

File System Performance

s5fs, ffs, lfs, etc.

CSCI3753

1

This lecture discussed file system performance, using three example file systems to illustrate how on-disk structure can lead to different performance behavior: s5fs, ffs, and lfs.

s5fs

- **First UNIX filesystem**
- **Simple (~1000 lines of C)**
- **Reasonable performance for the time**

CSCI3753

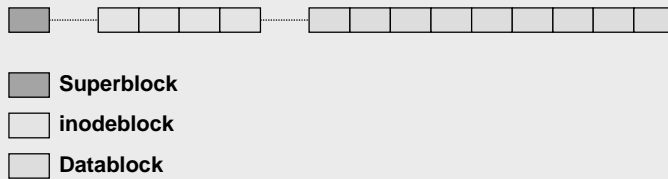
2

s5fs, or the System V filesystem, is the first UNIX filesystem. UNIX was when initially developed, it was designed for simplicity. It exhibited reasonable performance characteristics. It's very slow in comparison to modern filesystems, but can be implemented in roughly a thousand lines of code, which is much much less than complex systems such as Veritas, which are in the tens of thousands.

. As the rest of the world moved to more complex systems, s5fs was left behind. At the time, it was very low in comparison to modern systems like Veritas, which are in the tens of thousands.

s5fsStructure

- Filesystem is a linear array of blocks
- Superblock
- inode blocks
- Datablocks



CSCI3753

3

s5fs makes a single simple abstraction of the disk, that it's a simple linear array of blocks. The first non-boot block is the superblock, which is followed by some number of inode blocks, then the remaining blocks are datablocks.

s5fs Superblock

- Magic number
- Number of inodes
- Free block list head
- Free inode list head
- Last modification time
- etc.

CSCI3753

4

The s5fs superblock contains filesystem-wide metadata. For example, how many inodes the filesystem has, the index of the first free inode, the index of the first free data block, and the last time the filesystem was modified. This data is critical to be able to mount the filesystem. Without it, there's no surefire way to tell some of these things.

```

struct filsys
{
    ushort s_isize;           /* size in blocks of i-list */
    daddr_t s_fsize;         /* size in blocks of volume */
    short s_nfree;           /* number of addresses in s_free */
    daddr_t s_free[NICFREE]; /* free block list */
    short s_ninode;          /* number of i-nodes in s_inode */
    s5ino_t s_inode[NICINOD]; /* free i-node list */
    char s_flock;            /* lock during free manip */
    char s_ilock;            /* lock during i-list manip */
    char s_fmod;             /* super block modified flag */
    char s_ronly;            /* mounted read-only flag */
    time_t s_time;           /* last super block update */
    short s_dinfo[4];        /* device information */
    daddr_t s_tfree;         /* total free blocks */
    s5ino_t s_tinode;        /* total free inodes */
    char s_fname[6];         /* file system name */
    char s_fpack[6];         /* file system pack name */
    long s_fill[13];         /* makes sizeof filsys be 512 */
    long s_magic;            /* magic number */
    long s_type;             /* type of new file system */
};

```

CSCI3753

5

Here's a C struct that represents the superblock of a 5fs. Note the locks for inode and free list manipulation. This struct could be used to read the superblock into memory, then casting it to a struct

filsys. The presence of populated by

s5fsinode

- **Filesize**
- **Lastmodification**
- **Owner,group,permissions**
- **Datablocks**

CSCI3753

6

Filespecificdataisstoredinans5fsinode,suchasthesize
timeoflastmodification.Italsocontainstheblockmap,thed
storeswhatblocksbelongtothefile.

ofthefileandthe
atastructurethat

s5fsDataBlocks

- **Directblocks**
- **Singleindirectblock**
- **Doubleindirectblock**
- **Tripleindirectblock**

