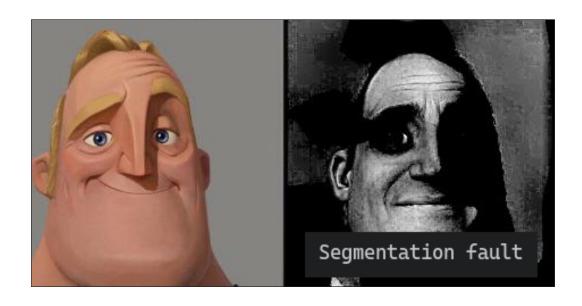
DM510 Operating systems Assignment 3 File system



Frederik List frlis21@student.sdu.dk

August 23, 2024

Introduction

When booted, one of the first things operating systems do is mount a block device, which probably has a file system on it. One could hack the kernel to support a new file system, but that's hard. FUSE lets us create a file system in userspace. In this report, I outline how I created the terrible file system (abbreviated TFS) with FUSE.

Design

TFS is very barebones. It can read and write (large) files and directories... and that's it. TFS is neither journalling nor log structured. The focus was on making TFS simple, which one might say TFS fails at, for reasons you will discover.

Layout

The layout of a TFS partition is as ordered:

- TFS Superblock, containing
 - the number of total blocks and nodes.
 - pointers to the heads of the free node and block linked lists.
- Nodes, where each node contains either
 - node metadata, such as node name, blocks, access time, etc.
 - a pointer to the next free node in the free node linked list.
- Blocks, where each block contains either
 - file data or directory entries.
 - indirect block pointers.
 - a single pointer to the next block in the free block linked list.

TFS allocates space for nodes near the beginning of the partition at format time.

Allocation scheme

TFS uses a combined index allocation scheme (exactly the one shown in figure 14.8 of the course textbook). Early in development, TFS used a file allocation table (FAT) for its simplicity, but the author suffers from occasional slight mental degradation and premature optimization syndrome, so he chose to switch to the more complicated and slightly faster indexed allocation scheme. (At least now he can open pictures of his dog a few milliseconds faster on a TFS mount.)

Random access is much faster for the indexed allocation scheme than for a FAT. This is because instead of having to go through the entire file to find a specific block, a formula can be used to calculate the block indices from the file offset.

If the required data cannot be found in the direct blocks, then for a file block offset x (starting at 0 regardless of the number of direct blocks), the indirect block level can be calculated with

$$d = \left\lfloor \log_b \left(x + 1 - \left\lfloor \frac{x+1}{b} \right\rfloor \right) \right\rfloor$$

where b = the number of block pointers that can fit in a block $(\frac{\text{block size}}{\text{pointer size}})$. Block level indices can then be calculated with

$$i_k = \left| \frac{x - b\frac{b^d - 1}{b - 1}}{b^k} \right| \mod b$$

where k = index level.

If the block size (and block pointer size, and hence max block pointers) is a power of 2, then bit twiddling can be used for fast logarithm and exponent calculation.

Free nodes and blocks

Free nodes contain a single pointer to the next free node, thus forming a linked list of free nodes.

Similarly, free blocks contain a single pointer to the next free block.

Implementation

TFS's code is somewhat littered with comments explaining most of the esoteric bits, this section serves to supplement those comments.

Backing store

TFS's backing store is in the form of a memory-mapped binary file. A TFS image may be created with the provided mktfs utility, much like Linux's mke2fs(8). As with mke2fs, a file needs to be allocated first, probably with the fallocate(1) utility. In the accompanying video a 16 MiB file is allocated with fallocate -l16M img.tfs, after which mktfs img.tfs is invoked.

Instead of reading and writing to the file with, for example, pread(2), TFS memory maps the entire backing store. The author feels some might not think this to be in the spirit of the project, but argues that without memory mapping, the code would just be polluted with some extra mallocs and preads. Memory mapping allows everyone to enjoy the internals of TFS without more cruft than is required!

Looking up entries

Entry lookup is implemented with the standard C library hash table (see hsearch(3)). Keys are paths and values are nodes. This makes entry lookup very fast, but the memory usage may be a concern for large file systems, as the hash table size is simply initialized to the number of nodes in the file system.

Reading and writing blocks

An iterator-like structure named the "block cursor" allows for random-access seeking and sequential iteration.

Reading and writing are much the same. A block cursor is simply set to the position to start from and iterated until no more bytes can or should be written or read.

Allocating and freeing blocks and nodes

Nodes are allocated and freed simply by popping from or pushing to the free node linked list.

Blocks are a little more tricky. This is because in addition to blocks storing actual data, index blocks also need to be allocated for sufficiently large files. The cursor iterator

presented earlier is able to take a callback "hook" which will be called for each block it iterates through, including index blocks which are transparent to iterator consumers.

Allocation can be easily implemented using this callback. A cursor is set to the end block of a file and iterated through with a function that pops free blocks from the free block linked list.

