

Extending ACID Semantics to the File System

CHARLES P. WRIGHT, RICHARD SPILLANE, GOPALAN SIVATHANU, and
EREZ ZADOK

An organization's data is often its most valuable asset, but today's file systems provide few facilities to ensure its safety. Databases, on the other hand, have long provided transactions. Transactions are useful because they provide atomicity, consistency, isolation, and durability (ACID). Many applications could make use of these semantics, but databases have a wide variety of non-standard interfaces. For example, applications like mail servers currently perform elaborate error handling to ensure atomicity and consistency, because it is easier than using a DBMS. A transaction-oriented programming model eliminates complex error-handling code, because failed operations can simply be aborted without side effects. We have designed a file system that exports ACID transactions to user-level applications, while preserving the ubiquitous and convenient POSIX interface. In our prototype ACID file system, called Amino, updated applications can protect arbitrary sequences of system calls within a transaction. Unmodified applications operate without any changes, but each system call is transaction protected. We also built a recoverable memory library with support for nested transactions to allow applications to keep their in-memory data structures consistent with the file system. Our performance evaluation shows that ACID semantics can be added to applications with acceptable overheads. When Amino adds atomicity, consistency, and isolation functionality to an application, it performs close to Ext3. Amino achieves durability up to 27% faster than Ext3, thanks to improved locality.

Categories and Subject Descriptors: D.2.5 [Software Engineering]: Testing and Debugging—*Error handling and recovery*; D.2.5 [Software Engineering]: Testing and Debugging—*Tracing*; D.4.3 [Operating Systems]: File Systems Management—*Access methods, Directory structures, File organization*; D.4.5 [Operating Systems]: Reliability—*Fault-tolerance*; H.2 [Database management]: Database Applications

General Terms: Design, Experimentation, Performance, Reliability

Additional Key Words and Phrases: File system transactions, Recoverable memory, Databases, File systems, `ptrace` monitors

1. INTRODUCTION

File systems offer a convenient and standard interface for user applications to store data, which is many organizations' most valuable asset. Computer hardware and software can be replaced, but lost or corrupted data can not. Providing reliable file system access is therefore an important goal of any operating system.

Database systems provide strong guarantees for the safety and consistency of data, but each database uses its own interface. Four key requirements define a transaction: atomicity, consistency, isolation, and durability—collectively known as the *ACID* properties. Despite their importance, most file systems have made no provisions to ensure that operations meet all four of these stringent requirements. Our goal is to combine the best part of databases, their reliability (embodied by the ACID properties)—with the best part of file systems, their common and easy-to-use POSIX API [16].

Next, we describe the ACID requirements, and how they relate to file systems.

Atomicity. Atomicity means that operations must complete or fail as a whole unit. Traditionally, file systems provided only limited atomicity (e.g., renaming a

file either fails or succeeds). Many applications undertake arduous procedures to try to perform atomic operations. For example, if Sendmail fails when attempting to append new mail messages to a mailbox, it then attempts to truncate the file to erase a partially written message [40]. Yet if the truncation fails, then the mailbox is left corrupted. To solve these problems, a file system should allow a sequence of operations to be encapsulated in a single atomic transaction. This has two key benefits: (1) error handling becomes easier, because transactions can simply be aborted, and (2) data corruption cannot occur, because no corrupted data ever reaches the file system. With this new functionality, Sendmail's append operations could be wrapped in a transaction. If they all succeeded, then Sendmail would commit the transaction. Otherwise, Sendmail would abort the transaction and the file-system state would not change.

Consistency. In the context of a database system, consistency means that the database enforces pre-defined integrity constraints. Examples of integrity constraints in a database system are that social security numbers must be unique or that a checking account must have a positive balance. File systems have similar constraints (e.g., inode numbers are unique and no directory entry points to a non-existent inode). By wrapping related operations in database transactions, a file system can maintain a consistent on-disk state.

Applications also have consistency requirements. For example, when committing files to CVS [2], lock files are created to protect against concurrent accesses. An integrity constraint in this example is that lock files only exist while an instance of CVS is updating the repository. In an unmodified CVS implementation, there are circumstances in which lock files are not properly deleted (e.g., on unexpected termination or occasionally when the user presses Control-C). Using transactions greatly improves error handling—with only four lines of code we were able to prevent CVS from leaving stale lock files. Additionally, we eliminated the possibility of some files being committed, and others not (e.g., if the process is terminated half-way through a commit). If CVS were to have used a transactional model from the start, then hundreds of lines of code through several source files could have been eliminated. Moreover, because the transactional interface does not commit data until all operations succeed, error-handling is much more robust than the several ad-hoc functions that are currently in use.

Isolation. Isolation (or serialization) means that one transaction will not affect the execution of another concurrently running transaction. This is not available in current file systems. For example, a set-UID program cannot use `access` to check whether a user has permission to create a file, because another process could create a symbolic link to a sensitive file between the `access` and the creation. This is known as a time-of-check-time-of-use (TOCTOU) security vulnerability. With a file system that maintains isolation, for example, `access` and file creation can safely be performed in a single transaction so that no other operations could be interleaved between the `access` and the creation; to improve performance, however, other operations may be interleaved, but the database management system ensures that the results are as if there was no interleaving.

Durability. Once a transaction is committed to disk, the data remains intact even across a software or a hardware crash. This is a desirable property for every

Table I. File system support for ACID. Current file systems cannot provide all ACID properties across multiple operations, but many do provide a subset of the ACID properties for a single operation (i.e., a system call or VFS-level operation). Amino provides all of the ACID properties for an arbitrary sequence of multiple operations.

* FFS-no-SU denotes FFS without SoftUpdates, and FFS+SU denotes FFS with SoftUpdates.

	Ext2 and FFS-no-SU*	Ext3	FFS+SU*	Amino
Atomicity	No	Single op	No	Multiple ops
Consistency	No	Multiple ops	Multiple ops, but resources may leak	Multiple ops
Isolation	Single op	Single op	Single op	Multiple ops
Durability	Only with O_SYNC	Only with O_SYNC	Only with O_SYNC	Legacy: each op. Enhanced: on commit.

application, but often operating systems (OSes) choose to sacrifice durability for better performance. OSs often make this choice because the synchronous I/O that is often required for durability can result in poor performance. Databases employ optimizations such as sequential logs, group commit, and ordered writes to provide durability efficiently.

As seen in Table I, current file systems do not support full ACID properties. Traditional file systems do not provide atomicity. For example, during `rename`, Ext2 and FFS can both create the file's new name, and then fail before the old name is removed. Journaling file systems like Ext3 provide atomicity for a single operation, so a rename operation cannot fail half-way through, but they do not provide atomicity for a sequence of multiple operations, which is vital for user applications. Many file systems do not provide consistency, which has resulted in the need to run a consistency checker before mounting them (`fsck`). Journaling file systems and SoftUpdates ensure that each operation is consistent, so the composition of many operations is also consistent [24]. Current file systems use VFS-level locking to provide isolation for a single operation. For example, a directory is locked before it is modified. However, there is no mechanism to isolate one sequence of operations from another operation (or sequence). To improve performance, current file systems do not provide durable writes unless the `O_SYNC` option is specified.

We believe that the ACID properties are desirable for many applications, especially applications like email that are expected to be highly reliable, or applications that require atomicity and isolation for security (e.g., updating a user's credentials). Therefore, we have designed a file system called *Amino* that extends ACID semantics to standard applications that use the POSIX interface. Legacy support is essential: unmodified applications and file systems continue to work as they have in the past. To exercise fine-grained control over transactions, existing applications need only slight modifications, and benefit from improved reliability.

It can be argued that databases are already taking over for the file system when reliable storage is required. For example, some commercial email systems store messages in databases instead of the file system [41], and it is becoming more common for revision-control systems to store information in a database [5]. However, we believe that writing applications that use the file system interface has inherent advantages over writing applications that use the database interface. When

an application is written to a database API, it severely limits its interoperability and adds to the burden of programmers and administrators. For example, with a mail server using a file system, an individual user's mail file can simply be copied to create a backup, or deleted to remove all of the user's messages (from personal experience working at an ISP, this is a not uncommon request). Moreover, any standard text processing package can be used to edit the file. When data is only accessible through a database interface, these types of convenient access are no longer possible. Instead, special applications must be written for each of these functionalities.

We have built Amino on top of the Berkeley Database (BDB) [38]. BDB is an embedded database package that provides efficient transaction-protected access to key-value pairs in hash tables or balanced trees. BDB provides the crucial database infrastructure such as logging, locking, and caching. However, BDB, does not provide or require the use of SQL, stored procedures, a specialized database server, or other heavyweight components often associated with a DBMS. This makes it ideal for use by other operating system components. Using BDB allows us to leverage almost 200,000 lines of time-tested industrial-strength code.

If we were to implement Amino as a traditional file system that interfaces with the VFS, we would be required to use the inode, dentry, and page caches. If a transaction aborted, then these caches would become stale with respect to the database. Therefore, we chose to implement Amino as a user-level monitor using the process-tracing facility (`ptrace`) provided by Linux. This interface allows us to intercept all system calls and use only the internal BDB caches. For internal Amino data structures, we developed a recoverable virtual memory (RVM) system on top of BDB. Our RVM system provides support for nested transactions and is transparent to applications.

We evaluated our prototype, and show that it can add atomicity, consistency, and isolation to existing applications with negligible performance overheads. To provide durability, existing file systems require an application to issue explicit `fsync` calls. Amino can implicitly provide durability, and is 20.2% faster than a traditional file system with `fsync` calls. If a programmer informs Amino when transactions begin and end, durable performance is 26.6% better than a traditional file system. Given that Amino is an unoptimized user-level prototype, we find these results encouraging and expect that performance can improve with more tuning.

The rest of this article is organized as follows. Section 2 provides an overview of our design. Section 3 describes our current Amino prototype. Section 4 evaluates Amino's performance. Section 5 describes related work. We conclude and discuss future work in Section 6.

2. DESIGN

The key decision to make when extending ACID semantics to a file system is whether to graft additional code to provide transactions onto an existing file system, or to build a file system on top of a system that already provides transactional semantics. The advantages of adding code to the file system is that you may end up with less overall code, which is more specialized to the task at hand. However, adding even a subset of the required code to an existing file system can take years.

For example, Ext3 shares most of its code with Ext2 and only adds atomicity to single file system operations, but it took more than two years to develop. To get a rough idea of how large a file system is versus a transactional processing system, we can compare the number of lines of code in Ext3 to the number of lines in version 4.1 of the open-source MySQL server [29] and version 4.2.52 of the Berkeley Database (BDB) [38; 42]. In Linux 2.6.11.12, Ext3 has 21,629 lines of code (including the block journaling layer, *jbd*, which is used only for Ext3). BDB has over 16,870 lines of code in just its transaction-related components, and BDB is a subset of MySQL’s overall transaction code (MySQL uses BDB to provide transactional tables). Aside from the transaction-related components, BDB provides efficient data access methods for key-value pairs (e.g., BDB’s balanced-tree implementation is 16,843 lines of code). We therefore chose to build our file system on top of BDB, because we can leverage the already existing transactions infrastructure and efficient access methods.

Once we decided to build the file system on top of a transaction-processing system, the next question was what transaction-processing system is an appropriate host for the file system. One option would have been to use an SQL server such as MySQL, PostgreSQL, or Oracle. We rejected using a full-fledged SQL server, because they require significant runtime resources. Moreover, each database update or query requires communication over a socket, adding extra context switches and data copies. These context switches and especially data copies could hurt file system performance. We therefore chose to use an embedded database, which runs directly in the address space of the client—thereby eliminating context switches and data copies. BDB fits our needs well. It is widely deployed, and has been thoroughly tested. BDB also scales both up and down: it can have a small memory footprint of less than 500KB, yet it also can be configured for databases as large as 256TB. BDB’s codebase is still tractable at about 200,000 lines of code. There are two key reasons that BDB’s codebase is manageable. First, BDB does not require or support SQL parsing, query planning, or other features often associated with a DBMS. As these features are not needed for a file system, having less code is a distinct advantage. Second, BDB has a highly modular design and the application designer can choose which components are required (e.g., the transaction subsystem can be used with normal files, or the access methods can be used without logging). Even though BDB is a relatively small DBMS, it still provides the key infrastructure for full ACID semantics: logging, locking, recovery, and a full-featured transactions API. It also provides four data access methods: a sorted balanced search tree, extended linear hashing, a fixed-length record queue, and access by logical record number.

The rest of this section is organized as follows. Section 2.1 provides an overview of BDB and its supported operations. Section 2.2 describes our database schema. Section 2.3 describes our internal use of transactions. Section 2.4 describes our use of transactional memory. Section 2.5 describes the transactions API that we expose to applications.

