094, Network Information Center, SRI International, Menlo Park, Calif., Mar.
- SWART, G., BIRRELL, A., HISGEN, A., JERIAN, C., AND MANN, T 1993. Availability in the Echo file system. Res. Rep. 112, Systems Research Center, Digital Equipment Corporation, Palo Alto, Calif., Sept.
- THACKER, C. P., STEWART, L. C., AND SATTERTHWAITE, E. H, JR. 1987. Firefly: A multiprocessor workstation. Res. Rep. 23, Systems Research Center, Digital Equipment Corporation, Palo Alto, Calif., Dec.
- WALKER, B., POPEK, G., ENGLISH, R., KLINE, C. AND THIEL, G 1983. The LOCUS distributed operating system. In *Proceedings 9th Symposium on Operating Systems Principles* (Bretton Woods, N.H., Oct. 12–13). ACM, New York, pp. 49–70.
- WELCH, B., AND OUSTERHOUT, J. 1986. Prefix tables: A simple mechanism for locating files in a distributed filesystem. In Proceedings 6th International Conference on Distributed Computing Systems (Cambridge, Mass., May 19–23). IEEE, New York, pp. 184–189

Received June 1993; revised January 1994; accepted February 1994

ACM Transactions on Computer Systems, Vol 12, No. 2, May 1994