

Caio Araújo Neponoceno de Lima

A Virtual Filesystem Layer implementation in the XV6 Operating System

Salvador - BA, Brazil

2016

Caio Araújo Neponoceno de Lima

A Virtual Filesystem Layer implementation in the XV6 Operating System

A monograph submitted to the Department of Computer Science of the Federal University of Bahia in partial fulfilment of the requirements for the Degree of Bachelor in Computer Science.

Federal University of Bahia – UFBA

Institute of Mathematics

Department of Computer Science

Advisor: Maurício Pamplona Segundo

Salvador - BA, Brazil

2016

Resumo

Este trabalho está apresentando uma implementação de um *Virtual File System* (VFS) que permite mais de um sistema de arquivos coexistir no XV6. Como uma prova de conceito, nós também apresentamos uma implementação básica do sistema de arquivos EXT2, um dos sistemas de arquivos mais populares entre usuários Linux na década de 90. O principal objetivo deste trabalho é documentar uma fonte de conhecimento simples e prática sobre desenvolvimento de sistemas de arquivos para ambientes Unix. Com a simplicidade do XV6, esta implementação do VFS torna possível adicionar suporte a novos sistemas de arquivos para o mesmo com o mínimo de esforço.

Keywords: VFS. Operating Systems. Filesystems. Unix.

Abstract

This work is presenting an implementation of a Virtual Filesystem Layer (VFS) to enable more than one filesystem coexist over the XV6. As a proof of concept, we also present a basic implementation of EXT2 filesystem, one of the most used filesystem by Linux users in 90's years. The main objective of this work is to document a simple and practical source of knowledge about filesystem development over an Unix-like environment. With the simplicity of XV6, this VFS implementation makes possible to add new filesystems support with minimal effort.

Keywords: VFS. Operating Systems. Filesystems. Unix.

Contents

1	INTRODUCTION	7
2	RELATED WORKS	9
2.1	File System Switch	9
2.2	The Sun VFS/Vnode Architecture	10
2.3	The Linux VFS Layer	13
3	XV6 VFS IMPLEMENTATION	17
3.1	XV6 VFS main structures	17
3.1.1	inode	17
3.1.2	superblock	18
3.1.3	<i>filesystem_type</i>	19
3.2	Filesystem-specific operations	20
3.2.1	<i>vfs_operations</i>	21
3.2.2	inode_operations	23
3.3	The <i>mount</i> system call	26
3.3.1	Mount table	26
3.3.2	Modifications over the XV6	30
3.3.3	The block device-filesystem mapping	34
4	IMPLEMENTING A NEW FILESYSTEM ON XV6	35
4.1	The EXT2 filesystem	35
4.1.1	EXT2 disk organization	35
4.1.2	Important EXT2 structures	37
4.1.2.1	<i>ext2_superblock</i>	37
4.1.2.2	<i>ext2_inode</i>	38
4.1.2.3	<i>ext2_dir_entry_2</i>	40
4.1.2.4	struct <i>ext2_block_group_desc</i>	40
4.2	Implementation strategy for new filesystems	41
4.3	<i>vfs_operations</i> for EXT2	41
4.4	<i>inode_operations</i> for EXT2	42
4.5	Configuring and registering the structure <i>filesystem_type</i>	42
4.6	Final remarks	43
5	XV6 VFS EVALUATION	45
5.1	Methodology	45

5.2	Experiments and results	45
6	CONCLUSION	49
6.1	Limitations and future works	49
	BIBLIOGRAPHY	51
	APPENDIX	53
	APPENDIX A – CONFIGURING AN ENVIRONMENT TO BUILD XV6	55
	APPENDIX B – SRC/S5.C	57
	APPENDIX C – SRC/EXT2.C	67

1 Introduction

Basic concepts of operating systems are extremely important to computer developers, mainly because the majority of software is developed to run as an application in an operating system. In this scenario, the operating system is a bottleneck in terms of performance, since it manages all computer resources, such as primary memory and storage access. Nowadays, we have access to cutting-edge commercial operating systems like GNU/Linux or BSD, but their source code are hard to understand. To solve this disparity between practical knowledge about operating system engineering and theoretical concepts, the MIT CSAIL Parallel and Distributed Operating Systems Group developed the XV6.

XV6 is a re-implementation of Dennis Ritchie and Ken Thompson's Unix Version 6 (v6) that runs on modern x86-based multiprocessors and is written in ANSI C. It was developed in the summer of 2006 for MIT's operating systems course, *6.828: Operating System Engineering*. Unlike Unix, Linux and other modern operating systems, XV6's source code is simple enough to be covered in one semester and also allows to cover the main concepts of Unix-like operating systems.

The XV6 is not the unique operating system for education puporses. The Minix, a microkernel-based operating system, was developed by [Tanenbaum \(2007\)](#) to apply the concepts covered in his book. In its third version, the primary objective shifted from education to the creation of a highly reliable and self-healing microkernel operating system, so it is possible to find modern concepts of operating systems since there are a lot of work and academic research being done on it. The biggest difference between XV6 and Minix 3 is in terms of design, since the XV6 is a monolithic kernel and Minix 3 is a microkernel. Unlike XV6, Minix 3 has an extensive code base and implements advanced concepts, such as network stack, which requires a detailed investigation before modifying the source code.

The filesystem¹ implemented in XV6 is a simpler version of original Unix v6 filesystem. However, its design is totally coupled with the rest of the XV6 code, difficulting tasks that require changes in this part of kernel, such as more than one filesystem mounted at the same time. Thus, it is not simple to introduce modern filesystem development techniques in this version of XV6.

This work presents an implementation of a Virtual Filesystem Layer (VFS) to enable having more than one filesystem coexisting over the XV6 and to discuss important concepts that influences modern filesystem design and implementation. As a proof of

¹ According to [Tanenbaum \(2007\)](#), a filesystem is an important kernel component used to control how data is stored and retrieved from disk or other medias. It is the component that organizes the access to data with "files" and "folders".

concept, we also present a basic implementation of the Second Extended Filesystem (EXT2), one of the most basic filesystems used by Linux users.

Our main objective is to document and present a simple and practical source of knowledge about filesystem development in Unix-like environments. Due to the simplicity of XV6, this VFS implementation makes possible to support new filesystems with minimal effort.

The remainder of this monograph is organized as follows. We will discuss some VFS implementations from different operating systems in Chapter 2. Chapter 3, we will present the necessary changes in XV6's original source code that were made to implement a VFS. In Chapter 4, we will show our EXT2 implementation for the XV6 VFS and a strategy to implement new filesystems. We show an experiment to verify the operability of our implementation in Chapter 5, followed by our conclusions and future works in Chapter 6.

2 Related Works

Before having multiple coexisting filesystems, the design of a filesystem was coupled with the kernel. This approach worked well in the past because there was only one filesystem type used by all Unix-like operating system implementation. When BSD Fast File System was created in the University of Berkeley, the necessity to support new filesystems finally became real. It was the first registry of a different filesystem layout since Unix became popular. After that, different filesystems were created, and, as consequence, an abstract layer to manipulate filesystems became a must have feature in commercial kernels. Some examples of this abstract layer are the File System Switch (FSS) in Unix System V Revision 3.0, the Sun VFS/Vnode in SunOS, and the Linux VFS layer, which are discussed below.

2.1 File System Switch

In order to solve the problem of multiple filesystems coexisting in the kernel, the File System Switch (FSS) was introduced in the Unix System V Revision 3.0 (PATE; BOSCH, 2003). As stated in this work, *"one of the main goals of SVR3 was to provide a framework under which multiple different filesystem types could coexist at the same time. Thus each time a call is made to mount, the caller could specify the filesystem type."*

The major difference between the filesystem-independent layer of the kernel and the filesystem-dependent layer is the inclusion of an in-core *inode*¹. The FSS implementation of the in-core inode contains fields that are abstract among different filesystem types (e.g. user and group IDs, file size and access permission), as well as the ability to reference data that is filesystem-specific (PATE; BOSCH, 2003). This way, each filesystem type can have a very different on-disk representation of its metadata/data and still be manipulated by the FSS layer.