Freeing blocks is not so easy. Ideally, a cursor would be set to the end block of a file and iterated through *backwards* while freeing blocks along the way. Instead, a hack was used that makes many assumptions about how the file system works: a cursor is set to the last required block of a file and iterated, adding any blocks along the way to a free block buffer. After an iteration is done, any blocks in the free block buffer are freed and the free block buffer is cleared. This is not so good, as free block and index block pointers are written to the same locations, and as the ancient proverb goes, "don't shit where you eat".

Crash recovery

Crash recovery in TFS is nonexistent; when the kernel decides to page out to the backing store is indeterminate. Disregarding memory mapping, it could be implemented with journalling: A portion of the TFS partition would be alloted to journalling, where dirty blocks (and their checksums) would periodically be sequentially written. After a journal is written, data can be moved to where it is supposed to go with slower random access without fear of losing data.

In case of a crash, the journal can be read and the rearranging of data can be resumed. Checksums are used to ensure only data that was successfully written to the journal are written to the disk.

Journalling can handle any number of crashes while recovering from a crash; the only data needed to recover is written to the nonvolatile disk. Of course, some small amount of data loss is always possible when a crash occurs while writing dirty blocks, the only remedy for which is to write often.

Tests

Mosts of the tests performed in the video are just the author trying to abuse TFS as hard as he can.

Test Time 00:00 Formatting File creation and deletion 00:10Directory creation and deletion 00:24 00:47Access and modification time using touch and cat Large files using seq, head, and tail 01:00 01:12 Directories with many entries 01:29 good boi File persistence 01:39

Figure 1: Tests with video timestamps

In more rigorous testing under development, an ad-hoc analysis concluded that TFS does not produce dead blocks or leak memory.

TFS breaks down a little bit when there is not much space left. Not all cases concerning lack of storage are well tested.

Conclusion

Though the author would not trust it for anything remotely important, it works. TFS is slow and probably riddled with nasty bugs, but it gets the job done.

Appendix A

#ifndef TFS_H

2

48

Source code

Listing A.1: tfs.h

```
#define TFS_H
    #include <stdlib.h>
   #include <sys/mman.h>
   typedef off_t blkoff_t;
   typedef off_t nodoff_t;
   // Block size must be a power of 2 for bit twiddlings
10
   #define BLOCK_SIZE 4096
#define BLOCKS_PER_NODE 4
13
   #define DIRECT_BLOCKS 12
   #define ILEVELS 3
14
15 #define NAME_LIMIT 64
   #define BLOCK_MAX_CHILDREN (BLOCK_SIZE / sizeof(nodoff_t))
   #define BLOCK_MAX_POINTERS (BLOCK_SIZE / sizeof(blkoff_t))
17
   #define END_BLOCKS -1
   #define END_NODES -1
19
20
   #define MIN(a, b) ((a) < (b) ? (a) : (b))
21
22 #define MAX(a, b) ((a) > (b) ? (a) : (b))
   // Absolute node size
24
   #define NODE_SIZE(node) ((node)->mode & S_IFDIR ? (node)->nlink * sizeof(nodoff_t) : (node)->size)
25
   // Number required blocks for data (not including indirect pointer blocks)
   #define NODE_NRBLOCKS(node) ((BLOCK_SIZE + NODE_SIZE(node) - 1) / BLOCK_SIZE)
27
28
29
30
    * TFS superblock:
    * Metadata needed to calculate everything.
31
32
    struct tfs_header {
            blkoff_t nblocks, free_block_head;
34
            nodoff_t nnodes, free_node_head;
    };
36
37
38
    * The TFS node, as it is represented in the image file.
39
40
    struct tfs_node {
41
            union {
42
                    struct {
43
                            mode_t mode;
44
                            char name[NAME_LIMIT];
                            // Direct blocks
46
                            blkoff_t blocks[DIRECT_BLOCKS];
                            // Indirect blocks
```

```
blkoff_t iblocks[ILEVELS];
49
 50
                              // Number of allocated blocks
                              fsblkcnt_t nblocks;
51
                              // Number of links of directory, file size otherwise
                              union {
53
                                      off_t size;
54
                                      nlink_t nlink;
55
                              };
56
                              struct timespec atim, mtim;
57
                      };
58
59
                      // Used for free node linked list.
                      nodoff_t next;
60
             };
61
     };
63
64
      * Collection of pointers to useful places and other info nice to have.
65
66
 67
     struct tfs_info {
             blkoff_t nblocks;
68
             nodoff_t nnodes;
             blkoff_t *free_block_head;
70
71
             nodoff_t *free_node_head;
             struct tfs_node *nodes;
72
73
             char (*data)[BLOCK_SIZE];
74
              /* no touchy */
             void *base;
75
             off_t filesize;
76
    };
77
78
79
     * Open a file as a TFS image.
80
 81
     int tfs_open(const char *filename);
82
83
 84
     * Open and initialize a TFS image.
85
 86
     int tfs_load(const char *filename);
87
88
 89
     * Format a TFS image.
90
91
     void tfs format():
92
93
94
     * Calculate pointers and other useful things.
95
96
     void tfs_init();
97
98
99
      * Write back any "queued" changes.
100
101
     int tfs_destroy();
102
103
104
      * (De)allocate necessary blocks for a node.
105
106
      * Must be called after changing the size or nlink of a node.
107
108
109
     int tfs_node_trim(struct tfs_node *node);
111
     * Get a node from the hash table given the path.
112
113
    struct tfs_node *get_node(const char *path);
114
115
     /**
116
```

```
* Get the directory of a node from the hash table given the path.
117
118
     struct tfs_node *get_directory(const char *path);
119
121
     * Collect children of a node into one contiguous array.
122
123
      * The array must be freed when you are done with it.
124
     struct tfs_node **tfs_node_children(struct tfs_node *node);
126
127
128
     * Read node data.
129
130
     int tfs_node_read(struct tfs_node *node, char *buf, size_t size, off_t offset);
131
132
133
     * Write node data.
134
135
     int tfs_node_write(struct tfs_node *node, const char *buf, size_t size, off_t offset);
136
137
138
     * Add a node.
139
140
     int tfs_add_node(const char *path, mode_t mode);
141
142
143
     * Remove a node.
144
145
     int tfs_remove_node(const char *path);
146
147
     #endif // TFS_H
148
                                        Listing A.2: fuse_tfs.c
     #define FUSE_USE_VERSION 29
 1
 2
    #include <assert.h>
 3
 4 #include <errno.h>
    #include <fuse.h>
 5
    #include <libgen.h>
    #include <math.h>
 8 #include <search.h> // hsearch
 9 #include <stdio.h>
    #include <stdlib.h>
 10
    #include <string.h>
 11
    #include <time.h>
 12
13
    #include "tfs.h"
 14
15
 16
     static int fuse_tfs_getattr(const char *path, struct stat *stbuf) {
             fprintf(stderr, "getattr %s\n", path);
17
             memset(stbuf, 0, sizeof(struct stat));
18
19
             struct tfs_node *node = get_node(path);
20
21
             if (!node)
                    return -ENOENT;
22
             stbuf->st_mode = node->mode;
24
             stbuf->st_nlink = (node->mode & S_IFDIR) ? node->nlink + 1 : 1;
25
26
             stbuf->st_size = NODE_SIZE(node);
             stbuf->st_atim = node->atim;
27
28
             stbuf->st_mtim = node->mtim;
29
            return 0;
30
```