2.1 BDB Overview

BDB provides a uniform API to access both hash tables and balanced search trees in a transactional manner. To open and use BDB databases, a database environ-

ment is opened first. The database environment provides caching, logging, and locking functionality for one or more databases (or even simple files). Transactions are associated with the environment, and they have three operations: begin, commit, and abort. Other database operations are protected by the transaction. If a transaction is committed, then all of the protected operations are applied to stable storage as a whole. If the transaction is aborted, then it has no effects. A single transaction can span multiple databases, but the databases must all belong to the same environment. Before a database is opened, a database handle is created and associated with an environment. Next, the handle's parameters are set (e.g., the page size, sorting or hashing function, etc.). Finally, the database is opened inside of a transaction using the fully configured handle. After the database or databases are opened, key-value pairs can be stored using a PUT operation and retrieved using a GET operation. These primitives take the database handle, a transaction, the key, and the value (for PUT) as arguments. Also, BDB provides support for *cursors*, which efficiently iterate through items in the database. The primary cursor operations we are concerned with are `DB_SET`, `DB_SET_RANGE`, and `DB_NEXT`, which find a given key, the first key that is greater than a given key, and the next key, respectively. There are many other BDB operations and parameters, which we omit here for brevity [42].

2.2 File System Schema

The database schema defines the format of our file system. The schema dictates the topology of the data, which in turn is directly related to what operations are possible, and how efficient each operation is. Our primary goal in developing our schema was to minimize the number of database accesses required for any given operation, because I/O operations are many orders of magnitude slower than in-memory operations. An organization that is appropriate for a normal disk-based file system is not necessarily appropriate for a database. For example, most FFS-like file systems use simple mappings of integers to disk blocks [25]. For example, to read a block from a file, first the root inode number is mapped to a disk block. After the root inode is read, the root directory's data blocks are scanned to find the inode number of the next pathname component. Reading each data block essentially maps a logical block in the file to a physical disk block using the inode's direct and indirect pointers. This procedure must be repeated for each pathname component, until the file is found.

BDB, on the other hand, provides more complex and efficient data structures. In BDB, the schema is defined by the set of databases and their key-value pairs. A file system can conceptually be divided into two halves: (1) a naming component and (2) a data storage component. For example, FreeBSD has a separate UFS component for naming and an FFS component for storage. Our schema, shown in Table II, has a similar division. We use a *Path* database to map pathnames to unique file identifiers, and a *Data* database to map unique file identifiers to file data. The *Orphan* database contains a list of identifiers that are not accessible through the name space, but is otherwise equivalent to the Path database.

In the rest of this section we describe the design considerations when developing our schema. First we discuss each database in turn: the Path database, the Data database, and then the Orphan database. We then describe path-local and data-

Table II. Our database schema. Directory-reading and lookup operations use the Path database, which maps full path names to path-local meta-data. Read, write, truncate, and other data-oriented operations use the Data database. The Data database has two types of keys: a file identifier points to its meta-data, and a file identifier concatenated with a page index point to the page's data. Files without any names are stored in the Orphan database.

Database	Key	Value
Path	Full Path	ID Path-local meta-data (e.g., stat information for a file without hardlinks)
Data	ID	Reference Count Data-local meta-data (e.g., stat information for a hard linked file)
	ID Page index	Page's data
Orphan	ID	Path-local meta-data (e.g., stat information for a file without hard links)

local meta-data.

The Path Database. The Path database is used for both lookup and directory-reading operations. Each file has a unique identifier, which is analogous to an inode number. In the Path database, the key is a full pathname and the value is a unique identifier. We designed our schema such that a given file can be looked up using a single database access. For any given path name we can quickly find the path's unique identifier, without the need to traverse each component's directory separately as is done in most Unix file systems. The Google file system uses a similar scheme [10]. When using a hash function, this yields constant time lookups. Using a balanced tree with a fan-out of 100 keys per page, four disk accesses are always sufficient to find any of 10^8 files.

The Path database is also suitable for the directory-reading operation. As the access method for the Path database, we selected a balanced tree structure using a customized sort function. In our database, pathnames are first sorted by depth (i.e., by an ascending number of pathname components) and then by standard lexicographic order. Using this sorting function means that for any given directory, every name is contiguous within the database. To read a directory, we use BDB's `DB.SET_RANGE` operator to position a cursor at the first path name within the directory. To read each subsequent entry we use the cursor's `DB.NEXT` operator until we encounter a path name in a different directory.

For the `lookup` operation, the sort function is not critical, as a name can be located correctly with any total ordering. However, our sorting function proves advantageous when reading a directory and performing `stat` operations on the entries. Because each path in the directory is located close to one another, fewer pages must be read in from disk. This type of operation is quite common (e.g, by `ls -l` or recursive tree scans), which is why NFSv3 introduced a single protocol primitive called `READDIRPLUS` for it [3].

The Data Database. To store the data pages, we use a balanced tree. If a file's unique identifier is stored in the tree, then the given file exists. We assign the identifier randomly, but as the tree is sorted, it is possible to influence data layout policies by modifying the identifier assignment and sort function. For each identifier, the database stores the file's reference counts and meta-data. There are two reference counts: one for the number of path names that reference it (a.k.a. a link count), and another for the number of open instances of the file.

The actual data associated with the file is also stored in the Data database. For a given page of the file, the key is the file's identifier concatenated with the page index. We first sort the tree by the file's identifier and then by the page index. This means that all of a file's data pages are allocated contiguously in the tree, thereby improving locality and allowing the use of database cursors.

Selecting database parameters properly is of the utmost importance for the Data database. In our experiments we found that there can be a factor of ten difference in performance based on page size, cursor use, and other database-tuning parameters. The page size is a particularly important parameter for data-intensive operations. BDB uses a configurable database page size of powers-of-two between 512 bytes and 64KB. It is often useful to make this page size the native page size of the underlying file system, so that BDB reads and writes pages that are compatible with the OS's native page size. The BDB page size also determines when and how *overflow pages* are used. For the Data database, most records are rather large, so they are stored in overflow pages, which means that they are not stored directly with the key. We have found that BDB will store only a single record within an overflow page. Therefore, if the database page size is larger than our file system's transfer unit (for the remainder of this paragraph we refer to our file systems page as a transfer unit to avoid confusion with BDB pages), then the remainder of the database's overflow page is wasted, reducing available disk space and imposing unnecessary I/O overheads. Similarly, if the overflow page size is less than or equal to the file system transfer unit, then BDB stores a small amount of internal meta-data in the beginning of the overflow page, and the first part of the actual data in the remainder of the first overflow page. Another complete overflow page is used for any remaining data, and the rest of it is wasted.

BDB's overflow page allocation behavior means that the file system transfer unit must be carefully selected to avoid performance conflicts with BDB. For example, with a file system transfer unit of 4,096 bytes and the default BDB page size of 16,384 bytes, only 4,122 bytes on each overflow page are used (4,096 for the data, and 26 bytes for BDB's internal meta-data), wasting the remaining $\frac{3}{4}$ of the page. This not only wastes space, but hurts performance because useless data needs to be sent to and from the disk. With a database page size of 4,096 and an equal transfer size, 26 bytes of meta-data are stored on the first overflow page and only 4,070 bytes of actual file-system data can be stored. On the second overflow page only the remaining 26 bytes of file-system data are stored—wasting nearly half of the space. Because of these considerations, we have chosen to use a transfer unit of 4,070 for our file system. Although this is a non-standard size, well-behaved applications should execute the `fstat` system call to find the optimal transfer unit stored in the `st_blksize` field. Poorly behaved applications work as expected, but with degraded performance. Our benchmarks show that when using a 4,096 block size, there is a 4% slow down for sequential reads, and an 18% slow down for sequential writes. Random operations have a greater performance penalty, because they do not benefit from locality as the sequential workloads do: reads are slowed by 48.4% and writes by 57.1%.

We have also found that using database cursors is essential for good sequential read performance. Simply iterating through the Data database using the GET

primitive without cursors can be twice as slow as sequentially reading the file with a cursor. Therefore, whenever possible we use cursor reads with the more efficient `DB_NEXT` flag instead of simple `GET` operations. We do not use write cursors as they are incompatible with transactions, and require locking an entire database environment.

The Orphan Database. Files that have been unlinked, but are still open, are stored in the Orphan database. The Orphan database is identical to the Path database, except that instead of storing the name, only the file's unique identifier is stored. In case of a system crash, we can quickly locate and remove all such orphaned files using a database cursor during the next mount.

Path-local and Data-local Meta-data. The `stat` system call returns vital information about a file, such as its size, owner, and access permissions. The performance of `stat` is quite important, as it constitutes a large portion of many workloads. Ellard's traces of NFS-mounted home directories show 24.6–72.4% of all calls were `GETATTR` and `ACCESS`, which both require `stat` information [8]. Because each file has a single set of attributes, the file's unique identifier determines the `stat` information even if the file has multiple pathnames. This means that the `stat` attributes are a *functional dependency* of the unique identifier. To avoid *logical redundancy*, or having the same data stored in two different places, and its associated pitfalls in a traditional SQL database, the `stat` information should be stored in a database with the unique identifier as the key [21]. In our schema, logical redundancy would introduce *update anomalies* in which one copy of the data could be updated, but the other might not. However, if `stat` information could be stored in the Path database, then performance would be improved because `stat` would require only one database access.

To solve this problem, we take advantage of the flexibility provided by BDB's key-value pair model to develop a more dynamic schema. Meta-data is divided into two classes: (1) *path-local* meta-data and (2) *data-local* meta-data. Path-local meta-data includes all meta-data that is specific to one path of a file. Data-local meta-data includes all meta-data that may refer to more than one path. For example, a newly created file's `stat` information is stored as path-local meta-data, because there is no other path name that references this `stat` information. However, if a hard link to the file is created, then the path-local meta-data is promoted to data-local meta-data, as both names could be used to reference the same underlying file. If one of the links is removed, then the data-local `stat` information could be demoted to path-local meta-data. Dividing meta-data into path-local and data-local components allows our schema to avoid the pitfalls associated with logical redundancy. Yet when the data has no logical redundancy, the `stat` information is stored right with the pathname to improve performance.

2.3 Internal File System Transactions

It is essential that each operation in an ACID file system be protected by a transaction. This is true even when the application that is executing that operation is not concerned with ACID semantics, because other applications must access a single consistent view of the database to ensure the isolation property. Also, to ensure that the file system is consistent, certain integrity constraints must be maintained.

We define our file system to be consistent, if and only if it meets the following seven integrity constraints:

UNIQID Each file identifier is unique.

REFCOUNT Each file's link reference count is equal to the number of path names that reference it.

NOORPHANEDFILES Each data-local meta-data block has a positive link count or open instance reference count. If the link count is zero, then an entry for this file must exist in the Orphan database.

NOORPHANEDBLOCKS Each data page in the Data database has an associated data-local meta-data block.

HARDLINKUSESDLMD If and only if a file has a link reference count greater than one, then it uses data-local meta-data.

PAGESMATCHSIZE A file has no data pages with an index greater than or equal to $\lceil \frac{FileSize}{TransferUnit} \rceil$.

LASTPAGEMATCHESIZE If there is page at index $\lfloor \frac{FileSize}{TransferUnit} \rfloor$, then it is no larger than $FileSize \bmod TransferUnit$ bytes.

Each of these integrity constraints is equivalent to a similar invariant in a standard file system and is also equivalent to common integrity constraints enforced by a database system. For example, **REFCOUNT** is equivalent to a foreign key constraint, and standard file systems verify the same when performing a **fsck**. In traditional file systems, constraints similar to **PAGESMATCHSIZE** and **LASTPAGEMATCHESIZE** are checked by **fsck** to ensure that no orphaned blocks exist, and that stale data does not reappear, respectively.

Our file system does not require a **fsck**, nor does it explicitly enforce the integrity constraints. Instead, each file system operation is designed to transition from one consistent file system state to another consistent file system state. Because each file system operation is surrounded by a transaction, it is atomically applied or it has no effect. Therefore, our file system is always consistent (because it meets the required integrity constraints). This strategy is different from *enforcement*, in that enforcement would require validating the constraints before committing every transaction. To recover the file system after a crash, it is enough to open the database with BDB's **DB_RECOVERY** flag, which replays the database log, and to remove any orphaned files (we efficiently locate these files using the Orphan database). BDB's internal support for recovery obviates the need for us to take complicated recovery steps in our file system code.

2.4 Transactional Memory

One major difficulty with any system that supports transactional semantics is how to deal with an abort operation. Transactional systems should be able to rewind to the state they were at just before the transaction began. This is part of supporting atomic behavior: the effects of a sequence of operations are realized if and only if the transaction containing that sequence is committed. If a transaction is aborted, the operations that were already performed must be reversed so that the state returns to how it was just before the transaction had begun. Of course, this is not restricted to the file system data: caches and other book-keeping memory regions

that describe the state of the file system also need to be reversible in this manner (e.g., the process's open file table).

