

The Brainix Manual

The Brainix Team

Copyright © Draft date September 4, 2009

Contents

I	Introductory matters	ii
II	Understanding the Brainix Operating System	1
1	On the Brainix File System	2
1.1	Introduction: How Servers Work (A Quick Gloss over it)	2
1.2	How File Systems Traditionally Work	3
1.2.1	The I-Node	3
1.2.2	The Directory	4
1.3	The File System Details	5
1.4	File System Initialization	6
1.4.1	block_init()	6
1.4.2	inode_init()	8
1.4.3	super_init()	9
1.4.4	dev_init()	11
1.4.5	descr_init()	12
1.5	File System Operations	13
1.5.1	REGISTER	13
III	Tutorials on Using the Brainix Operating System	27
2	How to run Brainix on Unix-like Operating Systems and Bochs	28
3	How to hack Brainix: Some Things to Bear In Mind	30
4	On the Implementation of the debugging for Brainix	32
4.1	FAQ	32
4.1.1	What is the format for debug?	32
4.1.2	I'm making a new driver, and I don't know what to do about this debugger...	33
IV	Appendices	34
A	GNU Free Documentation License	35
A.1	APPLICABILITY AND DEFINITIONS	35
A.2	VERBATIM COPYING	37
A.3	COPYING IN QUANTITY	37
A.4	MODIFICATIONS	38

A.5	COMBINING DOCUMENTS	39
A.6	COLLECTIONS OF DOCUMENTS	40
A.7	AGGREGATION WITH INDEPENDENT WORKS	40
A.8	TRANSLATION	40
A.9	TERMINATION	40
A.10	FUTURE REVISIONS OF THIS LICENSE	41
	Bibliography	42

Part I

Introductory matters

Part II

Understanding the Brainix Operating System

Chapter 1

On the Brainix File System

“ ‘Where shall I begin, please your Majesty?’

‘Begin at the beginning,’ the King said gravely, ‘and go on till you come to the end: then stop.’ ”¹

Note that this is released under the GNU Free Documentation License version 1.2. See the file fdl.tex for details of the license. This, like the rest of the Brainix Project, is a work in progress.

1.1 Introduction: How Servers Work (A Quick Gloss over it)

The structure for the file system is simple, it is structured like all servers for the micro-kernel:

```
/* typical_server.pseudo_c */

int main(void) {
    init(); //This starts up the server and initializes values
           //registers it with the kernel and file system
           //if necessary, etc.
    msg* m; //this is the message buffer

    //You can tell I didn't program this otherwise
    //SHUT_DOWN would be GO_TO_HELL
    while((&m = msg_receive(ANYONE))->op != SHUT_DOWN)
    {
        switch(m->op) {
            case OP_ONE: /* ... */ break;
            /* other op cases supported by the server */
            default: panic("server", "unrecognized message!");
        }
        //The following deals with the reply
        switch(m->op)
        {
            case OP_ONE: /* ... */ break;
            /* other replies that require modifications */
            default: msg_reply(m);
        }
    }
}
```

¹ *Alice's Adventures in Wonderland* by Lewis Carroll, chapter 12 “Alice's Evidence”

```

    deinit(); //this is called to de-initialize the server
              //to prepare for shut down
    shut_down(); //I would've named it "buy_the_farm()"
                //or "go_to_hell()"
    return 0;
}

```

With most servers, this is the entirety of the `main.c` file. The actual implementation of the methods (i.e. “the dirty work is carried out through”) auxiliary files.

The “op” field of the message refers to the operation; which is a sort of parallel to the monolithic kernel system call. The system call is merely handled in user space.

1.2 How File Systems Traditionally Work

There are probably a number of introductory texts and tutorials on Unix-like file systems. I will mention a few worthy of note [1] [2] [4] [3]. I will **attempt** to briefly explain how the Unix file system works, and explain its implementation in operating systems such as Linux and maybe FreeBSD.

File systems deal with long term information storage. There are three essential requirements for long-term information storage that Tanenbaum and Woodhull recognize [2]:

1. It must be possible to store a very large amount of information.
2. The information must survive the termination of the process using it.
3. Multiple processes must be able to access the information concurrently.

With the exception of the GNU-Hurd solution to these problems, the answer is usually to store information on hard disks in units called **files**. The management of these units is done by a program called the **file system**. (What’s so interesting and exciting about Unix and Unix-like operating systems is that it’s object oriented: everything “is-a” file!)

Some few notes on the geometry of the structure of hard disks. There are sectors, which consist of 512 bytes. There are blocks, which consist of 2^n sectors (where n is usually 3, but varies between 1 and 5). That is a block is 1024 to 16384 bytes. Typically it is 4096 bytes per block.

1.2.1 The I-Node

The file in Unix² is represented by something called an **inode (index-node)**. This lists the attributes and disk addresses of the file’s blocks. The skelix code³ shall be used (with permission of course) as an example of the simplest inode:

```

1 /* from /skelix07/include/fs.h */
2 /* Skelix by Xiaoming Mo (xiaoming.mo@skelix.org)
3  * Licence: GPLv2 */
4 #ifndef FS_H
5 #define FS_H
6
7 #define FT_NML    1

```

²Out of sheer laziness, “Unix” should be read as “Unix and Unix-like operating systems”.

³Specifically from here <http://skelix.org/download/07.rar>

```

8 #define FT_DIR    2
9
10 struct INODE {
11     unsigned int i_mode;        /* file mode */
12     unsigned int i_size;        /* size in bytes */
13     unsigned int i_block[8];
14 };

```

Note that the different types of inodes there are is defined in lines 06 and 07. The permissions and type of the inode is on line 10. The actual addresses to the blocks that hold the data for the file are stored in the array on line 12. At first you look and think “Huh, only 8 blocks per file? That’s only, what, 32768 bytes?!” Since it is incredibly unlikely that all the information you’d ever need could be held in 32 kilobytes, the last two addresses refers to *indirect* addresses. That is the seventh address refers to a sector that contains (512 bytes per sector)(1 address per 4 bytes) = 128 addresses. The seventh entry is called a **indirect block** (although because Skelix is so small, it’s an indirect sector). The last entry refers to an indirect block, for this reason it is called a **double indirect block**. The indirect block holds 128 addresses, each address refers to a 512 byte sector (in other operating systems they refer to blocks), so each indirect block refers to $128 \times 512 = 65536$ bytes or 64 kilobytes. The last double indirect block contains 128 single indirect blocks, or $128 \times 64 = 8192$ kilobytes or 8 Megabytes.

In bigger operating systems, there are triple indirect blocks, which if we implemented it in skelix we would get $128 \times 8192 = 1048576$ kilobytes or 1024 megabytes or 1 gigabyte. “Surely there must be quadruple indirect blocks, as I have a file that’s several gigabytes on my computer!” Well, the way it is implemented on Linux is that rather than refer to sectors, there are groups of sectors called **block groups**. Instead of accessing *only* 512 byte atoms, we are accessing **4 kilobyte atoms**! Indeed, if I am not mistaken, the Minix 3 file system refers to blocks instead of sectors too.

1.2.2 The Directory

So what about the directory? Well, in Unix file systems, the general idea is to have a file that contains **directory entries**. Directory entries basically hold at least two things: the file name, and the inode number of the entry. There are other things that are desirable like the name length of the entry, the type of file the entry is, or the offset to be added to the starting address of the directory entry to get the starting address of the next directory entry (the “rectangular length”). Consider the implementation in Skelix:

```

15 /* from /skelix07/include/fs.h */
16 extern struct INODE iroot;
17
18 #define MAX_NAME_LEN 11
19
20 struct DIR_ENTRY {
21     char de_name[MAX_NAME_LEN];
22     int de_inode;
23 };

```

The directory entry is, like the skelix inode, extremely simplistic. It consists of the address to the entry, and the entry’s name. Suppose one had the following directory:

inode number	name
1	.
1	..
4	bin
7	dev

One wants to run a program, so one looks up the program `/bin/pwd`. The look-up process then goes to the directory and looks up `/bin/`, it sees the inode number is 4, so the look up process goes to inode 4. It finds:

```
( I-Node 4 is for /bin/ )
Mode
Size
```

...

I-node 4 says that `/bin/` is in block 132. It goes to block 132:

6	.
1	..
19	bash
30	gcc
51	man
26	ls
45	pwd

The look up process goes to the last entry and finds `pwd` - the program we're looking for! The look up process goes to block 45 and finds the inode that refers to the blocks necessary to execute the file. That's how the directory system works in Unix file systems.

Every directory has two directory entries when they are made: 1) `.` which refers to "this" directory, 2) `..` which refers to the parent of "this" directory. In this sense, the directories are a sort of doubly linked lists.

1.3 The File System Details

"The rabbit-hole went straight on like a tunnel for some time, and then dipped suddenly down, so suddenly that Alice had not a moment to think about stopping herself before she found herself falling down a very deep well."⁴

So if you actually go and look at the file system directory, there are a number of ops that are implemented. Some of them are obvious, like `read()`, `write()`, etc. Others are not really intuitively clear why they're there, like `execve()`. The reason for this is because Brainix attempts to be POSIX-Compliant, and POSIX really wasn't made with Microkernels in mind. So we're stuck having an odd design like this; but the advantage is that we can eventually use a package manager like Portage⁵. The advantages really outweigh the cost of odd design.

So this section will inspect the various operations, and follow the code "down the rabbit hole". Yes we shall inspect the nitty-gritty details and analyze as much as

⁴ *Alice's Adventures in Wonderland* by Lewis Carroll, chapter 1 "Down the Rabbit-Hole"

⁵For those that do not know, Portage is the package manager for the Gentoo distribution of Linux. As far as I know it has been ported to FreeBSD, Open-BSD, Net-BSD, Darwin, and other operating systems because Portage is distributed via its source code. It works by downloading and compiling source code auto-magically and optimizing it as much as possible with the GCC.

possible. That is my duty as the file system hacker to explain as much as possible, using code snippets where appropriate. So we begin with the initialization of the file system.

1.4 File System Initialization

Looking in the file `/brainix/src/fs/main.c` one finds:

```
void fs_main(void)
{
    /* Initialize the file system. */
    block_init(); /* Initialize the block cache. */
    inode_init(); /* Initialize the inode table. */
    super_init(); /* Initialize the superblock table. */
    dev_init();   /* Initialize the device driver PID table. */
    descr_init(); /* Init the file ptr and proc-specific info tables. */
}
```

This is the initialization code that we are interested in. Let's analyze it line by line. First there is a call to the function `block_init()`. So let us inspect this function's code.

1.4.1 `block_init()`

There is the matter of the data structure that is involved here extensively that we ought to investigate first: `block_t`.