The FSS design implements the abstract behavior over filesystem operations through a structure called *fstypsw*, where all filesystem-specific operations are defined. To support multiple filesystems, the kernel maintains a global array of this structure, where each entry represents a possible filesystem. Then, when a system call to access a file is made, the *inode* representing it in memory only needs to have an index of the correct filesystem to access the list of adequate operations. Filesystem-specific operations used by FSS are described in Table 1. As may be observed, the operations in this list are related to a group

¹ According to Tanenbaum (2007), an inode is the representation of a filesystem file used internally by Unix-like kernels to manipulate its information. There are two major types of inode: the on-disk inode, which is the inode stored in disk, and the in-core inode, which is the representation of an on-disk inode in memory.

Table 1 – List of operations for the File System Switch, as presented by [Rodeh, Bacik and Mason \(2013\)](#)

FSS Operation	Description
<code>fs_init</code>	Each filesystem can specify a function that is called during kernel initialization allowing the filesystem to perform any initialization tasks prior the first <i>mount</i> call
<code>fs_iread</code>	Read the inode (during pathname resolution)
<code>fs_iput</code>	Release the inode
<code>fs_iupdat</code>	Update the inode timestamps
<code>fs_readi</code>	Called to read data from file
<code>fs_writei</code>	Called to write data to a file
<code>fs_itrunc</code>	Truncate a file
<code>fs_statf</code>	Return file information required by <i>stat()</i>
<code>fs_namei</code>	Called during pathname transversal
<code>fs_mount</code>	Called to mount a filesystem
<code>fs_umount</code>	Called to unmount a filesystem
<code>fs_getinode</code>	Allocate a file for a pipe
<code>fs_openi</code>	Call the device open routine
<code>fs_closei</code>	Call the device close routine
<code>fs_update</code>	Sync the superblock to disk
<code>fs_statfs</code>	Used by <i>statfs()</i> and <i>ustat()</i>
<code>fs_access</code>	Check access permissions
<code>fs_getdents</code>	Read directory entries
<code>fs_allocmap</code>	Build a block list map for demanding page
<code>fs_freemap</code>	Frees the demand paging block list map
<code>fs_readmap</code>	Read a page using the block list map
<code>fs_setattr</code>	Set file attributes
<code>fs_notify</code>	Notify the filesystem when the file attributes change
<code>fs_fcntl</code>	Handle the <i>fcntl()</i> system call
<code>fs_fsinfo</code>	Return filesystem-specific information
<code>fs_ioctl</code>	Called in response to <i>ioctl()</i> system call

of Unix-like system calls (e.g *mount*, *umount*, *statf* or *ioctl*) and their sub-procedures (e.g *fs_namei* or *fs_access*).

The FSS architecture is the first practical work to support multiple filesystems on Unix. The following architectures presented in this work were widely influenced by it, with some variations on layer organization or filesystem-dependent functions.

2.2 The Sun VFS/Vnode Architecture

This VFS implementation was developed for the Sun Microsystem's SunOS ([KLEIMAN, 1986](#)) and is grounded on the four following goals:

1. *The filesystem implementation should be clearly split into a filesystem independent*

and filesystem-dependent layer. The interface between the two should be well defined.

2. It should support local disk filesystems such as the 4.2BSD Fast File System (FSS), non-UNIX like filesystems such as MS-DOS, stateless filesystems such as NFS, and stateful filesystems such as RFS.
3. It should be able to support the server side of remote filesystems such as NFS and RFS.
4. Filesystem operations across the interface should be atomic such that several operations do not need to be encompassed by locks.

The main difference between FSS and VFS/Vnode is that the VFS/Vnode architecture is based on two major structures, the *vfsops*, the interface for abstract filesystem operations, and the *vnops*, the interface to enable individual file operations. Figure 1 shows a diagram of the VFS/Vnode architecture, which describes a high level interaction between SunOS' components. As may be seen, the VFS/VOP/veneer layer, which is the VFS/Vnode layer, separates other kernel components from supported filesystems. In addition, this diagram shows that the abstract layer can handle diskless filesystems (e.g Network File System (NFS), a network-based filesystem).

Since one goal of this architecture is to support diskless and non-Unix filesystems, the in-core *inode*, also named as *vnode*, stores only the common data among all filesystems, like the FSS in-core *inode* presented before. This structure stores all information required by the filesystem-independent layer and stores a pointer to a private data used by the filesystem-dependent layer.

The VFS layer (represented by the *vfsops*) is responsible to store all operations related to the specific filesystem. The set of operations stored by the VFS layer is described in Table 2.

The Vnode Layer (represented by *vnops*) is where all operations over files are stored. This set of operations is described in Table 3.

Table 3 – List of *vnops* operations, as presented by Pate and Bosch (2003).

VNode Operation	Description
vop_open	This function is only applicable to device special files, files in the namespace that represent hardware devices. It is called once the vnode has been returned from a prior call to voplookup

vop_close	This function is only applicable to device special files. It is called once the vnode has been returned from a prior call to voplookup.
vop_rdwr	Called to read from or write to a file. The information about the I/O is passed through the uio structure.
vop_ioctl	This call invokes an ioctl on the file, a function that can be passed to device drivers.
vop_select	This vnodeop implements select().
vop_getattr	Called in response to system calls such as stat(), this vnodeop fills in a vattr structure, which can be returned to the caller via the stat structure.
vop_setattr	Also using the vattr structure, this vnodeop allows the caller to set various file attributes such as the file size, mode, user ID, group ID, and file times.
vop_access	This vnodeop allows the caller to check the file for read, write, and execute permissions. A cred structure that is passed to this function holds the credentials of the caller.
voplookup	This function replaces part of the old namei() implementation. It takes a directory vnode and a component name and returns the vnode for the component within the directory.
vop_create	This function creates a new file in the specified directory vnode. The file properties are passed in a vattr structure.
vop_remove	This function removes a directory entry.
vop_link	This function implements the link() system call.
vop_rename	This function implements the rename () system call.
vop_mkdir	This function implements the mkdir() system call.
vop_rmdir	This function implements the rmdir() system call.
vop_readdir	This function reads directory entries from the specified directory vnode. It is called in response to the getdents() system call.
vop_symlink	This function implements the symlink() system call.
vop_readlink	This function reads the contents of the symbolic link.
vop_fsync	This function flushes any modified file data in memory to disk. It is called in response to an fsync() system call.
vop_inactive	This function is called when the filesystem-independent layer of the kernel releases its last hold on the vnode. The filesystem can then free the vnode.

vop_bmap	This function is used for demand paging so that the virtual memory (VM) subsystem can map logical file offsets to physical disk offsets.
vop_strategy	This vnodeop is used by the VM and buffer cache layers to read blocks of a file into memory following a previous call to vop_bmap().
vop_bread	This function reads a logical block from the specified vnode and returns a buffer from the buffer cache that references the data.
vop_brelse	This function releases the buffer returned by a previous call to vop_bread.

Our XV6 VFS architecture was strongly influenced by this implementation because of its simplicity and flexibility to support non-Unix and diskless filesystems.

2.3 The Linux VFS Layer

The Linux VFS implementation is one of the most successful in terms of design, mainly because of its complex requirements of portability and performance. A family of data structures represents the abstract file model. The Linux implementation also follows the design of generic structures for the data required by a filesystem-independent layer and maintain data and pointers to filesystem-dependent information.

The four primary structures types of Linux VFS are:

- *superblock*: This structure stores a superblock of a mounted filesystem and contains the global filesystem's metadata;
- *inode*: This structure stores a file and its metadata (e.g. access permissions);
- *dentry*: This structure represents a directory entry, which is a single component of a file;
- *file*: This structure represents an open file attached to a running process.