Through the use of BDB's support for application-specific recovery, we built a recoverable virtual memory (RVM) library. Our library supports rolling back allocation, deallocation, and writes to a recoverable region. Because one of our requirements was to allow applications to use nested transactions, our RVM library supports nested transactions. By allocating memory regions related to the file system state with our recoverable memory routines, we can easily rewind our state to the proper one upon abort. Our library internally uses `mmap`, `mprotect`, and signal handlers to protect memory regions transparently.

After catching the page faults, we log the memory's content. This allows us to access recoverable memory transparently using traditional memory references, without the need for error-prone explicit logging functions. This is especially important if the library is to be used to retrofit transactional semantics onto existing applications or infrastructure. It is relatively easy to locate all of the points where data structures are allocated and deallocated, whereas locating each access to a data structure can be very difficult.

2.5 Transactions API for Applications

Legacy applications need no changes to enjoy the benefits of a consistent file system, which uses transactions for each individual operation (as applications do today with a journaling file system). However, some applications require more stringent atomicity, consistency, isolation, and durability properties. For example, a mail server must append large messages to the end of a mailbox, and a password update system must consistently update `/etc/passwd` and `/etc/shadow` together. Importantly, both legacy and enhanced applications can coexist and use the same data—without the need to access a data store using a specialized interface.

For these types of applications, our file system exports a transactions API to user applications. Our primary design goal for the API was to avoid any changes to existing system calls, which means that we could not add a transaction argument to each call. To begin a transaction, an application issues a new system call that associates a *current transaction* with the process (or thread in multi-threaded applications). Each file system operation after that point is protected by the current transaction. The application can then commit or abort the transaction, with the expected semantics: an aborted transaction has no effect on the file system, and a committed transaction is safely written to stable storage. Aborting a transaction can greatly simplify error handling code, but developers still must take care not to persistently change state during an aborted transaction (e.g., internal application data structures). One simple way to ensure this property is to exit after an abort (many programs already exit on unexpected failures). A better option is to use our RVM facilities to rewind data structures transparently. We believe that one reason many applications are structured such that error handling consists of shutting down the current process or thread is that ad-hoc error recovery is so difficult, hard to debug, and error-prone that fault-tolerant applications, despite their benefits, are often impractical to develop on current systems. We believe that if transactional semantics for the file system and data structures were provided, then programmers may structure their programs to be more robust in the face of failures rather than

coding their programs to exit upon failure.

Using BDB's support for nested transactions, each of the file system's internal transactions is started as a child of the current transaction. This simplifies error handling in the file system, because a transaction for a failed system call can just be aborted. If the child transaction is committed, then it is committed to stable storage only if the parent transaction is committed as well. If a child transaction is aborted, then its effects are undone, but the parent transaction can continue. Our design makes use of this, by wrapping each individual system call in a transaction. In this way, our file system can abort transactions, even if the application is wrapping a set of system calls into a transaction. This functionality is also exposed to user applications. If a process already has a current transaction, and a new transaction is created, then a new current transaction is created as a child of the existing current transaction. This creates a stack of nested current transactions associated with the process.

We employ a simple shared-memory like API to allow processes to share transactions, and we support multiple concurrent transactions without changing the existing system call API. When a transaction begins, it is assigned a unique identifier that the process can then use to manipulate the transaction. A process with sufficient permissions can set its current transaction by attaching to the unique identifier. In this way, two processes can share the same transaction. Similarly, a process can detach from its current transaction, so that future operations are not transaction protected. If all processes have detached from a transaction, then it is automatically aborted (this policy ensures that no transaction-protected data reaches the file system if it was not explicitly committed). If a process temporarily wants to stop using a transaction, but not abort it, then it may suspend the transaction (e.g., to temporarily switch between transactions). The suspend and detach primitives allow processes to switch between transactions without adding system call arguments. For example, a network server may concurrently service many separate clients. Each client's data should be protected by separate transactions. On exit, all uncommitted transactions are automatically detached.

Transactions can be automatically inserted into an existing application's system call stream using pre-defined *profiles*. For example, a profile can protect an entire application by inserting a begin-transaction call on `exec`, and a commit-transaction call on `exit`. Another profile could use file sessions to insert transactions [34]: on the first `open` system call, a transaction is begun; on each subsequent successful open, a counter is incremented; and decremented on close. When the counter reaches zero, then the transaction is committed. Other transaction profiles can be designed and developed, either for a general class of applications or even for the behavior of a specific application.

3. IMPLEMENTATION

We developed a prototype ACID file system on Linux, called *Amino*. The key implementation question for our file system is how to intercept calls and direct them to the database transparently. We evaluated six techniques with respect to the following two criteria:

- Legacy applications should not be modified. In the best case, unmodified binaries

can run without recompiling or relinking. We also considered techniques in which the application must be recompiled or relinked, but its source code is unmodified.

- The interception technique should not insert caches between the application’s system calls and the database. This is because any caches that are not managed by the database suffer from two problems. First, if a transaction that spans multiple operations is aborted, then the cache becomes stale. Second, if the caches are accessed without consulting the database, then the isolation property is violated.

Finally, we considered the implementation effort and attempted to minimize changes to existing infrastructure. We considered six choices.

In-kernel file system. The most direct approach would be to write a standard in-kernel file system. In-kernel file systems do not require relinking of binaries, and such file systems fit into the existing kernel architecture. They also have the advantage of running in kernel mode, so they can minimize data copies and context switches.

In-kernel file systems, however, have two key disadvantages. The first is that standard in-kernel file systems are intimately tied together with caches. This means that substantial code changes would be required to ensure coherency between the internal database caches and the external VFS caches. The second disadvantage is that all of the database code would need to be ported to the kernel, and then execute within the kernel address space. Although this is not an insurmountable problem, it would introduce a code base into the kernel that is ten times larger than most existing file systems.

FUSE. FUSE or Filesystem in Userspace is a hybrid user-kernel approach [43]. Like a standard kernel-level file system, no application modifications are required. A standard kernel file system is used to interface with the VFS, but VFS calls are sent to a user-space demon via a device. The user-space demon executes the call and returns the data and status codes to the kernel-level file system, which in turn passes them on to the user. This means that the database code need not run within the kernel, eliminating one concern about developing an in-kernel file system. Unfortunately, this approach still suffers from the same caching problems, as a standard kernel level file system, in that cached accesses do not consult the DBMS. As FUSE file systems run outside of the kernel, and have less control over the VFS than a standard file system, these problems would be more difficult to solve than with a standard kernel-level file system.

User-level NFS server toolkits. A user-level NFS server toolkit, like the SFS-toolkit [23], has many of the same advantages and disadvantages as FUSE: applications need not be modified and the database can run in user level, but the kernel caches information inside of the NFS client, thereby violating the isolation property and creating coherence problems with the database caches. Additionally, user-level NFS servers require additional data copies through the network stack, as well as context switches.

LD_PRELOAD library. Another option is to run our file system directly in the address space of user processes and intercept system-call wrappers using the

LD_PRELOAD runtime-linker mechanism. This approach has three main advantages. First, as file-system calls are intercepted at the highest possible level, there are no cache coherency or isolation issues to contend with. Second, the database does not need to run in the kernel. Third, data copies between the process and the kernel are not required. There are, however, three disadvantages. First, statically linked binaries cannot use the file system, so they must be recompiled. Second, the C library itself continues to use the existing calls, so every call of interest must be intercepted (e.g., `fprintf` must be intercepted because applications use it to write to the file system). Third, system calls that do not use the library wrappers are not intercepted, so not all code would work with this approach.

Modified C library. Directly modifying the C library is another option to extend new file-system functionality to applications [20]. The advantages and disadvantages are similar to the LD_PRELOAD mechanism, but high-level calls like `fprintf` do not need to be modified if the corresponding low-level library calls like `write` are handled correctly. Three additional disadvantages of using a modified C library instead of an LD_PRELOAD are that all applications must be relinked with the new C library, modifying the C library requires significant implementation effort, and that circular dependencies would exist between the BDB library and the C library. For example, BDB needs the `fwrite` library call, but that call in turn would depend on BDB.

ptrace. The final option we considered was using the process-tracing facility, `ptrace` [14]. The process-tracing facility allows a *monitor* to intercept system calls and modify the calls and their arguments. From the perspective of the application, the monitor is equivalent to the OS, so no application modifications are required. As shown in Figure 1, the monitor runs in user-level, so BDB does not need to execute within the kernel. Unlike the library approach, a single instance of the monitor can handle multiple processes, so it is simpler to share data, caches, and other resources.

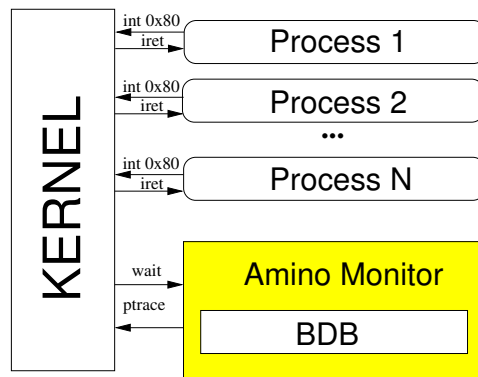


Fig. 1. The Amino monitor can trace an arbitrary number of processes. At system call entry, the kernel signals the monitor via the `wait` system call. Amino manipulates the monitored processes' state with `ptrace` primitives. BDB executes within the monitor's address space, and uses standard system calls.

The major disadvantage of the `ptrace` approach is that performance may suffer for system-call-intensive programs, as more context switches are required for each system call. However, we felt that ease of development and cache consistency outweighed performance concerns.

In Section 3.1 we describe the process tracing primitives. In Section 3.2 we describe the structure of the Amino monitor. In Section 3.3 we describe Amino's process control blocks, and in Section 3.4 we describe Amino's path resolution and mount framework. In Section 3.5 we discuss address space issues.

3.1 Process Tracing Primitives

The `ptrace` framework provides three primitives to establish tracing: the monitor can issue `PTRACE_ATTACH` to begin tracing a currently running process, the monitor can issue `PTRACE_DETACH` to stop tracing, and one of the monitor's children can issue `PTRACE_TRACEME` to be traced by the monitor. Our monitor begins by forking a new child, issuing `PTRACE_TRACEME`, and then executing the to-be-traced executable. From this point onward, the monitor is notified via the `wait` system call whenever the child needs attention.

The monitor uses three primitives to control the execution of the child process. (1) `PTRACE_SYSCALL` continues execution until the next entry or exit from a system call. If the child is in user-mode, then the child process is stopped before the kernel enters the system call handler, so that the monitor can change the arguments, or even the system call to be executed. If the child process is in the midst of executing a system call, then the kernel completes the routine and the monitor can examine and change any return values. (2) `PTRACE_CONT` continues execution until the child receives a signal. (3) `PTRACE_SINGLESTEP` continues execution until the next instruction.

When the child is in the stopped state, the monitor uses four primitives to observe and manipulate the child process: `PTRACE_GETREGS`, `PTRACE_SETREGS`, `PTRACE_PEEKDATA`, and `PTRACE_POKEDATA`.

`PTRACE_GETREGS` retrieves the values of the registers saved during a context switch from the kernel's process control block. On the Intel 80x86 architecture, the `eip` register contains the program counter, the `eax` register indicates what system call the process wants to execute, and the remaining general purpose registers contain the system call's arguments. Our current implementation is tied to the 80x86 architecture, because it references these registers, but it would not be difficult to add support for other architectures as the ABI is similar on all Linux platforms. In our prototype, only 353 out of 12,187 lines of code reference 80x86 specific registers.

The monitor can also manipulate the registers with the `PTRACE_SETREGS` primitive. Before a system call, the call to execute can be changed by setting `eax`, and the arguments can be changed by updating the corresponding registers. After a system call is executed, the return value can be set by updating the value of `eax`. At any point in time, the execution flow of the program can be changed by modifying `eip`. This is required when a single system call must be implemented in terms of several other system calls.

Finally, there are two primitives to examine and update a word in the child process's memory: `PTRACE_PEEKDATA` and `PTRACE_POKEDATA`. These primitives are used when the system call takes pointer arguments (e.g., file names are passed as

strings, and `stat` fills in a user-supplied buffer).

