

Virtual Log Based File Systems for a Programmable Disk*

Randolph Y. Wang[†] Thomas E. Anderson[‡] David A. Patterson[†]

Abstract

In this paper, we study how to minimize the latency of small transactional writes to disks. The basic approach is to write to free sectors that are near the current disk head location by leveraging the embedded processor core inside the disk. We develop a number of analytical models to demonstrate the performance potential of this approach. We then present the design of a variation of a log-structured file system based on the concept of a *virtual log*, which supports fast small transactional writes without extra hardware support. We compare our approach against traditional update-in-place and logging systems by modifying the Solaris kernel to serve as a simulation engine. Our evaluations show that random updates on an unmodified UFS execute up to an order of magnitude faster on a virtual log than on a conventional disk. The virtual log can also significantly improve LFS in cases where delaying small writes is not an option or on-line cleaning would degrade performance. If the current trends of disk technology continue, we expect the performance advantage of this approach to become even more pronounced in the future.

1 Introduction

In this paper, we set out to answer a simple question: how do we minimize the latency of small transactional writes to disk?

The performance of small synchronous disk writes impact the performance of important appli-

cations such as recoverable virtual memory [29], persistent object stores [2, 19], and database applications [34, 35]. These systems have become more complex in order to deal with the increasing relative cost of small writes [20].

Similarly, most existing file systems are carefully structured to avoid small synchronous disk writes. UFS by default delays data writes to the disk. More recent research has shown that it is possible to delay metadata writes as well if they are carefully ordered [10]. The Log-structured File System (LFS) [27] claims to be write optimized, but does so by batching small writes. While the structural integrity of file systems can be maintained using these delayed write techniques, none of them provides data reliability. Write-ahead logging systems [4, 7, 13, 33] accumulate small updates in a log and replay the modifications later by updating in place. Databases often place the log on a separate disk to avoid having the small updates to the log conflict with reads to other parts of the disk. Our interest is in the limits to small write performance for a single disk. Understanding the single disk case also may enable us to improve the performance of log truncation.

Because of the limitations imposed by disks, non-volatile RAM (NVRAM) or an uninterruptable power supply (UPS) is often used to provide fast transactional writes [3, 16, 17, 20]. However, when write locality exceeds buffer capacity, performance degrades. There are also applications that demand stricter guarantees of reliability and integrity than that of either NVRAM or UPS. Fast small disk writes can provide a cost effective complement to NVRAM solutions in these respects.

Unlike most of the existing file systems, we do not seek to avoid small disk writes. Instead, our basic approach is to write to a disk location that is closest to the current head location, while retaining transactional behavior. We call this *eager writing*. Eager writing requires the file system to be aware of the precise disk head location and disk geometry. One way to satisfy this requirement is to enhance the disk interface to the host so that the host file

*This work was supported in part by the Defense Advanced Research Projects Agency (DABT63-96-C-0056), the National Science Foundation (CDA 9401156), California MICRO, the AT&T Foundation, Digital Equipment Corporation, Hewlett Packard, IBM, Intel, Sun Microsystems, and Xerox Corporation. Anderson was also supported by a National Science Foundation Presidential Faculty Fellowship.

[†]Computer Science Division, University of California, Berkeley, {rywang,pattsrn}@cs.berkeley.edu

[‡]Department of Computer Science and Engineering, University of Washington, Seattle, tom@cs.washington.edu

system can have precise knowledge of the disk state. A second solution is to migrate into the disk some of the file system responsibilities which are traditionally executed on the host. In the rest of this paper, we will assume this second approach, although the techniques that we will describe do not necessarily depend on the ability to run file systems inside disks and can work equally well using the first solution.

Several technology trends have simultaneously enabled and necessitated the approach of migrating file system responsibility into the disk. First, Moore's Law has driven down the relative cost of CPU power to disk bandwidth, enabling powerful systems to be embedded on disk devices [1]. For example, the 233 Mhz StrongARM 110 by Digital costs under \$20 [6]. As this trend continues, it will soon be possible to run the entire file system on the disk.

Second, disk bandwidth has been scaling faster than other aspects of the disk system. Disk bandwidth has been growing at 40% per year [12]. I/O bus performance has been scaling less quickly [22]. The ability of the file system to communicate with the disk (to reorganize the disk, for example) without consuming valuable I/O bus bandwidth has become increasingly important. Seek time, half-rotation delay, head switch time, and track switch time are improving at even slower rates. Typical disk latency has improved at an annual rate of only 10% in the past decade [22]. A file system whose small write latency is largely decided by the disk bandwidth instead of any other parameters will continue to perform well.

Third, the increasing complexity of the modern disk drives and the fast product cycles make it increasingly difficult for operating system vendors to incorporate useful device heuristics into their file systems to improve performance. The traditional host file system and the disk firmware do not share perfect knowledge of each other and the approach of increasing complexity of the file system to match the complexity of the modern disks cannot work. In contrast, by running file system code inside the disk, we can combine the precise knowledge of the file system semantics and detailed disk mechanism to perform optimizations that are otherwise impossible.

The basic concept of performing writes near the current disk head position is by no means a new one [5, 9, 11, 14, 24]. But these systems either do not provide transactional behavior, or have poor failure recovery times, or require NVRAM for transactions. In this work, we present the design of a *virtual log*, a logging strategy based on eager writing with these

unusual features:

- Virtual logging exploits fast writes to provide transactions without special hardware support by writing the commit record to the next free block after the data write(s).
- The virtual log allows space occupied by obsolete entries to be reused without recopying live entries.
- The virtual log boot straps its recovery from the log tail pointer, which can be stored at a fixed location on disk as part of the firmware power down sequence, allowing efficient normal operations.

We discuss two designs in which the virtual log can be used to improve file system performance. The first is to use it to implement a logical disk interface. This design, called a VLD, does not alter the existing disk interface and can deliver the performance advantage of eager writing to an unmodified file system. In the second approach, which we have not implemented, we seek a tighter integration of the virtual log into the file system; we present the design of a variation of LFS, called VLFS, which should provide even better performance than the VLD. While our work focuses on a traditional architecture in which the user of the transactional file system executes on the host, the technique can also be extended to move the transaction code itself into the disk to speed up operations such as read-modify-write [25]. We develop analytical models and algorithms to answer a number of fundamental questions about eager writing:

- What is the theoretical limit of the performance of this approach?
- How should we re-engineer the host/disk interface to take advantage of the knowledge of the disk mechanism and file system semantics?
- How can we ensure open space under the disk head?
- How does this approach fare as different components of the disk mechanism improve at different rates?

We evaluate our approach against update-in-place and logging by modifying the Solaris kernel to serve as a simulation engine. Our evaluations show that an *unmodified* UFS on an eager writing disk runs about ten times as fast as an update-in-place system for random small updates. Eager writing's economical use of bandwidth also allows it to significantly improve LFS in cases where delaying small writes is not an option or on-line cleaning would degrade performance. As disk bandwidth continues to improve faster than disk latency, the performance advantage of eager writing should become more pro-

found in the future.

Of course, these benefits may come at a price of potentially reducing read performance. Like LFS, the virtual log does not necessarily place data in the optimal position for future reads. But as increasingly large file caches are employed, modern file systems such as the Network Appliance file system report predominantly write traffic [17]. Large caches also provide opportunity for read reorganization before the reads happen [23].

Although our evaluations show significant performance promise of the virtual log, this paper remains a preliminary study. A full evaluation of the approach would require solutions to algorithmic and policy questions of the data reorganizer as part of a complete VLFS implementation. These issues, as well as questions such as how to extend the virtual logging technique to multiple disks are subjects of our ongoing research.

The remainder of the paper is organized as follows. Section 2 presents the eager writing analytical models. Section 3 presents the design of the virtual log and the virtual log based LFS. Section 4 describes the experimental platform that is used to evaluate the update-in-place, logging, and eager writing strategies. Section 5 evaluates these alternatives with a series of experiments. Section 6 describes some of the related work. Section 7 concludes.

2 Limits to Low Latency Writes

The principle of writing data near the current disk head position is most effective when the disk head is always on a free sector when a write request arrives. This is not always possible in reality. In this section, we develop a number of analytical models to estimate the amount of time needed to locate free sectors under various utilizations. These models will help us evaluate whether eager writing is a sound strategy, set performance targets for real implementations, and predict future improvements as disks improve. We also use the models to motivate new file system allocation and reorganization algorithms. Because we are interested in the theoretical limits of eager writing latency, the models are for the smallest addressable unit: a disk sector (although the validity of the formulas do not depend on the actual sector size and can apply equally well to larger blocks).