There is a pointer to a group of operations on each structure presented, which stores the filesystem-dependent functions. These structures describe methods called by the kernel when the filesystem-dependent function need to be executed:

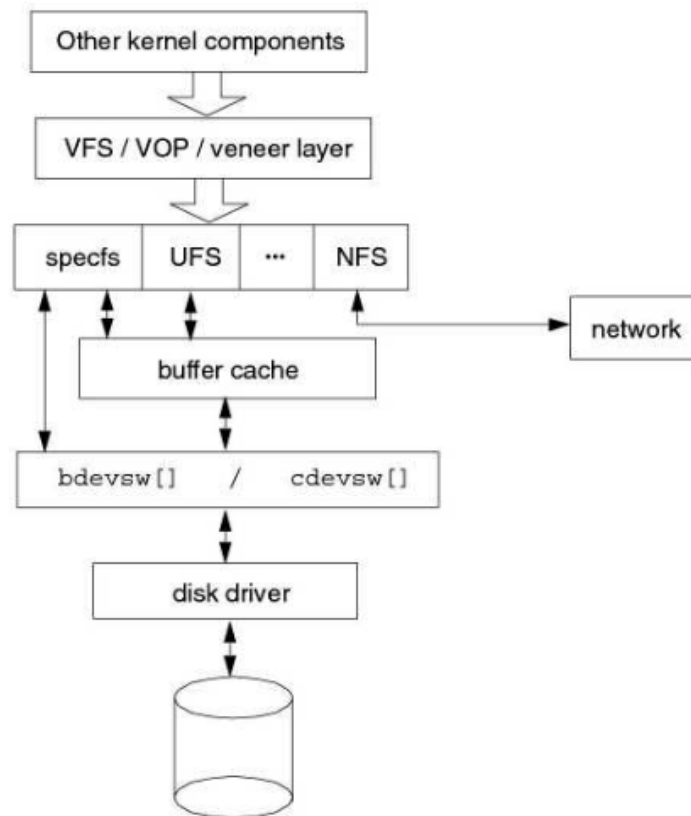


Figure 1 – SunOS VFS/Vnode architecture, extracted from [Pate and Bosch \(2003\)](#)

- *super_operations*: contains methods the kernel invokes on a specific filesystem, such as `write_inode()`.
- *inode_operations*: contains methods to operate over filesystem *inodes*, such as `create()` or `mkdir()`.
- *dentry_operations*: contains operations to manipulate a specific directory entry, such as `d_compare()`.
- *file_operations*: contains operations to manipulate opened files, such as `open()` and `close()`.

These structures and their operations are much more complex than previous presented VFSs. If you are interested in seeing what abstract information is covered by the Linux VFS, we recommend reading the Robert Love's work([LOVE, 2010](#)).

The major advantage of Linux VFS is its constant evolution, because it is the biggest open source project in the world. Thus, there are some features presented in this VFS implementation that were not implemented in others such as *page cache* utilization.

Table 2 – List of *vfsops* operations table, as presented by [Pate and Bosch \(2003\)](#)

VFS Operation	Description
<code>vfs_mount</code>	This function is called to mount the filesystem correctly.
<code>vfs_unmount</code>	This function is called to unmount a filesystem.
<code>vfs_root</code>	This function returns the root vnode for this filesystem and is called during pathname resolution.
<code>vfs_statfs</code>	This function returns filesystem-specific information in response to the <code>statfs()</code> system call. This is used by commands such as <code>df</code> .
<code>vfs_sync</code>	This function flushes file data and filesystem structural data to disk, which provides a level of filesystem hardening by minimizing data loss in the event of a system crash.
<code>vfs_fid</code>	This function is used by NFS to construct a file handle for a specified vnode.
<code>vfs_vget</code>	This function is used by NFS to convert a file handle returned by a previous call to <code>vfs_fid</code> into a vnode on which further operations can be performed.

3 XV6 VFS Implementation

In this section, we present our VFS implementation over XV6. Figure 2 shows the XV6 VFS architecture and it is possible to notice the influence of SunOS VFS/Vnode architecture. This is main contribution of this work, and almost everything that will be discussed later will refer the concepts and designs presented here.

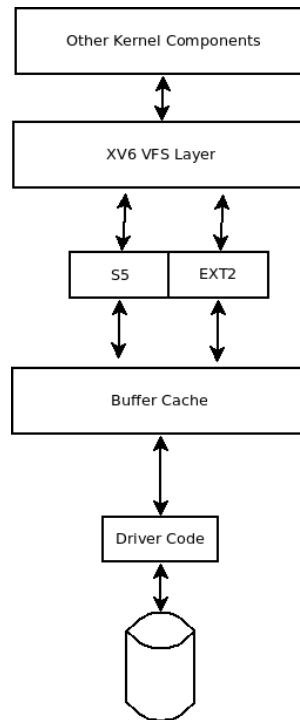


Figure 2 – XV6 VFS architecture

3.1 XV6 VFS main structures

Our design for the XV6 VFS was strongly influenced by the SunOS VFS/Vnode (KLEIMAN, 1986) architecture presented in Section 2.2. There are three major structures to handle filesystem-independent information: *inode*, *superblock* and *filesystem_type*. Details of the implementation of these structures are given below.

3.1.1 inode

The *inode* structure is responsible for representing files (e.g. regular files and directories). Filesystem-related system calls that manipulate the *inode* structure has no knowledge of the filesystem type in use. Our *inode* structure is presented on Listing 3.1.

Listing 3.1 – struct inode

```

1 // in-memory copy of an inode
2 struct inode {
3     uint dev;           // Minor Device number
4     uint inum;          // Inode number
5     int ref;            // Reference count
6     int flags;          // I_BUSY, I_VALID
7     struct filesystem_type *fs_t;
8     struct inode_operations *iops; // Specific inode operations
9     void *i_private;    // Filesystem-dependent information
10
11     short type;         // File type
12     short major;        // Major device number (T_DEV only)
13     short minor;        // Minor device number (T_DEV only)
14     short nlink;        // Number of links to inode in file system
15     uint size;          // Size of file (bytes)
16 };

```

The meaning of each variable in Listing 3.1 is presented on Table 4.

3.1.2 superblock

The *superblock* is the structure responsible for storing filesystem metadata information like total filesystem size, available amount of storage and block size. Filesystem-related system calls manipulate the *superblock* as an abstract representation of the filesystem-independent structure. Different from the Linux *superblock* implementation, that contains much more fields, XV6's *superblock* stores only necessary information to manipulate basic filesystem operations like inode and block allocation, as presented in Listing 3.2.

Listing 3.2 – struct superblock

```

1 struct superblock {
2     // Driver major number of block device this superblock is stored in.
3     int major;
4     // Driver major number of block device this superblock is stored in.
5     int minor;
6     // Block size of this superblock
7     uint blocksize;
8     // Filesystem-specific info
9     void *fs_info;
10    unsigned char s_blocksize_bits;
11    // Superblock Falgs to map its usage
12    int flags;

```

Table 4 – Description of *inode* variables.

Inode field	Description
dev	Stores the minor device number. In commercial operating systems, like Linux, it is prudent to store a structure representing the block device where this <i>inode</i> is stored on. We decided to store the minor device number because the XV6 does not have a block device structure, since it is planned to support only one type of block device (i.e. IDE hard disk).
inum	Stores the <i>inode</i> number, which is the identity of the <i>inode</i> used by all operations that manage it. It is heavily used by the <i>inode</i> cache.
ref	Tracks if this <i>inode</i> is in use. If its value is 0, the <i>inode</i> cache will reuse the space to store information of a new <i>inode</i> .
flags	Stores flags used internally by the VFS algorithms.
fs_t	Is a pointer to fast access filesystem-dependent operations.
iops	Points to filesystem-specific operations used by this <i>inode</i> .
i_private	Stores filesystem-specific data, and usually points to an filesystem-specific <i>inode</i> representation (e.g. "struct ext2_inode*" for an EXT2 filesystem). Its data is used on filesystem-specific functions (e.g. "int ext2_writei(struct inode *ip, char *src, uint off, uint n)" to write data into an <i>inode</i>).
type	Stores the file type, which can be T_DIR, T_FILE, T_DEV and T_MOUNT
major	Stores the device major number when file type equals T_DEV.
minor	Stores the device minor number when file type equals T_DEV.
nlink	Number of links pointing to this <i>inode</i> .
size	Stores the file size.