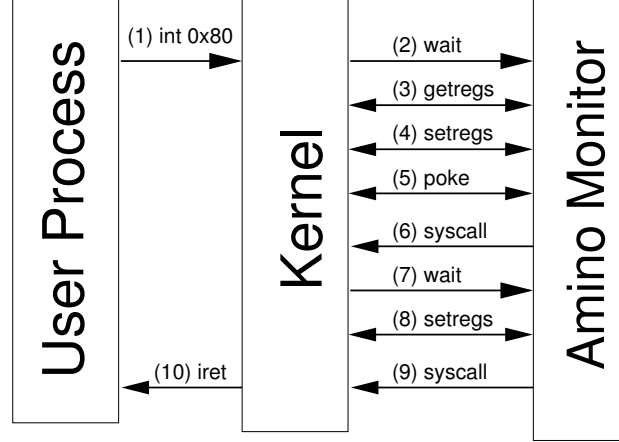


Fig. 2. `ptrace` primitives used to handle a `read` system call. Arrows indicate control transfer. Double arrows indicate that the function was called and returned immediately.

Figure 2 shows an example of how the Amino monitor handles a `read` system call destined for the database file system on behalf of a user process. There are ten steps involved in this call:

- (1) The user process issues a system call using `int 0x80`. The system call to execute is stored in `eax`.
- (2) The `wait` system call in the monitor returns the process ID of the user process.
- (3) The monitor issues a `PTRACE_GETREGS` call to retrieve the value of `eax`. Based on `eax` and the call’s arguments, Amino determines whether this call is destined for the database. If the call is not destined for the database, then Amino allows the process to continue with no further intervention.
- (4) If this call is destined for the database, then Amino changes the registers to prevent the kernel from handling the call. In the case of `read`, Amino sets `eax` to `-1`, thus the kernel essentially ignores the call because no handler is associated with `-1`.
- (5) Amino performs the database `read` operation, and uses the `PTRACE_POKEDATA` primitive to write the returned data into the user process’s address space (we also have an optimized mechanism described in Section 3.5).
- (6) Amino instructs the kernel to continue execution until the end of the call and calls `wait` (in this case the call returns immediately without performing any service, because `eax` was set to `-1` in step 5).
- (7) The kernel executes the system call, and returns from `wait`.
- (8) Amino uses the `PTRACE_SETREGS` primitive to store the return value of the previously executed `read` in `eax`.
- (9) Amino uses the `PTRACE_SYSCALL` primitive to allow the user process to continue executing.

- (10) The kernel issues an `iret` instruction to return control to the user process. The user process reads the return value from `eax`, and it is as if the system call were serviced by the kernel.

3.2 Amino Structure

The Amino monitor begins by forking a child process to trace. After the fork, the child executes the program to be monitored. All of the process's descendants are also monitored, and each monitored process is assigned a state. The two most common states are `INUSER` and `INCALL`, which indicate that the process is executing user-level code or it is executing a system call, respectively. To service requests, Amino calls the `wait` system call. When a process requires attention, usually because it is entering or exiting a system call, the kernel returns its process ID as the result of the `wait` system call (`wait` also returns when a signal is delivered or a process exits).

After returning from `wait`, Amino retrieves the current process's state and performs an appropriate action. There are currently 19 states (including `INUSER` and `INCALL`). Most of the states indicate that the user process is in the midst of a specific call, for example `clone`, `exec`, `chdir`, or `dup`. One of the most important states is `INFORCERET`, which indicates that the return value of the presently executing system call should be overridden by a given value. This state is used by most database calls to pass back status information. In the example in Section 3.1, the return value of the `read` is determined in step 5, but is not yet returned. When the return value is determined in step 5, the monitor sets the state to `INFORCERET`. After step 7, Amino looks up the state and because it is `INFORCERET`, Amino sets the value of `eax` to the proper return value. Two other states of note are `REDOCALL` which indicates that the current system call should be repeated, and `RESTOREREGS` which indicates that the process's registers should be set to their original values. `REDOCALL` allows us to insert a new system call into the stream (e.g., to create shared memory regions), and `RESTOREREGS` is used when we need to change system call arguments (e.g., when rewriting file names).

3.3 Process Control Blocks

The monitor maintains each process's state in a private *process control block* (PCB). The monitor's PCB is independent of the OS PCB, and contains the process ID to use as a search key, a copy of the process's registers, the current state of the process (e.g., `INFORCERET`), and all state-specific information (e.g., the return value to be passed back to the application). Encapsulating all of this information in a single structure allows the monitor to handle concurrent processes.

Like an OS PCB, the monitor's PCB contains an open-file table and present working directory (PWD). The open-file table is a simple array with a slot for each possible file descriptor. If a given file descriptor is connected to an Amino file, then its slot contains a pointer to a structure describing the file; otherwise it is empty (`NULL`). If a system call uses a file descriptor as an argument, it is looked up in the open-file table. If the file descriptor's slot is empty, then the system call proceeds with no further intervention. Otherwise, Amino extracts the schema data (i.e., the database and environment handles) and the unique file identifier from the open-file table and directs the call to BDB.

Amino cannot arbitrarily assign file descriptors to the user-level process, because the kernel would not know that a given file descriptor is in use. To handle this situation, Amino uses *shadow descriptors*. When opening a file in the database, Amino changes the path name to “/” before letting the system call proceed. The resulting file descriptor (in the child process) is used as a place holder, and no system calls are issued against it. The kernel does not assign the descriptor to any other file, so Amino can correctly identify the calls that it handles; in case of a software error, most calls on this file descriptor fail with `EISDIR` (because “/” is a directory). For efficiency, Amino reuses this file descriptor with `dup` on subsequent `open` calls.

3.4 Mount Subsystem

The Amino monitor must maintain a mount table to associate pathnames with database schemas. On startup, an Amino configuration file provides a list of paths to manage, and for each path, the mount type and data (the configuration file is essentially equivalent to `/etc/fstab`). Currently, Amino supports BDB mounts that take the database pathname as an argument. When Amino encounters a system call that references one of these paths, Amino passes it to the appropriate routine.

Pathnames passed to system calls can be rather complex. If they are relative path names, then they depend on the process’s context. Any path can use the “.” operator to move one level up the directory tree. We store paths as stacks, with the root path represented as an empty stack, and a path such as `/usr/local/bin` is represented by a stack containing `usr`, `local`, and `bin`. If a path is managed by Amino, then it is a child of one of the mount-table entries described in the configuration file. To rapidly determine if one path is a child of another, the path structure also contains a depth, and a length for each path component.

Each PCB contains a path stack for the PWD. When a `chdir` or `fchdir` system call is issued, the new PWD is stored as a candidate. If the system call is successful, then the candidate becomes the PWD. The mount table also uses a path stack to identify the path for each mount.

To resolve a path that is passed to a system call, first the process’s PWD is copied to a new stack. If the path begins with a “/,” then the stack is emptied. Each subsequent component is pushed onto the stack, unless it is “.” in which case an element is popped off (unless of course the stack is already empty). After converting the string pathname into a path stack, the monitor searches the mount table for any mount that contains this path. The path structure is optimized for this purpose: if the path has a lower depth than the mount, then it cannot be a child; and the length is stored with each component so the component names only need to be compared if they have equal length. If one is found, then the path components after the root of the mount are extracted (e.g., if the path is `/usr/local/src/amino` and the mount is rooted at `/usr/local`, then `src/amino` is extracted). The mount private data containing the database handles and the extracted path are then passed to the BDB call. If the path name is not contained in a mount, then Amino allows the system call to go through without any changes.

3.5 Address Spaces

There are two distinct address spaces involved in executing the Amino monitor: (1) the address space of the monitor and (2) the address space of the user process. The `ptrace` primitives to access the user process's address space are rather limited—they can only examine or change one word at a time. Thankfully, Linux provides a more powerful interface to it through the `/proc` file system. A process with permission to `ptrace` another process may read from the traced process's memory using the `/proc/pid/mem` file, where `pid` is the PID of the traced process. This allows the transfer of up to a page (1,024 words on the 80x86) in a single system call. Linux also has support to write to `/proc/pid/mem`, but it is disabled by default. For our prototype, we have enabled a writable `/proc/pid/mem` to allow bi-directional bulk transfers. If the `/proc/pid/mem` interface is not available for reading or writing, then Amino falls back to `PTRACE_PEEKDATA` and `PTRACE_POKEDATA`. An improved interface would be to allow regions of the child's address space to be memory-mapped into the monitor. This would provide a zero-copy transfer method.

All system call arguments must be in the user processes' address space. For example, the first argument to `open` is a pointer to a string. If Amino needs to update these values, then it must manipulate the child's address space. It is not always possible to manipulate the file name in place, because the new file name may be longer than the existing file name, and the memory segment may be read only. To address this issue, previous `ptrace` monitors have modified either the stack, or the first writable segment. In Amino, we establish a System-V shared-memory region between each user process and the monitor. When the first system call is issued with an argument that needs to be updated, the monitor creates a shared memory region. Next, the monitor inserts a shared-memory attach operation in to the child's system call stream. At this point, Amino writes the new file name into its own address space, and updates the child's registers to point to the shared memory in the child's address space. After the call, the child's original registers are restored. Subsequent arguments can be rewritten by simply updating the local region, and the child's registers. This approach has the advantage of requiring no data copies, and the child's existing memory is not modified, therefore the child's memory does not need to be restored after the call.

3.6 `ptrace` Enhancements

The standard `ptrace` interface requires at least six context switches for each system call: (1) the traced process traps into the kernel; (2) the kernel transfers control to the monitor; (3) the monitor transfers control to the kernel; (4) after executing the system call the kernel transfers control back to the monitor so that the return value can be manipulated; (5) the monitor transfers control back to the kernel; and finally, (6) the kernel transfers control back to the traced process. In reality, more context switches are required as the monitor must retrieve the values of traced process's registers, issue system calls to provide OS-like services, etc.

Clearly, reducing the number of times that the monitor is called improves performance. For most calls the monitor only needs to be notified on entry. If the call is not destined for an Amino file system, the monitor does not need to examine the return value so the call could execute without further intervention by the monitor.

If the call will be handled by the Amino file system, the return value could be set and the monitor need not be notified. Unfortunately, these two modes of operations are not possible under the current `ptrace` interface.

We created two new `ptrace` operations: `PTRACE_CHECKEMU` and `PTRACE_SYSSKIP`. The `PTRACE_CHECKEMU` operation is similar to the `PTRACE_SYSEMU` operation that was recently introduced to improve the performance of User Mode Linux [6]. The primitive `PTRACE_SYSEMU` allows all of a process's system calls to be emulated, but it is not suitable for the Amino monitor, because we only emulate a subset of the system calls. Our `PTRACE_CHECKEMU` interface allows the monitor to determine whether emulation is required after examining the registers. The UML developers agree that our more general `PTRACE_CHECKEMU` interface is an improvement over the existing `PTRACE_SYSEMU` [12]. The corollary to `PTRACE_CHECKEMU` is `PTRACE_SYSSKIP`. When the Amino monitor does not implement a call, then it issues `PTRACE_SYSSKIP` instead of `PTRACE_SYSCALL` to bypass notification of this system calls return value and go directly to the start of the next system call. Together these primitives reduced traps into the monitor by 30.8% during an OpenSSH compile.

Finally, there are also many non-file-system system calls that the monitor need not intercept at all (e.g., `time` or `getpid`). To reduce the number of extraneous calls into the monitor, we added an optional bitmap of system calls to the task structure. By using a new `ptrace` primitive, `PTRACE_SELECT`, the monitor selects precisely the set of calls that need to be traced. This method reduced the number of traps to the monitor by an additional 12.8% during an OpenSSH compilation. Overall, these techniques reduced the number of traps to the monitor by 43.7%.

These three improvements can benefit a wide variety of `ptrace` monitors. For example, the `PTRACE_CHECKEMU` grew out of work for User Mode Linux, but provides a more flexible interface that can be used by a monitor that emulates a subset of system calls. Many security-oriented monitors only need to examine which system calls are being executed and their arguments, but not their return value. For these types of monitors, `PTRACE_SYSSKIP` would greatly improve their performance. The `strace` program provides support for filtering the set of system calls to display (e.g., file system, process, or IPC related calls), but this filtering is done in user-space. By using `PTRACE_SELECT`, `strace` could have the kernel perform this filtering.

4. EVALUATION

We evaluated the performance of our system by running several general-purpose workloads and micro-benchmarks. We chose three general-purpose benchmarks to evaluate our system. In Section 4.1, we present results for the Postmark benchmark [18]. In Section 4.2 we present results for an OpenSSH compile, and we present results for a Sendmail benchmark in Section 4.3. We also ran two sets of micro-benchmarks. In Section 4.4 we present results for meta-data-intensive micro-benchmarks, and in Section 4.5 we present results for data-intensive micro-benchmarks. The testbed ran Fedora Core 4 with all updates as of July 19, 2005. All experiments were located on a dedicated 40GB Maxtor IDE disk. We compared Ext3 to Amino using BDB databases stored on Ext2. We used Ext2 as the underlying file system for Amino, because BDB provides ACID semantics even without a journaling file system. We chose to use Ext3 as a basis for comparison, because

it provides a limited subset of the ACID properties, whereas Ext2 does not. To ensure a cold cache, we remounted the file systems between each iteration of a benchmark. For all tests, we computed the 95% confidence intervals for the mean elapsed, system, and user time using the Student- t distribution. In each case, the half-widths of the intervals for the elapsed and system times were less than 5% of the mean.