2.1 A Single Track Model

Suppose a disk track contains n sectors, its utilization is $(1 - p)$, and the free space is randomly distributed. On average, the number of sectors the disk head must skip before arriving at any free sector is:

$$\frac{(1 - p)n}{1 + pn} \quad (1)$$

Proof. Suppose n is the track size, k is the number of free sectors in the track, and $E(n, k)$ is the expected number of used sectors we encounter before reaching the first free sector. With a probability of k/n , the first sector is free and the number of sectors the disk head must skip is 0. Otherwise, with a probability of $(n - k)/n$, the first sector is used and we must continue searching in the remaining $(n - 1)$ sectors, of which k are free. Therefore, the expected delay in this case is $[1 + E(n - 1, k)]$. This yields the recurrence:

$$E(n, k) = \frac{n - k}{n}[1 + E(n - 1, k)] \quad (2)$$

By induction on n , it is easy to prove that the following is the unique solution to (2):

$$E(n, k) = \frac{n - k}{1 + k} \quad (3)$$

(1) follows if we substitute k in (3) with pn . \square

For example, when there is only one free sector left in the track ($pn = 1$), (1) degenerates to a half rotation. More generally, (1) is roughly the ratio between occupied sectors and free ones. This is a promising result for the eager writing approach because, for example, even at a relatively high utilization of 80%, we can expect to incur only a four-sector rotation delay to locate a free sector. For today's disks, this translates to less than 100 μs . In six to seven years, this delay should improve by another order of magnitude because it scales with platter bandwidth. In contrast, it is difficult for an update-in-place system to avoid at least a half-rotation delay. With today's technology, this is at best 3 ms , and it improves at a slower rate than bandwidth. This difference is the fundamental reason why eager writing has the potential to outperform update-in-place.

2.2 A Single Cylinder Model

We now extend (1) to cover a single cylinder. We compare the time needed to locate the nearest sector in the current track against the time needed to find one in other tracks in the same cylinder and take

the minimum of the two. Therefore the expected latency can be expressed as:

$$\sum_x \sum_y \min(x, y) \cdot f_x(p, x) \cdot f_y(p, y) \quad (4)$$

where x is the delay (in units of sectors) experienced to locate the closest sector in the current track, y is the delay experienced to locate the closest sector in other tracks in the current cylinder, and $f_x(p, x)$ and $f_y(p, y)$ are the probability functions of x and y , respectively, under the assumption of a free space percentage of p . Suppose a head switch costs s and there are t tracks in a cylinder, then the probability functions can be expressed as:

$$f_x(p, x) = p(1 - p)^x \quad (5)$$

$$f_y(p, y) = f_x(1 - (1 - p)^{t-1}, y - s) \quad (6)$$

In other words, f_x is the probability that there are x occupied sectors followed by a free sector in the current track, and f_y is the probability that the first $(y - s)$ rotational positions in all $(t - 1)$ tracks are occupied and there is one free sector at the next rotational position in at least one of these tracks.

Figure 1 validates the model of (4) with a simulation of the HP97560 and the Seagate ST19101 disks whose relevant parameters are detailed in Table 1. The HP97560 is a relatively old disk; we chose it because its characteristics are well documented and its existing simulators are thoroughly validated [18, 28]. The Seagate is state-of-art technology. It is more complex and the simulator is at best a coarse approximation. For example, the models are for a single density zone while the actual disk has multiple density zones. The simulated eager writing algorithm in Figure 1 is not restricted to the current cylinder and always seeks to the nearest sector. (We will show a modified eager writing algorithm in the presence of a free space compactor in the next section.) The figure shows that the single cylinder model is in fact a good approximation for an *entire* zone. This is a simple consequence of the relative seek cost – nearby cylinders are not much more likely than the current cylinder to have a free sector in a rotationally good position. This is also because the head switch time and single cylinder seek time are close to each other.

Compared to the half-rotation delays of 7 ms for the HP and 3 ms for the Seagate that an update-in-place system may incur, Figure 1 promises significant performance improvement, especially at lower utilizations. Indeed, the figure shows that the eager writing latency has improved by nearly an order of magnitude on the newer Seagate compared to

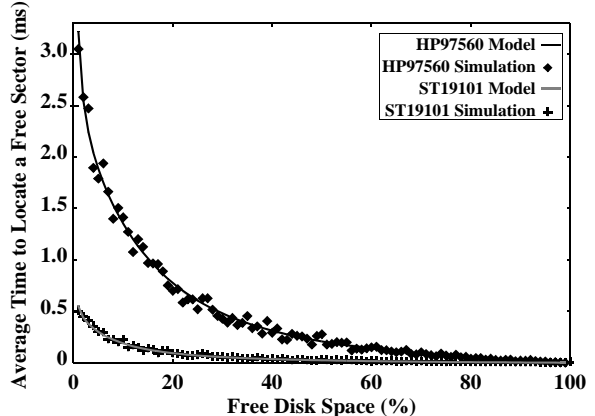


Figure 1: Amount of time to locate the first free sector as a function of the disk utilization. The model is for a single cylinder while the simulation results are for the entire zone.

	HP97560	ST19101
Sectors per Track (n)	72	256
Tracks per Cylinder (t)	19	16
Head Switch (s)	2.5 ms	0.5 ms
Minimum Seek	3.6 ms	0.5 ms
Rotation Speed (RPM)	4002	10000
SCSI Overhead (o)	2.3 ms	0.1 ms

Table 1: Parameters of the HP97560 and the Seagate ST19101 disks.

the HP disk. As the disk utilization increases, however, it takes longer for an eager writing disk head to find a nearby free sector and the performance degrades. When this occurs, one solution is to compact the free space under the disk head to ensure low latency. We next develop a model assuming the presence of a compactor.

2.3 A Model Assuming a Compactor

Compacting free space and reorganizing data are responsibilities of the storage subsystem and it is natural to dedicate these tasks to the embedded processor in the disk. A self-reorganizing disk can take advantage of the “free” bandwidth between the disk head and the platters during idle periods without consuming valuable I/O bus bandwidth, polluting host cache and memory, or interfering with host CPU operation.

We now develop a performance model of eager writing assuming the existence of a free space compactor that constantly generates empty tracks. We start writing to an empty track and continue to fill the same track with newly arrived write requests until the utilization of the current track reaches a

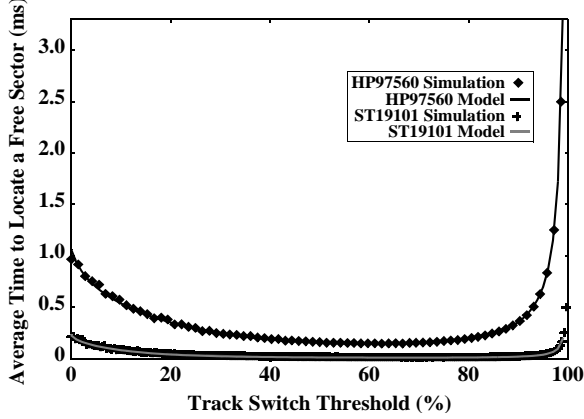


Figure 2: Average latency to locate free sectors for all writes performed to an initially empty track as a function of the track switch threshold. The track switch threshold is the percentage of free sectors reserved per track before a switch occurs. A high threshold corresponds to frequent switches.

threshold. Then we switch to the next empty track and repeat the process. If writes arrive with no delay, then it would be trivial to fill the track. If writes arrive randomly and we assume the free space distribution is random at the time of the arrivals, then writes between successive track switches follow the model of (1). Therefore, substituting pn in (1) with i , the number of free sectors, the total number of sectors to skip between track switches can be expressed as:

$$\sum_{i=k+1}^n \frac{n-i}{1+i} \quad (7)$$

where n is the total number of sectors in a track and k is the number of free sectors reserved per track before switching tracks. Suppose each track switch costs s and the rotation delay of one sector is r ; because we pay one track switch cost per $(n-k)$ writes, the average latency per sector of this strategy can be expressed as:

$$\frac{s + r \cdot \sum_{i=k+1}^n \frac{n-i}{1+i}}{n-k} \quad (8)$$

So far, we have assumed that the free space distribution is always random. This is not the case with the approach of filling the track up to a threshold because, for example, the first free sector immediately following a number of used sectors is more likely to be selected for eager writing than one following a number of free sectors. In general, this non-randomness increases latency. Although the precise modeling of the non-randomness is quite complex, through approximations and empirical trials, we have found that adding the following ϵ function

to (7) works well for a wide range of disk parameters in practice:

$$\epsilon(n, k) = \frac{(n-k-0.5)^{p+2}}{(8 - \frac{n}{96}) \cdot (p+2) \cdot n^p} \quad (9)$$

where $p = 1 + n/36$. By approximating the summation in (7) with an integral and adding the correction factor of (9) to account for the non-randomness, we arrive at the final latency model:

$$\frac{s + r \cdot [(n+1) \ln \frac{n+2}{k+2} - (n-k) + \epsilon(n, k)]}{n-k} \quad (10)$$

Figure 2 validates the model of (10). When we switch tracks too frequently, although the free sectors in any particular track between switches are plentiful, we are penalized by the high switch cost. When we switch tracks too infrequently, it becomes increasingly difficult to locate free sectors in an increasingly crowded track and the performance is also non-optimal. In general, the model aids the selection of an optimal threshold for a particular set of disk parameters. More importantly, the model also reassures us that we need not suffer the performance degradation seen at high utilizations in Figure 1 if we judiciously choose a track switch threshold. As long as there is sufficient idle time to compact free space, we can expect to stay comfortably at the lower part of the curve in Figure 2, which again promises significant improvement over an update-in-place approach. Figure 1 and Figure 2 also show why the demands on the compactor are less than the demands on the LFS cleaner. This is because the models indicate that compacting is only necessary for high utilizations and the compactor can move data at small granularities.