31 }

```
static int fuse_tfs_mknod(const char *path, mode_t mode, dev_t rdev) {
33
34
             fprintf(stderr, "mknod %s\n", path);
             return tfs_add_node(path, mode);
35
36
     }
37
     static int fuse_tfs_mkdir(const char *path, mode_t mode) {
38
             fprintf(stderr, "mkdir %s\n", path);
39
             // Bitwise OR with S_IFDIR because the documentation says to.
40
             return tfs_add_node(path, mode | S_IFDIR);
41
    }
42
43
     static int fuse_tfs_unlink(const char *path) {
44
             fprintf(stderr, "unlink %s\n", path);
45
             struct tfs_node *node = get_node(path);
46
47
             if (!node)
                     return -ENOENT;
48
             if (node->mode & S_IFDIR)
49
                     return -EISDIR;
50
51
             return tfs_remove_node(path);
52
53
     }
54
     static int fuse_tfs_rmdir(const char *path) {
55
             fprintf(stderr, "rmdir %s\n", path);
56
             struct tfs_node *node = get_node(path);
57
             if (!node)
58
                     return -ENOENT;
59
             if (!(node->mode & S_IFDIR))
60
                     return -ENOTDIR;
61
             if (node->size > 0)
62
                     return -ENOTEMPTY;
63
64
             return tfs_remove_node(path);
    }
66
67
     static int fuse_tfs_truncate(const char *path, off_t size) {
68
             fprintf(stderr, "truncate %s\n", path);
69
70
             struct tfs_node *node = get_node(path);
             if (!node)
71
72
                     return -ENOENT;
             if (node->mode & S_IFDIR)
73
                     return -EISDIR;
74
             node->size = size:
76
77
             return tfs_node_trim(node);
78
    }
79
     static int fuse_tfs_open(const char *path, struct fuse_file_info *fi) {
81
             fprintf(stderr, "open %s\n", path);
82
83
             if (!get_node(path))
                     return -ENOENT;
84
85
             return 0;
86
    }
87
88
     static int fuse_tfs_read(const char *path, char *buf, size_t size, off_t offset, struct fuse_file_info *fi) {
89
90
             fprintf(stderr, "read %s\n", path);
             struct tfs_node *node = get_node(path);
91
92
             if (!node)
                     return -ENOENT;
93
             if (node->mode & S_IFDIR)
                     return -EISDIR;
95
96
97
             return tfs_node_read(node, buf, size, offset);
    }
98
99
    static int fuse_tfs_write(const char *path, const char *buf, size_t size, off_t offset, struct fuse_file_info *fi) {
100
```

```
fprintf(stderr, "write %s\n", path);
101
102
              struct tfs_node *node = get_node(path);
              if (!node)
103
104
                      return -ENOENT;
              if (node->mode & S_IFDIR)
105
                     return -EISDIR;
106
107
             return tfs_node_write(node, buf, size, offset);
108
109
     }
110
     static int fuse_tfs_release(const char *path, struct fuse_file_info *fi) {
111
              fprintf(stderr, "release %s\n", path);
112
              return 0; // no-op
113
     }
114
115
     static int fuse_tfs_readdir(const char *path, void *buf, fuse_fill_dir_t filler, off_t offset,
116
                                   struct fuse_file_info *fi) {
117
              fprintf(stderr, "readdir %s\n", path);
118
119
              struct tfs_node *node = get_node(path);
              if (!node)
120
121
                      return -ENOENT;
              if (node->mode & S_IFREG)
122
                      return -ENOTDIR;
123
124
              filler(buf, ".", NULL, 0);
125
             filler(buf, "..", NULL, 0);
126
127
              struct tfs_node **children = tfs_node_children(node);
128
129
             for (int i = 0; i < node->nlink; i++) {
130
                      struct tfs_node *child = children[i];
131
132
                      filler(buf, child->name, NULL, 0);
133
              }
134
135
             free(children);
136
137
138
              return 0;
     }
139
140
     static void fuse_tfs_destroy(void *data) {
141
              tfs_destroy();
142
             hdestroy();
143
     }
144
145
     static int fuse_tfs_utimens(const char *path, const struct timespec tv[2]) {
146
             fprintf(stderr, "utimens %s\n", path);
147
148
              struct tfs_node *node = get_node(path);
              if (!node)
149
                      return -ENOENT;
150
151
             node->atim = tv[0];
152
              node->mtim = tv[1];
153
154
155
              return 0;
     }
156
157
158
     struct tfs_config {
              char *tfs_file_path;
159
160
     };
161
     enum {
162
             KEY_HELP,
163
     };
164
165
     // We intercept the help flag.
166
     static struct fuse_opt tfs_opts[] = {FUSE_OPT_KEY("-h", KEY_HELP), FUSE_OPT_KEY("-help", KEY_HELP), FUSE_OPT_END};
167