13 };

Variables in Listing 3.2 are described with more details on Table 5.

3.1.3 *filesystem_type*

The XV6 VFS layer keeps a registration list containing all supported filesystems. To store it into the kernel in an organized way, there is a structure named *filesystem_type* that is responsible for storing important data about a filesystem, such as name, *inode* operations and global operations. This structure is defined in *src/vfs.h* and is shown in Listing 3.3.

Table 5 – Description of *superblock* variables.

Superblock field	Description
major	It is the block device major identifier. It is used internally by the kernel to correctly map what driver is used to access the block device this superblock is stored in.
minor	It is the block device minor identifier. It is used internally by the block device driver to correctly map the correct device this superblock is stored in.
blocksize	It is used to inform the kernel what is the size of the block for this filesystem.
fs_info	It is a generic pointer used to store the filesystem-specific information.
s_blocksize_bits	It is also used to inform the kernel what is the size of the block for this filesystem, but it is used by bitwise operations.
flags	It is used to store flags to internal kernel control.

Listing 3.3 – struct filesystem_type

```

1 struct filesystem_type {
2     char *name;
3     struct vfs_operations *ops;
4     struct inode_operations *iops;
5     struct list_head fs_list;
6 };

```

There is a list pointing to other registered filesystem. To manipulate this list, there are two helper functions used internally by the kernel, both shown in Listing 3.4.

Listing 3.4 – List of helper functions for filesystem_type

```

1 int register_fs(struct filesystem_type *fs);
2 struct filesystem_type* getfs(const char *fs_name);

```

The *register_fs()* is used to install a new filesystem into the internal kernel filesystem list. The *getfs()* is used to retrieve a filesystem type named as *fs_name* (see the *mount* system call in Listing 3.10).

3.2 Filesystem-specific operations

Our XV6 VFS implementation offers an interface between the implemented filesystem and the kernel code through the use of two structures, named *vfs_operations* and *inode_operations*.

3.2.1 *vfs_operations*

This structure stores operations that affect the entire filesystem, and after almost every operations, there is a change over the state of the filesystem. This structure is a list of function pointers and is one of the main components that make the kernel use filesystem operations with a satisfactory abstraction level. It is stored in a structure that manages general filesystem information, such as the *filesystem_type* in Section 3.1.3 or the *inode* in Section 3.1.1.

The operations stored by this structure are presented in Listing 3.5 and detailed in Table 6.

Listing 3.5 – struct *vfs_operations*

```

1
2 struct vfs_operations {
3     int          (*fs_init)(void);
4     int          (*mount)(struct inode *devi, struct inode *ip);
5     int          (*unmount)(struct inode *);
6     struct inode* (*getroot)(int, int);
7     void         (*readsb)(int dev, struct superblock *sb);
8     struct inode* (*ialloc)(uint dev, short type);
9     uint         (*balloc)(uint dev);
10    void         (*bzero)(int dev, int bno);
11    void         (*bfree)(int dev, uint b);
12    void         (*brelse)(struct buf *b);
13    void         (*bwrite)(struct buf *b);
14    struct buf*   (*bread)(uint dev, uint blockno);
15    int          (*namecmp)(const char *s, const char *t);
16 };

```

Table 6 – Description of *vfs_operations* functions.

vfs_operations field	Description
<code>fs_init(void)</code>	This operation is called when the filesystem is being loaded by the kernel (i.e. when the kernel is bootstrapping itself). Its main purpose is to provide an interface where filesystem developers can initialize data structures internally used by their code.

mount(struct inode *devi, struct inode *ip)	This operation is almost self-explained. This function is called when a new instance of the filesystem is being mounted. Usually, a lot of operations are handled in a mount operation, such as read the <i>superblock</i> , initialize device structures like filesystem logger, and add an entry on <i>mount_table</i> . Function parameters are the device inode being mounted (<i>devi</i>) and the directory where the filesystem is going to be mounted (<i>ip</i>).
unmount(struct inode *ip)	This operation is called when the filesystem is going to be unmounted. The parameter <i>ip</i> is the directory being unmounted.
getroot(int major, int minor)	This operation is responsible to allocate, read, fill and return the necessary information of the root <i>inode</i> , including the filesystem-specific information. The parameters <i>major</i> and <i>minor</i> are used to read the root <i>inode</i> from the correct device.
readsb(int dev, struct superblock *sb)	This operation reads the <i>superblock</i> from the device and store it on memory to be manipulated by the kernel. The parameter <i>dev</i> indicate the correct device and <i>sb</i> is a pointer to an already allocated superblock structure. In this version, the implementation is using only the <i>dev</i> as parameter to identify what is the block device the <i>superblock</i> is being read, because we are considering that the XV6 is hardcoded to support only one type of block device. To support more than one type of block device, this parameter shall be a <i>struct block_device</i> .
ialloc(uint dev, short type)	This operation allocates new <i>inodes</i> for this filesystem. This function searches for a free inode in the filesystem (sometimes there is an <i>inode</i> table) and return it to the kernel. When the <i>inode</i> is allocated, the disk is updated to avoid double allocation and data losses. The parameter <i>dev</i> indicates which device this <i>inode</i> is located and <i>type</i> is the <i>inode</i> type (e.g. directory, regular file, device file).
balloc(uint dev)	This operation allocates a block. This function searches for an available block in the filesystem, set the information to avoid double allocation, and return the block number. The parameter <i>dev</i> is the device identifier.

<code>bzero(uint dev, int bno)</code>	This operation is used to fill and write the block <i>bno</i> from <i>dev</i> device with zero and persist this information.
<code>bfree(uint dev, uint b)</code>	This function is the opposite of <i>balloc</i> , where a block <i>b</i> from the device <i>dev</i> is released to be reused in the future. This function is usually called by system calls that delete files or directories such as <i>unlink</i> and <i>rmdir</i> .
<code>brelse(struct buf *b)</code>	This operation is used by the kernel to release the buffer <i>b</i> and make its space available to other process that are trying to get a buffer cache. This function is used to let the filesystem layer handle buffer release operations, but it is usually set as the default internal <i>brelse</i> function.
<code>bwrite(struct buf *b)</code>	This operation writes the content stored on buffer <i>b</i> in the block device. The purpose of this operation is to let filesystem developers handle the <i>bwrite</i> operation and implement custom behaviors. However, the default internal <i>bwrite</i> function is usually used.
<code>bread(uint dev, uint blockno)</code>	This operation is used by the kernel to read the block with number <i>blockno</i> from the device <i>dev</i> . The purpose of this operation is to enable filesystem developers handle the <i>bread</i> operation and implement custom behaviors. Usually, the default internal <i>bwrite</i> function is used.
<code>namecmp(const char *s, const char *t)</code>	This operation defines how the way directory entries' names will be compared. For example, there are some filesystem with case-sensitive and non-case-sensitive rules.

3.2.2 inode_operations

This structure is the main glue between the filesystem code and *inode* operations performed by the kernel. Its operations manipulate the filesystem-dependent *inode* information and return messages of success or failure to the abstract kernel code (i.e. normally filesystem-related system calls). Operations like reading or writing to a file, reading directory entries from a folder and directory lookup are stored in *inode_operations*. This structure is presented in Listing 3.6 and its operations are detailed in Table 7.