`block_t`

```
43 /* from /brainix/inc/fs/block.h */
44 /* A cached block is a copy in RAM of a block on a device: */
45 typedef struct block
46 {
47     /* The following field resides on the device: */
48     char data[BLOCK_SIZE]; /* Block data. */
49
50     /* The following fields do not reside on the device: */
51     dev_t dev; /* Device the block is on. */
52     blkcnt_t blk; /* Block number on its device. */
53     unsigned char count; /* Number of times the block is used. */
54     bool dirty; /* Block changed since read. */
55     struct block *prev; /* Previous block in the list. */
56     struct block *next; /* Next block in the list. */
57 } block_t;
```

This is all rather straight forward. The `dev_t` field tells us what device we are dealing with, rather what device the file system is dealing with. To be more precise about what exactly `dev_t` is we look to the code:

`dev_t`

```
48 /* from /brainix/inc/lib/sys/type.h */
49 /* Used for device IDs: */
50 #ifndef _DEV_T
51 #define _DEV_T
52 typedef unsigned long dev_t;
53 #endif
```

which is pretty self-explanatory that `dev_t` is little more than an unsigned long. The `blkcnt_t blk` field gives more precision with what we are dealing with, which is a rather odd field because I don't know what the `blkcnt_t` type is off hand so I doubt that you would either. Let us shift our attention to this type!

`blkcnt_t`

```

30 /* from /brainix/inc/lib/sys/type.h */
31 /* Used for file block counts: */
32 #ifndef _BLKCNT_T
33 #define _BLKCNT_T
34 typedef long blkcnt_t;
35 #endif

```

So this is a rather straight forward type that needs no explanation it seems. We can continue our analysis of the `block_t` struct. The `unsigned char count;` is little more than a simple counter it seems, and the `bool dirty;` tells us whether the block has changed since last read or not. The last two entries tells us this `block_t` data structure is a doubly linked list. This is common, the use of doubly linked lists that is, because it is common to lose things at such a low level.

Now we may proceed to analyze the `block_init()` function defined in the `block.c` file:

block_init()

```

32 /* /brainix/src/fs/block.c */
33 void block_init(void)
34 {
35
36     /* Initialize the block cache. */
37
38     block_t *block_ptr;
39
40     /* Initialize each block in the cache. */
41     for (block_ptr = &block[0]; block_ptr < &block[NUM_BLOCKS]; block_ptr++)
42     {
43         block_ptr->dev = NO_DEV;
44         block_ptr->blk = 0;
45         block_ptr->count = 0;
46         block_ptr->dirty = false;
47         block_ptr->prev = block_ptr - 1;
48         block_ptr->next = block_ptr + 1;
49     }
50
51     /* Make the cache linked list circular. */
52     block[0].prev = &block[NUM_BLOCKS - 1];
53     block[NUM_BLOCKS - 1].next = &block[0];
54
55     /* Initialize the least recently used position in the cache. */
56     lru = &block[0];
57 }

```

Line 37 simply initializes a block pointer that is used to initialize the blocks. Lines 40 to 48 (the for-loop) uniformly sets all the blocks to be identical with the exact same fields. The fields are self explanatory; the device number is set to no device (line 42), the number of times the block has been used is set to zero (line 43), the block has not changed since it's last been read (line 44), the previous block and next block are rather elementarily defined.

At first one would think looking up until line 47 that there would have to be a negative block, and that block would require another, and so on *ad infinitum*. But lines 50 to 52 make the block a circularly doubly linked list. Line 51 makes the zeroeth block's previous block `prev` refer to the last block, and line 52 makes the last block's `next` field refers to the zeroeth block's address.

What's the significance of line 55? Well, I don't know. It does not seem to relevant at the moment, though undoubtedly we shall have to come back to it in the future.

1.4.2 inode_init()

Just as we had the `block_init()` we have a `inode_init()`. If you are new to this whole Unix-like file system idea, it is highly recommended that you read [5] [6] [7] [8] [9] [10]. Perhaps in a future version of this documentation it will be explained in further detail. The original motivation I suspect (yes, this is a baseless conjecture I made up from my own observations that is probably not true at all) was to have something similar to a hybrid of Linux and Minix 3, and this is somewhat reflected by the choice of attempting to support the ext2 file system (the file system from the earlier Linux distributions). The inode data structure is identical to its description in the third edition of *Understanding the Linux Kernel*. However I am making this an independent, stand-alone type of reference...so that means I am going to inspect the data structure, line by line.

inode_t

```

42 /* from /brainix/inc/fs/inode.h */
43 /* An inode represents an object in the file system: */
44 typedef struct
45 {
46     /* The following fields reside on the device: */
47     unsigned short i_mode;          /* File format / access rights. */
48     unsigned short i_uid;          /* User owning file. */
49     unsigned long i_size;          /* File size in bytes. */
50     unsigned long i_atime;         /* Access time. */
51     unsigned long i_ctime;         /* Creation time. */
52     unsigned long i_mtime;         /* Modification time. */
53     unsigned long i_dtime;         /* Deletion time (0 if file exists). */
54     unsigned short i_gid;          /* Group owning file. */
55     unsigned short i_links_count; /* Links count. */
56     unsigned long i_blocks;        /* 512-byte blocks reserved for file. */
57     unsigned long i_flags;         /* How to treat file. */
58     unsigned long i_osd1;          /* OS dependent value. */
59     unsigned long i_block[15];     /* File data blocks. */
60     unsigned long i_generation;    /* File version (used by NFS). */
61     unsigned long i_file_acl;      /* File ACL. */
62     unsigned long i_dir_acl;       /* Directory ACL. */
63     unsigned long i_faddr;         /* Fragment address. */
64     unsigned long i_osd2[3];       /* OS dependent structure. */
65
66     /* The following fields do not reside on the device: */
67     dev_t dev;                    /* Device the inode is on. */
68     ino_t ino;                    /* Inode number on its device. */
69     unsigned char count;          /* Number of times the inode is used. */
70     bool mounted;                 /* Inode is mounted on. */
71     bool dirty;                   /* Inode changed since read. */
72 } inode_t;

```

A lot of this code is seemingly unused. All that really matters is that the `inode_t` data type is a wrapper for the addresses (line 58), with some constraints for permissions and so forth (lines 46 to 57), and some device specific fields (lines 66 to 70). This data structure is nearly identical to the ext2 file system's `inode` struct. As stated previously, the motivation was to incorporate the ext2 file system into Brainix. This proved too difficult since the ext2 file system is intimately related to the Linux virtual file system. It seems that the most appropriate description for the Brainix file system is a fork of the ext2 one.

Now on to the `inode_init()` code itself:

inode_init()

```

32 void inode_init(void)
33 {

```

```

34
35 /* Initialize the inode table. */
36
37     inode_t *inode_ptr;
38
39     /* Initialize each slot in the table. */
40     for (inode_ptr = &inode[0]; inode_ptr < &inode[NUM_INODES]; inode_ptr++)
41     {
42         inode_ptr->dev = NO_DEV;
43         inode_ptr->ino = 0;
44         inode_ptr->count = 0;
45         inode_ptr->mounted = false;
46         inode_ptr->dirty = false;
47     }
48 }

```

End of /brainix/src/fs/inode.c

Line 37 tells us there is a dummy inode pointer that is used later on, more specifically it is used in lines 40 to 47 when the inode table is initialized. The for-loop, as stated, initializes the inode-table. Line 42 sets the device that the inode is on to NO_DEV, line 43 sets the inode number to zero, the next line (line 44) sets the number of times the inode is used to zero, line 45 sets the boolean checking whether the inode is mounted or not to false (the inode is initialized to be not mounted), and line 46 tells us that the inode has not changed since we last dealt with it.

Now that the inode table has been initialized, we now look to the initialization of the super block.

1.4.3 super_init()

To inspect the inner workings of the `super_init()` method we need to first investigate the `super` struct representing the super block.

typedef struct {...} super

```

35 /* /brainix/inc/fs/super.h */
36 /* The superblock describes the configuration of the file system: */
37 typedef struct
38 {
39     /* The following fields reside on the device: */
40     unsigned long s_inodes_count;    /* Total number of inodes.      */
41     unsigned long s_blocks_count;    /* Total number of blocks.      */
42     unsigned long s_r_blocks_count;  /* Number of reserved blocks.   */
43     unsigned long s_free_blocks_count; /* Number of free blocks.      */
44     unsigned long s_free_inodes_count; /* Number of free inodes.      */
45     unsigned long s_first_data_block; /* Block containing superblock. */
46     unsigned long s_log_block_size;  /* Used to compute block size.  */
47     long s_log_frag_size;            /* Used to compute fragment size. */
48     unsigned long s_blocks_per_group; /* Blocks per group.            */
49     unsigned long s_frags_per_group;  /* Fragments per group.         */
50     unsigned long s_inodes_per_group; /* Inodes per group.            */
51     unsigned long s_mtime;           /* Time of last mount.          */
52     unsigned long s_wtime;           /* Time of last write.          */
53     unsigned short s_mnt_count;       /* Mounts since last fsck.      */
54     unsigned short s_max_mnt_count;  /* Mounts permitted between fscks. */
55     unsigned short s_magic;           /* Identifies as ext2.          */
56     unsigned short s_state;           /* Cleanly unmounted?          */
57     unsigned short s_errors;          /* What to do on error.         */
58     unsigned short s_minor_rev_level; /* Minor revision level.        */
59     unsigned long s_lastcheck;        /* Time of last fsck.           */
60     unsigned long s_checkinterval;    /* Time permitted between fscks. */
61     unsigned long s_creator_os;       /* OS that created file system. */
62     unsigned long s_rev_level;        /* Revision level.              */
63     unsigned short s_def_resuid;      /* UID for reserved blocks.     */

```

```

64     unsigned short s_def_resgid;          /* GID for reserved blocks.      */
65     unsigned long s_first_ino;           /* First usable inode.          */
66     unsigned short s_inode_size;         /* Size of inode struct.        */
67     unsigned short s_block_group_nr;     /* Block group of this superblock. */
68     unsigned long s_feature_compat;      /* Compatible features.          */
69     unsigned long s_feature_incompat;     /* Incompatible features.        */
70     unsigned long s_feature_ro_compat;    /* Read-only features.           */
71     char s_uuid[16];                     /* Volume ID.                    */
72     char s_volume_name[16];               /* Volume name.                  */
73     char s_last_mounted[64];              /* Path where last mounted.      */
74     unsigned long s_algo_bitmap;          /* Compression methods.          */
75
76     /* The following fields do not reside on the device: */
77     dev_t dev;                            /* Device containing file system. */
78     blksize_t block_size;                 /* Block size.                   */
79     unsigned long frag_size;              /* Fragment size.                */
80     inode_t *mount_point_inode_ptr;       /* Inode mounted on.             */
81     inode_t *root_dir_inode_ptr;          /* Inode of root directory.      */
82     bool dirty;                           /* Superblock changed since read. */
83 } super_t;

```

This is the super block, and - as previously iterated a number of times - this is from the ext2 file system. The superblock should have the magic number `s_magic` which tells us this is indeed the ext2 file system. Line 61 tells us the revision level which allows the mounting code to determine whether or not this file system supports features available to particular revisions. When a new file is created, the values of the `s_free_inodes_count` field in the Ext2 superblock and of the `bg_free_inodes_count` field in the proper group descriptor must be decremented. If the kernel appends some data to an existing file so that the number of data blocks allocated for it increases, the values of the `s_free_blocks_count` field in the Ext2 superblock and of the `bg_free_blocks_count` field in the group descriptor must be modified. Even just rewriting a portion of an existing file involves an update of the `s_wtime` field of the Ext2 superblock. For a more in-depth analysis of the ext2 file system's `super_block` data structure which was forked for the Brainix file system, see Chapter 18 [3] or [7] [11] [12].

super_init()

```

32 /* /brainix/src/fs/super.c */
33 void super_init(void)
34 {
35
36     /* Initialize the superblock table. */
37
38     super_t *super_ptr;
39
40     /* Initialize each slot in the table. */
41     for (super_ptr = &super[0]; super_ptr < &super[NUM_SUPERS]; super_ptr++)
42     {
43         super_ptr->dev = NO_DEV;
44         super_ptr->block_size = 0;
45         super_ptr->frag_size = 0;
46         super_ptr->mount_point_inode_ptr = NULL;
47         super_ptr->root_dir_inode_ptr = NULL;
48         super_ptr->dirty = false;
49     }
50 }

```

As previously stated, the brainix file system is perhaps more properly thought of as a fork (rather than an implementation) of the ext2 file system. In the ext2 file system, each block group has a super block (as a sort of back up), and this feature

has been inherited in the brainix file system. This `init()` method is pretty much identical to the other ones. There is a pointer struct (line 37) that's used in a for-loop to set all the super blocks to be the same (lines 40 to 48).