168
```

```
static int tfs_opt_proc(void *data, const char *arg, int key, struct fuse_args *outargs) {
169
170
             struct tfs_config *config = data;
171
172
             switch (key) {
             case KEY_HELP:
173
                      fprintf(stderr,
174
                              "usage: %s file mountpoint [fuse options]\n"
175
                              "\n"
176
                              "`file` must exist and must be initialized with `mktfs`."
177
                              "\n"
178
179
                              "See fuse(8) for more options.\n",
                              outargs->argv[0]);
180
                      return -1;
181
             case FUSE_OPT_KEY_NONOPT:
182
                      if (config->tfs_file_path)
183
                              return 1;
184
185
                      config->tfs_file_path = strdup(arg);
186
187
                      return 0;
             }
188
189
             return 1;
190
     }
191
192
     static struct fuse_operations fuse_tfs_oper = {.getattr = fuse_tfs_getattr,
193
                                                       .mknod = fuse_tfs_mknod,
194
                                                       .mkdir = fuse_tfs_mkdir,
195
                                                       .unlink = fuse_tfs_unlink,
196
                                                       .rmdir = fuse_tfs_rmdir,
197
                                                       .truncate = fuse_tfs_truncate,
198
199
                                                       .open = fuse_tfs_open,
                                                       .read = fuse_tfs_read,
200
                                                       .write = fuse_tfs_write,
201
                                                       .release = fuse_tfs_release,
202
                                                       .readdir = fuse_tfs_readdir,
203
                                                       .destroy = fuse_tfs_destroy,
204
                                                       .utimens = fuse_tfs_utimens};
205
206
     int main(int argc, char *argv[]) {
207
208
             struct fuse_args args = FUSE_ARGS_INIT(argc, argv);
             struct tfs_config config = {.tfs_file_path = NULL};
209
210
211
             if (fuse_opt_parse(&args, &config, tfs_opts, tfs_opt_proc) == -1)
                      return 1:
212
213
             if (!config.tfs_file_path) {
214
                      fprintf(stderr, "tfs: missing file to mount\n");
215
216
                      return 1;
             }
217
218
             int ret = tfs_load(config.tfs_file_path);
219
             if (ret)
220
221
                      return ret;
222
223
             return fuse_main(args.argc, args.argv, &fuse_tfs_oper, NULL);
     }
224
                                            Listing A.3: tfs.c
    #include "tfs.h"
    #include <assert.h>
    #include <errno.h>
    #include <fcntl.h>
     #include <libgen.h>
    #include <math.h>
    #include <search.h>
    #include <stdio.h>
```

```
#include <string.h>
9
    #include <time.h>
10
    #include <unistd.h>
11
    static struct tfs_info tfs_info;
13
14
    // Node number from pointer.
15
    #define NODENO(node) ((node)-tfs_info.nodes)
16
    // Block number from pointer.
    #define BLOCKNO(block) ((block)-tfs_info.data)
18
    // Cast block data to an array of node offsets.
19
    #define BLOCK_NODES(block) ((nodoff_t *)(tfs_info.data[block]))
20
    // Cast block data to an array of block offsets.
21
    #define BLOCK_POINTERS(block) ((blkoff_t *)(tfs_info.data[block]))
23
    // Cast block data to the next member in the free block linked list.
    #define NEXT_FREE_BLOCK(block) BLOCK_POINTERS(block)[0]
24
25
26
27
     * Iterating through indirect levels is painful,
     * so the process is abstracted away with the help of this iterator-like thingy.
28
29
    struct block_cursor {
30
            struct tfs_node *node;
31
            blkoff_t i;
32
             int level;
33
             blkoff_t pos[ILEVELS];
34
            blkoff_t block[ILEVELS];
35
36
37
    // Get what block an iterator is currently on.
38
    #define CURRENT_BLOCK(cursor)
39
             ((cursor)->i < DIRECT_BLOCKS ? (cursor)->node->blocks[(cursor)->i]
40
                                          : BLOCK_POINTERS((cursor)->block[(cursor)->level])[(cursor)->pos[(cursor)->level]])
41
42
    // Iterator definition helper
43
    #define DEFINE_BLOCK_CURSOR(var, nodeptr)
44
            struct block_cursor var = {
45
                 .node = nodeptr,
                 .level = -1,
47
48
             }:
49
    // Most significant bit
50
    static int msb(unsigned int n) {
            // I'm pretty sure gcc -0>1 compiles this down to a single BSR instruction.
52