3 A Virtual Log

Eager writing allows a high degree of location independence of data. It also presents two challenges: how to locate data as its physical location can constantly change, and how to recover this location information after a failure. In this section, we explain the rationale behind the design of the virtual log and the file systems built on it.

3.1 The Indirection Map

To support location independence, we introduce a level of indirection, similar to the technique found in a number of existing file systems [5, 8, 9]. The host file system manipulates *logical addresses*. The

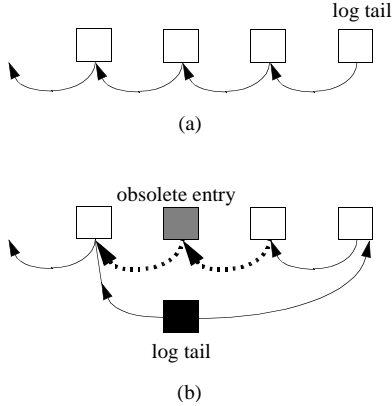


Figure 3: Maintaining the virtual log. (a) New map entry sectors are appended to a backward chain. (b) Implementing the backward chain as a tree to support map entry overwriting.

disk system maintains a *logical-to-physical address map*. As new blocks are allocated, new map entries are created; as blocks are overwritten, their map entries are updated; as blocks are read, their map entries are queried for their physical locations; and as blocks are freed, so are their map entries. Most modern disk drives already maintain a similar map so that they can transparently remap failed sectors. Updates to such a map, however, occur rarely. In the rest of this section, we will describe how we keep the map persistent, how we perform updates, how we recover the map after a failure, and how file systems can be built using these primitives. Our goal is to avoid some of the inefficiencies and inflexibilities of the previous approaches. These pitfalls include:

- lengthy scanning of large portions of the disk to recover the map,
- reliance on NVRAM, which is not always available (and even when it is, NVRAM can usually be put to more productive uses than storing map entries for recovery),
- excessive overhead in terms of space and extra I/O's needed to maintain the map, and
- altering the physical block format on disks to include, for example, a self-identifying header for each disk block.

3.2 Maintaining and Recovering the Virtual Log

We implement the indirection map as a table. To keep the map persistent, we leverage the low latency offered by eager writing. Whenever an update takes place, we write the piece of the table that contains the new map entry to a free sector near the disk head. Suppose the file system addresses the disk

at sector granularity and each map entry consumes a word, the indirection map will consume a storage overhead of less than 1% of the disk capacity. We will discuss how to further reduce this storage overhead to enable the entire indirection map to be stored in disk memory in the next section. If the transaction includes multiple data blocks and their map entries do not fall in the same map sector, then multiple map sectors may need written. Although the alternative of logging the multiple map entries using a single sector may better amortize the cost of map entry updates, it requires garbage collecting the obsolete map entries and is not used in our current design. We will discuss the issues of how to minimize map entry update I/O and how to avoid garbage collection later in the paper.

To be able to recover the map after a failure, we must be able to identify the locations of the map entry sectors which are scattered throughout the disk due to eager writing. One way to accomplish this is to thread all the map entry sectors together to form a log. We term this a *virtual log* because the components of the log are not necessarily physically contiguous. Because eager writing prevents us from predicting the location of the next on-disk map entry, we cannot maintain forward pointers in the map entries. Instead, we chain the map entries backward by appending new tail entries to the virtual log as shown in Figure 3a. Note that we can adapt this technique to a disk that supports a small writable header per block, in which case a map entry including the backward pointer will be placed in the header. We do not assume this support in the rest of the discussion.

As map entries are overwritten, if we continuously append to the virtual log, the backward chain will accumulate obsolete sectors over time. We cannot simply reuse these obsolete sectors because doing so will break the backward chain. Our solution is to implement the backward linked list as a tree as shown in Figure 3b. Whenever an existing map entry is overwritten, a new log tail is introduced as the new tree root. One branch of the root points to the previous root; the other points to the map sector following the overwritten map entry. The overwritten sector can be recycled without breaking the virtual log. As Section 3.3 will show, we can keep the entire virtual log in disk memory during normal operation. Consequently, overwriting a map entry requires only one disk I/O to create the new log tail.

In order to be able to recover the virtual log without scanning the disk, we must remember the location of the log tail. Modern disk drives use residual power left in the drive to park their heads in a land-

ing zone at the outer perimeter of the disks prior to spinning down the drives. It is easy to modify the firmware so that the drive records the current log tail location at a fixed location on the disk before it parks the actuator[21, 32]. To be sure of the validity of the content stored at this fixed location, we can protect it with a checksum and clear it after recovery. In the extremely rare case when the power down mechanism fails, we can detect the failure by computing the checksum and resort to scanning the disk to retrieve the log tail.

With a stable log tail, recovering the virtual log is straightforward. We start at the log tail as the root and traverse the tree on the frontier based on age. Obsolete log entries can be recognized as such because their updated versions are younger and traversed earlier.

3.3 Implementing LFS on the Virtual Log

So far, we have described 1) a generic logging strategy that can support transactional behavior, and 2) an indirection map built on the log which implements a logical disk interface that supports location independence. One advantage of this approach is that we can implement eager writing behind a logical disk interface and deliver its performance advantage to an unmodified file system. We now describe another application of the virtual log: implementing a variant of the log-structured file system (VLFS). Unlike the logical disk approach of the previous section, the VLFS requires modifying the disk interface to the host file system. As a result of seeking a tighter integration between the file system and the programmable disk, however, this allows a number of optimizations impossible with an unmodified UFS. Currently, we have not directly implemented VLFS. Instead, the experiments in Section 5 are based on file systems running on the virtual log via the logical disk interface as described in the last section. We indirectly deduce the VLFS performance by evaluating these file systems.

One disadvantage of the indirection map as described in Section 3.1 is the amount of storage space and the extra I/O's needed to maintain and query the map. To solve this inefficiency, the file system can store physical addresses of the data blocks in the inodes (and indirection blocks). This is the same approach taken by LFS and is shown in Figure 4. As file blocks are written, the data blocks, the inode blocks that contain physical addresses of the data blocks, and inode maps that contain physical addresses of the inodes are all appended to the log.

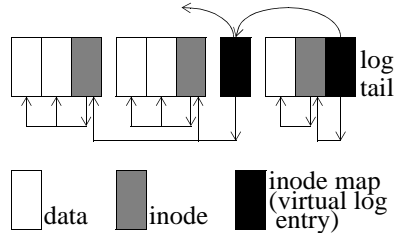


Figure 4: Implementing LFS on the virtual log. VLFS uses eager writing to write to the disk so the log is not necessarily physically contiguous. The virtual log contains only the inode map blocks.

What is different in the virtual log based implementation (VLFS) is that the log need not be physically contiguous, and only the inode map blocks logically belong to the virtual log. This is essentially adding a level of indirection to the indirection map. The advantage is that the inode map, which is the sole content of the virtual log, is now compact enough to be stored in memory; it also reduces the number of I/O's needed to maintain the indirection map because VLFS simply takes advantage of the existing indirection data structures in the file system without introducing its own.

Another LFS optimization that can also be applied to VLFS is checkpointing for recovery. Periodically, we write the entire inode map to the disk contiguously and truncate the virtual log. At recovery time, while LFS reads a checkpoint at a known disk location and rolls forward, VLFS traverses the virtual log backwards from the log tail towards the checkpoint.

In order to support VLFS, the disk needs to export a new interface to the host. In deciding how to partition the responsibility between the disk and the host, we face at least two conflicting goals. On the one hand, we need to move enough file system operations into the disk so that they can take advantage of the knowledge of the disk mechanism which is unavailable on the host. On the other hand, we would like the host to retain certain operations so that we do not place undue stress on the I/O bus. Two key questions are where to place the data cache and where to place the inode cache. To minimize the crossing of the I/O bus, one design is to retain both caches on the host. Under this design, on a write, in addition to the data blocks, the disk is also given the inode. The disk writes the data blocks, inserts their physical addresses into the inode, writes the inode, writes the inode map blocks, and returns the modified inode to the host. This design complicates the effort of keeping the inode cache coherent if the disk is to modify the inodes later during the

process of free space compacting or data reorganization. A simpler alternative is to keep the inode cache on the disk. This alternative drastically simplifies the interface and the cache coherence issue at the expense of more communication across the I/O bus. The precise placement of these caches and the associated interfaces are subjects of our ongoing research.