We used the following seven configurations for our tests:

- VANILLA The benchmark is run on Ext3.
- VANSYNC The benchmark is run on Ext3, but the file system is mounted with the `sync` mount option to provide durability.
- STRACE The benchmark is run on Ext3, but is monitored by `strace -cf`. This configuration shows the overhead of the `ptrace` facilities, but does not modify any system calls or produce any output during execution.
- AMINONULL The benchmark is run on Ext3, but is monitored by the Amino monitor. This configuration shows the overhead of `ptrace` and our path-name resolution infrastructure.
- AMINOACI The benchmark is run through the Amino monitor with a BDB database stored on an Ext2 file system. BDB is configured to provide atomicity, consistency, and isolation, but not durability.
- AMINOACID This configuration is the same as AMINOACI, but durability is also provided because BDB flushes the log to disk on each commit.
- AMINOTXN This configuration is the same as AMINOACID, but the benchmark is modified to insert calls to begin and commit Amino transactions. This improves performance, because data needs to be flushed to disk only after the transaction is committed, rather than after every system call.

4.1 Postmark

Postmark 1.5 is an I/O-intensive benchmark that stresses the file system by performing a series of file system operations such as directory look ups, creations, and deletions on small files [18]. The first Postmark configuration we chose is to create 5,000 files ranging from 512 bytes to 10KB, and perform 20,000 *transactions* (this is Postmark's term for an operation, and is distinct from Amino transactions). We used the `read` and `write` system calls (as opposed to Unix buffered I/O), and a transfer size of 4,096 bytes for Ext3 and 4,070 bytes for Amino as these are the optimal transfer sizes reported by `fstat`. Because Amino has a smaller block size, it never requires fewer system calls than Ext3, but occasionally requires more.

The Postmark results are shown in Figure 3. The VANILLA configuration took 30.6 seconds to execute. The VANSYNC configuration synchronously writes data and meta-data to disk to provide durability. The VANSYNC configuration was slower than VANILLA by a factor of 10.6, due to additional synchronous disk writes. The STRACE and AMINONULL have overheads of 62.3% and 68.9% over VANILLA, respectively. This shows the overhead of the process-tracing facilities. Amino uses slightly less system time (2.6%), because it accesses all process registers using a single system call instead of one system call for each register. However, Amino uses more user time (0.51 seconds or 40.93%) because it resolves each path name. AMINOACI

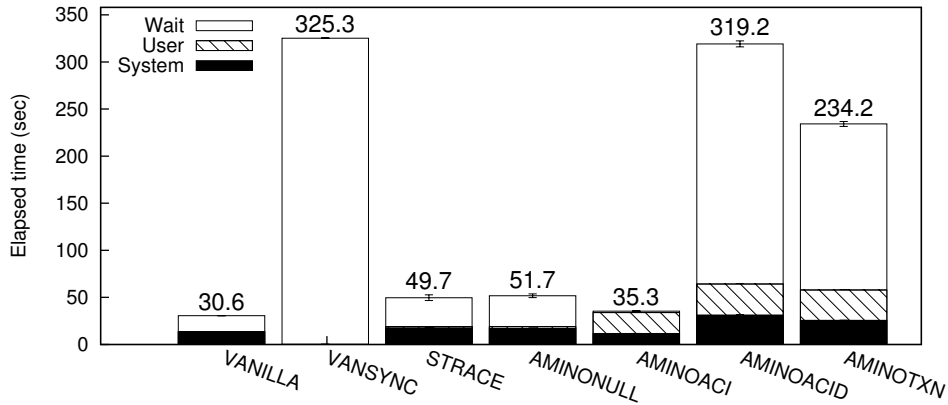


Fig. 3. Postmark: 5,000 files and 20,000 transactions.

provides atomicity, consistency, and isolation using BDB. It is 15.4% slower than VANILLA, and 31.7% faster than AMINONULL. This shows that a file system built on a database can provide atomicity, consistency and isolation with good performance, even for I/O-intensive applications, because we can quickly access files and directories with our schema and BDB efficiently writes data to the log.

AMINOACID provides all four ACID properties: atomicity, consistency, isolation, and durability. To provide durability, the database log must be synchronously written to disk after each transaction. This leads to an expected overhead of a factor of 10.4 over VANILLA, but AMINOACID provides semantics closer to VANSYNC. When compared to VANSYNC, AMINOACID improves performance by a slight 1.9%, but AMINOACID also provides stronger guarantees—the contents of the file system are always defined. In VANSYNC, the contents of the file system are undefined because there are no explicit commits.

In AMINOACID, each individual system call flushes the log to disk. In AMINOTXN, we modified Postmark to begin and end Amino transactions before each high-level operation (i.e., create, remove, read, or write a file) that Postmark refers to as a transaction. This required 33 lines of code: one to include a header file, four to begin transactions, four to commit transactions, and the remaining 24 were simple checks of the begin and commit return values. In this configuration, Amino requires fewer synchronous writes because individual system calls are grouped into a larger transaction, and hence performance improved by 26.6% over AMINOACID. In sum, we show that Amino provides performance comparable to Ext3. Moreover, with only small modifications, applications can improve durable performance and benefit from full ACID semantics.

Alternate Configurations. We also ran two alternate Postmark configurations that are slight modifications of the first configuration. The first alternative configuration has ten times larger files: from 5,120–102,400 bytes. The second configuration has more files: we increased the number of initial files to 25,000 and

the number of transactions to 100,000 (a factor of five for each). For this configuration we introduced 100 subdirectories, so that Ext3 would not be required to perform linear scans over 25,000 files for some operations. For each of these two configurations we ran VANILLA, VANSYNC, AMINOACI, AMINOACID, and AMINOTXN.

For the first alternative configuration (large files), the VANILLA and AMINOACI results were quite similar to the original configuration. The VANILLA configuration took 113.7 seconds and AMINOACI took 130.6 seconds, for an overhead of 14.9%. The VANSYNC configuration only had an overhead of 3.7 times over VANILLA, because more I/O operations were done by VANILLA during this longer running configuration. The AMINOACID configuration had a 93.2% overhead over VANSYNC. The reason that AMINOACID performed poorer with large files is that it had to write 2.8 times more sectors to disk than VANSYNC did (as reported by `/proc/diskstats`). We plan to investigate alternative database layouts to reduce the number of sectors written. The AMINOTXN configuration's elapsed time was indistinguishable from that of VANSYNC. However, AMINOTXN provides stronger guarantees than VANSYNC does.

For the second alternative configuration (more files), VANILLA took 600.8 seconds. For this configuration, AMINOACI outperformed VANILLA by 70.3%. The reason is that VANILLA spreads the files through many cylinder groups, but AMINOACI stores them together in a balanced tree, improving locality thereby reducing wait time. However, AMINOACI did use 6.1 times more CPU time. As expected, the synchronous configurations were slower. VANSYNC had a 5.7 times slow down, and AMINOACID had a 3.0 times slowdown. Again, AMINOTXN was the most efficient with only a 2.5 times slowdown compared to VANILLA.

4.2 OpenSSH Compile

To simulate a more CPU-intensive typical user workload, we adapted the SSH build workload [39], but used OpenSSH 4.2p1 as it builds cleanly on our systems whereas SSH 1.2.26 does not. This workload stresses the Amino monitor, as it requires significant amounts of additional CPU time in order to intercept system calls. The compile benchmark is divided into three phases: (1) unpack, (2) configuration, and (3) build. We measured each of the phases separately to isolate their different characteristics. In the unpack phase, the package is uncompressed and new files are created by `tar`. In the system-call-intensive configuration phase, the `configure` shell script performs many small configuration tests, which involve a fair mix of file-system operations. The build phase is more CPU-intensive and builds 157 object files, two libraries, eleven executables, and sixteen man pages.

Figure 4 shows the results of each phase of the OpenSSH compile benchmark. The unpack phase (shown in Figure 4(a)) took 0.20 seconds on VANILLA. The STRACE configuration added an overhead of 53.9%, and AMINONULL had a similar overhead of 48.0%. For all three of these Ext3 configurations, the benchmark completed quickly, because no disk writes were performed during program execution due to the buffer cache. The AMINOACI configuration took 0.86 seconds to complete, which is a factor of 4.1 slower than Ext3. The reason that AMINOACI is slower than Ext3 is that the CPU time used increased by 0.48 seconds from 0.17 seconds to 0.65 seconds. Of this 0.48 second increase, one quarter of it (0.12 seconds) can be attributed to the monitoring infrastructure. The remaining increase of 0.36 seconds is due to

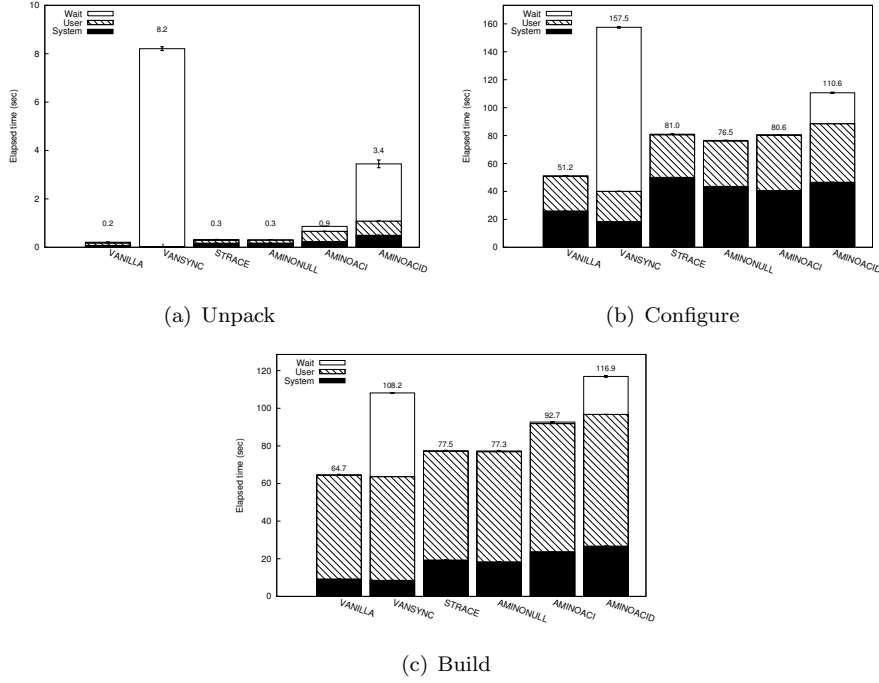


Fig. 4. OpenSSH Compile Results.

performing BDB operations and data copying between the monitor and `tar` process. The last two configurations we tested were VANSYNC and AMINOACID. The VANSYNC configuration is 39.2 times slower than the VANILLA configuration, because changes are written to the disk synchronously. The AMINOACID configuration provides the same functionality, it is only 16.5 times slower than VANILLA, because BDB is optimized for durable performance.

The second phase of benchmark, configuration, is shown in Figure 4(b). On VANILLA, this phase took 51.2 seconds. This phase of the benchmark is CPU and system call intensive, so the STRACE and AMINONULL configurations had overheads over VANILLA of 58.2% and 49.4%, respectively. The AMINOACI configuration has an overhead over VANILLA of 57.5%. When compared with AMINONULL, the overhead of AMINOACI is only 13%. This demonstrates that our file system is relatively efficient, though the CPU intensive nature of this workload causes the context switches and data-copying induced by the monitoring infrastructure to degrade performance. When durability is added, VANSYNC is 3.0 times slower than VANILLA, and AMINOACID is 2.1 times slower than VANILLA. Again, this demonstrates that Amino efficiently provides durable performance.

The build phase (shown in Figure 4(c)) took 64.7 seconds on VANILLA. Even though this phase is the most CPU intensive phase of all, this is the least system call intensive. Therefore, the monitoring infrastructure has a lower overhead than in the configuration phase: 19.7% for STRACE and 19.5% for AMINONULL. The AMINOACI configuration had an overhead of 43.3%. Most of this was due to a 42.8% increase in

CPU time from 64.3 seconds to 91.8 seconds, caused by BDB operations, additional data copying, and context switches. The VANSYNC configuration was 67.2% slower than VANILLA, and AMINOACID was 80.7% slower than VANILLA. The reason that AMINOACID was slower than VANSYNC is that the build phase is remained more CPU intensive than the other phases when durability was added. For VANSYNC, CPU utilization was 59% and for AMINOACID it was 82%. As CPU was a bottleneck in this configuration, the extra context switches and data-copying hurt AMINOACID more than the improved durable performance helped it.