CSCI3753

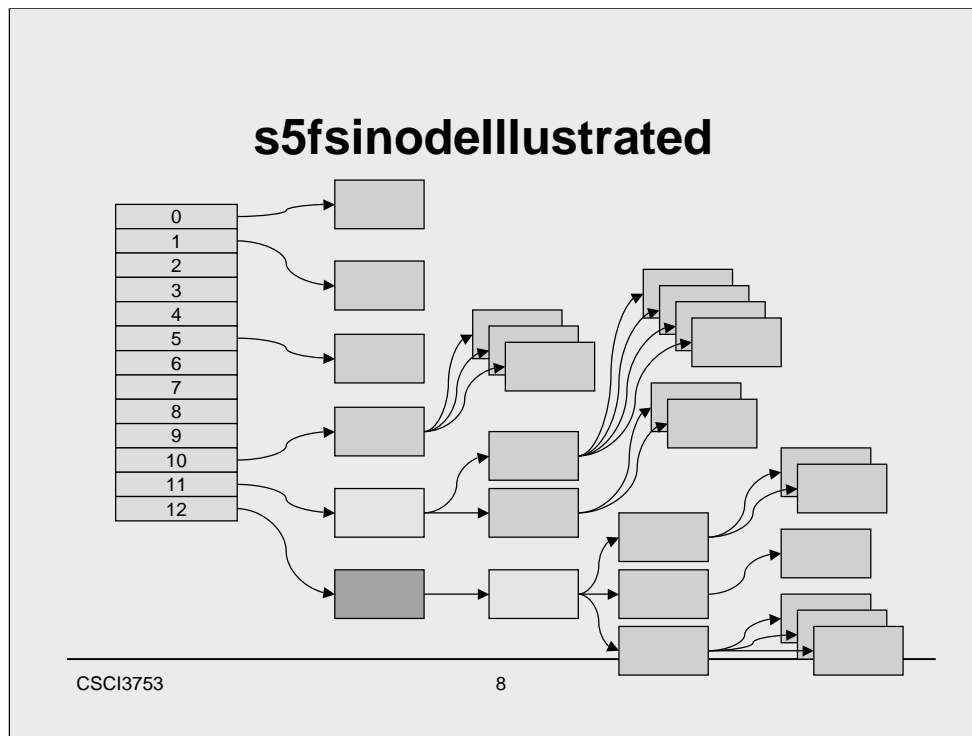
7

s5fs has three types of data blocks. First are direct blocks, which are single indirect blocks, which store block numbers for direct way, instead of having a single large table of direct blocks in on the disk can be allocated for storing the disk map of large files. Next are double indirect blocks, which store the block numbers of single indirect blocks. Last are triple indirect blocks, which store the block numbers of double indirect blocks.

The hierarchy of single, double, and triple indirect blocks allows for large files while keeping the inode disk map fairly small. It does mean that file might require multiple block reads, but the worst case is 5 reads (inode, triple, double, single, data).

ich stored data. Next are blocks. This the inode, blocks are files. Next are indirect blocks. fdouble indirect

ws for large files later blocks in a reads (inode,



Here's a pictorial representation of the inode block map. An s5fs total of 13 block entries in its block map. The first ten are direct blocks. The eleventh is a single indirect block. The twelfth is a double indirect block, and the thirteenth is a triple indirect block.

sinode has a rect blocks. The irect block, and

Block numbers are ordered in indirect blocks. The eleventh block is the first block number in the single indirect block, while the twelfth is the second.

k of a file is welf this the

s5fs supports sparse files. The file size provides the kernel with the last valid block in a file. However, if large (blocks sized all zeros), then having a large number of zero-filled blocks is a waste. Instead, the file system only allocates blocks that are needed. If a process then seeks far into it before writing a single byte, of course a block must be allocated for that byte, but blocks need not be allocated for all of the space that was sought over.

th knowledge of reas) of the file are waste. Instead, the opens a new file, block must be lof the space that

Indirect blocks are wasted disk space, as they're not actually storing any file data. A degenerate series of files might require up to four times the expected storage. This is a pretty degenerate case, though.

toring any file stheir expected

Maximum File Size

- 4 byte block index
- 1024 byte blocks
- 256 block indices per block
- $10 + 256 + 256^2 + 256^3$
- 16GB files (actually, 2GB, but they)

CSCI3753

9

Given a few pieces of data, it's clear why there are triple indirect blocks. Assuming a 4 byte block index and a 1 K block, each block can contain 256 block entries. Doing a little arithmetic, it's easy to determine that a triple indirect block allows for files that are about 16GB in length. Only having a double indirect block would make the maximum file size roughly 64MB. Of course, although it's possible to have a 16GB file, a 32 bit architecture could never reference the latter 14GB of the file. (Quick quiz: Why is a 32 bit architecture limited to 2GB instead of 4GB? Hint: what's the type of an offset pointer?)

s5fsFreeList

- Linkedlistoffreeblocks
- Linkedlistoffreeinodes

CSCI3753

10

Freeinodesandblocksarekepttrackofbythefreelist. Theheadofthelistis onthesuperblock. Freeinodesandblockscanthenreferencelaterones,forming alinkedlistofsorts. Whenanewblockorinodeisfreed, itcanbeputonthe headofthelist(you **really** don'twanttowalktotheendofthelist);whenoneis needed, itcanberemoved.