Listing 3.6 – struct *inode_operations*

```
1 struct inode_operations {
```

```

2  struct inode* (*dirlookup)(struct inode *dp, char *name, uint *off);
3  void (*iupdate)(struct inode *ip);
4  void (*itrunc)(struct inode *ip);
5  void (*cleanup)(struct inode *ip);
6  uint (*bmap)(struct inode *ip, uint bn);
7  void (*ilock)(struct inode* ip);
8  void (*iunlock)(struct inode* ip);
9  void (*stati)(struct inode *ip, struct stat *st);
10 int (*readi)(struct inode *ip, char *dst, uint off, uint n);
11 int (*writei)(struct inode *ip, char *src, uint off, uint n);
12 int (*dirlink)(struct inode *dp, char *name, uint inum, uint type);
13 int (*unlink)(struct inode *dp, uint off);
14 int (*isdirempty)(struct inode *dp);
15 };

```

Table 7 – Description of *vfs_operations* fuctions.

inode_operations field	Description
dirlookup(struct inode *dp, char *name, uint *off)	This operation has a big role in making the way kernel manages <i>inodes</i> highly abstract. This function verifies if a directory entry named <i>name</i> is a child of the folder <i>dp</i> . If this entry exists, this function reads and returns the corresponding <i>inode</i> and set the parameter <i>off</i> with the byte offset of the entry in the <i>dp</i> data. The internal path to <i>inode</i> translator kernel function <i>namex()</i> , presented in Listing 3.13, is the main caller of this operation.
iupdate(struct inode *ip)	This function is the implementation of the <i>inode</i> update operation performed by kernel in arbitrary system calls (e.g. <i>mkdir</i>). It updates the in-disk <i>inode</i> with the data stored in the in-memory <i>inode ip</i> . The pointer <i>inode->i_private</i> is used to store the filesystem-dependent information to be written on disk.
itrunc(struct inode *ip)	This operation cleans all the information stored in an <i>inode ip</i> .
cleanup(struct inode *ip)	This operation is called by the kernel when an <i>inode</i> is being released, because there is no reference to it. Filesystems should free the <i>inode</i> and its blocks, and optionally, but strongly recommended, erase its content.

bmap(struct inode *ip, uint bn)	This is one of the most important functions used by the VFS layer, because it translates the block <i>bn</i> of the <i>inode ip</i> into the block number in the filesystem. If there is no allocated block for <i>bn</i> , the <i>bmap</i> function allocates one block and returns it. The return value is the filesystem block number.
ilock(struct inode* ip)	This function is called by the kernel VFS layer when it is necessary to lock the access to the <i>inode ip</i> . If the <i>inode ip</i> is already locked, the caller process sleeps and waits until <i>ip</i> becomes available. It is important to the kernel synchronization mechanism.
iunlock(struct inode* ip)	It is the opposite of <i>ilock</i> . This function unlocks the <i>inode ip</i> and wakes up all process waiting for <i>ip</i> . As almost all unlock functions do the same operations, there is an internal kernel function called <i>generic_iunlock</i> that can be used instead.
stati(struct inode *ip, struct stat *st))	This operation is called by the <i>stati</i> system call and is responsible for filling the parameter <i>st</i> with information from the <i>inode ip</i> .
readi(struct inode *ip, char *dst, uint off, uint n)	This function implements the <i>read</i> system call and is internally used by the kernel when is necessary to read a directory entry. The <i>read</i> operation tranfers <i>n</i> bytes from <i>inode ip</i> starting from the byte offset <i>off</i> to the <i>dst</i> buffer. The function returns the number of bytes read.
writel(struct inode *ip, char *src, uint off, uint n)	This function implements the <i>write</i> system call and is internally used by the kernel when is necessary to write to a directory entry. The write operation is performed on <i>inode ip</i> starting from the byte offset <i>off</i> , where <i>n</i> bytes from the buffer <i>src</i> will be written. The function returns the number of bytes written.
dirlink(struct inode *dp, char *name, uint inum, uint type)	This function is called by <i>mkdir</i> and <i>creat</i> system calls. Its purpose is to add the <i>inode inum</i> of type <i>type</i> to the directories map <i>dp</i> with name <i>name</i> . It keeps the filesystem's hierarchical structure updated.
unlink(struct inode *dp, uint off)	This function is called by <i>rmdir</i> and <i>rm</i> system calls. Its objective is to remove the directory entry located starting in the offset byte <i>off</i> and freeing it from the <i>inode dp</i> .

<code>isdirempty(struct</code>	<code>inode</code>	This function checks if the directory <i>dp</i> does not contain
<code>*dp)</code>		directory entries.

3.3 The *mount* system call

The filesystem hierarchy is the interface that a process uses to access files. This abstraction is powerful because application developers do not need to think or even know how its application data will be stored in a block device.

The XV6, as an Unix-like operating system, implements file access through filesystem hierarchy and originally supports only one block device. Being able attach more than one block device is not a VFS feature, but without it VFS would not be useful enough. So, it was necessary to implement a prototype *mount* system call to support multiple block devices.

The *mount* is the operation used to attach a new block device to the filesystem hierarchy. This way, application developers have a high level of abstraction to manipulate data from one block device to another. The XV6 *mount* system call is defined as:

```
int mount(char *special_device_file, char *mount_point_directory, char *filesystem_type)
```

where *special_device_file* is the file specifying which disk will be mounted, *mount_point_directory* is the directory where the new filesystem will be mounted in, and *filesystem_type* is a valid and supported type of the filesystem to be mounted. Different from our implementation, commercial operating systems also offer an interface to pass options to *mount* operations, such as *read-only* or *no-recovery* flags.

To exemplify the *mount* operation, lets consider we have a block device called */dev/hdc* with an EXT2 filesystem, and we want to mount it on */mnt* directory. To do that, we should make the following call:

```
mount("/dev/hdc", "/mnt", "ext2")
```

After the *mount* call, it is possible to access the */dev/hdc* filesystem tree through */mnt*. Figure 3 shows how the filesystem hierarchy will look like after the mount operation.

3.3.1 Mount table

When the operating system supports the mount operation, some data-structures and changes over methods that translate path names into *inodes* (i.e. *namex* function in Listing 3.13) are required because there are cases where the path translation needs to cross *mount* points.

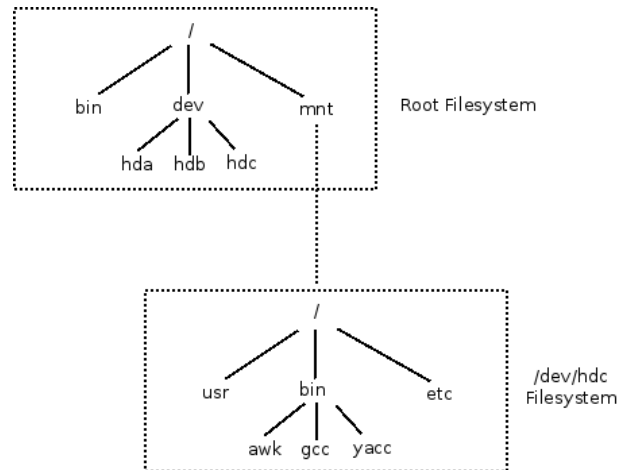


Figure 3 – Filesystem tree with `/dev/hdc` device mounted on `/mnt`. Based on [Bach \(1986\)](#)

The *mount table* is responsible for storing information about mounted filesystems. Our implementation over the XV6 is based on the work of [Bach \(1986\)](#) and represents the *mount table* using two major structures: *mntentry*, representing each table entry; and the global structure *mtable*, representing the table itself. Both *mntentry* and *mtable* are respectively shown in Listings 3.7 and 3.8.

Listing 3.7 – struct mntentry

```

1 // Mount Table Entry
2 struct mntentry {
3     struct inode *m_inode;
4     struct inode *m_rtinode;
5     void *pdata;
6     int dev;
7     int flag;
8 };

```

Listing 3.8 – struct mtable

```

1 // Mount Table Structure
2 struct {
3     struct spinlock lock;
4     struct mntentry mpoint[MOUNTSIZE];
5 } mtable;

```

Each entry in the mount table stores the information presented in Table 8.

The global mount table representation stores an array of mount entries with size *MOUNTSIZE* and a spin lock to be used internally by the kernel to control concurrent access. There are some utility functions defined to manipulate this table, as presented in Listing 3.9.

Table 8 – Description of *mntentry* variables.

mntentry field	Description
<code>m_inode</code>	A pointer to the inode that is the mount point named <i>m_inode</i> . ("/mnt" of the root filesystem in Figure 3).
<code>m_rtinode</code>	A pointer to the inode that is the root of the mounted filesystem. ("/" of the "/dev/hdc" filesystem in Figure 3).
<code>pdata</code>	A pointer to entry's private data (normally it is the superblock).
<code>dev</code>	The block device identifier.
<code>flags</code>	A flag member for internal kernel manipulation.