4.3 Sendmail

Using an identically configured machine as a client, we ran a Sendmail 8.13.4 server and varied the backing store for the `/var/mail` directory, where user mailboxes are stored. We did not run Sendmail through the Amino monitor, because it does not access the mail files. Instead, it delegates that task to the local mailer. We used the default local mailer for the VANILLA configuration. To provide isolation and an approximation of atomicity, the local mailer performs locking and complex checks (e.g., repeatedly calling `stat` to ensure that the file does not change). To ensure that mail is not lost (i.e., provide durability), the local mailer calls `fsync` after writing the message. These checks are unnecessary under Amino, as our file system transparently provides isolation to applications, without the need for explicit locking calls or repeated checks. Instead of using the default local mailer, we wrote a simple replacement that uses an Amino transaction to provide ACID properties for the AMINOTXN configuration. The Amino mailer is only 78 lines long, which is less than one fifth of the local mailer's 450 line delivery function, because the transaction abort primitive removes the need for specialized error handling code.

For our benchmark, we developed a Perl script that stress tests the mail server by continuously sending mail. We created a pool of 100 users to receive the mail, and each message had a randomly selected recipient. The messages sizes were normally distributed with a mean of 5,993 bytes and a standard deviation of 4,166. We chose the size parameters based a 2.5%-trimmed mean of our non-spam email for the past year. The test begins with a 60 second warmup period, in which the test runs without measurement to avoid startup effects. After the warmup, messages are sent for five minutes, and we record the mean achieved rate.

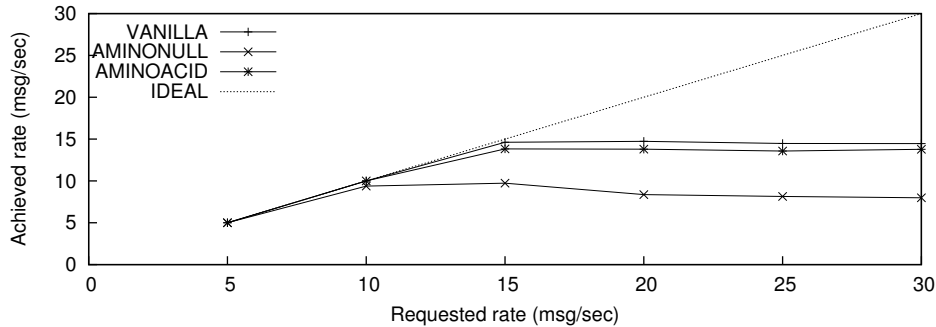


Fig. 5. Local mailer: requested vs. achieved load.

We ran the test for requested rates of 5–30 messages per second (MPS), and plotted the requested rate against the achieved rate in Figure 5. The `VANILLA` configuration handled 5 and 10 MPS well, but achieved only 14.6 MPS when 15 MPS were requested. For requested loads of 20 MPS, `VANILLA` only achieved 14.7–14.5 MPS. `AMINONULL` had significantly degraded performance. For a request rate of 10 MPS, only 9.4 MPS was achieved, and for request rates of 15–30 MPS, the achieved rate declined from 9.7 MPS to 7.9 MPS. `AMINOTXN` fared better than `AMINONULL`, and was able to handle 5 and 10 MPS, but for 15–30 requested MPS, it only achieved 13.6–13.8 MPS. This is 6.5% below `VANILLA`, but is 42.0% better than `AMINONULL`. This shows that our file system code is quite competitive with `VANILLA`, even though its performance is reduced by `ptrace` operations. The performance of Amino decreases compared to the Postmark results, because this benchmark is more data intensive than Postmark. Additionally, as Sendmail uses multiple processes to deliver mail, there is increased lock contention to provide the isolation. We plan to investigate alternative schema designs that may yield a higher degree of concurrency. Even though `AMINOTXN` is 6.5% slower, the code is significantly smaller and simpler, which means that fewer bugs and security flaws are possible, and the system is more reliable.

4.4 Meta-data Micro-benchmarks

We ran several micro-benchmarks on Amino to evaluate the overheads of certain primitive file system functions like file creation and deletion as well as reading and writing data. We broadly classify our micro-benchmarks into metadata and data benchmarks. We describe the meta-data micro-benchmarks in this section, and the data benchmarks in Section 4.5. The meta-data operations we evaluated are `create` (and `mkdir`), `unlink` (and `rmdir`), `stat`, and `readdir`. We chose these meta-data operations because they are a broad cross-section of file system operations, and together with data operations account for the vast majority of operations [8].

To generate metadata operations, we developed a C program that operates on several directories each containing a fixed number of files. We used this method rather than a generic data set (e.g., the source of a package), because when evaluating one specific method we did not want to use directory reading operations or lookups to determine which files must be operated upon. For all the metadata workloads, we disabled `atime` updates on both in Ext3 and in Amino so to isolate the overheads of the metadata operation to be tested.

Create. To evaluate the overhead of the `create` and `mkdir` operations, we used our C program to create 500 directories with 1,000 files each. As seen in Figure 6(a), the `AMINONULL` configuration had an elapsed time overhead of 61% compared to `VANILLA`. This is primarily because of the context switches caused by the `ptrace` monitor, as evidenced by the system time overhead of 54%.

The `AMINOACI` configuration ran 3.2 times slower than `VANILLA` in terms of elapsed time. There was significant increase in the user time (0.7 seconds vs. 159 seconds) and system time (59 seconds vs. 108 seconds), caused mainly due to data copies, comparisons traversing B-trees, and locking overheads. The `AMINOACID` configuration ran 20% faster than the synchronous mode of Ext3. This is because of a 42% decrease in wait time of `AMINOACID`. Synchronous mode Ext3 incurs more

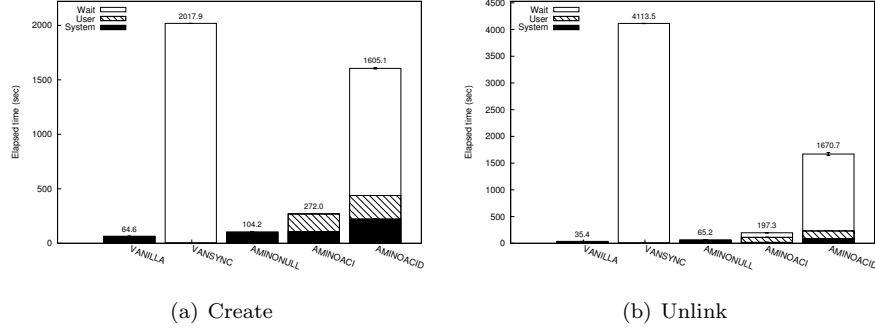


Fig. 6. Creation and deletion micro-benchmark results.

seeks as it has to synchronously commit several structures (directory files, inodes, and inode bitmaps), whereas AMINOACID often just needs to update two leaf nodes of the tree.

Unlink. To evaluate the performance of `unlink` and `rmdir`, we removed the files and directories created by the create workload described above. We unmounted and remounted the file system to ensure cold cache between the create and unlink workloads. Figure 6(b) shows that the AMINONULL configuration had an overhead of 84% compared to VANILLA, mostly because of the context switches of `ptrace`. AMINOACI ran 4.5 times slower than VANILLA. The break up of the overhead is similar to that of the `create` workload. The AMINOACID configuration ran 59% faster than Ext3 in its synchronous mode. This is because of a 64% decrease in the wait time, caused by the fewer number and smaller seeks in AMINOACID.

Stat. Directory lookups are one of the most common operations, because they are a precursor for almost every meta-data operation (e.g., opening a file, creating an entry, deleting an entry, etc.). To evaluate the lookup operation, we ran `stat` on 5,000 directories with 100 files each. After unmounting and remounting the file system, we performed a `stat` system call on each of the files. Figure 7(a) shows the results for this workload. The AMINONULL configuration had an elapsed time overhead of 75% compared to vanilla Ext3. This is because of two reasons: first the monitor context switches contribute to the increased system time (17 seconds vs. 45 seconds). This is because of the overheads of context switches from user to kernel caused by the `ptrace` monitor. Second, the increase in user time (0.5 seconds vs. 7.9 seconds) is caused by resolving each path to determine if it is destined for a BDB mount in the monitor.

The AMINOACI configuration ran 2.9 times slower than vanilla Ext3 in terms of elapsed time. This is because of a 43% increase in system time, a 51.5 second increase in user time, and 2.3 times increase in the wait time. The increase in system time is caused by the additional context switches and data copies due to `ptrace`. The increase in user time is because of locking overhead, path resolution, and B-tree comparison overheads. These operations are counted against user time, because our monitor executes in user-space, whereas a kernel file system would perform many of the same operations and charge them to the process’s system time. The AMINOACID

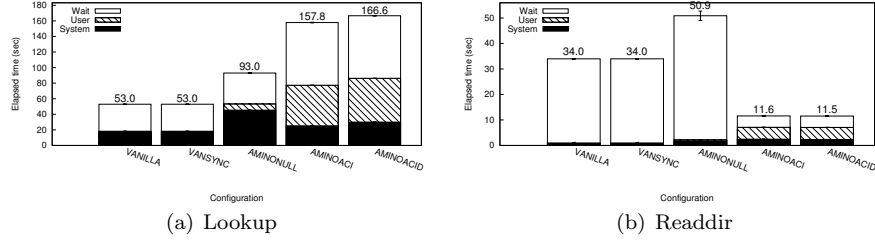


Fig. 7. Directory operation micro-benchmark results.

configuration also had an overhead of 2.9 times compared to VANSYNC, for much the same reasons that AMINOACI is slower than VANILLA. Because there were no writes in this workload, the durable performance improvements usually realized by the more efficient database log are not seen. AMINOACID has an overhead of 5.5% as enabling durable transactions executes additional database code.

Readdir. We used the same working set as in the `stat` micro-benchmark to evaluate the performance of the `readdir` operation. We performed a `readdir` on each of the 5,000 directories in sequence. The results are shown in Figure 7(b). The AMINONULL configuration had an elapsed time overhead of 34% compared to VANILLA. The overheads associated with the AMINONULL configuration are due to the same factors described in the lookup micro-benchmark. However, as the total number of operations performed in the `readdir` workload is smaller than the lookup workload, AMINONULL incurs fewer context switches.

The AMINOACI configuration ran 65% faster than vanilla Ext3 for this workload. The improvement is mainly due to an 86% decrease in wait time. Wait time is reduced because Ext3 requires seeks to read each directory, as it does not place directories close to each other on the disk. AMINOACI stores the path names in a B-tree and hence has better spatial locality. Therefore it requires fewer and shorter seeks to read directories. The use of B-trees to store metadata and data makes Amino suitable for metadata-intensive workloads which benefit from this locality. The difference in overheads between VANSYNC and AMINOACID is similar to AMINOACI as this is a read-only workloads, so we omit further discussion for brevity.

4.5 Data Micro-benchmarks

To evaluate the performance of data operations we ran a random read, random write, sequential read, and sequential write micro-benchmark. For each benchmark we performed 20,000 operations on a single 1GB file. For sequential operations, we operated (read or write) on consecutive pages of the file in sequence. For random operations, we generated a pre-populated pattern by randomly shuffling a sequential list of page numbers, and operated on the file using the shuffled list. This method ensures that there are no repeated pages so that caching does not affect the results. We used each file system’s native page size: 4,096-byte pages for Ext3 and 4,070-byte pages for Amino.

Random Read. The results of the random read benchmark are shown in Figure 8(a). The AMINONULL configuration had a marginal 0.7% elapsed time overhead compared to vanilla Ext3, mostly due to the context switches for `ptrace`. The AMINOACI ran 10% slower than VANILLA in terms of elapsed time. This slow down is due to a 6.9% increase in wait time and several times increase in the user time (0.06 for VANILLA vs. 3.2 seconds for AMINOACI) and system time (0.8 vs. 3.0 seconds). The increase in wait time is because of the additional read I/O required for traversal of the B-tree while reading nodes in random order. The increases in user and system time are caused by context switches, locking overhead, and additional data copies. The overheads of AMINOACID compared to VANSYNC are similar, as this is a read-only workload.