s5fsDirectoryStructure

Filename(14bytes)	inode
passwd	34
shadow	12
inetd.conf	123
hosts	753
DataBlock	

CSCI3753

11

Directories in s5fs have a fixed structure. They have a data block file, but have a special bit set in the inode to denote that the data block has a specific format, broken up into directory entries. A directory entry has a 14-byte filename and an inode index. 14 bytes allow a reasonable length without wasting a great deal of space. Nowadays, a 14-byte length limit is unheard of, but back then, it was OK, and a larger significant waste or a more complex structure. The size of a directory file (in the inode) allows the kernel to know how many entries are in the directory.

Just like any other file, a directory is a file. A directory file (in the inode) allows the kernel to know how many entries are in the directory.

s5fs is great, but...

- Disjoint data and meta information
- No locality in data
- Ignores disk characteristics
- Small block size
- Superblock fragility

CSCI3753

12

s5fs was great for its time and is very simple, but has a few problems that came to light as UNIX grew in popularity. As UNIX moved from a research to a production system with the BSD project, s5fs was no longer adequate.

fsignificant transitioned swas no longer

The foremost problem with s5fs is its performance. Metadata (inode) and data are disjoint and distant on the disk. All of the inodes are next to one another, while all of the data blocks are next to one another. The structure of s5fs ignores the structure and behavior of a disk.

des) and data to one another, ure of s5fs ignores

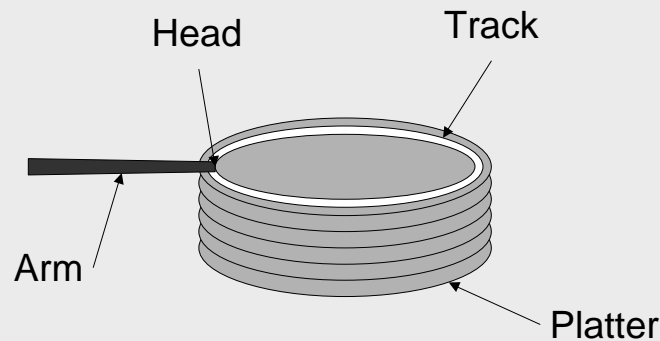
Additionally, s5fs has a very small block size. This means that sizes start requiring indirect blocks. Since there's no promise that a block is close to its data blocks, this can require a lot of time across the disk.

files over 10K in hat an indirect espent seeking

One more major problem is s5fs' fragility in the face of crashes. A single copy of the superblock. If it's lost due to a crash, the way of getting your file system back. Some sort of redundancy would be nice, such that a single crash is unlikely to destroy your entire file system.

.Notably, there's ere's really no good uld be nice, system.

Disk Structure: A Crash Course



CSCI3753

13

Here's a simple pictorial representation of a disk. The disk has each of which has its own disk head. The platters are separated which is the area of the disk that an unmoving head covers as the

as the platters, into tracks, as the disk spins.

Disks are really slow in comparison to main memory. This latency fact that physical parts have to move. Disk latency can be broken into rotational latency and seek latency. The latter is much worse than the former. Rotational latency refers to the time it takes a block to be rotated under an unmoving disk head; it's determined by how quickly the disk spins. Seek latency is how long it takes the head to move from track to track. For a 5,000 RPM disk (slow by today's standards), the rotational latency is on the order of 12 milliseconds. A corresponding seek latency might be 10 milliseconds. While these might seem like short periods of time, they're not in terms of modern CPUs. A 500 MHz processor could execute 100,000 instructions during the rotational latency, or 5 million during the seek latency.

This latency is due to the fact that physical parts have to move. Disk latency can be broken into two parts, rotational latency and seek latency. The latter is much worse than the former. Rotational latency refers to the time it takes a block to be rotated under an unmoving disk head; it's determined by how quickly the disk spins. Seek latency is how long it takes the head to move from track to track. For a 5,000 RPM disk (slow by today's standards), the rotational latency is on the order of 12 milliseconds. A corresponding seek latency might be 10 milliseconds. While these might seem like short periods of time, they're not in terms of modern CPUs. A 500 MHz processor could execute 100,000 instructions during the rotational latency, or 5 million during the seek latency.