Listing 3.9 – List of helper functions for the mount table.

```

1 // Utility functions
2 struct inode* mtablertinode(struct inode * ip);
3 struct inode* mtablemntinode(struct inode * ip);
4 int isinoderoot(struct inode* ip);
5 void mountinit(void);

```

The function *mtablertinode()* returns the root *inode* of the mounted filesystem in which *ip* is a mount point. The function *mtablemntinode()* returns the *inode* of the mount point where the root *inode ip* is mounted. Both functions are used by the modified implementation of the *namex* function, shown in Listing 3.13.

The *mount* system call was implemented in *src/sysfile.c*. We present the main idea of this system call in Listing 3.10, omitting error checking to improve readability.

Listing 3.10 – Simplified version of the mount system call.

```

1 int sys_mount(void) {
2     char *devf; char *path; char *fstype;
3     struct inode *ip, *evi;
4     // Handle syscall arguments
5     if (argstr(0, &devf) < 0 ||
6         argstr(1, &path) < 0 ||
7         argstr(2, &fstype) < 0) {
8         return -1;
9     }
10    // Get inodes
11    if ((ip = namei(path)) == 0 ||
12        (devi = namei(devf)) == 0) {
13        return -1;
14    }
15
16    struct filesystem_type *fs_t = getfs(fstype);

```

```

17 // Open the device and check if everything is ok.
18 bdev_open(devi);
19 // Add this to a list of filesystem type
20 putvfsonlist(devi->major, devi->minor, fs_t);
21 // Call specific fs mount operation.
22 fs_t->ops->mount(devi, ip);
23 //Turn the current ip into a mount point
24 ip->type = T_MOUNT;
25 return 0;
26 }

```

From Line 2 to Line 14, this function is setting up the local variables by parsing the system call arguments and getting the necessary information to perform the *mount* operation. Line 16 checks if the *fstype* is supported by the kernel. Line 18 opens the device to be mounted and checks if it is possible to access the hardware without error. There is a registration of this mount operation indicating the device identifiers (i.e. major and minor numbers) and the filesystem type in Line 20. Line 22 is using the VFS layer to call the filesystem-specific *mount* operation. This operation will read the device *superblock* and request a *mount entry* on the *mount table*. To complete the operation, the inode representing the mount point is marked as *T_MOUNT*.

Filesystem-specific *mount* operations change for different filesystem types. To help the comprehension of the *mount* system call, a simple version of the filesystem-specific *mount* operation for the XV6's default filesystem, that we named S5, is shown in Listing 3.11. You can find the complete implementation in *src/s5.c* (see Appendix B).

Listing 3.11 – S5's mount operation handler.

```

1 int s5_mount(struct inode *devi, struct inode *ip) {
2     struct mntentry *mp;
3     s5_ops.readsb(devi->minor, &sb[devi->minor]);
4     struct inode *devrtip = s5_ops.getroot(devi->major, devi->minor);
5     for (mp = &mtable.mpoint[0]; mp < &mtable.mpoint[MOUNTSIZE]; mp++) {
6         // This slot is available
7         if (mp->flag == 0) {
8             found_slot:
9                 mp->dev = devi->minor;
10                mp->m_inode = ip;
11                mp->pdata = &sb[devi->minor];
12                mp->flag |= M_USED;
13                mp->m_rtinode = devrtip;
14                initlog(devi->minor);
15                return 0;

```

```

16     } else {
17         // The disk is already mounted
18         if (mp->dev == devi->minor) {
19             return -1;
20         }
21         if (ip->dev == mp->m_inode->dev &&
22             ip->inum == mp->m_inode->inum)
23             goto found_slot;
24     }
25 }
26 return -1;
27 }

```

Almost every operation performed in *s5_mount* can be shared with other filesystems. The *superblock* is read in Line 3 through the *s5_ops* structure, which is global in *src/s5.c*. Line 4 reads the root inode of the device being mounted. The loop between Lines 5 and 25 searches for a empty entry on mount table and, when it finds one, the *mntentry* is set. Line 14 initializes the log system for this filesystem. In addition, this implementation does not enable a device be mounted twice, as may be seen in Line 18. Finally, Line 21 checks if the mount point is already an entry and updates it to point to the new device to be mounted.

To exemplify the cases discussed in Section 3.3, Figure 4 illustrates the relationship between *Inode Table* and *Mount Table*. The XV6 implement this diagram with the support of the *mtabmntinode()* and *mtablertinode()* functions, both presented in Listing 3.9.

3.3.2 Modifications over the XV6

The *mount* operation required changes over the XV6 code to be implemented. The first change was over the IDE driver code, because it was hard coded to support just two IDE disk and both slots were already in use as boot disk and root XV6's filesystem. We changed the driver to use the Slave Bus (TECHNOLOGY, 1993) and now it is possible to attach 4 IDE devices on XV6. Also, *namex* and *iget* functions were changed to support path translation with crossing mount points.

The updated *iget* function checks if the required inode is a mount point (i.e. its type is *T_MOUNT*). If it is true, it finds the *mount table entry* for this *inode*, then get the root *inode* of the mounted filesystem using the *mtablertinode()*, and return it as the requested *inode*. This algorithm ensures that a path translation crossing mount points follows the direction from the mount point to the mounted filesystem correctly. Listing 3.12 shows the updated version of *iget()* function.

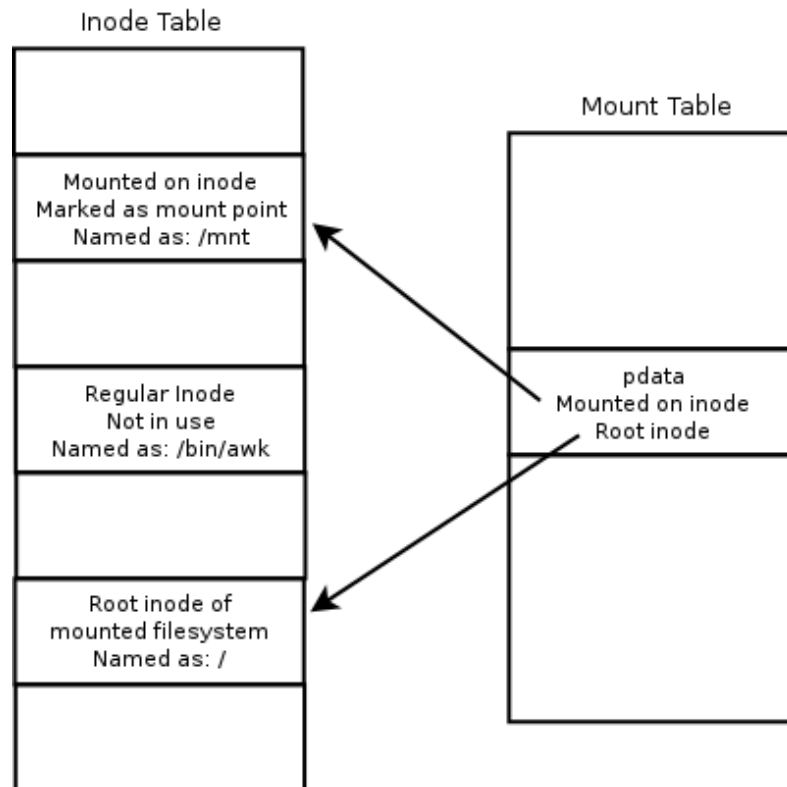


Figure 4 – Diagram showing the relation between *Inode Table* and *Mount Table* based on [Bach \(1986\)](#).