To support the VLFS free space compactor, which is equivalent to an LFS-style cleaner, we introduce the equivalent of the LFS segment summary, a *track summary*. The track summary details the content of a track, recording the inode number and block offset of each block. Like all other data structures of VLFS, the track summary has no permanent location and can be written anywhere in the track. In fact, although the track summary is a logically distinct data structure, it can be piggy-backed onto a log map entry to avoid extra disk writes¹.

During VLFS compacting, a whole track is read and its track summary is examined to determine the locations of the inodes that point to the live blocks. If the inodes in question are in the current track, then we can simply copy the data to its destination track, update the corresponding inodes, and rewrite the inodes. If the inodes in question are not found in the current track, then we have the choice of either reading the inodes from other tracks and completing the cleaning or “cleaning around” these blocks by leaving them where they are. Placing the only copy of the inode cache on the disk instead of on the host can simplify the implementation of the VLFS compactor by making the disk the sole entity that can read or write inodes. This can also allow concurrent write operations by both the host and the compactor using techniques similar to LFS “optimistic cleaning” [15, 30] that avoids expensive file locks.

3.4 Comparing VLFS with LFS

VLFS and LFS share a number of common advantages. Both can benefit from an asynchronous memory buffer by preventing short-lived data from ever reaching the disk. Both can benefit from disk reorganization during idle time: the LFS cleaner can generate empty segments while the VLFS compactor can compact free space under the disk head, so that during busy time both file systems can work with large extents of free space to ensure good performance.

¹We can avoid inadvertent freeing of the the piggy-backed track summaries by storing the track summary locations in memory.

Due to eager writing, VLFS possesses a number of unique advantages. First, small synchronous writes perform well on VLFS whereas the LFS performance suffers if an application requires frequent “fsync” operations². This is because LFS can only efficiently write large segments. Eager writing, however, can efficiently perform small writes at the most convenient disk locations. Second, while the free space compactor is only an optimization for VLFS, the cleaner is a necessity for LFS. In cases where idle time is scarce or disk utilization is high, LFS performance may degrade due to the repeated copying of live data by the cleaner [23, 30, 31]. Eager writing avoids this bandwidth waste and fills the available space in the most efficient way possible. Third, cleaning under LFS is a relatively heavy weight operation because it must move data at segment granularity, typically 0.5 MB or larger. As a result, LFS needs large idle intervals to mask the cleaning overhead. The VLFS compactor has no restriction in terms of the amount of data it can move so it can take advantage of relatively short idle intervals. Finally, reads can interfere with LFS writes by forcing the disk head away from free space and/or disturbing the track buffer (which can be sometimes used to absorb writes without accessing the disk platters). VLFS can gracefully accommodate reads by performing intervening writes near the data being read.

VLFS and LFS also share some common disadvantages. For example, data written randomly may have poor sequential read performance. In some of these situations, reorganization techniques that can improve LFS performance [23] should be equally applicable to VLFS. In some other situations, sequential reads after random updating is nothing more than an artifact of the I/O interface. For example, an application that streams through all the records of a database usually does not particularly care about the ordering of the records. An I/O interface that hints to the storage system to deliver the records in *any* order it deems convenient can perform well for both read and write operations. Indeed, an interface similar to that used by “Informed Prefetching” [26] can serve the virtual log well.

4 Experimental Platform

To study the effectiveness of eager writing, we evaluate the following four combinations of file sys-

²“fsync” is the system call that forces dirty data to the disk.

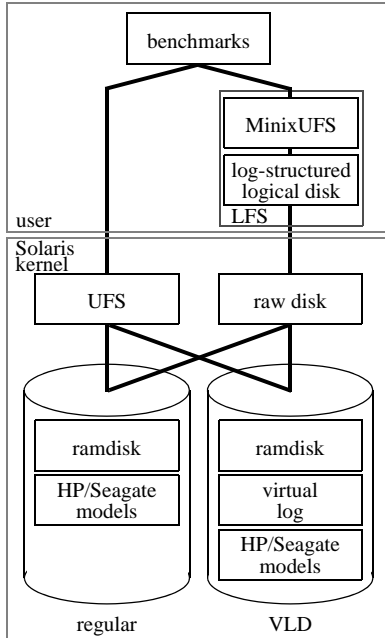


Figure 5: Architecture of the experimental platform.

tems and simulated disks: a UFS on a regular disk, a UFS on a Virtual Log Disk (VLD), an LFS on a regular disk, and an LFS on a VLD. These combinations are illustrated in Figure 5. Although we do recognize that a complete evaluation of the VLFS would require a complete implementation of VLFS, with its associated data reorganizer and free space compactor, in this paper, we take the first step of deducing the behavior of VLFS by examining the file systems running on the VLD.

Two host machines are used: one is a 50 Mhz SUN SPARCstation-10 which is equipped with 64 MB of memory and runs Solaris 2.6; the other is a similarly configured UltraSPARC-170 workstation that runs at 167 Mhz. The SPARCstation-10 supports both 4 KB and 8 KB file blocks while the UltraSPARC-170 only supports 8 KB file blocks. Because our focus is small write performance, our experiments are run on the SPARCstation-10 unless explicitly stated to be otherwise. We perform some additional experiments on the UltraSPARC-170 only to study the impact of host processing speed. In the rest of this section, we briefly describe each of the disk modules and file systems shown in Figure 5.

4.1 The Regular Disk

The regular disk module simulates a portion of the HP97560 disk or the Seagate ST19101 disk. Their relevant parameters are shown in Table 1. A

ramdisk driver is used to store file data using 24 MB of kernel memory. The Dartmouth simulator [18] is ported into the kernel to ensure realistic timing behavior of the HP disk. We simply adjust the parameters of the Dartmouth model to coincide with those of the Seagate disk to simulate the faster disk. Although the Seagate simulator is not as precise as the HP model, it gives a clear indication of how the improvements in disk technology are affecting the different disk allocation strategies. We have run all our experiments on both the the HP model and the Seagate model. Unless stated otherwise, however, the results presented are those obtained on the fast Seagate model.

One advantage of the ramdisk simulator is that it allows two possible modes of simulation. In one mode, the simulator sleeps the right amount of time reported by the Dartmouth model to imitate the behavior of the physical disks and we can conduct the evaluations by directly timing the application. In the second mode, the simulator does not sleep in the kernel; it runs at memory speed instead and only reports statistics. This allows us to speed up certain phases of the experiments whose actual elapsed time is not important. The disadvantage of the ramdisk simulator is its small size due to the limited available kernel memory. We only simulate 36 cylinders of the HP97560 and 11 cylinders of the Seagate. The results reported in Section 5, however, should be applicable to a much larger zone. We plan to use a real disk to store the file data in a future version of the simulator, although this modification would slightly complicate timing measurements because the time spent in the real disk will need to be carefully accounted for in the final analysis.

4.2 The Virtual Log Disk

The VLD adds the virtual log to the disk simulator described above. It exports the same block device driver interface so it can be used by any existing file system. The VLD maintains an in memory indirection map. As the map entries are updated, the VLD appends new map entries to the virtual log as described in Section 3.2. To avoid having the disk head trapped in regions of high utilization during eager writing, the VLD simply performs cylinder seeks only in one direction until it reaches the last cylinder, at which point it starts from the first cylinder again.

One challenge of implementing the VLD while preserving the existing disk interface is handling deletes, which are not visible at the block device driver interface. This is a common problem faced by

logical disks. Our solution is to monitor overwrites: when a logical address is re-used by a write, the VLD detects that the old logical-to-physical mapping can be freed and a new one needs to be established. The disadvantage of the approach is that it does not capture the blocks that are freed by the file system but not yet overwritten.

Another question that arose during the implementation of the VLD is the choice of the physical block size. For example, suppose the host file system writes 4 KB blocks. If the VLD chooses a physical block size of 512 B, it may need to locate eight separate free sectors to complete eager writing one logical block. If the VLD chooses a physical block size of 4 KB instead, although it takes longer to locate the first free 4 KB-aligned sector, it is guaranteed to have eight consecutive free sectors afterwards. To understand this tradeoff, we can extend the model of (1) in Section 2.1. Suppose the file system logical block size is B and the VLD physical block size is b ($b \leq B$), then the average amount of time (expressed in the numbers of sectors skipped) needed to locate all the free sectors for a logical block is:

$$\frac{(1-p)n}{b+pn} \cdot B \quad (11)$$

(11) indicates that the latency is lowest when the physical block size matches the logical block size. In all our experiments, we have used a physical block size of 4 KB. Each physical block requires one map entry; so the entire map consumes 24 KB and is stored in memory in our simulations.