MoreDiskStructure

- **Cylinder:** set of tracks that have same location on different platters

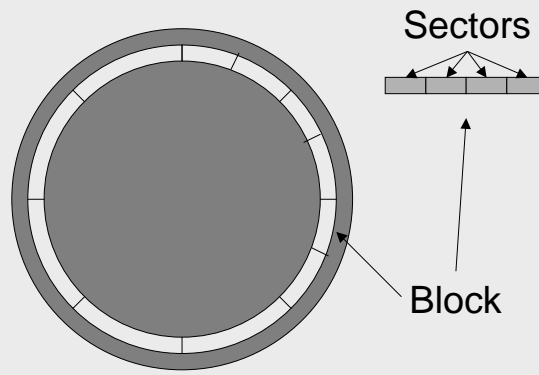
CSCI3753

14

Generally, all of the heads move as one unit. Therefore, an additional abstraction is useful, that of a cylinder. A cylinder is the set of the heads. You can think of this as the same track of each of the platters. This distinction of a cylinder is useful because if the head is on a particular track, then any block in the cylinder can be accessed within the rotational latency of the disk.

tion
of tracks under all
the different
heads don't seek,
on all latency of

EvenMoreDiskStructure



CSCI3753

15

The smallest unit of reading and writing on a disk is the sector. Sectors are quite small (512 or 1 Kbytes), so file systems often make an abstraction on top of the sector known as a block. A block is the smallest unit of reading and writing for the file system.

. Usually, sectors are an abstraction on top of reading and writing.

Blocks and Sectors

- A sector is a unit of disk allocation
- A block is a unit of filesystem allocation
- Block size is a design choice: fragmentation

CSCI3753

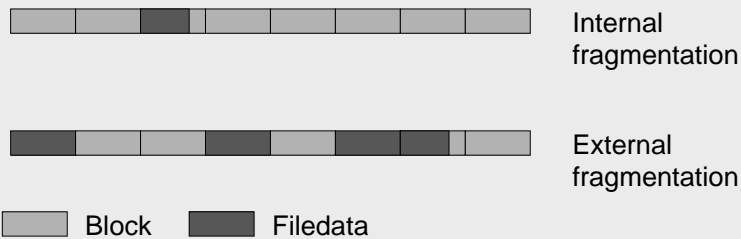
16

Choosing the block size properly can lead to markedly increased performance for a filesystem. However, there's an inherent tradeoff present in choosing the block size, one tradeoff between different forms of fragmentation.

performance
in deciding the
ion.

Fragmentation

- Internal Fragmentation (wasted space)
- External Fragmentation (low locality)



CSCI3753

17

Here's a picture that demonstrates why block size presents a tradeoff between forms of fragmentation. Internal fragmentation is when not the entire block is used; space is wasted. External fragmentation is when a single piece of data is spread out across distant blocks. This results in lower throughput as the disk has to seek or wait for platters to rotate before it can service requests.

de off between
ntire block is
iece of data is
ut as the disk has
uests.

Tradeoffs

- **Large block size**
 - Low external fragmentation
 - High internal fragmentation
- **Small block size**
 - High external fragmentation
 - Low internal fragmentation

CSCI3753

18

A large block size reduces external fragmentation but increases internal fragmentation. A small block size does the opposite. Hopefully, there's some perfect middle value that results in little latency and not too much waste.

FFS

- **Filesystem organized in cylinder groups**
- **Blocks can have fragments**
- **Superblock redundancy**
- **Data locality**
- **Many of these are now assumed techniques**

CSCI3753

19

The Berkeley Fast File System (known as `ufs` or `ffs`) was developed in the 70s to deal with these limitations. It kept the basic abstraction of `fs` (inodes, data blocks, etc.) while using a completely different on-disk structure.

The most critical abstraction used in `ffs` is that of a cylinder group, which is a small number of adjacent cylinders. `ffs` divides the disk up into cylinder groups, each of which has its own inodes and data blocks. Hopefully, the data blocks in the cylinder group are associated with its inodes, therefore requiring few seeks. Also, blocks of a file are usually close to one another.

`FFS` deals with the issue of fragmentation by dividing blocks up into fragments. Usually, entire blocks are allocated, but in special circumstances fragments can be used, in order to decrease internal fragmentation.

In addition, `ffs` has multiple copies of the superblock, to prevent its loss resulting from a crash.

Many of the techniques used in `ffs` are generally assumed when constructing modern file systems.