Listing 3.12 – `iget()` - updated function supporting path translation with crossing mount points.

```

1 struct inode* iget(uint dev, uint inum,
2                     int (*fill_inode)(struct inode *)){
3     struct inode *ip, *empty;
4     struct filesystem_type *fs_t;
5     acquire(&icache.lock);
6     empty = 0;
7
8     // Is the inode already cached?
9     for(ip = &icache.inode[0]; ip < &icache.inode[NINODE]; ip++){
10         if(ip->ref > 0 && ip->dev == dev && ip->inum == inum){
11             // If the current inode is an mount point
12             if (ip->type == T_MOUNT) {
13                 struct inode *rinode = mtablertinode(ip);
14                 if (rinode == 0) {
15                     panic("Invalid_Inode_on_Mount_Table");
16                 }
17                 rinode->ref++;
18                 release(&icache.lock);

```

```

19         return rinode;
20     }
21     ip->ref++;
22     release(&icache.lock);
23     return ip;
24 }
25 if(empty == 0 && ip->ref == 0)    // Remember empty slot.
26     empty = ip;
27 }
28 // Recycle an inode cache entry.
29 if(empty == 0)
30     panic("iget: no inodes");
31 fs_t = getvfsentry(IDEMAJOR, dev)->fs_t;
32 ip = empty;
33 ip->dev = dev;
34 ip->inum = inum;
35 ip->ref = 1;
36 ip->flags = 0;
37 ip->fs_t = fs_t;
38 ip->iops = fs_t->iops;
39 release(&icache.lock);
40 if (!fill_inode(ip)) {
41     panic("Error on fill_inode");
42 }
43 return ip;
44 }

```

This implementation of the *iget()*, however, does not handle path translation in the opposite direction (i.e. coming from the mounted device to the mount point direction). This case can happen when you go back in the filesystem tree (i.e. the path has `..`). To handle this problem, the *namex* function had to be modified, as presented in Listing 3.13.

Listing 3.13 – *namex()* - updated function supporting path translation with crossing mount points.

```

1 static struct inode*
2 namex(char *path, int nameiparent, char *name) {
3     struct inode *ip, *next, *ir;
4     if(*path == '/')
5         ip = rootfs->fs_t->ops->getroot(IDEMAJOR, ROOTDEV);
6     else
7         ip = idup(proc->cwd);
8     while((path = skipelem(path, name)) != 0){

```

```

9      ip->iops->ilock(ip);
10     if(ip->type != T_DIR){
11         iunlockput(ip);
12         return 0;
13     }
14     if(nameiparent && *path == '\0'){
15         // Stop one level early.
16         ip->iops->iunlock(ip);
17         return ip;
18     }
19     component_search:
20     if((next = ip->iops->dirlookup(ip, name, 0)) == 0){
21         iunlockput(ip);
22         return 0;
23     }
24     ir = next->fs_t->ops->getroot(IDEMAJOR, next->dev);
25     if (next->inum == ir->inum &&
26         isinoderoot(ip) &&
27         (strncmp(name, "..", 2) == 0)) {
28         struct inode *mntinode = mtablemntinode(ip);
29         iunlockput(ip);
30         ip = mntinode;
31         ip->iops->ilock(ip);
32         ip->ref++;
33         goto component_search;
34     }
35     iunlockput(ip);
36     ip = next;
37 }
38 if(nameiparent){
39     iput(ip);
40     return 0;
41 }
42 return ip;
43 }
```

The modification to support crossing mount points in the *namei* function was made between Lines 20 and 35. It checks if the reached *inode* is a root *inode* and if the previous path component is *".."*. If it is true, it is a root *inode* from a mounted filesystem and we have to find the mount point *inode* into the mount table (Line 29) to look for the *".."* directory entry from it. This is the main reason why the mount point is required to be a folder, once all folders contain at least *"."* and *".."* directory entries. This change ensures

that a path translation crossing a mount point from the mounted filesystem to the mount point will be performed correctly.

3.3.3 The block device-filesystem mapping

Some kernel operations must be aware about what is the filesystem type for the block device it is manipulating (see Listing 3.12 which requires this information to setup the *inode* correctly). To avoid wasting time looking for this information in the mount table, XV6 VFS contains an object called *vfsmlist* that is responsible for storing a list of *vfs*. This structure is defined on *src/vfs.h* and is shown in Listing ??.

Listing 3.14 – struct vfs

```

1 struct vfs {
2     int major;
3     int minor;
4     int flag;
5     struct filesystem_type *fs_t;
6     struct list_head fs_next; // Next mounted on vfs
7 };

```

With this structure, it is possible to map a device through its major and minor identifiers into a filesystem type. It is important to say that this mapping should be implemented using hash, but as the performance is not the major concern of this work, we used a linked list instead.

Listing 3.15 – List of helper functions for vfs.

```

1 struct vfs* getvfscopy(int major, int minor);
2 int putvfsonlist(int major, int minor, struct filesystem_type *fs_t);

```

To abstract the use of this list, XV6 VFS offers two helpers function: *getvfscopy()*, used to retrieve a *vfs* reference for a device and *putvfsonlist()*, used to link device to its filesystem type.

4 Implementing a new filesystem on XV6

In this chapter, we describe the necessary steps to implement a new filesystem using our XV6 VFS. To illustrate this process, we use a basic version of EXT2 ([CARD; TS'O; TWEEDIE, 2010](#)). All listings presented in this chapter can be found in *src/ext2.c* and *src/ext2.h* (see Appendix C). Before discussing implementation, we present an overview of EXT2's concepts. To a complete documentation of EXT2's design and implementation, please read the official documentation ([POIRIER, 2011](#)).

4.1 The EXT2 filesystem

EXT2 is a block-based filesystem implemented by Rémy Card, Theodore Ts'o and Stephen Tweedie to substitute the Extended Filesystem, maintaining old internal structures while providing new functionalities. It was first released on January 1993 as part of the Linux kernel and was further used as the standard filesystem in different major Linux distributions. Structures from EXT4 and EXT3 were strongly influenced by the EXT2 internal structure. Unlike recent filesystems, EXT2 does not support any optimization feature, such as journaling, journal checksums or extents. However, due its simplicity, it is a good start point for filesystem developers.

4.1.1 EXT2 disk organization

The EXT2 disk organization is strongly based on the layout of the BSD filesystem ([MCKUSICK et al., 1984](#)). Unlike previous filesystems, EXT2 is physically divided in block groups to improve sequential access, since it allows to allocate related data, such as directories and files, physically near each other. The physical structure of an EXT2 filesystem is illustrated in Figure 5.



Figure 5 – EXT2 filesystem architecture.

There is a boot sector in the first 1024 bytes. Then, follows N block groups, where each block group is divided in blocks with 1 up to 8KB. The number of block groups, inodes and blocks is varies depending on the partition size and the block size. These parameters can be configured when the filesystem is being installed in a device using the *mkfs.ext2* utility.

The first block group contains a copy of important filesystem control information, such as superblock and filesystem descriptors (e.g. block group descriptor table) as well as part of the filesystem itself, with a block bitmap, an inode bitmap, a piece of the inode table and data blocks, as shown in Figure 6.

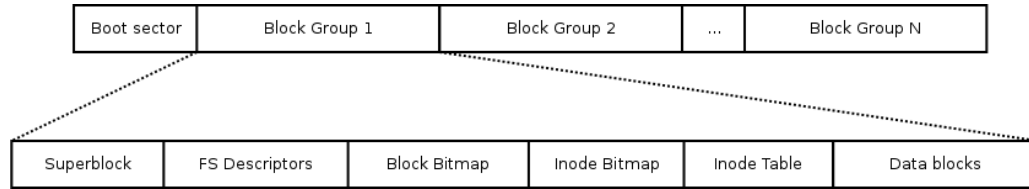


Figure 6 – Layout of the first EXT2 Block Group.

Table 9 shows a layout of a 20MB EXT2 filesystem with block size of 1KB and block group size of 8MB. As may be observed, there are backups for *superblock* and *filesystem descriptors*. They increase the reliability, since they make it possible to recover the original data in case of corruption.