More specifically, in more detail, line 42 sets each super block's device to `NO_DEV`. The block size for the super block is initialized to be zero as well, with no fragments either (lines 43 and 44). The inode holding the mount point information is set to be `NULL` as is the root directory inode pointer. Since we just initialized the super blocks, they haven't changed since we last used them, so we tell that to the super blocks with line 47.

1.4.4 `dev_init()`

The `pid_t` data structure

We should first inspect the vital data structure relevant to discussion here: `pid_t` which is defined in `/brainix/inc/lib/unistd.h`:

pid_t

```
112 /* /brainix/inc/lib/unistd.h */
113 #ifndef _PID_T
114 #define _PID_T
115 typedef long pid_t;
116 #endif
```

That's pretty much the only new data structure (or type, rather) that's relevant for discussion here.

Back to the `dev_init()` method

The next function called in the `init()` section of the file system server is the `dev_init()`. This is defined in the `/brainix/src/fs/device.c` file:

dev_init()

```
55 /* /brainix/src/fs/device.c */
56 void dev_init(void)
57 {
58
59 /* Initialize the device driver PID table. */
60
61     unsigned char maj;
62
63     for (maj = 0; maj < NUM_DRIVERS; maj++)
64         driver_pid[BLOCK][maj] =
65         driver_pid[CHAR][maj] = NO_PID;
66 }
```

This is the `dev_init()` code, that basically initializes the device driver part of the PID⁶ table. At first looking at the for-loop, one says "This won't work!" But upon further inspection, the line 63 doesn't have a semicolon, so the compiler continues to the next line (line 64). It sets the `driver_pid[BLOCK][maj]` to be `NO_PID`. It does this for every major device (more precisely, for the number of drivers `NUM_DRIVERS`). Note that the first index of the matrix that represents the device driver PID table is capable of having values 0 and 1, represented by `BLOCK` and `CHAR` respectively.

⁶"PID" stands for "Process identification".

1.4.5 descr_init()

For this method, there are global variables defined in the headers:

```

54 /* /brainix/inc/fs/fildes.h */
55 /* Global variables: */
56 file_ptr_t file_ptr[NUM_FILE_PTRS]; /* File pointer table. */
57 fs_proc_t fs_proc[NUM_PROCS];      /* Process-specific information table. */

```

This allows us to introduce the data structures `fs_proc_t` and `file_ptr_t`.

file_ptr_t

```

35 /* from /brainix/inc/fs/fildes.h */
36 /* A file pointer is an intermediary between a file descriptor and an inode: */
37 typedef struct
38 {
39     inode_t *inode_ptr; /* Inode pointer. */
40     unsigned char count; /* Number of references. */
41     off_t offset; /* File position. */
42     int status; /* File status. */
43     mode_t mode; /* File mode. */
44 } file_ptr_t;

```

This is self explanatory thanks to the comments. The `file_ptr_t` is an intermediary between a file descriptor and an i-node. It consists of the inode it intermediates for (line 38), the number of references made to the inode in the file descriptor (line 39), the file position's offset (line 40), the status of the file (41), and the mode of the file (42). The other important data structure is:

fs_proc_t

```

45 /* from /brainix/inc/fs/fildes.h */
46 /* Process-specific file system information: */
47 typedef struct
48 {
49     inode_t *root_dir; /* Root directory. */
50     inode_t *work_dir; /* Current working directory. */
51     mode_t cmask; /* File mode creation mask. */
52     file_ptr_t *open_descr[OPEN_MAX]; /* File descriptor table. */
53 } fs_proc_t;

```

Which gives us information about the file system which is process-specific, as the comment suggests. More to the point, the root directory inode, the current directory inode, the “file mode creation mask” which is little more than telling the file what you **DON’T** want (“Setting a mask is the opposite of setting the permissions themselves; when you set a mask, you are telling the computer the permissions you do not want, rather than the permissions you do” [13]), and more importantly the file descriptor table.

This is the last step in the file system initialization. It essentially initializes a few other tables that we are going to use.

descr_init()

```

32 /* from /brainix/src/fs/fildes.c */
33 void descr_init(void)
34 {
35
36     /* Initialize the file pointer table and the process-specific file system
37      * information table. */
38
39     int ptr_index;
40     pid_t pid;

```

```

41     int descr_index;
42
43     /* Initialize the file pointer table. */
44     for (ptr_index = 0; ptr_index < NUM_FILE_PTRS; ptr_index++)
45     {
46         file_ptr[ptr_index].inode_ptr = NULL;
47         file_ptr[ptr_index].count = 0;
48         file_ptr[ptr_index].offset = 0;
49         file_ptr[ptr_index].status = 0;
50         file_ptr[ptr_index].mode = 0;
51     }
52
53     /* Initialize the process-specific file system information table. */
54     for (pid = 0; pid < NUM_PROCS; pid++)
55     {
56         fs_proc[pid].root_dir = NULL;
57         fs_proc[pid].work_dir = NULL;
58         fs_proc[pid].cmask = 0;
59         for (descr_index = 0; descr_index < OPEN_MAX; descr_index++)
60             fs_proc[pid].open_descr[descr_index] = NULL;
61     }
62 }

```

There is nothing new here, only two for-loops rather than one to initialize two (rather than one) tables. But where are these tables defined? They seem to fall from thin air into our laps!

1.5 File System Operations

This section, unlike the previous, is in a seemingly random order. It does not logically follow the structure of the code as it would appear to a new comer to the Brainix kernel. Instead, it inspects the more important operations first, discussing them at length. We shall begin with the most recently inspected method: `REGISTER`.

1.5.1 REGISTER

This method was long thought to be a problem child, until some clever debugging proved it to be little more than a nuisance. It is a function in the `device.c` file:

fs_register()

```

----- Beginning of /brainix/src/fs/super.c -----
70 void fs_register(bool block, unsigned char maj, pid_t pid)
71 {
72
73     /* Register a device driver with the file system - map a device's major number
74      * to its driver's PID. If the driver for the device containing the root file
75      * system is being registered, mount the root file system and initialize the
76      * root and current working directories. */
77
78     dev_t dev;
79
80     /* Register the device driver with the file system. */
81     driver_pid[block][maj] = pid;
82
83     if (block && maj == ROOT_MAJ)
84     {
85         /* The driver for the device containing the root file system is
86          * being registered. */
87         mount_root();
88         dev = maj_min_to_dev(ROOT_MAJ, ROOT_MIN);
89         fs_proc[FS_PID].root_dir = inode_get(dev, EXT2_ROOT_INO);

```

```

90     fs_proc[FS_PID].work_dir = inode_get(dev, EXT2_ROOT_INO);
91 }
92 }

```

----- End of /brainix/src/fs/device.c -----

This register method is rather straightforward: it adds the device driver to the device driver PID table, then it checks to see if this is a root device we are mounting. If it is, then it calls some additional functions to mount the root file system on the device (line 87), it creates the `dev_t` from the major and minor numbers of the device, assigns the inodes to the root directory and working directory. We shall investigate each of these components of the function in turn.

First, the `mount_root()` method which unsurprisingly mounts the root file system.

mount_root()

```

----- Beginning of /brainix/src/fs/mount.c -----
130 void mount_root(void)
131 {
132
133     /* Mount the root file system. */
134
135     super_t *super_ptr;
136
137     /* Open the device. */
138     dev_open_close(ROOT_DEV, BLOCK, OPEN);
139     if (err_code)
140         /* The device could not be opened. */
141         panic("mount_root", strerror(err_code));
142 }
----- End of /brainix/src/fs/mount.c -----

```

So let us explore this far and say to ourselves “Aha! So, it calls this function ‘`dev_open_close()`’, let’s see what that does exactly!” We look for this method and find it:

dev_open_close()

```

----- Beginning of /brainix/src/fs/super.c -----
146 int dev_open_close(dev_t dev, bool block, bool open)
147 {
148
149     /* If open is true, open a device. Otherwise, close a device. */
150
151     unsigned char maj, min;
152     pid_t pid;
153     msg_t *m;
154     int ret_val;
155 }
----- End of /brainix/src/fs/device.c -----

```

So far several variables are initialized as dummy variables, that is “local variables” which are not used permanently. They are used only temporarily, like a counter.

```

----- Beginning of /brainix/src/fs/super.c -----
156 /* Find the device driver's PID. */
157 dev_to_maj_min(dev, &maj, &min);
158 pid = driver_pid[block][maj];
159 if (pid == NO_PID)
160     return -(err_code = ENXIO);
----- End of /brainix/src/fs/device.c -----

```

We have all ready seen the driver PID table before, but line 157 is completely foreign. We have yet to see exactly what `dev_to_maj_min` does. We can easily locate it however:

dev_to_maj_min()

```

----- Beginning of /brainix/src/fs/super.c -----
32 void dev_to_maj_min(dev_t dev, unsigned char *maj, unsigned char *min)
33 {
34

```

```

35  /* Extract the minor number and the minor number from a device number. */
36
37      *maj = (dev & 0xFF00) >> 8;
38      *min = (dev & 0x00FF) >> 0;
39  }

```

Which is pretty self explanatory code. Line 37 uses bitwise operators to set the Major number to be a modification of the last 8 bits of the `dev`, and line 38 uses bitwise operators to set the Minor number to be a modification of the first 8 bits of the `dev` variable. This code is solid and has been tested, it probably shouldn't need to be changed. At any rate, back to the `dev_open_close()` method:

```

Beginning of /brainix/src/fs/super.c
162  /* Send a message to the device driver. */
163  m = msg_alloc(pid, open ? SYS_OPEN : SYS_CLOSE);
164  m->args.open_close.min = min;
165  msg_send(m);
End of /brainix/src/fs/device.c

```

The code explains itself quite readily. Line 163 allocates a message, line 164 sets the minor number, and line 165 sends the message to the device.

```

Beginning of /brainix/src/fs/super.c
167  /* Await the device driver's reply. */
168  m = msg_receive(pid);
169  ret_val = m->args.open_close.ret_val;
170  msg_free(m);
171  if (ret_val < 0)
172      err_code = -ret_val;
173  return ret_val;
174  }
End of /brainix/src/fs/device.c

```

This code is also self-explanatory. Line 168 waits for the message from the driver, presumably in reply to the message sent from line 165 though this may or may not be the case; line 169 analyzes the message's return value. Now that the return value has been extracted, we can delete the message to free up space (line 170) assuming this wasn't a different message sent by the device driver asking to do some other method, the code - as you can tell - does not check. It does return the return value from the message however (line 173) and catches any possible errors (lines 171-2).