The third issue concerns the interaction between eager writing and the disk track buffer read-ahead algorithm. When reading, the HP97560 (or at least the Dartmouth simulator) keeps in cache only the sectors from the beginning of the current request through the current read-ahead point and discards the data whose addresses are lower than that of the current request. This algorithm makes sense for sequential reads of data whose physical addresses increase monotonically. This is the case for file systems which directly manipulates physical addresses. For VLD, however, the combination of eager writing and the logical-to-physical address translation means that the sequentially read data may not necessarily have monotonically increasing physical addresses. As a result, the HP97560 tends to purge data prematurely from its read-ahead buffer under VLD. The solution is to aggressively prefetch the entire track as soon as the head reaches the target track and not discard data until it is delivered to the host during sequential reads on the VLD. The measurements of sequential reads on the VLD in

Section 5 were taken with this modification.

Lastly, the VLD module also implements a free space compactor. Although the eager writing strategy should allow the compactor to take advantage of idle intervals of arbitrary length, for simplicity, our compactor compacts free space at the granularity of tracks. During idle periods, the compactor reads the current track and uses eager writing to copy the live data to other tracks, in a way similar to hole-plugging under LFS [23, 36]. Currently, we choose compaction targets randomly and plan to investigate optimal VLD compaction algorithms (e.g., ones that preserve or enhance read and write locality) in the future. Applying the lessons learned from the models in Section 2.3, the VLD fills empty tracks to a certain threshold (75% in the experiments). After exhausting empty tracks generated by the compactor, the VLD reverts to the greedy algorithm modeled in Section 2.2.

4.3 UFS

Because both the regular disk and the VLD export the standard block device driver interface, we can run the Solaris UFS (subsequently also labeled as UFS) unmodified on these disks. We configure UFS on both disks with a block size of 4 KB and a fragment size of 1 KB. Like most other Unix local file systems, the Solaris UFS updates metadata synchronously while the user can specify whether the data writes are synchronous. It also performs aggressive prefetching after several sequential reads are detected.

4.4 LFS

We have ported the MIT Log-Structured Logical Disk [8], a user level implementation of LFS (subsequently also labeled as LFS). It consists of two modules: the MinixUFS and the log-structured logical disk, both running at user level. MinixUFS accesses the logical disk with a block interface while the logical disk interfaces with the raw disk using segments. The block size is 4 KB and the segment size is 0.5 MB. Read-ahead in MinixUFS is disabled. A number of other issues also impact the LFS performance. First, MinixUFS employs a file buffer cache of 6.1 MB. Unless “sync” operations are issued, all writes are asynchronous. In some of the experiments in Section 5, we assume this buffer to be made of NVRAM so that the LFS configuration can have a similar reliability guarantee as that of the synchronous systems. In this context, we will examine the effectiveness of NVRAM.

Second, the logical disk’s response to a “sync” operation is determined by a tunable parameter called *partial segment threshold*. If the current segment is filled above the threshold at the time of the “sync”, the current segment is flushed to the disk as if it were full. If it is filled below the threshold, the current segment is written to the disk but the memory copy is retained to receive more writes. The partial segment threshold in the experiments is set to 75%.

Third, the original version of the logical disk only invokes the cleaner when it runs out of empty segments. We have modified the cleaner so that it can be invoked during idle periods before it runs out of free space.

5 Experimental Results

In this section, we compare the performance of eager writing against that of update-in-place and logging with a variety of micro-benchmarks. We first run the small file and large file benchmarks that are commonly used by similar file system studies. Then we use a benchmark that best demonstrates the strength of eager writing: small random synchronous updates with no idle time. This benchmark also illustrates the effect of disk utilization. Next we examine the effect of technology trends as disks improve. Last, we examine how the availability of idle time impacts the performance of eager writing and logging. Unless explicitly stated to be otherwise, the experimental platform is based on the SPARCstation-10 running on the simulated Seagate disk.

5.1 Small File Performance

We first examine two benchmarks similar to the ones used by both the original LFS study [27] and the Logical Disk study [8]. In the first benchmark, we create 1500 1 KB files, read them back after a cache flush, and delete them. The benchmark is run on empty disks. The results of this benchmark are shown in Figure 6. Under LFS, updates are flushed to disk only if the memory buffer is filled. (Otherwise, LFS pays a heavy price by continuously flushing partial segments to the disk.) Under UFS, updates are synchronous.

As expected, UFS on the VLD is able to significantly outperform the same file system on the regular disk for the create and delete phases of the benchmark. This is due to eager writing’s ability to complete small writes much faster than update-

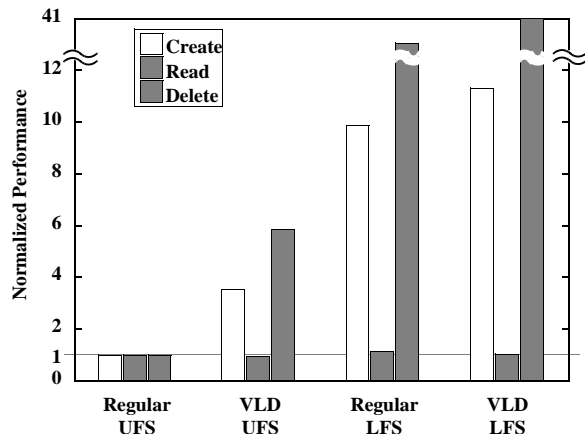


Figure 6: Small file performance. The benchmark creates, reads, and deletes 1 KB files. All performance is normalized to that of UFS running on a regular disk.

in-place. The performance improvement, however, is not as great as the analytical models would suggest. We will see why this is the case in Section 5.4. The read performance on the VLD is slightly worse than that on the regular disk because of the overhead introduced by the indirection map and the fact that read-ahead inside the disk is not as effective. We will see the same pattern in read performance of other experiments as well.

LFS comfortably out-performs UFS by batching small writes and preventing overwritten directories from reaching the disk. The performance of LFS on the regular disk and the VLD are close because eager writing can support large segment writes as efficiently as update-in-place does. If we do not consider the cost of LFS cleaning (which is not triggered in this experiment), the performance of a UFS that employs delayed write techniques such as those proposed in [10] should be between that of the unmodified UFS and that of LFS. Just like under LFS, however, delayed writes under UFS do not offer data reliability.

From this benchmark, we speculate that by integrating LFS with the virtual log, the VLFS (which we have not implemented) should be able to approximate the performance of UFS on the VLD when we must write synchronously to the disk, while retaining the benefits of LFS when asynchronous buffering is acceptable.

5.2 Large File Performance

In the second benchmark, we write a 10 MB file sequentially, read it back sequentially, write 10 MB of data randomly to the same file, read it back sequentially again, and finally read 10 MB of ran-

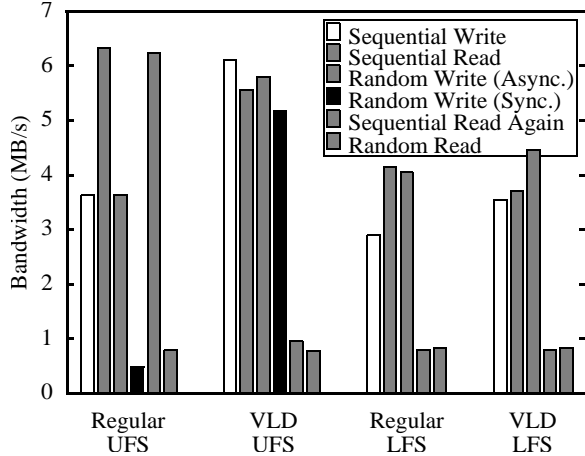


Figure 7: Large file performance. The benchmark sequentially writes a 10 MB file, reads it back sequentially, writes it again randomly (both asynchronously and synchronously for the UFS runs), reads it again sequentially, and finally reads it randomly.

dom data from the file. The performance of random I/O can also be an indication of the effect of interleaving a large number of independent concurrent workloads. The benchmark is again run on empty disks. Figure 7 shows the results. The writes are asynchronous with the exception of the two random write runs on UFS that are labeled as “Sync”. Neither the LFS cleaner nor the VLD compactor is run in the benchmark.

We first point out a few characteristics that are results of implementation artifacts. The first two phases of the LFS performance are not as good as those of UFS. This is because LFS relies on Solaris for the raw disk interface and is less efficient than the native in-kernel UFS. Furthermore, LFS disables prefetching, which explains its low sequential read bandwidth. The UFS sequential read performance is significantly better than its write performance due to aggressive prefetching both at the file system level and inside the disk. With these artifacts explained, we now examine a number of interesting performance characteristics.