            // Godbolt says it does, in which case this happens to be fast enough.
53
            // In any case, this is faster than log2().
54
            unsigned r = 0;
55
             while (n >>= 1)
                    r++;
57
58
59
             return r;
60
61
    // msb(x) == floor(log2(x))
62
    #define MAX_POINTERS_NBITS msb(BLOCK_MAX_POINTERS)
63
    #define BLOCK_SIZE_NBITS msb(BLOCK_SIZE)
64
    // floor-log base MAX_POINTERS
65
    #define FLOG(x) (msb(x) / MAX_POINTERS_NBITS)
66
    // This bit-twiddling is why BLOCK_SIZE must be a power of 2.
67
    #define MAX_POINTERS_POW(e) (e == 0 ? 1 : BLOCK_MAX_POINTERS << (MAX_POINTERS_NBITS * (e - 1)))
68
69
     * Set the position of a cursor for random access.
71
72
73
    static blkoff_t block_seek(struct block_cursor *cursor, blkoff_t pos) {
            blkoff_t nblocks = cursor->node->nblocks;
74
            cursor->i = pos;
             cursor \rightarrow level = -1;
76
```

```
77
78
              if (pos > nblocks)
                     return -1:
79
80
              if (pos < DIRECT_BLOCKS)</pre>
                      return cursor->node->blocks[pos];
81
82
              // Start from 0
 83
             pos -= DIRECT_BLOCKS;
84
              cursor->level = FLOG(pos + 1 - (pos + 1) / BLOCK_MAX_POINTERS);
             blkoff_t accum = BLOCK_MAX_POINTERS * (MAX_POINTERS_POW(cursor->level) - 1) / (BLOCK_MAX_POINTERS - 1);
86
87
              blkoff_t offset = pos - accum;
88
             cursor->block[0] = cursor->node->iblocks[cursor->level];
89
             for (int i = 0; i < cursor->level; i++) {
91
                      blkcnt_t N = MAX_POINTERS_POW(cursor->level - i);
92
                      cursor->pos[i] = offset / N;
93
                      offset %= N;
94
                      cursor->block[i + 1] = BLOCK_POINTERS((cursor)->block[i])[(cursor)->pos[i]];
95
             }
96
97
              cursor->pos[cursor->level] = offset;
98
99
              for (int i = cursor->level + 1; i < ILEVELS; i++)</pre>
100
                      cursor->pos[i] = 0;
101
102
              return CURRENT_BLOCK(cursor);
103
     }
104
105
106
107
      * Iterate a cursor through a callback.
108
      * Useful for inspecting all blocks we iterate through, including between levels.
109
110
     static blkoff_t iter_through(struct block_cursor *cursor,
111
                                    blkoff_t (*callback)(struct block_cursor *cursor, int level)) {
112
             if (++cursor->i < DIRECT_BLOCKS)</pre>
113
114
                      return callback(cursor, -1);
115
116
              int level = cursor->level;
              while (level >= 0 && ++cursor->pos[level] >= BLOCK_MAX_POINTERS)
117
                      cursor->pos[level--] = 0;
118
              if (level == -1)
120
                      // Done with this level, go to next level.
121
                      cursor->level += 1;
122
123
             for (; level < cursor->level; level++)
124
                      cursor->block[level + 1] = callback(cursor, level);
125
126
127
              return callback(cursor, cursor->level);
     }
128
129
130
131
      * Callback for simply getting the block a cursor is sitting on.
132
     static blkoff_t _next_block_callback(struct block_cursor *cursor, int level) {
133
134
             if (cursor->i >= cursor->node->nblocks)
                      return END_BLOCKS;
135
             if (cursor->i < DIRECT_BLOCKS)</pre>
136
                      return cursor->node->blocks[cursor->i];
137
              if (level == -1)
138
                     return cursor->node->iblocks[cursor->level];
139
140
141
             return BLOCK_POINTERS((cursor)->block[level])[(cursor)->pos[level]];
    }
142
143
    /**
144
```

```
* Sequential access.
145
146
     #define next_block(cursor) iter_through(cursor, _next_block_callback)
147
148
     int tfs_open(const char *filename) {
149
             int fd = open(filename, O_RDWR);
150
             if (fd == -1)
151
                      return -errno;
152
153
             // Memory map the file.
154
             // This will work for most x86-64 machines, but I'm not so sure about much else...
155
             tfs_info.filesize = lseek(fd, 0, SEEK_END);
156
             tfs_info.base = mmap(NULL, tfs_info.filesize, PROT_READ | PROT_WRITE, MAP_SHARED, fd, 0);
157
             close(fd);
158
159
             if (tfs_info.base == MAP_FAILED)
160
                      return -errno;
161
162
             return 0;
163
     }
164
165
     void tfs_format() {
166
             struct tfs_header *header = tfs_info.base;
167
             // Allocate the FAT (implicitly), blocks, and nodes.
168
             header->nblocks = tfs_info.filesize / (BLOCK_SIZE + (sizeof(struct tfs_node) / BLOCKS_PER_NODE));
169
             header->nnodes = header->nblocks / BLOCKS_PER_NODE;
170
             // Root takes 1 node.