We can assume that the device driver coder knows what he's doing, so we won't investigate the interactions of this message with regards to the device driver. The interested reader can look up the appropriate code in the driver documentation (or supposing that it has yet to be written, which implies the Brainix operating system is still early in development, look at the `/brainix/src/driver/floppy.c`).

We continue our investigation of the `mount_root()` method, after finding out quite a bit about the `dev_open_close()` method.

```

Beginning of /brainix/src/fs/mount.c
143  /* Read the superblock. */
144  super_ptr = super_read(ROOT_DEV);
145  if (err_code)
146      /* The superblock could not be read. */
147      panic("mount_root", strerror(err_code));
End of /brainix/src/fs/mount.c

```

The `mount_root()` reads in the super block from the root directory on the device (line 144). If it could not have been read in, there is a kernel panic (line 147).

We shall now shift our focus onto the `super_read()` method:

super_read()

```

_____ Beginning of /brainix/src/fs/super.c _____
75 super_t *super_read(dev_t dev)
76 {
77
78 /* Read a superblock from its block into the superblock table, and return a
79 * pointer to it. */
80
81     super_t *super_ptr;
82     block_t *block_ptr;
_____ End of /brainix/src/fs/super.c _____

```

Again, as is the style of Brainix, the dummy variables are defined first (lines 81-2) and a brief comment description of the method is given (lines 78-9).

```

_____ Beginning of /brainix/src/fs/super.c _____
84 /* Find a free slot in the table. */
85 if ((super_ptr = super_get(NO_DEV)) == NULL)
86     /* There are no free slots in the table --- too many mounted
87      * file systems. */
88     return NULL;
_____ End of /brainix/src/fs/super.c _____

```

This segment of code from the `super_read()` calls the `super_get()` method. Indeed it is a bit convoluted, but it's the easiest way to program it. Let us now analyze the `super_get()` method:

super_get()

```

_____ Beginning of /brainix/src/fs/super.c _____
54 super_t *super_get(dev_t dev)
55 {
56
57 /* Search the superblock table for a superblock. If it is found, return a
58 * pointer to it. Otherwise, return NULL. */
59
60     super_t *super_ptr;
61
62     /* Search the table for the superblock. */
63     for (super_ptr = &super[0]; super_ptr < &super[NUM_SUPERS]; super_ptr++)
64         if (super_ptr->dev == dev)
65             /* Found the superblock. Return a pointer to it. */
66             return super_ptr;
67
68     /* The superblock is not in the table. */
69     return NULL;
70 }
_____ End of /brainix/src/fs/super.c _____

```

Line 60 initializes the dummy variable, lines 63-66 is a simple, linear for-loop search through the superblock table that is looking for a superblock on the device `dev`. If it is found (line 64), then it returns a pointer to that super block struct (line 66). Supposing that the for loop has run out of places to look on the table, it returns `NULL` indicating that the superblock is not in the table.

Returning our focus to the `super_read()` method:

```

_____ Beginning of /brainix/src/fs/super.c _____
90 /* Copy the superblock from its block into the free slot. */
91 block_ptr = block_get(dev, SUPER_BLOCK);
92 memcpy(super_ptr, block_ptr->data, offsetof(super_t, dev));
_____ End of /brainix/src/fs/super.c _____

```

Lines 91 and 92 introduce two new functions that we will need to investigate: `block_get()` and `memcpy()`. Essentially, line 91 looks on the device `dev` for the block `SUPER_BLOCK`. Get some caffeine in your system, because the `block_get()` method requires more focus and attention:

block_get()

```

166 block_t *block_get(dev_t dev, blkcnt_t blk)
167 {
168
169 /* Search the cache for a block. If it is found, return a pointer to it.
170  * Otherwise, evict a free block, cache the block, and return a pointer to
171  * it. */
172
173     block_t *block_ptr;
174
175     /* Search the cache for the block. */
176     for (block_ptr = lru->prev; ; )
177         if (block_ptr->dev == dev && block_ptr->blk == blk)
178         {
179             /* Found the block. Increment the number of times it is
180              * used, mark it recently used, and return a pointer to
181              * it. */
182             block_ptr->count++;
183             recently_used(block_ptr, MOST);
184             return block_ptr;
185         }
186         else if ((block_ptr = block_ptr->prev) == lru->prev)
187             /* Oops - we've searched the entire cache already. */
188             break;

```

The file system keeps a block cache. Whenever you change anything, it changes the file system's block cache. `Block_put()` writes these changes to the disk, that's the whole point of `block_put()`. Right now, however, we are interested in looking through this cache for a specific block. It may be a little inelegant by most standards, but we are absolved by virtue of this being an operating system ("breaks" do not exist in the Queen's C!).

The cache is searched through until either the specific block in question is found (lines 176-183) or we've run out of cache (we're baroque, that is out of Monet, by line 187). If we do run out of cache, that means the requested block is not cached. Which means there is more to this method than meets the eye:

```

189 /* The requested block is not cached. Search the cache for the least
190  * recently used free block. */
191 for (block_ptr = lru; ; )
192     if (block_ptr->count == 0)
193     {
194         /* Found the least recently used free block. Evict it.
195          * Cache the requested block, mark it recently used, and
196          * return a pointer to it. */
197         block_rw(block_ptr, WRITE);
198         block_ptr->dev = dev;
199         block_ptr->blk = blk;
200         block_ptr->count = 1;
201         block_ptr->dirty = true;
202         block_rw(block_ptr, READ);
203         recently_used(block_ptr, MOST);
204         return block_ptr;
205     }
206     else if ((block_ptr = block_ptr->next) == lru)
207         /* Oops - we've searched the entire cache already. */
208         break;

```

What happens is that the cache is searched through again (this may be a source of inefficiency to search the cache twice, just an aside) for a block that has a `count` of 0. Upon finding it we write the `block_ptr` to that block location (line 197),

and set some new values for our `block_ptr` (lines 198 to 201). We indicate that we have changed the block since it has last been read (that is what the dirty flag indicates...and how long it's been since the block had a bath). This seems intuitively circular to change this only to have it completely ignored by line 202 when `block_rw()` essentially sets the fields of `block_ptr` to whatever the `block_ptr` is. Just as before, there is a method to break out of the for-loop using pragmatic C coding.

Let us now shift our attention to the method `block_rw()`:

block_rw()

```

_____ Beginning of /brainix/src/fs/block.c _____
86 void block_rw(block_t *block_ptr, bool read)
87 {
88
89 /* If read is true, read a block from its device into the cache. Otherwise,
90  * write a block from the cache to its device. */
91
92     dev_t dev = block_ptr->dev;
93     off_t off = block_ptr->blk * BLOCK_SIZE;
94     void *buf = block_ptr->data;
95     super_t *super_ptr;
_____ End of /brainix/src/fs/block.c _____

```

Lines 92-95 are the dummy variables that are used throughout the method, as is usual in the Brainix coding style. Note the comment that tells us what exactly this method does.

```

_____ Beginning of /brainix/src/fs/block.c _____
97     if (!block_ptr->dirty)
98         /* The cached block is already synchronized with the block on
99          * its device. No reason to read or write anything. */
100         return;
_____ End of /brainix/src/fs/block.c _____

```

If the block pointer is dirty, that means that it has changed since last inspected, then `block_ptr->dirty=TRUE`. We are hoping that the `block_ptr` is dirty, that's the entire point of this method. If it's not dirty, we leave the method right here and now.

```

_____ Beginning of /brainix/src/fs/block.c _____
101 /* Read the block from its device into the cache, or write the block
102  * from the cache to its device. */
103     dev_rw(dev, BLOCK, read, off, BLOCK_SIZE, buf);
_____ End of /brainix/src/fs/block.c _____

```

Unfortunately, we have not yet had the good fortune to investigate the `dev_rw()` method, so let us do so now!

dev_rw()

```

_____ Beginning of /brainix/src/fs/super.c _____
179 ssize_t dev_rw(dev_t dev, bool block, bool read, off_t off, size_t size,
180               void *buf)
181 {
182
183 /* If read is true, read from a device. Otherwise, write to a device. */
184
185     unsigned char maj, min;
186     pid_t pid;
187     msg_t *m;
188     ssize_t ret_val;
189
190     /* Find the device driver's PID. */
191     dev_to_maj_min(dev, &maj, &min);
192     pid = driver_pid[block][maj];
193     if (pid == NO_PID)
194         return -(err_code = ENXIO);
_____ End of /brainix/src/fs/device.c _____

```

Lines 185-188 initialize the dummy variables that hold the values for this method. The really interesting part begins at line 190, wherein the device driver's PID is found. If the PID is NO_PID, then an error is returned (line 194).

```

196      Beginning of /brainix/src/fs/super.c
197      /* Send a message to the device driver. */
198      m = msg_alloc(pid, read ? SYS_READ : SYS_WRITE);
199      m->args.read_write.min = min;
200      m->args.read_write.off = off;
201      m->args.read_write.size = size;
202      m->args.read_write.buf = buf;
203      msg_send(m);
204      End of /brainix/src/fs/device.c

```

We hope that the driver can read or write for us, so assuming the driver coder did his or her homework, then we have no worries. We simply allocate a message (line 197), give the message the minor number of the device (198), the offset to read/write (199), the size of the buffer (200), and the buffer to read to or write from (201). The message is then sent.

```

205      Beginning of /brainix/src/fs/super.c
206      /* Await the device driver's reply. */
207      m = msg_receive(pid);
208      ret_val = m->args.read_write.ret_val;
209      msg_free(m);
210      if (ret_val < 0)
211          err_code = -ret_val;
212      return ret_val;
213      End of /brainix/src/fs/device.c

```

We wait for a reply. We assume, and I can't stress this enough, assume that the message from the process with PID `pid` is in response to the message sent. The return value is extracted from the reply (line 206), and we free up the message (207). We check to see if there is an error, and then we return the `ret_val`. It's pretty simple.