First, sequential read after random write performs poorly in all LFS and VLD systems. This is because both logging and eager writing destroy spatial locality. This, fortunately, is a problem that can be solved by caching, data reorganization [23], or I/O interface changes as explained in Section 3.4.

Second, the LFS random write bandwidth is higher than that of sequential write. This is because during the random write phase, some of the blocks are written multiple times and the overall number of bytes that reach the disk is smaller than that of

the sequential case. This is an important benefit of delayed writes.

Third, on a UFS, while it is not surprising that the synchronous random writes do well on the VLD, it is interesting to note that even sequential writes perform better on the VLD. This is because of the occasional inadvertent miss of disk rotations on the regular disk. Interestingly enough, this phenomenon does not occur if we run the benchmark on the slower HP97560 disk. This evidence supports our earlier contention that tuning the host operating system to match changing technologies in both the disk and the host is indeed a difficult and error-prone task. The approach of running the file system inside the disk in general and the concept of a virtual log in particular can simplify such efforts.

Fourth, although our regular disk simulator does not implement disk queue sorting, UFS does sort the asynchronous random writes when flushing to disk. Therefore, the performance of this phase of the benchmark, which is also worse on the regular disk than on the VLD due to the reason described above, is a best case scenario of what disk queue sorting can accomplish. In general, disk queue sorting is likely to be much less effective when the disk queue length is short compared to the working set size during random updates. Similarly, if the size of a write-ahead log is small compared to the size of the database, the throughput of log truncation will be limited. The VLD based systems need not suffer from these limitations.

Finally, LFS running on the VLD delivers slightly better write bandwidth than LFS on the regular disk. In addition to the potential occasional miss of rotations on the regular disk, this is also because even with large segment-sized writes, LFS on a regular disk still needs to perform the occasional long-distance seeks, which the VLD diligently avoids with eager writing. In summary, the benchmark further demonstrates the power of combining lazy writing by the file system with eager writing by the disk.

5.3 Effect of Disk Utilization

There are a number of questions that are still unanswered by the first two benchmarks. First, the VLD always has plenty of free space to work with in the previous benchmarks. How much performance degradation can we expect when the disk is fuller? Second, the LFS cleaner is not invoked under the previous benchmarks. What then is the impact on performance when the cleaner is forced to run under various disk utilizations? Third, we know that LFS

performs poorly with frequent flushes to the disk. How much then can NVRAM help? We attempt to answer these questions with the third benchmark.

In this benchmark, we create a single file of certain size. Then we repeatedly choose a random 4 KB block to update. There is no idle time between writes. For UFS, the file is opened with the flag that indicates to the file system that all writes are to be performed synchronously. In other words, the “write” system call does not return until the block is written to the disk surface. For LFS, we do not flush to the disk until the 6.1 MB file buffer cache is full; we assume in this case that the buffer cache is made of NVRAM. We measure the steady state bandwidth of UFS on the regular disk, UFS on the VLD, and LFS on the regular disk. We repeat the experiment under different disk utilizations by varying the size of the file we update.

Figure 8 plots the average latency experienced per write. UFS on the regular disk suffers from excessive disk head movement due to the update-in-place policy. The latency increases slightly as the updated file grows because the disk head needs to travel a greater distance between successive writes. This increase may have been larger had we simulated the entire disk.

LFS provides excellent performance when the entire file fits in the NVRAM buffer. As soon as the file outgrows the available NVRAM, writes are forced to the disk; as this happens and as disk utilization increases, the cleaner must run, quickly dominating performance. The plateau between roughly 60% and 85% disk utilization is due to the fact that the LFS cleaner tends to choose less utilized segments to clean; with certain distribution patterns of free space, the number of segments to clean in order to generate the same amount of free segments may be the same as (or even larger than) that that required under a higher utilization.

With eager writing, the VLD suffers from neither the excessive disk head movements, nor the bandwidth waste as a result of cleaning. As the disk utilization increases, the VLD latency also rises. The rise, however, is not significant compared to the various overheads in the system. We examine the effects of system overheads on the VLD in the next section.

5.4 Effect of Technology Trends

Despite the advantage that eager writing has demonstrated over update-in-place so far, the performance difference is not as great as the analytical models might suggest. To see why this is the

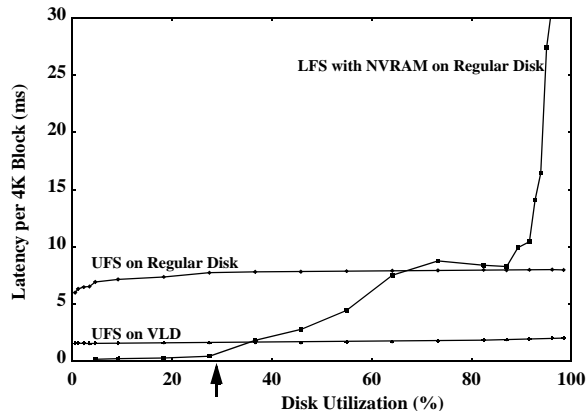


Figure 8: Performance of random small synchronous updates under various disk utilizations. The disk utilization is obtained from the Unix “df” utility and includes about 12% of reserved free space that is not usable. The arrow on the x-axis points to the size of the NVRAM used by LFS.

case, we now provide a more detailed breakdown of the benchmark times reported in the last section. We also examine how the technology trends impact the relative ratios of the individual components that make up the overall latency.

We repeat the experiment of the last section on three different platforms under the same disk utilization (80%)³. Figure 9 shows the result. The first two bars compare the performance of update-in-place and virtual logging on a SPARCstation-10 and HP disk combination. In the next run, we replace the older HP disk with the new Seagate disk. In the third run, we replace the older SPARCstation-10 with the newer UltraSPARC-170. We see that the performance gap between update-in-place and virtual logging widens from less than three-fold to almost an order of magnitude as both disks and host processors improve.

Figure 10 reveals the underlying source of this performance difference by providing the detailed breakdown of the latencies shown in the preceding figure. The component labeled as “SCSI overhead” is the fixed amount of time that the embedded processor in the disk spends processing each SCSI command. The component labeled as “transfer” is the time it takes to move the bits to or from the media after the head has been positioned over the target sector. The component labeled as “locate sectors” is the amount of time the disk spends positioning the disk head. It includes seek, rotation, and head switch times. The component labeled as “other” includes the operating system processing overhead,

³The VLD latency in this case is taken immediately after running a compactor, as explained in Section 5.5.

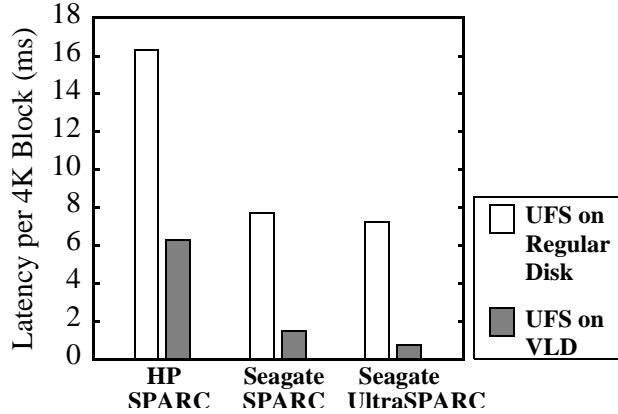


Figure 9: Performance gap between update-in-place and virtual-logging widens as disks and host processors improve.

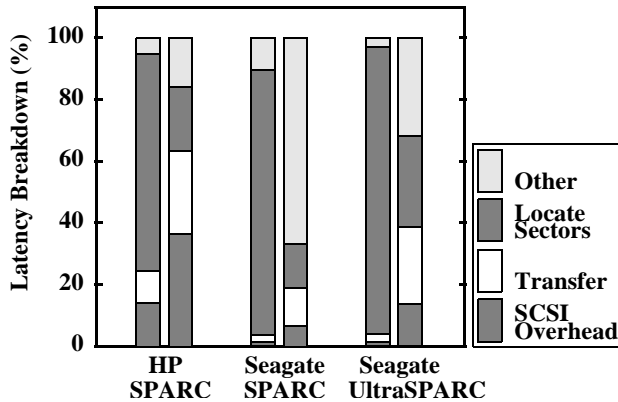


Figure 10: Breaking down the total latency into SCSI overhead, transfer time, time required to locate free sectors, and other processing time. The latency breakdowns show the underlying reason between the performance difference. Update-in-place performance becomes increasingly dominated by mechanical delays of disk while virtual logging achieves a balance between processor and disk improvements.

which includes the time to run the virtual log algorithm because the disk simulator is part of the host kernel. The time consumed by simulating the disk mechanism itself, however, is less than 5% of this component.