171
             header->free_node_head = 1;
172
             header->free_block_head = 0;
173
174
             // Now (re)calculate pointers to FAT n' stuff.
175
             tfs_init();
176
             // Initialize root node:
178
             struct tfs_node *root = &tfs_info.nodes[0];
179
             root->mode = S_IFDIR | 644;
180
             root->name[0] = '\0'; // Root has no name.
181
182
             root->nblocks = 0;
             root->nlink = 0;
183
184
             clock_gettime(CLOCK_REALTIME, &root->atim);
             root->mtim = root->atim;
185
186
              // Initialize free blocks:
187
             for (int i = *tfs_info.free_block_head; i < tfs_info.nblocks - 1; i++)
188
                      NEXT_FREE_BLOCK(i) = i + 1;
189
             NEXT_FREE_BLOCK(tfs_info.nblocks - 1) = END_BLOCKS;
190
191
              // Initialize free nodes:
192
             for (int i = *tfs_info.free_node_head; i < tfs_info.nnodes - 1; i++)</pre>
193
                      tfs_info.nodes[i].next = i + 1;
194
             tfs_info.nodes[tfs_info.nnodes - 1].next = END_NODES;
195
     }
196
197
198
      * Walk the entire filesystem, adding all nodes to the hash table.
199
200
     static void init_htable(const char *path, struct tfs_node *node) {
201
202
             ENTRY entry;
203
204
             if (path) {
                      entry.key = malloc(strlen(path) + 1 + strlen(node->name) + 1);
205
                      sprintf(entry.key, "%s/%s", path, node->name);
             } else {
207
                      // Special case for root
208
209
                      entry.key = strdup("/");
             }
210
211
             entry.data = node;
212
```

```
hsearch(entry, ENTER);
213
214
              fprintf(stderr, "found %s\n", entry.key);
215
216
              if (!(node->mode & S_IFDIR))
217
                      return;
218
219
              // Recurse through directory:
220
221
             struct tfs_node **children = tfs_node_children(node);
222
223
             for (int i = 0; i < node->nlink; i++)
                      init_htable(path ? entry.key : "", children[i]);
224
225
             free(children);
226
     }
227
228
     void tfs_init() {
229
             struct tfs_header *header = tfs_info.base;
230
231
              tfs_info.nblocks = header->nblocks;
232
233
              tfs_info.nnodes = header->nnodes;
              tfs_info.free_block_head = &header->free_block_head;
234
              tfs_info.free_node_head = &header->free_node_head;
235
             tfs_info.nodes = tfs_info.base + sizeof(struct tfs_header);
236
237
              tfs_info.data = (void *)tfs_info.nodes + sizeof(struct tfs_node) * tfs_info.nnodes;
238
              fprintf(stderr, "nblocks: %ld\n", tfs_info.nblocks);
239
             fprintf(stderr, "nnodes: %ld\n", tfs_info.nnodes);
240
             fprintf(stderr, "free_node_head: %ld\n", *tfs_info.free_node_head);
241
              fprintf(stderr, "free_block_head: %ld\n", *tfs_info.free_block_head);
242
243
     }
244
     int tfs_load(const char *filename) {
              int ret = tfs_open(filename);
246
              if (ret)
247
248
                      return ret;
249
250
              tfs_init();
251
252
              hcreate(tfs_info.nnodes); // Initialize hash table, see hsearch(3)
             init_htable(NULL, &tfs_info.nodes[0]);
253
254
              return ret;
255
     }
256
257
     struct tfs_node *get_node(const char *path) {
258
             ENTRY entry = {
259
260
                  .key = strdup(path),
             }:
261
              ENTRY *result = hsearch(entry, FIND);
262
263
             free(entry.key);
264
              if (!result)
265
                      return NULL;
266
267
              return result->data;
268
     }
269
270
     struct tfs_node *get_directory(const char *path) {
271
272
              char *pathc = strdup(path);
             char *parent_path = dirname(pathc);
273
             struct tfs_node *parent_node = get_node(parent_path);
274
              free(pathc);
275
              return parent_node;
276
277
     }
278
279
     * Set or update a node give a path.
280
```

```
281
      * Especially useful for unsetting a node by passing NULL.
282
283
284
     static void set_node(const char *path, struct tfs_node *node) {
              ENTRY entry = {
285
                  .key = strdup(path),
286
                  .data = node,
287
              };
288
              ENTRY *result = hsearch(entry, FIND);
290
              if (!result) {
291
                      hsearch(entry, ENTER);
292
                      return;
293
              }
295
              result->data = node;
296
297
              free(entry.key);
     }
298
299
300
301
      * Iterator callback for freeing blocks we iterate through.
302