Back to our discussion on `block_rw()`, recall the code we left off at was:

```

101      Beginning of /brainix/src/fs/block.c
102      /* Read the block from its device into the cache, or write the block
103       * from the cache to its device. */
104      dev_rw(dev, BLOCK, read, off, BLOCK_SIZE, buf);
105      End of /brainix/src/fs/block.c

```

So now we know what exactly the `dev_rw()` method is, we can understand that line 103 is really asking to read the device `dev`, this is indeed a `BLOCK` device that we are reading from rather than a character device, we are indeed `read`-ing from it with an offset of `off`, we are reading exactly 1 `BLOCK_SIZE` into the buffer `buf` by the method we just inspected above.

```

106      Beginning of /brainix/src/fs/block.c
107      /* The cached block is now synchronized with the block on its device. */
108      block_ptr->dirty = false;
109      if (!read)
110      {
111          super_ptr = super_get(block_ptr->dev);
112          super_ptr->s_wtime = do_time(NULL);
113          super_ptr->dirty = true;
114      }
115      End of /brainix/src/fs/block.c

```

We have not updated the inspection with the `block_ptr` so we may set its `dirty` flag to be `false` (clean as a whistle).

Now, we go on to investigate if we are writing to the file (that is checking that the boolean `read` is false, if it is that means we are of course writing to the block, and we simply follow out lines 109 to 111; however, we are not really interested in that at the moment so we will not really inspect those lines of code here).

Back on track to our analysis of `block_get()`:

```

_____ Beginning of /brainix/src/fs/block.c _____
210  /* There are no free blocks in the cache.  Vomit. */
211  panic("block_get", "no free blocks");
212  return NULL;
213 }
_____ End of /brainix/src/fs/block.c _____

```

Which are the final lines of `block_get()`. It basically calls a kernel panic (line 211) and returns `NULL`, there's nothing elegant needing explanation here.

Back to the `super_read()` method:

```

_____ Beginning of /brainix/src/fs/super.c _____
93  block_put(block_ptr, IMPORTANT);
94
95  return super_ptr;
96 }
_____ End of /brainix/src/fs/super.c _____

```

We have not seen `block_put()` although we have seen `block_get()`. There is a difference between the two, and now we shall investigate `block_put()`:

block_put()

```

_____ Beginning of /brainix/src/fs/block.c _____
218 void block_put(block_t *block_ptr, bool important)
219 {
220
221  /* Decrement the number of times a block is used.  If no one is using it, write
222   * it to its device (if necessary). */
223
224  if (block_ptr == NULL || --block_ptr->count > 0)
225  {
226      return;
227  }
228 }
_____ End of /brainix/src/fs/block.c _____

```

If the `block_ptr` is null, or if the `block_ptr`'s count is greater than 1, then we exit this method. I think this code needs to be modified, as I think line 224 is supposed to be `--block_ptr->count < 0`. This code works however, so I wouldn't touch it just yet (although don't feel discouraged or intimidated when meddling with code!).

```

_____ Beginning of /brainix/src/fs/block.c _____
226  switch (ROBUST)
227  {
228      case PARANOID:
229          block_rw(block_ptr, WRITE);
230          return;
231      case SANE:
232          if (important)
233              block_rw(block_ptr, WRITE);
234          return;
235      case SLOPPY:
236          return;
237  }
238 }
_____ End of /brainix/src/fs/block.c _____

```

Basically, this tells us *when* `block_rw()` should be called. As previously mentioned, the file system has its own block cache, and this method (`block_put()`) essentially writes the changes in this cache. How often it happens depend on the ROBUST-ness

of the configuration of the Brainix kernel. Basically, **SLOPPY** optimizes performance, **PARANOID** always writes blocks when the cache changes slightly, and **SANE** is a center between the two where the cache is written if and only if it is **IMPORTANT**.

Back to the `mount_root()` method which invoked this long aside:

```

149  _____ Beginning of /brainix/src/fs/mount.c _____
150  /* Perform the mount. Fill in the superblock's fields. */
151  super_ptr->s_mtime = do_time(NULL);
152  super_ptr->s_mnt_count++;
153  super_ptr->s_state = EXT2_ERROR_FS;
154  _____ End of /brainix/src/fs/mount.c _____

```

Now that we are mounting the root file system, we have to fill in the super block's fields. We start by setting the time of last mount to be `do_time(NULL)` which is, as far as we know so far, an unknown function. We continue our pattern of looking functions up and find `do_time()`. However, it is a messy system call rather than some ordinary function, so we will not exactly look at the code line by line. It simply gets the number of seconds since January 1, 1970. That is the mount time for the super block (line 150).

Recall that `s_mnt_count` is the number of Mounts since last `fsck`. Line 151 simply tells the super block that it is getting mounted one more time.

Note that in the Brainix `/brainix/inc/fs/super.h` header, we define:

```

_____ Beginning of /brainix/inc/fs/super.h _____
89 #define EXT2_ERROR_FS 2 /* Mounted or uncleanly unmounted. */
_____ End of /brainix/inc/fs/super.h _____

```

so really line 152 of the `mount_root()` method is telling the super block that it's mounted, and nothing more.

```

153 _____ Beginning of /brainix/src/fs/mount.c _____
154 memcpy(super_ptr->s_last_mounted, "\0", 2);
155 _____ End of /brainix/src/fs/mount.c _____

```

We do not know the `memcpy()` method, so allow us to look it up as usual:

memcpy()

```

_____ Beginning of /brainix/src/lib/string.c _____
75 void *memcpy(void *s1, const void *s2, size_t n)
76 {
77
78  /* Copy bytes in memory. */
79
80  char *p1 = s1;
81  const char *p2 = s2;
82
83  for (; n; n--, p1++, p2++)
84    *p1 = *p2;
85  return s1;
86 }
_____ End of /brainix/src/lib/string.c _____

```

This function is pretty straightforward. There is a pair of dummy variables declared (lines 80-81), and then the addresses are switched in a for-loop (lines 83-84). The buffer `void *s1` is returned after the copying (rather, switching of addresses) has occurred.

Back to the line we left off at in the `mount_root()` method:

```

153 _____ Beginning of /brainix/src/fs/mount.c _____
154 memcpy(super_ptr->s_last_mounted, "\0", 2);
155 _____ End of /brainix/src/fs/mount.c _____

```

This basically tells us that we copy to the `s_last_mounted` component of the super block the null character `\0`, we only copy 2 bytes.

```

_____ Beginning of /brainix/src/fs/mount.c _____
154     super_ptr->dev = ROOT_DEV;
155     super_ptr->block_size = 1024 << super_ptr->s_log_block_size;
156     super_ptr->frag_size = super_ptr->s_log_frag_size >= 0 ?
157         1024 << super_ptr->s_log_frag_size :
158         1024 >> -super_ptr->s_log_frag_size;
159     super_ptr->mount_point_inode_ptr = NULL;
_____ End of /brainix/src/fs/mount.c _____

```

We set the `dev` device for the super block to be the `ROOT_DEV` device. The block size is set to be 2^{10} shifted to the right by `s_log_block_size`, and the fragment is variable depending on whether `s_log_frag_size` is less than zero or not. If it is not less than zero, then you shift 2^{10} to the left by `s_log_frag_size`, otherwise you shift to the right by negative one times `s_log_frag_size`.

```

_____ Beginning of /brainix/src/fs/mount.c _____
160     super_ptr->root_dir_inode_ptr = inode_get(ROOT_DEV, EXT2_ROOT_INO);
161     super_ptr->dirty = true;
162 }
_____ End of /brainix/src/fs/mount.c _____

```

Well, we have yet to cover what exactly the `inode_get()` method is, but we know line 161 is telling us that the super block has changed since we last inspected it, which makes sense since we just obtained it and set its values. Let us now have a more intelligible investigation of the `inode_get()` method (besides simply guessing “Well, it has the words ‘get inode’ so I’m guessing it gets an inode...”):

inode_get()

```

_____ Beginning of /brainix/src/fs/inode.c _____
95     inode_t *inode_get(dev_t dev, ino_t ino)
96     {
97
98     /* Search the inode table for an inode.  If it is found, return a pointer to it.
99     * Otherwise, read the inode into the table, and return a pointer to it. */
100
101     inode_t *inode_ptr;
102
103     /* Search the table for the inode. */
104     for (inode_ptr = &inode[0]; inode_ptr < &inode[NUM_INODES]; inode_ptr++)
105         if (inode_ptr->dev == dev && inode_ptr->ino == ino)
106         {
107             /* Found the inode.  Increment the number of times it is
108             * used, and return a pointer to it. */
109             inode_ptr->count++;
110             return inode_ptr;
111         }
_____ End of /brainix/src/fs/inode.c _____

```

The inode table which was initialized previously in section (4.2), we now search the very same table for an inode with the device field equal to `dev` and inode number equal to `ino`. If it is found, the field indicating how many times its been used `count` is incremented (line 109), and the inode pointer to it is returned.

Then there is the case where the inode is not in the table:

```

_____ Beginning of /brainix/src/fs/inode.c _____
113     /* The inode is not in the table.  Find a free slot. */
114     for (inode_ptr = &inode[0]; inode_ptr < &inode[NUM_INODES]; inode_ptr++)
115         if (inode_ptr->count == 0)
116         {
117             /* Found a free slot.  Read the inode into it, and
118             * return a pointer into it. */
119             inode_ptr->dev = dev;
120             inode_ptr->ino = ino;
121             inode_ptr->count = 1;

```

```

122         inode_ptr->mounted = false;
123         inode_ptr->dirty = true;
124         inode_rw(inode_ptr, READ);
125         return inode_ptr;
126     }

```

End of /brainix/src/fs/inode.c

We look for the first inode that has absolutely no references whatsoever (line 115). If a free slot is found, then we simply write the inode into the slot, and return a pointer to it.

However, we also have not seen the method `inode_rw()` before either (line 124). We shall investigate that method now:

inode_rw()

```

Beginning of /brainix/src/fs/inode.c
53 void inode_rw(inode_t *inode_ptr, bool read)
54 {
55
56     /* If read is true, read an inode from its block into the inode table.
57      * Otherwise, write an inode from the table to its block. */
58
59     blkcnt_t blk;
60     size_t offset;
61     block_t *block_ptr;
62     super_t *super_ptr = super_get(inode_ptr->dev);
63
64     if (!inode_ptr->dirty)
65         /* The inode in the table is already synchronized with the inode
66          * on its block. No reason to read or write anything. */
67         return;

```

End of /brainix/src/fs/inode.c

The basic approach is the same as always, check to see if the inode is not dirty (if it isn't, there's no point in writing what's all ready there).

```

Beginning of /brainix/src/fs/inode.c
69 /* Get the block on which the inode resides. */
70 group_find(inode_ptr, &blk, &offset);
71 block_ptr = block_get(inode_ptr->dev, blk);

```

End of /brainix/src/fs/inode.c

The `group_find()` method is completely foreign to us! So, you know the drill, let's look its data structure up first:

```

Beginning of /brainix/inc/fs/group.h
34 typedef struct
35 {
36     unsigned long bg_block_bitmap;    /* First block of block bitmap. */
37     unsigned long bg_inode_bitmap;    /* First block of inode bitmap. */
38     unsigned long bg_inode_table;     /* First block of inode table. */
39     unsigned short bg_free_blocks_count; /* Number of free blocks. */
40     unsigned short bg_free_inodes_count; /* Number of free inodes. */
41     unsigned short bg_used_dirs_count; /* Number of directories. */
42     unsigned short bg_pad;             /* Padding. */
43     unsigned long bg_reserved[3];      /* Reserved. */
44 } group_t;

```

End of /brainix/inc/fs/group.h

I reiterate a main point: this is exactly identical to the ext2 block group data structure. I won't really go in depth into any analysis of it, but merely presented it to point out what it is and where you can find it.