FFShasblocksandfragments

- Everyblockiscomprisedof2ⁿ fragments
- Filesareallocatedintermsofblocks, *except for thelastblock*
- Lastpartoffilecanbeoneormorefragments
- Bestofbothworlds

CSCI3753

20

Everyblockinffsisdividedintoasmallnumberoffragments(2,4,8,etc.). Filesareallocatedintermsofblocks,exceptforitsend,which canbeallocated intermsofafragments. Thisway,smallfileswithinopportune sizes(blocksize +1bytes)won'twasteentireblocks,butthefileisn'tbroken upintopiecesas smallasfragments. ThisallowsFFStokeepexternalfragmentation downwitha largeblocksize,butalsokeepinternalfragmentationdownatt heendofafile withfragments. Ofcourse,asparsedfilestillwastesdiskspace basedonthe blocksize,butthisislesssofaconcernthantheendofafile ;usersaregenerally angrieriftheir1Kfileactuallytakes4Kthaniftheir200Kfi lewithonlyafew bytesinitactuallytakes200K.

s5fssuperblock

- **s5fshasonesuperblock**
- **Containsfilesystemmeta -meta-data**
 - **blocksize**
 - **numberofinodes**
 - **numberofblocks**
- **Losethesuperblock,goodluckrecoveringthe filesystem!**

CSCI3753

21

The s5fssuperblock is a big problem due to its fragility. As it contains the file system meta -metadata that exists nowhere else, reconstructing the filesystem without it can be impossible. Having only one copy of the superb lock (in a very conspicuous place) is asking for trouble.

Superblock redundancy

- **FFS keeps redundant copies of the superblock**
- **Different cylinder groups, platters, rotational positions, etc.**
- **No localized physical failures should be able to destroy the superblock**
- **Superblock consistency can be resolved at crash recovery**

CSCI3753

22

ffs solves the superblock fragility problem by having multiple copies of it scattered across the disk. They're all on different tracks, platters, and rotational positions, so that if a single hardware failure were to destroy them all, you'd be pretty sure you'd lost all of your data as well. If the system crashes in the midst of updating the superblock, inconsistency between them can be solved when the system reboots, as you're probably going to run fsck.

s5fs'BadLocality

- s5fsfreelistbecomesrandomizedafteruse(two successivelyallocatedblockscanbedistant)
- s5fsinodesanddataseparate

CSCI3753

23

Thefreelistsofs5fsareasimplemechanism,butleadtohorribleperformance. They'reperfectlyorderedwhenthefilesystemisfirstcreated, butquickly become randomized. Two adjacent blocks in the freelist can be very distant on the disk. Understandably, this is a problem. Chances are if you access block 1 of a file, you'll soon access block 2. Placing them close together on disks saves costly seeks.

bleperformance. butquickly erydistanton accessblock1of ondissaves

Additionally, theseparationofinodesandtheirdataisareal problem. Opening afileandreadingthefirstbytecaninvolve two seeks across theentiredisk.

problem. Opening heentiredisk.

TheFFSSolution

- **FFSdata blocks are allocated on a cylinder group basis, not a global basis**
- **FFS has inodes and data in each cylinder group**
- **FFS tries to keep related files (same directory) together**

CSCI3753

24

ffs solves the data locality problem through its use of cylinder groups. Instead of maintaining a global freelist, the cylinder group has a bit field of blocks and files in the same logical file blocks fragments in use. ffs tries to keep most of the blocks of most files in the same logical file blocks cylinder group past their inode. This way, not only are adjacent logical file blocks actually nearby physical blocks, but the inode is close as well.

In addition, ffs tries to keep all of the files in a given directory in the same cylinder group.

DataBlockAllocation

- If next rotationally available block is free, use it
- Else if a block on cylinder is free, use it
- Else if a block in cylinder group is free, use it
- Else if a block in new cylinder group is free, use it
 - New cylinder group selected through quadratic hashing
- Else try searching all cylinder groups
- Exceptions for when file reaches certain size (don't want to use all of one cylinder group)

Here's the basic data block allocation algorithm for FFs.

RelatedFileLocality

- **ls** isn't a problem...
- **ls -li**s...

```
itsydev /tmp_mnt/home/cia/levis/public_html -10-> ls -li
total 1016
-rw-r--r--  1 levis  student      4876 Apr 18 20:02 first.txt
drwxr-xr-x  2 levis  student      4096 Aug 16 1999 img
-rw-r--r--  1 levis  student      3310 Sep  5 22:39 index.html
-rw-r--r--  1 levis  student        393 Aug 17 1999 links.html
drwxr-xr-x  2 levis  student      4096 Sep  5 12:15 os
-rw-r--r--  1 levis  disk        994571 Aug 30 22:47 paper.pdf
-rw-r--r--  1 levis  disk        9987 Nov 30 1999 plans.txt
-rw-r--r--  1 levis  student      1974 Aug 17 1999 techlinks.html
```

CSCI3753

26

Recalling the 5fs directory layout, **ls** doesn't pose a large problem, as all of the directory entries are in the data block of the directory file. More complex operations, however, do cause problems. For example, the command **ls -li** lists a lot of metadata on each file, all of which resides in the inodes. This requires going to the inode of each of the files, reading it into memory, and using some of the data. Although all of the inodes are in one section of the disk, it can still require a good number of costly seeks or rotation to read all of them.