       * WARNING extreme hacks, detailed in report.
303
304
     static blkoff_t free_block_buffer[ILEVELS + 1];
305
     static blkoff_t _free_callback(struct block_cursor *cursor, int level) {
306
              return free_block_buffer[level + 1] = _next_block_callback(cursor, level);
307
     }
308
309
310
      * Iterator callback for allocating and setting blocks.
311
312
      * Also kind of hacky...
313
314
     static blkoff_t _alloc_callback(struct block_cursor *cursor, int level) {
315
              blkoff_t block = *tfs_info.free_block_head;
316
317
              if (block == END_BLOCKS)
318
                      return -1;
319
320
              *tfs_info.free_block_head = NEXT_FREE_BLOCK(block);
321
322
              if (cursor->i < DIRECT_BLOCKS)</pre>
                      return cursor->node->blocks[cursor->i] = block;
324
              if (level == -1)
325
                      return cursor->node->iblocks[cursor->level] = block;
326
327
              return BLOCK_POINTERS((cursor)->block[level])[(cursor)->pos[level]] = block;
328
     }
329
     int tfs_node_trim(struct tfs_node *node) {
331
              blkoff_t nrblocks = NODE_NRBLOCKS(node);
332
              blkoff_t dblocks = nrblocks - node->nblocks;
333
334
              DEFINE_BLOCK_CURSOR(cursor, node);
335
336
              if (dblocks < 0) {
337
338
                      block_seek(&cursor, nrblocks - 1);
                      while (dblocks && iter_through(&cursor, _free_callback) != END_BLOCKS) {
339
                               dblocks += 1:
340
341
                               // Clear out free block buffer:
                               for (int i = 0; i <= cursor.level + 1; i++) {</pre>
343
                                       if (free_block_buffer[i] < 0)</pre>
344
345
                                               continue:
                                       NEXT_FREE_BLOCK(free_block_buffer[i]) = *tfs_info.free_block_head;
346
347
                                       *tfs_info.free_block_head = free_block_buffer[i];
                                       free_block_buffer[i] = -1;
348
```

```
}
349
350
                      node->nblocks = nrblocks;
351
352
             } else {
                      block_seek(&cursor, node->nblocks - 1);
353
                      while (dblocks && iter_through(&cursor, _alloc_callback) != END_BLOCKS)
354
355
                              dblocks -= 1:
                      node->nblocks = nrblocks - dblocks;
356
             }
357
358
              // In case we couldn't allocate enough blocks, set sizes correctly.
359
             if (node->mode & S IFDIR)
360
                      node->nlink = MIN(node->nlink, node->nblocks * BLOCK_MAX_CHILDREN);
361
             else
                      node->size = MIN(node->size, node->nblocks * BLOCK SIZE):
363
364
             return dblocks > 0 ? -ENOSPC : 0;
365
     }
366
367
     int tfs_node_read(struct tfs_node *node, char *buf, size_t size, off_t offset) {
368
369
             DEFINE_BLOCK_CURSOR(cursor, node);
             size_t chunk, to_read = size;
370
             blkoff_t block = block_seek(&cursor, offset / BLOCK_SIZE);
371
372
             while (offset < NODE_SIZE(node) && (chunk = MIN(to_read, BLOCK_SIZE - (offset % BLOCK_SIZE)))) {
373
                      memcpy(buf, &tfs_info.data[block][offset % BLOCK_SIZE], MIN(chunk, NODE_SIZE(node) - offset));
374
                      block = next_block(&cursor);
375
                      to_read -= chunk;
376
                      offset += chunk;
377
                      buf += chunk;
378
             }
379
380
             clock_gettime(CLOCK_REALTIME, &node->atim);
381
382
             return size - to_read;
383
384
     }
385
386
     int tfs_node_write(struct tfs_node *node, const char *buf, size_t size, off_t offset) {
             node->size = MAX(node->size, offset + size);
387
388
             int ret = tfs_node_trim(node);
389
             DEFINE_BLOCK_CURSOR(cursor, node);
390
             size_t chunk, to_write = size;
             blkoff_t block = block_seek(&cursor, offset / BLOCK_SIZE);
392
393
             while (offset < NODE_SIZE(node) && (chunk = MIN(to_write, BLOCK_SIZE - (offset % BLOCK_SIZE)))) {
394
                      memcpy(&tfs_info.data[block][offset % BLOCK_SIZE], buf, chunk);
395
                      block = next_block(&cursor);
396
                      to_write -= chunk;
397
                      offset += chunk;
398
                      buf += chunk;
399
400
401
             clock_gettime(CLOCK_REALTIME, &node->mtim);
402
403
             return ret < 0 ? ret : size - to_write;</pre>
404
405
406
     struct tfs_node **tfs_node_children(struct tfs_node *node) {
407
             nodoff_t *children = malloc(NODE_SIZE(node));
408
             if (!children)
409