Meanwhile, the `group_find()` method we still need to analyze:

group_find()

```

Beginning of /brainix/src/fs/group.c
32 void group_find(inode_t *inode_ptr, blkcnt_t *blk_ptr, size_t *offset_ptr)
33 {

```

```

34
35 /* Find where an inode resides on its device - its block number and offset
36  * within that block. */
37
38     super_t *super_ptr;
39     unsigned long group;
40     unsigned long index;
41     block_t *block_ptr;
42     group_t *group_ptr;
43     unsigned long inodes_per_block;
44
45     /* From the superblock, calculate the inode's block group and index
46      * within that block group. */
47     super_ptr = super_get(inode_ptr->dev);
48     group = (inode_ptr->ino - 1) / super_ptr->s_inodes_per_group;
49     index = (inode_ptr->ino - 1) \% super_ptr->s_inodes_per_group;

```

Again, the brainix style has the dummy variables defined first (lines 38-43). We then find the superblock in order to calculate out the inode's block group and index therein. The super block extraction occurs on line 47, recall from Code fragment 29 the `super_get()` method.

Lines 48-49 uses bitwise operator black magic to actually come up with the group and index therein, but it is valid black magic.

Continuing our analysis of `group_find()`:

```

Beginning of /brainix/src/fs/group.c
51 /* From the group descriptor, find the first block of the inode
52  * table. */
53     block_ptr = block_get(inode_ptr->dev, GROUP_BLOCK);
54     group_ptr = (group_t *) block_ptr->data;
55     *blk_ptr = group_ptr[group].bg_inode_table;
56     block_put(block_ptr, IMPORTANT);
End of /brainix/src/fs/group.c

```

We see that line 53 uses `block_get()` which was explored in code fragment 31, it gets the block with the inode's dev component, and the `GROUP_BLOCK blkcnt_t`. If you remember we explored briefly the `blkcnt_t` data structure back in code fragment 7...it's basically a (signed) long.

Line 54 casts the data of the `block_ptr` as a `group_t`, which is used later on...the next line as a matter of fact, to find out what the `*blk_ptr` is (or more precisely, assign an address to the `*blk_ptr`).

Then, regardless of the ROBUST-ness of the operating system (except for SLOPPY of course), line 56 writes the block cache to the disk.

```

Beginning of /brainix/src/fs/group.c
58 /* Finally, calculate the block on which an inode resides (it may or may
59  * not be the first block of the inode table) and its offset within that
60  * block. */
61     inodes_per_block = super_ptr->block_size / super_ptr->s_inode_size;
62     *blk_ptr += index / inodes_per_block;
63     *offset_ptr = (index \% inodes_per_block) * super_ptr->s_inode_size;
64 }
End of /brainix/src/fs/group.c

```

The last thing that occurs is the calculation of the block which an inode resides, and its offset therein. This occurs using the same old bitwise black magic that was seen previously.

Back to the `inode_rw()` method that we left off at:

```

Beginning of /brainix/src/fs/inode.c

```

```

69      /* Get the block on which the inode resides. */
70      group_find(inode_ptr, &blk, &offset);
71      block_ptr = block_get(inode_ptr->dev, blk);
      _____ End of /brainix/src/fs/inode.c _____

```

We assign addresses to `blk` and `offset` by use of `group_find()` on line 70. We then invoke `block_get()` (recall from code fragment 31 what exactly `block_get()` is) to assign an address to `block_ptr`.

```

      _____ Beginning of /brainix/src/fs/inode.c _____
73      if (read)
74          /* Read the inode from its block into the table. */
75          memcpy(inode_ptr, &block_ptr->data[offset],
76                super_ptr->s_inode_size);
      _____ End of /brainix/src/fs/inode.c _____

```

This basically reads the inode from the block into the table, by means of copying the address using `memcpy()`.

```

      _____ Beginning of /brainix/src/fs/inode.c _____
77      else
78      {
79          /* Write the inode from the table to its block, and mark the
80           * block dirty. */
81          memcpy(&block_ptr->data[offset], inode_ptr,
82                super_ptr->s_inode_size);
83          block_ptr->dirty = true;
84      }
      _____ End of /brainix/src/fs/inode.c _____

```

This is the alternative case where we are writing to the block from the inode table, and mark the block as dirty. Note which fields of `memcpy()` have changed compared to lines 75-76.

Regardless of which path was taken, we conclude:

```

      _____ Beginning of /brainix/src/fs/inode.c _____
86      /* Put the block on which the inode resides, and mark the inode in the
87       * table as no longer dirty. */
88      block_put(block_ptr, IMPORTANT);
89      inode_ptr->dirty = false;
90  }
      _____ End of /brainix/src/fs/inode.c _____

```

By putting the block wherein the inode resides, and mark the inode in the inode table as clean as a whistle. Recall, again, code fragment 31 for `block_put()`.

We continue on with our analysis of `inode_get()`:

```

      _____ Beginning of /brainix/src/fs/inode.c _____
128     /* There are no free slots in the table. Vomit. */
129     panic("inode_get", "no free inodes");
130     return NULL;
131 }
      _____ End of /brainix/src/fs/inode.c _____

```

Basically, the file system cries if there are no free slots. It calls a kernel panic, and returns empty handed. Such is life I guess.

We have concluded our investigation of the nitty-gritty details of `mount_root()` but we left off at line 88 of `fs_register()` in fragment 19. We shall return to it here:

```

      _____ Beginning of /brainix/src/fs/super.c _____
83      if (block && maj == ROOT_MAJ)
84      {
85          /* The driver for the device containing the root file system is

```

```

86         * being registered. */
87     mount_root();
88     dev = maj_min_to_dev(ROOT_MAJ, ROOT_MIN);
89     fs_proc[FS_PID].root_dir = inode_get(dev, EXT2_ROOT_INO);
90     fs_proc[FS_PID].work_dir = inode_get(dev, EXT2_ROOT_INO);
91 }
92 }

```

Line 87 we just investigated thoroughly, so let us investigate starting with line 88. The `dev` device is assigned based on the device Major and Minor numbers. We have seen `dev_to_maj_min()` but we have not seen `maj_min_to_dev()`, let us try investigating it here:

```

44 dev_t maj_min_to_dev(unsigned char maj, unsigned char min)
45 {
46
47     /* Build the device number from a major number and a minor number. */
48
49     return ((maj & 0xFF) << 8) | ((min & 0xFF) << 0);
50 }

```

This is basically bitwise black magic that undoes the `dev_to_maj_min()`. Note that if we were to look at the `dev` variable as bits (that is, a string of 1s and 0s) and compare it to the `maj` and `min` as bits, we should in theory see that about half way through the `dev` one half looks similar to the `maj` and the other half resembles `min`. That is no coincidence, that is how the `dev` is assigned a value.

Back to the `fs_register()` code:

```

83     if (block && maj == ROOT_MAJ)
84     {
85         /* The driver for the device containing the root file system is
86          * being registered. */
87         mount_root();
88         dev = maj_min_to_dev(ROOT_MAJ, ROOT_MIN);
89         fs_proc[FS_PID].root_dir = inode_get(dev, EXT2_ROOT_INO);
90         fs_proc[FS_PID].work_dir = inode_get(dev, EXT2_ROOT_INO);
91     }
92 }

```

So line 88 basically constructs a device number out of the `ROOT_MAJ` and `ROOT_MIN` numbers. Lines 89-90 assign the `fs_proc` entries of `FS_PID` (the file system process ID) to have a root directory (line 89). Recall `inode_get()` from code fragment 52 takes in as arguments the device number `dev` and the inode to search the inode table for (`EXT2_ROOT_INO` in our case).

Upon completion of those functions, the `fs_register()` method is complete.

Part III

Tutorials on Using the Brainix Operating System

Chapter 2

How to run Brainix on Unix-like Operating Systems and Bochs

OK, you have downloaded your copy of brainix by first installing subversion (you need to install subversion!), then type in:

```
$ svn checkout http://brainix.googlecode.com/svn/trunk/ brainix
```

It should start downloading. You wait until its downloading is complete. I assume that the working directory is `/home/your_user_name_here/`. You type into the command line:

```
$ mkdir mnt
```

Then go into the Brainix directory. Create several scripts:

```
#!/usr/bin/bash
# from brainix/compile.sh

make clean
make
```

The `compile.sh` script.

```
#!/usr/bin/bash
# brainix/update_bootdisk.sh

sudo mount -o loop Bootdisk.img ../mnt/
sudo cp ../bin/Brainix ../mnt
sudo rm ../mnt/Brainix.gz
sudo gzip ../mnt/Brainix
sudo umount ../mnt/
```

This is the `update_bootdisk` script that updates the boot disk with the Brainix binary image (the Brainix binary image is made after compilation and put in `/brainix/bin`). Now, to make a bochs `bochsrc` file. I have included a `dot-bochsrc` file in this documentation (its external to this file), I use it to run Brainix in bochs.

Now to put it all together:

```
#!/usr/bin/bash
# brainix/run.sh
```

```
sh compile.sh
sudo sh update-bootdisk.sh
bochs -f dot-bochsrc
```

This is the script you invoke in order to compile and run the compiled image in bochs.

I have included sample scripts **THAT YOU NEED TO MOVE TO THE BRAINIX FOLDER**, i.e. type into the command line:

```
$ mv compile.sh ../
$ mv update-bootdisk.sh ../
$ mv dot-bochsrc ../
$ mv run.sh ../
```

Chapter 3

How to hack Brainix: Some Things to Bear In Mind

So you're an eager young (or not so young) programmer who wants to start fiddling around with Brainix's code and you'd like to submit what you've done. There are some things that you've got to bear in mind with regard to the documentation and the license.

First and foremost, if you modify pre-existing code, logically it changes how the code works. So you need to modify the documentation. It's important that other people know what the hell is going on, and no one knows that better than the person who changed it!¹ It helps no one if the explanation is in your noodle!

To add a chapter to the Brainix manual, be sure to make your new file a .tex file and make it begin:

```
\chapter{My Chapter Name}
```

If you would like to add code, simply type:

```
\begin{code}[numbers=left,
             firstnumber=...,
             label={[[Beginning of /path/to/code]End of /path/to/code]}]
...
\end{code}
```

Then you are golden...provided you add the code and have a unix-like file path.