We see that the mechanical delay becomes a dominant factor of update-in-place latency. We also see that eager writing has indeed succeeded in significantly reducing the disk head positioning times. The overall performance improvement on the older disk or on the older host, however, is low due to the high overheads. After we replace the older disk, the performance of virtual logging becomes host limited as the component labeled as “other” dominates. After we replace the older host, however, the latency components again become more balanced. This in-

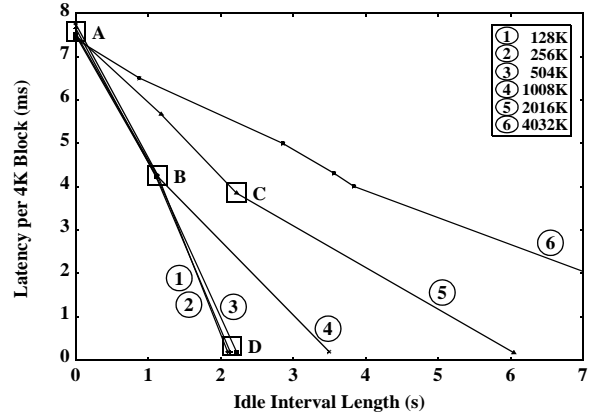


Figure 11: Performance of LFS (with NVRAM) as a function of available idle time.

dicates that the virtual log is able to ride the impressive disk bandwidth growth, achieving a balance between processor and disk improvements.

5.5 Effect of Available Idle time

The benchmark in Section 5.3 assumes zero idle time. This severely stresses LFS because the time consumed by the cleaner is not masked by idle periods. It also penalizes the VLD by disallowing free space compacting. In this section, we examine how LFS on a regular disk and how UFS on a VLD behave as more idle time becomes available. We modify the benchmark of Section 5.3 to perform a burst of random updates, pause, and repeat. The disk utilization is kept at 80%.

Figure 11 shows how the LFS performance responds to the increase of idle interval length. Each curve represents a different burst size. At point A, no idle time is available. LFS fills up the NVRAM with updates, and then proceeds to flush the buffered segments to the regular disk, invoking the cleaner when necessary.

At point B, enough idle time is available to clean one segment. If the burst size is less than or equal to 1 MB, at the time when the NVRAM becomes full, the cleaner has already generated the maximum number of empty segments possible. Consequently, it takes a constant amount of time to flush the NVRAM. This is why the first four curves meet at point B⁴. A similar latency is achieved at point C where cleaning two segments per idle interval is sufficient to compact all the free space.

⁴In this case, because the NVRAM size is larger than the available free space, the cleaner still needs to run during flushing to reclaim the free space created by the overwrites buffered in the NVRAM.

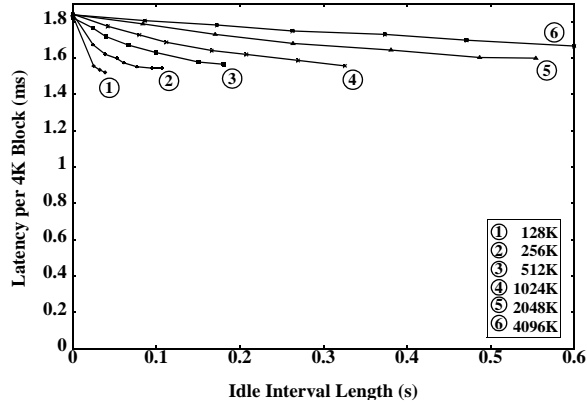


Figure 12: Performance of UFS on VLD as a function of available idle time.

Point D corresponds to sufficient idle time to flush the entire burst from NVRAM to the disk. The benchmark is able to run at NVRAM speed because the flushing time, which includes occasional cleaning, is entirely masked by the long idle periods. The first three curves coincide at this point because they correspond to burst sizes that can fit in a single segment.

We do not pretend to know the optimal LFS cleaning algorithm in the presence of NVRAM, which is outside the scope of this paper. Nevertheless, this experiment reveals two characteristics in terms of the impact of available idle time on LFS performance. First, because the cleaner moves segment-sized data, LFS can only benefit from relatively coarse-grained idle intervals. Second, unless there is sufficient idle time to mask the flushing of the burst buffered in the NVRAM, the LFS performance remains poor.

Figure 12 shows the result of repeating the same benchmark on UFS running on a VLD. Unlike the LFS cleaner, the VLD compactor has no restriction in terms of the granularity of the data it moves. For convenience, our implementation of the compactor moves data at track granularity, which is still much smaller than the LFS segment granularity. The performance on the VLD improves along a continuum of idle interval lengths. Also note the smaller scale of the axes in Figure 12 than those of Figure 11. Furthermore, the VLD performance is also more predictable, whereas LFS experiences large variances in performance depending on an array of factors including whether the NVRAM is full and whether disk cleaning is necessary. The experiments of this section run on the SPARCstation-10. As we have seen in the last section, a more powerful host processor such as the UltraSPARC-170 can easily cut the latency in Figure 12 in half.

The disadvantage of UFS on the VLD compared to LFS with NVRAM is the limiting performance with infinite amount of idle time. Because each write is synchronously written to the disk surface, the VLD experiences the overheads detailed in Section 5.4. The overheads also render the impact of the compactor less important. Fortunately, as explained in that section, as disks continue to improve and operating system implementors pay more careful attention to streamlining the I/O path, we expect this gap to narrow.

Furthermore, eager writing does not dictate the use of a UFS, nor does it preclude the use of NVRAM. A lazy writing file system that employs both an NVRAM and an eager writing VLD can enjoy 1) the low latency of and the filtering of the short-lived data by the NVRAM, and 2) the effective use of the available bandwidth by the eager writing strategy when idle intervals are short or disk utilization is high.

6 Related Work

The work presented in this paper builds upon a number of existing techniques including reducing latency by writing data near the disk head, a transactional log, file systems that support data location independence, and log-structured file systems. The goal of this study is not to demonstrate the effectiveness of any of these individual ideas. Rather, the goals are 1) provide a theoretical foundation of eager writing with the analytical models, 2) show that the integration of these ideas at the disk level can provide a number of unique benefits to both UFS and LFS, 3) demonstrate implementations that realize the benefits of eager writing without the semantic compromise of delayed writes or extra hardware support such as NVRAM, and 4) conduct a series of systematic experiments to quantify the differences of the alternatives. We are not aware of existing studies that aim for these goals.

Simply writing data near the disk head is not a new idea. Many efforts have focused on improving the performance of the write-ahead log. This is motivated by the observation that appending to the log may incur extra rotational delay even when no seek is required. The IBM IMS Write Ahead Data Set (WADS) system [11] addresses this issue for drums (fixed-head disks) by keeping some tracks completely empty. Once each track is filled with a single block, it is not re-used until the data is copied out of the track into its normal location.

Likewise, Hagmann places the write-ahead log in

its own logging disk[14], where each log append can fill any open block in the cylinder until its utilization reaches a threshold. Eager writing in our system, while retaining good logging performance, assumes no dedicated logging disk and does not require copying of data from the log into its permanent location. The virtual log *is* the file system.

A number of file systems have also explored the idea of lowering latency of small writes by writing near the disk head location. In [24], Menon proposes to use this technique to speed up parity updates in disk arrays. This is similar to some write-ahead logging systems in that it requires multiple disks and rewriting of data.

Mime [5], the extension of Loge [9], is the closest in spirit to our system. Using an indirection map, it also writes near disk head locations to improve small write performance. There are a number of differences between Mime and the virtual log. First, Mime relies on self-identifying disk blocks. Second, Mime scans free segments to recover its indirection map. As disk capacity increases, potentially reaching a terabyte by the year 2002 [22], this scanning may become a time consuming process. Third, the virtual log incorporates a free space compactor, whose importance is shown in Section 2.3.

The Network Appliance file system, WAFL [16, 17], checkpoints the disk to a consistent state periodically, uses NVRAM for fast writes between checkpoints, and can write data and metadata anywhere on the disk. An exception of the write-anywhere policy is the root inodes, which are written for each checkpoint and must be at fixed locations. Unlike Mime, WAFL supports fast recovery by rolling forward from a checkpoint using the log in the NVRAM. One goal of the virtual log is to support fast transactions and fast recovery without NVRAM, with its capacity, reliability, and cost limitations. Another difference between WAFL and the virtual log is that the WAFL write allocation decisions are made at the RAID controller level, so the opportunity to optimize for rotational delay is limited.

Our idea of fast transactional writes using the virtual log originated as a generalization of the AutoRAID technique of hole-plugging [36], to improve LFS performance at high disk utilizations without AutoRAID hardware support for self-describing disk sectors. In hole-plugging, partially empty segments are freed by writing their live blocks into the holes found in other segments. This outperforms traditional cleaning at high disk utilizations by avoiding reading and writing a large number of nearly full segments just to produce a few empty

segments[23]. AutoRAID requires an initial log-structured write of a physically contiguous segment, after which it is free to copy the live data in a segment into any empty block on the disk. Compared to AutoRAID, our approach eliminates the initial segment write and more efficiently schedules the individual writes.