                      return NULL;
410
411
             tfs_node_read(node, (void *)children, NODE_SIZE(node), 0);
412
413
             struct tfs_node **children_nodes = malloc(node->nlink * sizeof(struct tfs_node *));
414
             if (!children_nodes)
415
                      return NULL;
416
```

```
417
              for (int i = 0; i < node->nlink; i++)
418
                      children_nodes[i] = &tfs_info.nodes[children[i]];
419
420
              free(children);
421
422
              return children_nodes;
423
     }
424
425
     int tfs_add_node(const char *path, mode_t mode) {
426
             if (get_node(path))
427
                      return -EEXIST;
428
              if (*tfs_info.free_node_head == END_NODES)
429
                      return -ENOSPC;
430
431
              char *basename = strrchr(path, '/') + 1;
432
              if (strlen(basename) + 1 > NAME_LIMIT)
433
                      return -ENAMETOOLONG;
434
435
              // Allocate node.
436
437
              nodoff_t nodei = *tfs_info.free_node_head;
             struct tfs_node *node = &tfs_info.nodes[nodei];
438
              *tfs_info.free_node_head = node->next;
439
440
              fprintf(stderr, "\tAllocated node %ld...\n", nodei);
441
              // Initialize node.
              strcpy(node->name, basename);
443
              node->mode = mode;
444
              if (node->mode & S_IFDIR)
445
                      node->nlink = 0;
446
              else
447
                      node->size = 0;
448
             node->nblocks = 0;
449
             clock_gettime(CLOCK_REALTIME, &node->atim);
450
              node->mtim = node->atim;
451
452
              // Add child to parent.
453
454
              struct tfs_node *parent_node = get_directory(path);
              parent_node->nlink += 1;
455
456
              tfs_node_trim(parent_node);
              DEFINE_BLOCK_CURSOR(cursor, parent_node);
457
              BLOCK_NODES(block_seek(&cursor, parent_node->nblocks - 1))
458
              [(parent_node->nlink - 1) % BLOCK_MAX_CHILDREN] = nodei;
459
460
              clock_gettime(CLOCK_REALTIME, &parent_node->mtim);
461
462
              // Update hash table.
463
464
              set_node(path, node);
465
              return 0;
466
     }
467
468
469
     int tfs_remove_node(const char *path) {
             struct tfs_node *node = get_node(path);
470
471
              struct tfs_node *parent_node = get_directory(path);
472
              // Don't rm -rf / -_-
473
474
              if (!parent_node)
                      return -ENOTSUP;
475
476
              // Remove from parent.
477
              DEFINE_BLOCK_CURSOR(cursor, parent_node);
              nodoff_t last_child =
479
                  BLOCK_NODES(block_seek(&cursor, parent_node->nblocks - 1))[(parent_node->nlink - 1) % BLOCK_MAX_CHILDREN];
480
481
              for (blkoff_t block = block_seek(&cursor, 0); block != END_BLOCKS; block = next_block(&cursor)) {
482
                      for (int i = 0; i < BLOCK_MAX_CHILDREN; i++) {</pre>
483
                               if (BLOCK_NODES(block)[i] == NODENO(node)) {
484
```

```
BLOCK_NODES(block)[i] = last_child;
485
486
                                       goto outer;
                              }
487
488
                      }
             }
489
     outer:
490
              parent_node->nlink -= 1;
491
             tfs_node_trim(parent_node);
492
              clock_gettime(CLOCK_REALTIME, &parent_node->mtim);
493
494
495
              // Deallocate blocks.
             node->size = 0;
496
              tfs_node_trim(node);
497
498
              // Deallocate node.
499
             node->next = *tfs_info.free_node_head;
500
              *tfs_info.free_node_head = NODENO(node);
501
502
503
              // Remove from hash table.
              set_node(path, NULL);
504
505
              return 0;
506
     }
507
508
     int tfs_destroy() {
509
              // Write back changes to disk.
510
              return munmap(tfs_info.base, tfs_info.filesize);
511
     }
512
                                           Listing A.4: mktfs.c
     #include "tfs.h"
     #include <math.h>
     #include <stdio.h>
     #include <stdlib.h>
     #include <string.h>
     int main(int argc, char *argv[]) {
             int ret = 0;
              if (argc < 2) {
 10
                      fprintf(stderr,
11
                               "usage: %s < file > n"
13
                               "Allocate space to a file using fallocate(1) first.\n",
 14
15
                               argv[0]);
                      return 1;
16
 17
             }
18
              ret = tfs_open(argv[1]);
              if (ret)
20
21
                      return ret;
22
23
              tfs_format();
              tfs_destroy();
24
25
26
              return 0;
    }
27
```