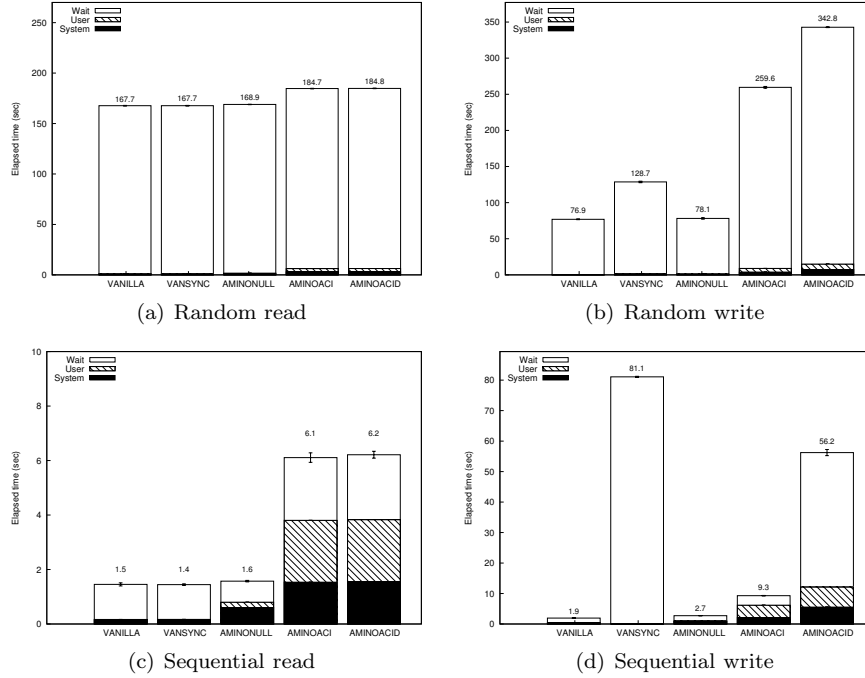


Fig. 8. Data micro-benchmark results: 20,000 operations.

Random Write. Figure 8(b) shows the overheads associated with Amino for the random write workload. The AMINONULL configuration had a small overhead of 1.5% elapsed time compared to vanilla Ext3. This is caused by the context switches done by `ptrace`. Under the AMINOACI configuration, Amino ran 2.3 times slower than vanilla Ext3 in terms of elapsed time. This is mainly because of the 2.2 times increase in the wait time caused by the additional reads performed by AMINOACI while traversing B-trees. The disk statistics revealed that the number of sectors read by AMINOACI for this benchmark was 177,650 whereas Ext3 performed just 2,064 sector reads. The AMINOACID configuration ran 1.6 times slower (elapsed

time) than Ext3 in its synchronous write mode. The largest component of the overhead was a 1.5 times increase in wait time. This is again due to the larger number of sectors read for traversing B-trees. Even though these results show that Amino has significant overheads for a data-intensive random write workload, as this pattern is rather rare [7].

Sequential Read. The overheads of Amino under the sequential read workload are shown in Figure 8(c). AMINONULL had a marginal overhead of 8.3% over VANILLA, mostly because of the context switches caused by **ptrace**. The AMINOACI configuration ran 3.2 times slower than VANILLA. The increase in user and system time is because of data copies and B-tree comparison overheads. The increase in wait time is because sequential reads require more seeks in Amino than in Ext3 as the data layout in Amino is a B-tree. As part of our future research, we plan to investigate more efficient data layouts, possibly including a hybrid model that combines a flat file and a database structure. The overheads of AMINOACID compared to VANSYNC are similar, as this is a read-only workload.

Sequential Write. Figure 8(d) shows the time taken for Amino for the sequential write workload. The AMINONULL configuration had an overhead of 40% over VANILLA due to **ptrace**. This is a similar absolute magnitude as read, but since writes can be asynchronous and thus take more CPU time, the increase in user time is more pronounced as a percentage. For AMINOACI the overhead was 4.7 times, primarily due to increased system and user CPU time. Amino in its AMINOACID configuration ran 30% faster than Ext3 in its synchronous mode of operation. The difference is primarily due to the 45% decrease in wait time because Ext3 requires more seeks to commit inodes and inode bitmaps synchronously, whereas Amino just needs to commit changes to a sequential log.

5. RELATED WORK

In this section we discuss four classes of related work. In Section 5.1 we describe systems that integrate databases with the file system. In Section 5.2, we discuss log-structured and journaling file systems. We discuss transactional memory systems in Section 5.3, and in Section 5.4 we describe various system-call interception methods.

5.1 Databases and File Systems

Previous simulations of transactions embedded in the file system showed that file system transactions can perform as well as a DBMS in disk-bound configurations [37]. The same simulations showed that for CPU-bound configurations, file system transactions usually have an overhead caused by system call costs of less than 20%.

The Inversion File System [30] is a user-level wrapper library with file-system-like functions that stores files in a POSTGRES database. Inversion uses POSTGRES to support transactions and fast crash recovery. Unfortunately, Inversion operates in its own namespace, separate from that of other file systems, and uses different functions from the traditional Unix API, making it unsuitable for integrating legacy and transactional applications.

OdeFS [9], Oracle's iFS [31], and DBFS [28] are user-level NFS servers that use databases as a backing store. This approach allows unmodified applications to

use the file system, but ACID properties cannot be extended to the application because the NFS client cache can serve requests without consulting the database system. Also performance suffers due to data copies and context switches related to the network stack. OdeFS is a file-system interface to objects already in the Ode object-oriented database. For each type of Ode object, new methods must be defined for read, write, and other file-system operations. iFS supports many access methods: NFS, FTP, WebDAV, Samba, SMTP, IMAP4, and POP3; it offers several useful features such as versioning, change notification, and indexing; iFS provides convenient access to files using a variety of network protocols, but most of these protocols are not compatible with many applications; and the clients may cache data. DBFS is a block-structured file system developed using BDB [28]. DBFS exceeds FFS's performance for meta-data operations. However, for data operations it is 50–80% slower. In our evaluation, we have observed a similar pattern: meta-data operations are often faster using the database, but raw I/O throughput is reduced.

WinFS is part of an upcoming version of Microsoft Windows [26] (originally WinFS was slated for Longhorn, but has been delayed to “some future date”). WinFS will integrate a full-fledged SQL DBMS into the OS. Using a heavyweight DBMS with SQL enables powerful queries, but could add significant code complexity. Additionally, overheads may be significant depending on schema design and query processing. WinFS uses the database as well as an NTFS file system as a backing store for all files. WinFS changes the basic unit of data storage from a file to an item (an object with attributes). The WinFS API supports explicit transactions for items, but since the API is so radically different, applications must change to take advantage of its new features.

QuickSilver is a distributed operating system developed by IBM research that makes use of transactional IPC [36]. QuickSilver was designed from the ground up using a microkernel architecture and IPC. Every IPC request has a transaction ID, and servers must be able to rollback requests on abort and write them to non-volatile storage on commit (assuming the server has non-volatile state). All resource management and notification in QuickSilver is handled by transactions. For example, on process termination (commit or abort) the window manager destroys all windows; the virtual terminal server closes the standard input and output file descriptors; and the task manager kills all of its children. The use of transactions removes the need to handle local vs. remote processes differently. Amino integrates transactions into the file system using simpler and more widely-used OS primitives than QuickSilver. Unlike QuickSilver, in which each OS component must provide specific transaction support for rollback and commit, Amino leverages BDB so that each OS component or application can use simple begin, commit, and abort calls, without managing its own rollback or commit.

5.2 Log-structured and Journaling File Systems

Log-structured and journaling file systems borrowed the technique of write-ahead logging from databases [33; 15; 19]. The key difference between a log-structured file system and a journaling file system is that in a log-structured file system the log is the permanent home of the data, whereas in a journaling file system the log is a temporary location until the data is checkpointed to a permanent location

on disk. In this respect, BDB, and hence Amino, is more similar to a journaling file system than a log-structured file system. When updates are made, they are first written to the database log file and then written to their permanent locations within the database file. In log-structured file systems, journaling file systems, and Amino, synchronous writes have improved performance because they are written sequentially to the log, obviating the need to seek to many locations of the disk for a single update.

Log-structured and journaling file systems write “transactions” to their log, but these transactions are completely controlled by the file system software—user applications cannot surround multiple file system operations in a single atomic transaction. Additionally, the transactions in a log-structured or journaling file system do not provide all of the elements of ACID. Instead, they provide atomicity and consistency for well-defined operations within the file system, and durability can be provided by flushing the log to disk. Notably, log-structured and journaling file systems do not include provisions for isolation apart from other facilities provided by the OS (e.g., directory-level semaphores). Amino provides atomicity, consistency, isolation, and durability for arbitrary sequences of file system operations.

5.3 Memory Transactions

Lightweight Recoverable Virtual Memory (LRVM) was developed to simplify Coda servers [35]. LRVM is designed to handle transactionally protected memory-mapping of a file into a process’s address space. To simplify LRVM’s design, the file should be a small portion of the total storage: the undo log was stored in memory. Durability was provided by writing a redo log to disk. This type of functionality is closer to Amino’s support for memory-mapped files than our in-memory transactions, as our memory-mapping essentially provides recoverable virtual memory. Our in-memory transactions, however, are designed to provide more efficient volatile transactions (i.e., without the need for any redo logging). LRVM was developed as a user-library and requires explicit calls to indicate that a given region of memory will be written to. We believe that our page fault handling mechanism for identifying writes is more convenient and robust. Indeed, the LRVM authors point out that the most common types of bugs were missing calls before manipulating a region, and suggest that language support for LRVM calls would be a good solution to these missing calls.

The Rio, or RAM I/O, project sought to bring persistence to standard memory [4]. If memory is persistent, then file systems can avoid writing data indefinitely, thereby improving performance by an order of magnitude. The key observation is that most data in memory is lost because of either power failures or software errors. Power failures can be solved through the use of UPSs. To cope with software errors, two approaches are taken. First, Rio memory uses page protection and checksums to prevent an errant instruction from writing to it. To update a page, it must be made writable, then the update is performed, and finally the page is made read-only again. Along with the update, checksums are stored along with the data so that errors can be detected. These two mechanisms raise the bar for updating memory, so that an errant instruction is unlikely to corrupt Rio memory, and even if Rio memory is corrupted, the change can be detected with a checksum. The second approach that Rio uses is saving memory across warm reboots. After a system

crash, the machine is rebooted, but the memory contents are preserved. Before the OS is fully booted, the memory is written to a swap partition. After the OS is booted, the contents of Rio memory are restored from the swap partition.

The authors implemented a file cache with Rio, and showed that it can be as reliable as a traditional disk-based file system under a variety of software faults. However, the two major problems with the Rio architecture are that not all architectures support warm reboot (e.g., an x86 cannot be rebooted without destroying RAM contents), and Rio also assumes that hardware and power failures are so rare as to be ignored. Unfortunately, hardware is becoming an increasingly large source of faults, as hardware components increase in number and complexity, and cost pressures force the use of less reliable components [27; 10].

Vista is a transactional RVM built on top of Rio [22]. Vista greatly improves the performance of an RVM system, because memory is assumed to survive a system crash—avoiding synchronous writes. Because Vista is built on top of Rio, it does not require a redo log, and the code complexity is much simpler than that of previous RVM systems, at around 700 lines. In our system, we provide a more rich transactional interface that includes nested transactions. Rather than implementing redo logging and its associated complexities, we leverage existing BDB code. Our completed transactional memory library is only 625 lines of code.

5.4 System Call Interception

The Ufo Global File system uses a similar interposition technique as our monitor [1]. Ufo provides transparent access to remote files via FTP or HTTP. Ufo’s monitor uses the Solaris `/proc` file system. The monitor operates on system calls such as `open`, `close`, and `stat`. When an access to a remote file is detected, the file is transparently fetched, and the system call is changed to open the local copy. Ufo does not implement other calls such as `read`, `write`, `getdents`, or `stat` internally, because the file’s local copy can be used without modifying the application. To implement `getdents` and `stat` properly, however, sparse files are used to create *stubs* for files that are not yet locally cached. Creating this hierarchy of stub files hurts performance. The `ptrace` interface was used by the Janus framework to sandbox untrusted applications [13]. Janus monitors file-system and network-related system call invocations, and applies configurable policies to allow or deny system call execution.

Several other interposition techniques operate at the same logical system-call level as Amino, but use a customized interface. SLIC [11] is an OS extensibility system that allows kernel-level extensions or user-level servers to register handlers for system calls, signals, and other OS entry points. SLIC has been used to patch security holes, encrypt files, and provide a restricted execution environment. SLIC extensions that run as user-level servers are quite similar to the `ptrace` interface. Interposition agents provide higher-level abstractions for system call interception [17]. The key insight for interposition agents is that system calls can be divided into classes that operate on independent sets of objects (e.g., path names or file descriptors).