7 Conclusion

In this paper, we have developed a number of analytical models which show the theoretical performance potential of eager writing, the technique of performing small transactional writes near the disk head position. We have presented the virtual log design which delivers fast transactional writes without the semantic compromise of delayed writes or extra hardware support such as NVRAM. This design requires careful log management for fast recovery. By conducting a systematic comparison of this approach against traditional update-in-place and logging approaches, we have demonstrated the benefits of applying the virtual log approach to both UFS and LFS. The availability of fast transactions can also simplify the design of file systems and sophisticated applications. As the current technology trends continue, we expect that the performance advantage of this approach will become increasingly important.

Acknowledgements

We would like to thank Arvind Krishnamurthy for many interesting discussions on the proofs of several analytical models, Marvin Solomn for discovering the simplest proof of the single track model, Daniel Stodolsky and Chris Malakapalli for helping us understand the working of Quantum and Seagate disks, and Jeanna Neefe Matthews for asking the question of how hole-plugging could be efficiently implemented without hardware support for self-describing disk sectors.

References

- [1] ADAMS, L., AND OU, M. Processor Integration in a Disk Controller. *IEEE Micro* 17, 4 (July 1997).
- [2] ATKINSON, M., CHISHOLM, K., COCKSHOTT, P., AND MARSHALL, R. Algorithms for a Persistent Heap. *Software - Practice and Experience* 13, 3 (March 1983), 259–271.

- [3] BAKER, M., ASAMI, S., DEPRIT, E., OUSTERHOUT, J., AND SELTZER, M. Non-Volatile Memory for Fast, Reliable File Systems. In *Proceedings of the Fifth International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS-V)* (Sept. 1992), pp. 10–22.
- [4] BIRRELL, A., HISGEN, A., JERIAN, C., MANN, T., AND SWART, G. The Echo Distributed File System. Technical Report 111, Digital Equipment Corp. Systems Research Center, Sept. 1993.
- [5] CHAO, C., ENGLISH, R., JACOBSON, D., STEPANOV, A., AND WILKES, J. Mime: a High Performance Parallel Storage Device with Strong Recovery Guarantees. Tech. Rep. HPL-CSP-92-9 rev 1, Hewlett-Packard Company, Palo Alto, CA, March 1992.
- [6] Chart Watch: Mobile Processors. Microprocessor Report, June 1997.
- [7] CHUTANI, S., ANDERSON, O., KAZAR, M., LEVERETT, B., MASON, W., AND SIEDBOTHAM, R. The Episode File System. In *Proc. of the 1992 Winter USENIX* (January 1992), pp. 43–60.
- [8] DE JONGE, W., KAASHOEK, M. F., AND HSIEH, W. C. The Logical Disk: A New Approach to Improving File Systems. In *Proc. of the 14th ACM Symposium on Operating Systems Principles* (December 1993), pp. 15–28.
- [9] ENGLISH, R. M., AND STEPANOV, A. A. Loge: a Self-Organizing Disk Controller. In *Proc. of the 1992 Winter USENIX* (January 1992).
- [10] GANGER, G. R., AND PATT, Y. N. Metadata Update Performance in File Systems. In *Proc. of the First Symposium on Operating Systems Design and Implementation* (November 1994), pp. 49–60.
- [11] GAWLICK, D., GRAY, J., LIMURA, W., AND OBERMARCK, R. Method and Apparatus for Logging Journal Data Using a Log Write Ahead Data Set. U.S. Patent 4507751 issued to IBM, March 1985.
- [12] GROCHOWSKI, E. G., AND HOYT, R. F. Future Trends in Hard Disk Drives. *IEEE Transactions on Magnetics* 32, 3 (May 1996).
- [13] HAGMANN, R. Reimplementing the Cedar File System Using Logging and Group Commit. In *Proc. of the 11th ACM Symposium on Operating Systems Principles* (October 1987), pp. 155–162.
- [14] HAGMANN, R. Low Latency Logging. Tech. Rep. CSL-91-1, Xerox Corporation, Palo Alto, CA, February 1991.
- [15] HARTMAN, J., AND OUSTERHOUT, J. The Zebra Striped Network File System. *ACM Transactions on Computer Systems* (Aug. 1995).
- [16] HITZ, D., LAU, J., AND MALCOLM, M. File System Design for an NFS File Server Appliance. In *Proc. of the 1994 Winter USENIX* (January 1994).
- [17] HITZ, D., LAU, J., AND MALCOLM, M. File System Design for an NFS File Server Appliance. Tech. Rep. 3002, Network Appliance, March 1995.
- [18] KOTZ, D., TOH, S., AND RADHAKRISHNAN, S. A Detailed Simulation Model of the HP 97560 Disk Drive. Tech. Rep. PCS-TR91-220, Dept. of Computer Science, Dartmouth College, July 1994.
- [19] LAMB, C., LANDIS, G., ORENSTEIN, J., AND WEINREB, D. The ObjectStore Database System. *Communications of the ACM* 34, 10 (October 1991), 50–63.
- [20] LOWELL, D. E., AND CHEN, P. M. Free Transactions with Rio Vista. In *Proc. of the 16th ACM Symposium on Operating Systems Principles* (October 1997), pp. 92–101.
- [21] MALAKAPALLI, C. Personal Communication, Seagate Technology, Inc., July 1998.
- [22] MASHEY, J. R. Big Data and the Next Wave of InfraStress. Computer Science Division Seminar, University of California, Berkeley, October 1997.
- [23] MATTHEWS, J. N., ROSELLI, D. S., COSTELLO, A. M., WANG, R. Y., AND ANDERSON, T. E. Improving the Performance of Log-Structured File Systems with Adaptive Methods. In *Proc. of the 16th ACM Symposium on Operating Systems Principles* (October 1997), pp. 238–251.
- [24] MENON, J., ROCHE, J., AND KASSON, J. Floating parity and data disk arrays. *Journal of Parallel and Distributed Computing* 17, 1 and 2 (January/February 1993), 129–139.
- [25] O'TOOLE, J., AND SHRIRA, L. Opportunistic Log: Efficient Installation Reads in a Reliable Storage Server. In *Proc. of the First Symposium on Operating Systems Design and Implementation* (November 1994), pp. 39–48.
- [26] PATTERSON, R. H., GIBSON, G. A., GINTING, E., STODOLSKY, D., AND ZELENKA, J. Informed Prefetching and Caching. In *Proceedings of the ACM Fifteenth Symposium on Operating Systems Principles* (December 1995).
- [27] ROSENBLUM, M., AND OUSTERHOUT, J. The Design and Implementation of a Log-Structured File System. In *Proc. of the 13th Symposium on Operating Systems Principles* (Oct. 1991), pp. 1–15.
- [28] RUEMLER, C., AND WILKES, J. An Introduction to Disk Drive Modeling. *IEEE Computer* 27, 3 (March 1994), 17–28.
- [29] SATYANARAYANAN, M., MASHBURN, H. H., KUMAR, P., STERE, D. C., AND KISTLER, J. J. Lightweight Recoverable Virtual Memory. In *Proc. of the 14th ACM Symposium on Operating Systems Principles* (December 1993), pp. 146–160.
- [30] SELTZER, M., BOSTIC, K., MCKUSICK, M., AND STAELIN, C. An Implementation of a Log-Structured File System for UNIX. In *Proc. of the 1993 Winter USENIX* (Jan. 1993), pp. 307–326.
- [31] SELTZER, M., SMITH, K., BALAKRISHNAN, H., CHANG, J., MCMAINS, S., AND PADMANABHAN, V. File System Logging Versus Clustering: A Performance Comparison. In *Proc. of the 1995 Winter USENIX* (Jan. 1995).
- [32] STODOLSKY, D. Personal Communication, Quantum Corp., July 1998.
- [33] SWEENEY, A., DOUCETTE, D., HU, W., ANDERSON, C., NISHIMOTO, M., AND PECK, G. Scalability in the XFS File System. In *Proc. of the 1996 Winter USENIX* (January 1996), pp. 1–14.
- [34] TRANSACTION PROCESSING PERFORMANCE COUNCIL. *TPC Benchmark B Standard Specification*. Waterside Associates, Fremont, CA, Aug. 1990.
- [35] TRANSACTION PROCESSING PERFORMANCE COUNCIL. *TPC Benchmark C Standard Specification*. Waterside Associates, Fremont, CA, August 1996.

- [36] WILKES, J., GOLDING, R., STAELIN, C., AND SULLIVAN, T. The HP AutoRAID Hierarchical Storage System. In *Proc. of the 15th ACM Symposium on Operating Systems Principles* (December 1995), pp. 96–108.