Also, don't forget that you **should** comment at the top the changelog. That is, the date, your name, and a one line explanation of what you did. Take for example:

```
/* June 2007 A. R. Hacker, modified the main()
 * function to include a better idle() implementation.
 */
```

There is enough information to know exactly where to look for the changed code. Be sure not to have too much information, e.g. "modified the main() function to

¹A funny anecdote, one day Brainix and Pqnelson were hacking on the kernel and Pqnelson randomly tinkered with some of the code hoping to solve a problem at the time. He succeeded but had no clue what he did. In rare events like this, make sure that the problem is solved and upload your modifications; then you can see what changes were made by the SVN interface on the command line. You can then go and change the documentation

include a modified `idle()` process fork such that an infinite for-loop is used rather than an infinite for-loop so the operating system won't go to hell after the file system registers the floppy disk." This is like listening to the crazy Aunt that rambles on about everything at the family reunion. On the other hand, don't include too little information, e.g. "modified `main()`". This tells us nothing about what you did! It's too cryptic and short, what happened with the `main()` function? Was it completely revamped or did you add a semicolon somewhere?

Don't forget, if you add a completely new file, to include the GPL copyright information. This is in addition to writing the documentation for it, of course.

Chapter 4

On the Implementation of the debugging for Brainix

It seems all too evident that a system requiring programming will require debugging, especially for something as important as an operating system! So it seemed all too obvious that, by virtue of Murphy's Law (if something can go wrong, it will), we need to implement debugging features quickly. The impromptu solution was to present a debug function:

```
/* from /brainix/src/drivers/video.c */
void debug(unsigned int priority, char *message, ...);
void dbug(char *message, ...);
```

Where the priority is ranked from 1 to some large number, 1 being the most important, and compared to `DUM_DEBUG` a constant defined in `kernel.h`. If you want to turn off debugging features, set it to -1. The higher the value for `priority`, the more esoteric the debugging results. Note that there is a shorter function `dbg` which essentially calls `debug` with a priority of 1.

4.1 FAQ

4.1.1 What is the format for debug?

The standard format should be

```
/* typical debug */
void debug(unsigned int priority-SYS_ESOTERIC, "file_name.method(): your message here", ...);
```

The `SYS_ESOTERIC` indicates the system (file system, driver, kernel, where-ever the debug is working in) debug emphasis. So if you have debug *ONLY* the kernel, or *ONLY* a part of the operating system, you go to its header (unless its a driver, then its defined in the driver's .c file) and you change the `SYS_ESOTERIC` to or something bigger. This will cause the debugger to print out only the messages from this system.

4.1.2 I'm making a new driver, and I don't know what to do about this debugger...

Well, suppose your driver file is `X.c`. Near the top, insert the following code:

```
/* from /brainix/src/drivers/X.c */  
#define X_ESOTERIC    0
```

and simply put debugging information that you might find useful as

```
/* from /brainix/src/drivers/X.c */  
    debug(1-X_ESOTERIC, "X.method(): debugging information...");
```

If you are debugging your dear code, then set `X_ESOTERIC` to be some large number like 10. And you don't have to set the value to `1-X_ESOTERIC` it could be any positive number minus `X_ESOTERIC`.

Part IV

Appendices

Appendix A

GNU Free Documentation License

Version 1.2, November 2002

Copyright © 2000,2001,2002 Free Software Foundation, Inc.

Franklin St, Fifth Floor, Boston, MA 02110-1301 USA

Everyone is permitted to copy and distribute verbatim copies of this license document, but changing it is not allowed.

Preamble

The purpose of this License is to make a manual, textbook, or other functional and useful document “free” in the sense of freedom: to assure everyone the effective freedom to copy and redistribute it, with or without modifying it, either commercially or noncommercially. Secondly, this License preserves for the author and publisher a way to get credit for their work, while not being considered responsible for modifications made by others.

This License is a kind of “copyleft”, which means that derivative works of the document must themselves be free in the same sense. It complements the GNU General Public License, which is a copyleft license designed for free software.

We have designed this License in order to use it for manuals for free software, because free software needs free documentation: a free program should come with manuals providing the same freedoms that the software does. But this License is not limited to software manuals; it can be used for any textual work, regardless of subject matter or whether it is published as a printed book. We recommend this License principally for works whose purpose is instruction or reference.

A.1 APPLICABILITY AND DEFINITIONS

This License applies to any manual or other work, in any medium, that contains a notice placed by the copyright holder saying it can be distributed under the terms of this License. Such a notice grants a world-wide, royalty-free license, unlimited in duration, to use that work under the conditions stated herein. The “**Document**”, below, refers to any such manual or work. Any member of the public is a licensee,

and is addressed as “**you**”. You accept the license if you copy, modify or distribute the work in a way requiring permission under copyright law.

A “**Modified Version**” of the Document means any work containing the Document or a portion of it, either copied verbatim, or with modifications and/or translated into another language.

A “**Secondary Section**” is a named appendix or a front-matter section of the Document that deals exclusively with the relationship of the publishers or authors of the Document to the Document’s overall subject (or to related matters) and contains nothing that could fall directly within that overall subject. (Thus, if the Document is in part a textbook of mathematics, a Secondary Section may not explain any mathematics.) The relationship could be a matter of historical connection with the subject or with related matters, or of legal, commercial, philosophical, ethical or political position regarding them.

The “**Invariant Sections**” are certain Secondary Sections whose titles are designated, as being those of Invariant Sections, in the notice that says that the Document is released under this License. If a section does not fit the above definition of Secondary then it is not allowed to be designated as Invariant. The Document may contain zero Invariant Sections. If the Document does not identify any Invariant Sections then there are none.

The “**Cover Texts**” are certain short passages of text that are listed, as Front-Cover Texts or Back-Cover Texts, in the notice that says that the Document is released under this License. A Front-Cover Text may be at most 5 words, and a Back-Cover Text may be at most words.

A “**Transparent**” copy of the Document means a machine-readable copy, represented in a format whose specification is available to the general public, that is suitable for revising the document straightforwardly with generic text editors or (for images composed of pixels) generic paint programs or (for drawings) some widely available drawing editor, and that is suitable for input to text formatters or for automatic translation to a variety of formats suitable for input to text formatters. A copy made in an otherwise Transparent file format whose markup, or absence of markup, has been arranged to thwart or discourage subsequent modification by readers is not Transparent. An image format is not Transparent if used for any substantial amount of text. A copy that is not “Transparent” is called “**Opaque**”.

Examples of suitable formats for Transparent copies include plain ASCII without markup, Texinfo input format, LaTeX input format, SGML or XML using a publicly available DTD, and standard-conforming simple HTML, PostScript or PDF designed for human modification. Examples of transparent image formats include PNG, XCF and JPG. Opaque formats include proprietary formats that can be read and edited only by proprietary word processors, SGML or XML for which the DTD and/or processing tools are not generally available, and the machine-generated HTML, PostScript or PDF produced by some word processors for output purposes only.

The “**Title Page**” means, for a printed book, the title page itself, plus such following pages as are needed to hold, legibly, the material this License requires to appear in the title page. For works in formats which do not have any title page as such, “Title Page” means the text near the most prominent appearance of the work’s title, preceding the beginning of the body of the text.

A section “**Entitled XYZ**” means a named subunit of the Document whose title either is precisely XYZ or contains XYZ in parentheses following text that translates XYZ in another language. (Here XYZ stands for a specific section name mentioned below, such as “**Acknowledgements**”, “**Dedications**”, “**Endorsements**”,

or “**History**”). To “**Preserve the Title**” of such a section when you modify the Document means that it remains a section “Entitled XYZ” according to this definition.

The Document may include Warranty Disclaimers next to the notice which states that this License applies to the Document. These Warranty Disclaimers are considered to be included by reference in this License, but only as regards disclaiming warranties: any other implication that these Warranty Disclaimers may have is void and has no effect on the meaning of this License.

A.2 VERBATIM COPYING

You may copy and distribute the Document in any medium, either commercially or noncommercially, provided that this License, the copyright notices, and the license notice saying this License applies to the Document are reproduced in all copies, and that you add no other conditions whatsoever to those of this License. You may not use technical measures to obstruct or control the reading or further copying of the copies you make or distribute. However, you may accept compensation in exchange for copies. If you distribute a large enough number of copies you must also follow the conditions in section 3.

You may also lend copies, under the same conditions stated above, and you may publicly display copies.

A.3 COPYING IN QUANTITY

If you publish printed copies (or copies in media that commonly have printed covers) of the Document, numbering more than 100, and the Document’s license notice requires Cover Texts, you must enclose the copies in covers that carry, clearly and legibly, all these Cover Texts: Front-Cover Texts on the front cover, and Back-Cover Texts on the back cover. Both covers must also clearly and legibly identify you as the publisher of these copies. The front cover must present the full title with all words of the title equally prominent and visible. You may add other material on the covers in addition. Copying with changes limited to the covers, as long as they preserve the title of the Document and satisfy these conditions, can be treated as verbatim copying in other respects.

If the required texts for either cover are too voluminous to fit legibly, you should put the first ones listed (as many as fit reasonably) on the actual cover, and continue the rest onto adjacent pages.

If you publish or distribute Opaque copies of the Document numbering more than 100, you must either include a machine-readable Transparent copy along with each Opaque copy, or state in or with each Opaque copy a computer-network location from which the general network-using public has access to download using public-standard network protocols a complete Transparent copy of the Document, free of added material. If you use the latter option, you must take reasonably prudent steps, when you begin distribution of Opaque copies in quantity, to ensure that this Transparent copy will remain thus accessible at the stated location until at least one year after the last time you distribute an Opaque copy (directly or through your agents or retailers) of that edition to the public.

It is requested, but not required, that you contact the authors of the Document well before redistributing any large number of copies, to give them a chance to provide you with an updated version of the Document.

A.4 MODIFICATIONS

You may copy and distribute a Modified Version of the Document under the conditions of sections 2 and 3 above, provided that you release the Modified Version under precisely this License, with the Modified Version filling the role of the Document, thus licensing distribution and modification of the Modified Version to whoever possesses a copy of it. In addition, you must do these things in the Modified Version:

- A. Use in the Title Page (and on the covers, if any) a title distinct from that of the Document, and from those of previous versions (which should, if there were any, be listed in the History section of the Document). You may use the same title as a previous version if the original publisher of that version gives permission.
- B. List on the Title Page, as authors, one or more persons or entities responsible for authorship of the modifications in the Modified Version, together with at least five of the principal authors of the Document (all of its principal authors, if it has fewer than five), unless they release you from this requirement.
- C. State on the Title page the name of the publisher of the Modified Version, as the publisher.
- D. Preserve all the copyright notices of the Document.
- E. Add an appropriate copyright notice for your modifications adjacent to the other copyright notices.
- F. Include, immediately after the copyright notices, a license notice giving the public permission to use the Modified Version under the terms of this License, in the form shown in the Addendum below.
- G. Preserve in that license notice the full lists of Invariant Sections and required Cover Texts given in the Document's license notice.
- H. Include an unaltered copy of this License.
- I. Preserve the section Entitled "History", Preserve its Title, and add to it an item stating at least the title, year, new authors, and publisher of the Modified Version as given on the Title Page. If there is no section Entitled "History" in the Document, create one stating the title, year, authors, and publisher of the Document as given on its Title Page, then add an item describing the Modified Version as stated in the previous sentence.
- J. Preserve the network location, if any, given in the Document for public access to a Transparent copy of the Document, and likewise the network locations given in the Document for previous versions it was based on. These may be placed in the "History" section. You may omit a network location for a work that was published at least four years before the Document itself, or if the original publisher of the version it refers to gives permission.
- K. For any section Entitled "Acknowledgements" or "Dedications", Preserve the Title of the section, and preserve in the section all the substance and tone of each of the contributor acknowledgements and/or dedications given therein.