6. CONCLUSIONS

Applications use an easy-to-use and standard POSIX API to access file systems, but file systems do not provide transactional semantics. Databases provide transactions, yet have differing and difficult-to-use APIs. Many applications can benefit from transactions, which can greatly improve error handling and provide atomic operations. For example, atomicity obviates the need for complex error handling, because a transaction can simply be aborted without any ill-effects. Atomicity can be used as a tool to ensure consistency, so that specialized and complex recovery code is not required. Isolation allows applications to provide race-free updates. Finally, durability ensures that data that was written actually reaches the persistent storage. Because transactions are so useful and the file system interface is convenient and ubiquitous, we therefore believe that file systems should provide transactional semantics to applications. Furthermore, we contend that the combination of file system transactions and recoverable memory will enable developers to use more robust and elegant error recovery methods than simply “giving up” and terminating an application upon a failure.

We have designed and developed *Amino*, a prototype file system with ACID semantics. Amino uses the Berkeley Database (BDB) as a backing store, with an efficient file system schema. Using BDB’s flexible key-value pair access model, meta-data properly migrates between the Path and Data databases—improving performance for common operations while avoiding the pitfalls of logical redundancy. Amino exports an easy-to-use, yet powerful, nested-transactions API to user space. Applications can begin, commit, and abort transactions. We designed a simple API to enable cooperating processes to share transactions. Using the same API, a single application can support multiple concurrent transactions. To provide powerful transactions to application data structures, we developed an RVM library with support for nested transactions and transparent logging.

We have evaluated our prototype, and have shown that it has acceptable performance. Amino configured for atomicity, consistency, and isolation is only 15.4% slower than Ext3 even though it runs in user space and has additional overheads due to `ptrace`. When durability is required, performance inevitably suffers because of synchronous disk writes. Whereas providing durability for unmodified applications on Ext3 degrades performance by a factor of 10, Amino provides modified applications durable performance with a slowdown of only 7.6 times. Moreover, Amino provides applications with the additional benefits of atomicity, consistency, and isolation. This validates our decision to build Amino on top of a database rather than an existing file system.

6.1 Future work

Applications are currently responsible for handling their own data structures during a transaction. If the application has internal state, and a transaction is aborted, then its state and the file system state are not coherent. We plan to create an API that will let applications use file system transactions to protect in-memory updates. This way, applications can safely update their in-memory structures together with an associated file. If the transaction aborts, then both the application’s memory and the file are restored.

We also plan to further improve the performance of Amino. There are two performance aspects we plan to focus on: (1) our data schema and (2) our `ptrace` monitoring interface. Currently our Data database uses a balanced tree. We plan to use a custom access method that will write pages to a file (or possibly raw disk) directly. This will give us more control over data placement, and allow us to align data properly with the native OS page size. BDB's modular design means that we can use the existing locking, logging, transaction, and caching components. This should help improve performance for data operations, which suffer compared to a standard disk based file system. We also plan to investigate changes to the Path database that would improve concurrent access. Because the database schema is so flexible compared to a file system layout, we also plan to explore application specific schemas (e.g., changing B-trees to hash tables, adding fields, or introducing indices).

To further improve the `ptrace` interface, we believe that we should reduce both the number of context switches and data copies between the kernel and the monitor. We believe that three key ways to do this are: (1) the kernel should use a shared-memory segment to manipulate the user process's registers so that data copies and context switches are reduced; (2) the monitor should be able to map regions from the user process's address space into its own; and (3) several `ptrace` operations could be bundled into a single system call (e.g., waiting for notification could be combined with retrieving registers) to reduce context switches [32]. Finally, we are also considering porting performance-sensitive subsets of the database and the monitor into the kernel (e.g., path name resolution and file-table lookups).

To obtain copies of the micro-benchmark programs used in this article go to www.fsl.cs.sunysb.edu/~cwright/benchmarks/.

REFERENCES

- A. D. Alexandrov, M. Ibel, K. E. Schauser, and C. J. Scheiman. Extending the Operating System at the User Level: the Ufo Global File System. In *Proceedings of the Annual USENIX Technical Conference*, pages 77–90, Anaheim, CA, January 1997. USENIX Association.
- B. Berliner and J. Polk. Concurrent Versions System (CVS). www.cvshome.org, 2001.
- B. Callaghan, B. Pawlowski, and P. Staubach. NFS Version 3 Protocol Specification. Technical Report RFC 1813, Network Working Group, June 1995.
- P. M. Chen, W. T. Ng, S. Chandra, C. Aycok, G. Rajmani, and D. Lowell. The Rio File Cache: Surviving Operating System Crashes. In *Proceedings of the Seventh International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS-VII)*, pages 74–83, Cambridge, MA, October 1996. ACM.
- CollabNet, Inc. Subversion. <http://subversion.tigris.org>, 2004.
- J. Dike. A user-mode port of the Linux kernel. In *Proceedings of the 4th Annual Linux Showcase and Conference*, pages 63–72, Atlanta, GA, October 2000. USENIX Association.
- D. Ellard, J. Ledlie, P. Malkani, and M. Seltzer. Everything You Always Wanted to Know about NFS Trace Analysis, but Were Afraid to Ask. Technical Report TR-06-02, Harvard University, Cambridge, MA, June 2002.
- D. Ellard and M. Seltzer. New NFS Tracing Tools and Techniques for System Analysis. In *Proceedings of the Annual USENIX Conference on Large Installation Systems Administration*, San Diego, CA, October 2003. USENIX Association.

- N. H. Gehani, H. V. Jagadish, and W. D. Roome. OdeFS: A File System Interface to an Object-Oriented Database. In *Proceedings of the Twentieth International Conference on Very Large Databases*, pages 249–260, Santiago, Chile, September 1994. Springer-Verlag Heidelberg.
- S. Ghemawat, H. Gobioff, and S. T. Leung. The Google File System. In *Proceedings of the 19th ACM Symposium on Operating Systems Principles (SOSP '03)*, pages 29–43, Bolton Landing, NY, October 2003. ACM SIGOPS.
- D. P. Ghormley, D. Petrou, S. H. Rodrigues, and T. E. Anderson. SLIC: An Extensibility System for Commodity Operating Systems. In *Proceedings of the Annual USENIX Technical Conference*, pages 39–52, Berkeley, CA, June 1998. ACM.
- P. Giarrusso. Fwd: Re: [patch 1/4] UML Support - Ptrace: adds the host SYSEMU support, for UML and general usage, July 2005. www.uwsg.iu.edu/hypermail/linux/kernel/0507.3/1992.html.
- I. Goldberg, D. Wagner, R. Thomas, and E. Brewer. A Secure Environment for Untrusted Helper Applications (Confining the Wily Hacker). In *Proceedings of the Sixth USENIX UNIX Security Symposium*, pages 1–13, San Jose, CA, July 1996. USENIX Association.
- M. Haardt and M. Coleman. *ptrace(2)*. Linux Programmer's Manual, Section 2, November 1999.
- R. Hagmann. Reimplementing the Cedar file system using logging and group commit. In *Proceedings of the 11th ACM Symposium on Operating Systems Principles (SOSP '87)*, pages 155–162, Austin, TX, October 1987. ACM Press.
- IEEE/ANSI. Information Technology–Portable Operating System Interface (POSIX)–Part 1: System Application: Program Interface (API) [C Language]. Technical Report STD-1003.1, ISO/IEC, 1996.
- M. B. Jones. Interposition Agents: Transparently Interposing User Code at the System Interface. In *Proceedings of the 14th Symposium on Operating Systems Principles (SOSP '93)*, pages 80–93, Asheville, NC, December 1993. ACM.
- J. Katcher. PostMark: A New Filesystem Benchmark. Technical Report TR3022, Network Appliance, 1997. www.netapp.com/tech_library/3022.html.
- M. L. Kazar, B. W. Leverett, O. T. Anderson, V. Apostolides, B. A. Bottos, S. Chutani, C. F. Everhart, W. A. Mason, S. Tu, and E. R. Zayas. DECORUM File System Architectural Overview. In *Proceedings of the Summer USENIX Technical Conference*, pages 151–164, Anaheim, OH, June 1990. USENIX Association.
- D. G. Korn and E. Krell. A New Dimension for the Unix File System. *Software-Practice and Experience*, 20(S1):19–34, June 1990.
- P. Lewis, A. Bernstein, and M. Kifer. *Databases and Transaction Processing: An Application-Oriented Approach*, chapter 8: Database Design II: Relational Normalization Theory, pages 211–260. Addison Wesley, Boston, MA, 2002.
- D. E. Lowell and P. M. Chen. Free transactions with Rio Vista. In *Proceedings of the 16th Symposium on Operating Systems Principles (SOSP '97)*, pages 92–101, Saint Malo, France, October 1997. ACM.
- D. Mazières. A Toolkit for User-Level File Systems. In *Proceedings of the Annual USENIX Technical Conference*, pages 261–274, Boston, MA, June 2001. USENIX Association.
- M. K. McKusick and G. R. Ganger. Soft Updates: A Technique for Eliminating Most Synchronous Writes in the Fast Filesystem. In *Proceedings of the Annual USENIX Technical Conference, FREENIX Track*, pages 1–18, Monterey, CA, JUNE 1999. USENIX Association.
- M. K. McKusick, W. N. Joy, S. J. Leffler, and R. S. Fabry. A fast file system for UNIX. *ACM Transactions on Computer Systems*, 2(3):181–197, August 1984.

- Microsoft Corporation. Microsoft MSDN WinFS Documentation. <http://msdn.microsoft.com/data/winfs/>, October 2004.
- Dejan Milojicic, Alan Messer, James Shau, Guangrui Fu, and Alberto Munoz. Increasing relevance of memory hardware errors: a case for recoverable programming models. In *Proceedings of the 9th ACM SIGOPS European workshop*, pages 97–102, Kolding, Denmark, 2000. ACM Press.
- N. Murphy, M. Tonkelowitz, and M. Vernal. The Design and Implementation of the Database File System. www.eecs.harvard.edu/~vernal/learn/cs261r/index.shtml, January 2002.
- MySQL AB. MySQL: The World’s Most Popular Open Source Database. www.mysql.org, July 2005.
- M. A. Olson. The Design and Implementation of the Inversion File System. In *Proceedings of the Winter 1993 USENIX Technical Conference*, pages 205–217, San Diego, CA, January 1993. USENIX.
- Oracle Corporation. Oracle Internet File System Archive Documentation. http://otn.oracle.com/documentation/ifs_arch.html, October 2000.
- A. Purohit, C. Wright, J. Spadavecchia, and E. Zadok. Develop in User-Land, Run in Kernel Mode. In *Proceedings of the 2003 ACM Workshop on Hot Topics in Operating Systems (HotOS IX)*, pages 109–114, Lihue, Hawaii, May 2003. USENIX Association.
- M. Rosenblum and J. K. Ousterhout. The design and implementation of a log-structured file system. In *Proceedings of 13th ACM Symposium on Operating Systems Principles*, pages 1–15, Asilomar Conference Center, Pacific Grove, CA, October 1991. Association for Computing Machinery SIGOPS.
- D. S. Santry, M. J. Feeley, N. C. Hutchinson, A. C. Veitch, R. W. Carton, and J. Ofir. Deciding When to Forget in the Elephant File System. In *Proceedings of the 17th ACM Symposium on Operating Systems Principles*, pages 110–123, Charleston, SC, December 1999. ACM.
- M. Satyanarayanan, H. H. Mashburn, P. Kumar, D. C. Steere, and J. J. Kistler. Lightweight recoverable virtual memory. *ACM Transactions on Computer Systems*, 12(1):33–57, 1994.
- F. Schmuck and J. Wylie. Experience with transactions in QuickSilver. In *Proceedings of the 13th ACM Symposium on Operating Systems Principles (SOSP ’91)*, pages 239–253, Pacific Grove, CA, October 1991. ACM Press.
- M. Seltzer and M. Stonebraker. Transaction Support in Read Optimized and Write Optimized File Systems. In *Proceedings of the Sixteenth International Conference on Very Large Databases*, pages 174–185, Brisbane, Australia, August 1990. Morgan Kaufmann.
- M. Seltzer and O. Yigit. A new hashing package for UNIX. In *Proceedings of the Winter USENIX Technical Conference*, pages 173–184, Dallas, TX, January 1991. USENIX Association.
- M. I. Seltzer, G. R. Ganger, M. K. McKusick, K. A. Smith, C. A. N. Soules, and C. A. Stein. Journaling Versus Soft Updates: Asynchronous Meta-data Protection in File Systems. In *Proc. of the Annual USENIX Technical Conference*, pages 71–84, San Diego, CA, June 2000. USENIX Association.
- Sendmail Consortium. Sendmail home page. www.sendmail.org, August 2004.
- Sendmail, Inc. Sendmail Advanced Message Server. www.sendmail.com/products/mailcenter/sams/, 2004.
- Sleepycat Software, Inc. *Berkeley DB Reference Guide*, 4.3.27 edition, December 2004. www.sleepycat.com.
- M. Szeredi. Filesystem in Userspace. <http://fuse.sourceforge.net>, February 2005.

Received December 2005; revised Month Year; accepted Month Year