RelatedFileLocality,cont.

- **The metadata resides on inodes**
- **Want inodes in a directory to be close**
- **Want data to be close to inodes**
- **Try to make all files in a directory reside on same cylinder group**

The ffs solution is to try to place all of the files in a directory in the same cylinder group. This way, only one long seek is needed.

or in a given cylinder

Results

- **Correct tuning can result in 900% increase in data throughput over 5fs**
- **Space loss to internal fragmentation roughly equivalent**

CSCI3753

28

ffs provides impressive improvement over 5fs. Notably, its throughput is 900% higher (10 to 1), while its fragment/block division keeps its internal fragmentation roughly equivalent to 5fs.

Steps Beyond FFS

- **Better block allocation (Smith and Seltzer)**
- **Larger block allocation - extents (Veritas, Episode, etc.)**
- **Better crash recovery - journaling, lfs**

CSCI3753

29

Many steps have been taken since FFS. For example, although FFS tries to keep adjacent logical file blocks close, they're often not actually adjacent, which is what's usually wanted. Work has been done to write new block allocation algorithms that keep the blocks closer. Enterprise level file systems (such as Veritas) use a different approach, that of extents. Extents are basically groups of contiguous blocks. Instead of allocating a single block at a time, the file system allocates an extent, part of which might be later freed.

In addition to disk throughput, a major concern for file systems is crash recovery. While FFS prevents the situation of crash recovery being impossible, it makes no promises for its rapidity. In modern computing environments, the time it can take to perform a crash recovery is critical. Take, for example, the recent CSEL crash recovery, which took several days. The basic issue in crash recovery is inconsistent state. For example, if the system crashes during a file write, the block may have been taken off the free list but not yet added to the file. Since we're generally unable to know to which file the block was supposed to be added, we could just return it to the free list. But first, we need to examine every inode to see if it belongs to one.

One way around this is journaling. In essence, one logs what operations are going to be performed just before it is. When recovering from a crash, all one has to do is examine the journaling log and compare it with the disk state. If the state matches the entry in the log, it's completed. If it isn't, then it didn't. The rest of the system is assured to be in a consistent state.

lfs:LogStructuredFileSystem

- Only one on-disk structure: the log
- Filesystem updates append to the log
- Metadata updates append to the log too



lfs, or the log structured filesystem, is an entirely different filesystem structure, and has very different performance characteristics. Instead of dividing the disk up into inode areas, data areas, and cylinder groups, the entire disk forms one large log. All operations result in an analog entry being appended to the log. For example, when a file is modified, the new version is written in its entirety onto the log.

proach to file
tics. Instead of
groups, the entire
y being appended
is written in

At first glance, this structure is very appealing. Crash recovery is merely a matter of looking at the end of the log. Writes are fast, as they are contiguous. Similarly, although a read might require a long seek, the file is entirely contiguous on the disk, so subsequent reads will

be easy, as one
they're always
to the file, the
be quick.

lfsBenefits

- **Fastwrites**
- **Filesalwayscontiguous**
- **Easycrashrecovery**

lfs Drawbacks

- Reads can be very slow
 - needs a big buffer cache
- Fragmentation of log
- The answer: *logcleaning*

CSCI3753

32

Of course, there are always problems. Although reading an entire file requires only one seek, rarely does a system just read in a single file at a time. Instead, the files are spread all across the disk, and so reads are generally fairly slow. A big buffer cache can solve this problem, though, as then the file can just be kept in memory.

In addition, there's the problem of log fragmentation. As files are updated and rewritten, then their old versions can be reclaimed as space for the log. Understandably, this can lead to external fragmentation, in which there's slots of free space, but none of it is very large. This prevents lfs from having good performance, as it can no longer just continue to simply append.

Log fragmentation is dealt with by a technique known as logcleaning, which basically reorders and restructures the log when the system is idle. This way, in use portions can be compacted, and large contiguous free areas can be made available.