- L. Preserve all the Invariant Sections of the Document, unaltered in their text and in their titles. Section numbers or the equivalent are not considered part of the section titles.
- M. Delete any section Entitled “Endorsements”. Such a section may not be included in the Modified Version.
- N. Do not retitle any existing section to be Entitled “Endorsements” or to conflict in title with any Invariant Section.
- O. Preserve any Warranty Disclaimers.

If the Modified Version includes new front-matter sections or appendices that qualify as Secondary Sections and contain no material copied from the Document, you may at your option designate some or all of these sections as invariant. To do this, add their titles to the list of Invariant Sections in the Modified Version’s license notice. These titles must be distinct from any other section titles.

You may add a section Entitled “Endorsements”, provided it contains nothing but endorsements of your Modified Version by various parties—for example, statements of peer review or that the text has been approved by an organization as the authoritative definition of a standard.

You may add a passage of up to five words as a Front-Cover Text, and a passage of up to words as a Back-Cover Text, to the end of the list of Cover Texts in the Modified Version. Only one passage of Front-Cover Text and one of Back-Cover Text may be added by (or through arrangements made by) any one entity. If the Document already includes a cover text for the same cover, previously added by you or by arrangement made by the same entity you are acting on behalf of, you may not add another; but you may replace the old one, on explicit permission from the previous publisher that added the old one.

The author(s) and publisher(s) of the Document do not by this License give permission to use their names for publicity for or to assert or imply endorsement of any Modified Version.

A.5 COMBINING DOCUMENTS

You may combine the Document with other documents released under this License, under the terms defined in section 4 above for modified versions, provided that you include in the combination all of the Invariant Sections of all of the original documents, unmodified, and list them all as Invariant Sections of your combined work in its license notice, and that you preserve all their Warranty Disclaimers.

The combined work need only contain one copy of this License, and multiple identical Invariant Sections may be replaced with a single copy. If there are multiple Invariant Sections with the same name but different contents, make the title of each such section unique by adding at the end of it, in parentheses, the name of the original author or publisher of that section if known, or else a unique number. Make the same adjustment to the section titles in the list of Invariant Sections in the license notice of the combined work.

In the combination, you must combine any sections Entitled “History” in the various original documents, forming one section Entitled “History”; likewise combine any sections Entitled “Acknowledgements”, and any sections Entitled “Dedications”. You must delete all sections Entitled “Endorsements”.

A.6 COLLECTIONS OF DOCUMENTS

You may make a collection consisting of the Document and other documents released under this License, and replace the individual copies of this License in the various documents with a single copy that is included in the collection, provided that you follow the rules of this License for verbatim copying of each of the documents in all other respects.

You may extract a single document from such a collection, and distribute it individually under this License, provided you insert a copy of this License into the extracted document, and follow this License in all other respects regarding verbatim copying of that document.

A.7 AGGREGATION WITH INDEPENDENT WORKS

A compilation of the Document or its derivatives with other separate and independent documents or works, in or on a volume of a storage or distribution medium, is called an “aggregate” if the copyright resulting from the compilation is not used to limit the legal rights of the compilation’s users beyond what the individual works permit. When the Document is included in an aggregate, this License does not apply to the other works in the aggregate which are not themselves derivative works of the Document.

If the Cover Text requirement of section 3 is applicable to these copies of the Document, then if the Document is less than one half of the entire aggregate, the Document’s Cover Texts may be placed on covers that bracket the Document within the aggregate, or the electronic equivalent of covers if the Document is in electronic form. Otherwise they must appear on printed covers that bracket the whole aggregate.

A.8 TRANSLATION

Translation is considered a kind of modification, so you may distribute translations of the Document under the terms of section 4. Replacing Invariant Sections with translations requires special permission from their copyright holders, but you may include translations of some or all Invariant Sections in addition to the original versions of these Invariant Sections. You may include a translation of this License, and all the license notices in the Document, and any Warranty Disclaimers, provided that you also include the original English version of this License and the original versions of those notices and disclaimers. In case of a disagreement between the translation and the original version of this License or a notice or disclaimer, the original version will prevail.

If a section in the Document is Entitled “Acknowledgements”, “Dedications”, or “History”, the requirement (section 4) to Preserve its Title (section 1) will typically require changing the actual title.

A.9 TERMINATION

You may not copy, modify, sublicense, or distribute the Document except as expressly provided for under this License. Any other attempt to copy, modify, sublicense or distribute the Document is void, and will automatically terminate your

rights under this License. However, parties who have received copies, or rights, from you under this License will not have their licenses terminated so long as such parties remain in full compliance.

A.10 FUTURE REVISIONS OF THIS LICENSE

The Free Software Foundation may publish new, revised versions of the GNU Free Documentation License from time to time. Such new versions will be similar in spirit to the present version, but may differ in detail to address new problems or concerns. See <http://www.gnu.org/copyleft/>.

Each version of the License is given a distinguishing version number. If the Document specifies that a particular numbered version of this License “or any later version” applies to it, you have the option of following the terms and conditions either of that specified version or of any later version that has been published (not as a draft) by the Free Software Foundation. If the Document does not specify a version number of this License, you may choose any version ever published (not as a draft) by the Free Software Foundation.

ADDENDUM: How to use this License for your documents

To use this License in a document you have written, include a copy of the License in the document and put the following copyright and license notices just after the title page:

Copyright © YEAR YOUR NAME. Permission is granted to copy, distribute and/or modify this document under the terms of the GNU Free Documentation License, Version 1.2 or any later version published by the Free Software Foundation; with no Invariant Sections, no Front-Cover Texts, and no Back-Cover Texts. A copy of the license is included in the section entitled “GNU Free Documentation License”.

If you have Invariant Sections, Front-Cover Texts and Back-Cover Texts, replace the “with ... Texts.” line with this:

with the Invariant Sections being LIST THEIR TITLES, with the Front-Cover Texts being LIST, and with the Back-Cover Texts being LIST.

If you have Invariant Sections without Cover Texts, or some other combination of the three, merge those two alternatives to suit the situation.

If your document contains nontrivial examples of program code, we recommend releasing these examples in parallel under your choice of free software license, such as the GNU General Public License, to permit their use in free software.

Bibliography

- [1] The Skelix Operating System Introduction to writing (an operating system including) a file system from scratch <http://en.skelix.org/skelixos/tutorial07.php>
- [2] Andrew Tanenbaum and Albert Woodhull, *Operating Systems: Design and Implementation* Third Edition, Upper Saddle River, New Jersey: Pearson Prentice Hall (2006).
- [3] Daniel P. Bovet and Marco Cesati, *Understanding the Linux Kernel* Sebastopol: O'Reilly Media Inc. (2006).
- [4] Marshall Kirk McKusick and George V. Neville-Neil *The Design and Implementation of the FreeBSD Operating System: FreeBSD Version 5.2* Upper Saddle River: Pearson Education, Inc. (2005).
- [5] Dr. Matloff's "File Systems in Unix" <http://heather.cs.ucdavis.edu/~matloff/UnixAndC/Unix/FileSyst.html>
- [6] Unix File System <http://unixhelp.ed.ac.uk/concepts/fssystem.html>
- [7] The Linux File System Explained, by Mayank Sarup <http://www.freeos.com/articles/3102/>
- [8] Chapter 9 of *The Linux Kernel* Available <http://www.tldp.org/LDP/tlk/fs/filesystem.html>
- [9] Inode Definition by the Linux Information Project <http://www.bellevuelinux.org/inode.html>
- [10] Understanding the Filesystem Inodes http://www.onlamp.com/pub/a/bsd/2001/03/07/FreeBSD_Basics.html
- [11] David A Rusling, The Second Extended File system (EXT2) <http://www.science.unitn.it/~fiorella/guidelinux/tlk/node97.html>
- [12] Understanding the Linux Kernel Second Edition Excerpt www.oreilly.com/catalog/linuxkernel2/chapter/ch17.pdf
- [13] "Setting your file creation mask" http://osr507doc.sco.com/en/OSTut/Setting_your_file_creation_mask.html
- [14] Kevin T. Van Maren *The Fluke Device Driver Framework* (1999). <http://www.cs.utah.edu/flux/papers/vanmaren-thesis.pdf>

- [15] David A Rusling *The Linux Kernel* Chapter 9 <http://www.tldp.org/LDP/tlk/fs/filesystem.html>
- [16] Balazs Gerofi *MINIX VFS: Design and implementation of the MINIX Virtual File system* http://www.minix3.org/doc/gerofi_thesis.pdf
- [17] S.R. Kleiman *Vnodes: An Architecture for Multiple File System Types in Sun UNIX* In Proc. of the Summer 1986 USENIX Technical Conference, June 1991, pp 238-247 <http://www.arl.wustl.edu/~fredk/Courses/cs523/fall01/Papers/kleiman86vnodes.pdf>
- [18] Neil Brown et al. *The Linux Virtual File-system Layer* (1999) <http://www.cse.unsw.edu.au/~neilb/oss/linux-commentary/vfs.html>
- [19] Richard Gooch et al. *Overview of the Linux Virtual File System* (2005) <http://www.gelato.unsw.edu.au/lxr/source/Documentation/kernel-docs.txt>
- [20] Max Bruning *A Comparison of Solaris, Linux, and FreeBSD Kernels* (2005) http://www.opensolaris.org/os/article/2005-10-14_a_comparison_of_solaris__linux__and_freebsd_kernels/
- [21] Vans Information *Virtual file system Part 1* (2001) <http://freeos.com/articles/3838/>
- [22] Vans Information *Virtual file system Part 2* (2001) <http://freeos.com/articles/3851/>
- [23] <http://fxr.watson.org/fxr/source/?v=RELENG62>
- [24] <http://fxr.watson.org/fxr/source/sys/?v=RELENG62>
- [25] Remy Card, Theodore Ts'o, and Stephen Tweedie, *The Design and Implementation of the Ext2 File System* (1994) <http://web.mit.edu/tytso/www/linux/ext2intro.html>
- [26] Anon. *The Ext2 Linux Documentation* (for kernel 2.6.20.1) <http://lxr.linux.no/source/Documentation/filesystems/ext2.txt>