Table 9 – Sample 20MB Ext2 filesystem using 1KiB block size. Based on Poirier (2011)

Block Offset	Length	Description
byte 0	512 bytes	boot record (if present)
byte 512	512 bytes	additional boot record data (if present)
– block group 1, blocks 1 to 8192 –		
byte 1024	1024 bytes	superblock
block 2	1 block	filesystem descriptor table
block 3	1 block	block bitmap
block 4	1 block	inode bitmap
block 5	214 blocks	inode table
block 219	7974 blocks	data blocks
– block group 2, blocks 8193 to 16384 –		
block 8193	1 block	superblock backup
block 8194	1 block	filesystem descriptor table backup
block 8195	1 block	block bitmap
block 8196	1 block	inode bitmap
block 8197	214 blocks	inode table
block 8408	7974 blocks	data blocks
– block group 3, blocks 16385 to 24576 –		
block 16385	1 block	block bitmap
block 16386	1 block	inode bitmap
block 16387	214 blocks	inode table
block 16601	3879 blocks	data blocks

The layout of a disk is predictable when block size, number of blocks per group, inodes per group are known as well. These information are located in the superblock structure, and the EXT2 implementations use these values to compute the correct offset of an inode entry on the *inode table*, to find a specific data block, and so on.

4.1.2 Important EXT2 structures

Every filesystem requires internal and specific structure representations to enable its data manipulation with a certain level of abstraction. These structures must strictly follow the documented layout, because they are used to access the raw data from block devices, and if wrong operations are performed, they may end up corrupting the filesystem.

The main structures in EXT2 are *ext2_superblock* for EXT2 *superblock* manipulation, *ext2_inode* for EXT2 *inode* manipulation, *ext2_dir_entry_2* for EXT2 directory entry manipulation, and *ext2_block_group_desc* for EXT2 block group manipulation.

4.1.2.1 *ext2_superblock*

The structure *ext2_superblock* is used to manage EXT2's physical superblock. The stored information is presented in Listing 4.1, and its complete documentation is available on the EXT2 documentation ([POIRIER, 2011](#)).

Listing 4.1 – struct *ext2_superblock*

```

1 struct ext2_superblock {
2     uint32 s_inodes_count;    /* Inodes count */
3     uint32 s_blocks_count;    /* Blocks count */
4     uint32 s_r_blocks_count;  /* Reserved blocks count */
5     uint32 s_free_blocks_count; /* Free blocks count */
6     uint32 s_free_inodes_count; /* Free inodes count */
7     uint32 s_first_data_block; /* First Data Block */
8     uint32 s_log_block_size;  /* Block size */
9     uint32 s_log_frag_size;   /* Fragment size */
10    uint32 s_blocks_per_group; /* # Blocks per group */
11    uint32 s_frags_per_group;  /* # Fragments per group */
12    uint32 s_inodes_per_group; /* # Inodes per group */
13    uint32 s_mtime;           /* Mount time */
14    uint32 s_wtime;           /* Write time */
15    uint16 s_mnt_count;        /* Mount count */
16    uint16 s_max_mnt_count;    /* Maximal mount count */
17    uint16 s_magic;            /* Magic signature */
18    uint16 s_state;            /* File system state */
19    uint16 s_errors;           /* Behaviour when detecting errors */
20    uint16 s_minor_rev_level;  /* minor revision level */
21    uint32 s_lastcheck;        /* time of last check */
22    uint32 s_checkinterval;    /* max. time between checks */
23    uint32 s_creator_os;       /* OS */
24    uint32 s_rev_level;        /* Revision level */
25    uint16 s_def_resuid;       /* Default uid for reserved blocks */

```

```

26  uint16 s_def_resgid;      /* Default gid for reserved blocks */
27  uint32 s_first_ino;      /* First non-reserved inode */
28  uint16 s_inode_size;     /* size of inode structure */
29  uint16 s_block_group_nr; /* block group # of this superblock */
30  uint32 s_feature_compat; /* compatible feature set */
31  uint32 s_feature_incompat; /* incompatible feature set */
32  uint32 s_feature_ro_compat; /* readonly-compatible feature set */
33  uint8 s_uuid[16];        /* 128-bit uuid for volume */
34  char s_volume_name[16];  /* volume name */
35  char s_last_mounted[64]; /* directory where last mounted */
36  uint32 s_algorithm_usage_bitmap; /* For compression */
37  uint8 s_prealloc_blocks; /* Nr of blocks to try to preallocate */
38  uint8 s_prealloc_dir_blocks; /* Nr to preallocate for dirs */
39  uint16 s_padding1;
40  uint8 s_journal_uuid[16]; /* uuid of journal superblock */
41  uint32 s_journal_inum;    /* inode number of journal file */
42  uint32 s_journal_dev;    /* device number of journal file */
43  uint32 s_last_orphan;    /* start of list of inodes to delete */
44  uint32 s_hash_seed[4];   /* HTREE hash seed */
45  uint8 s_def_hash_version; /* Default hash version to use */
46  uint8 s_reserved_char_pad;
47  uint16 s_reserved_word_pad;
48  uint32 s_default_mount_opts;
49  uint32 s_first_meta_bg;  /* First metablock block group */
50  uint32 s_reserved[190]; /* Padding to the end of the block */
51  };

```

4.1.2.2 ext2_inode

The structure `ext2_inode` keeps track of every directory, regular file, symbolic link or special file stored in the filesystem. It stores their location, size, type and access rights. Filenames are not stored in the inode itself, since this information is contained in directory entries. Listing 4.2 presents this structure.

Listing 4.2 – struct ext2_inode

```

1  struct ext2_inode {
2      uint16 i_mode; /* File mode */
3      uint16 i_uid;  /* Low 16 bits of Owner Uid */
4      uint32 i_size; /* Size in bytes */
5      uint32 i_atime; /* Access time */
6      uint32 i_ctime; /* Creation time */
7      uint32 i_mtime; /* Modification time */

```



```

8     uint32 i_dtime; /* Deletion Time */
9     uint16 i_gid;   /* Low 16 bits of Group Id */
10    uint16 i_links_count; /* Links count */
11    uint32 i_blocks; /* Blocks count */
12    uint32 i_flags;  /* File flags */
13    union {
14        struct {
15            uint32 l_i_reserved1;
16        } linux1;
17        struct {
18            uint32 h_i_translator;
19        } hurd1;
20        struct {
21            uint32 m_i_reserved1;
22        } masix1;
23    } osd1; /* OS dependent 1 */
24    uint32 i_block[EXT2_N_BLOCKS]; /* Pointers to blocks */
25    uint32 i_generation; /* File version (for NFS) */
26    uint32 i_file_acl; /* File ACL */
27    uint32 i_dir_acl; /* Directory ACL */
28    uint32 i_faddr; /* Fragment address */
29    union {
30        struct {
31            uint8 l_i_frag; /* Fragment number */
32            uint8 l_i_fsize; /* Fragment size */
33            uint16 i_pad1;
34            uint16 l_i_uid_high; /* these 2 fields */
35            uint16 l_i_gid_high; /* were reserved2[0] */
36            uint32 l_i_reserved2;
37        } linux2;
38        struct {
39            uint8 h_i_frag; /* Fragment number */
40            uint8 h_i_fsize; /* Fragment size */
41            uint16 h_i_mode_high;
42            uint16 h_i_uid_high;
43            uint16 h_i_gid_high;
44            uint32 h_i_author;
45        } hurd2;
46        struct {
47            uint8 m_i_frag; /* Fragment number */
48            uint8 m_i_fsize; /* Fragment size */
49            uint16 m_pad1;

```

```

50     uint32 m_i_reserved2[2];
51 } masix2;
52 } osd2;    /* OS dependent 2 */
53 };

```

It is important to say that, even if the EXT2 implementation does not use some members of this structure, it is necessary to keep them to avoid corrupting the filesystem meta-data.

4.1.2.3 ext2_dir_entry_2

EXT2's directory entries are stored by a linked list, and each entry contains the *inode* number, total entry length, name length, file type and filename. Listing 4.3 presents this structure.

Listing 4.3 – struct ext2_dir_entry_2

```

1 struct ext2_dir_entry_2 {
2     uint32 inode;        /* Inode number */
3     uint16 rec_len;      /* Directory entry length */
4     uint8  name_len;     /* Name length */
5     uint8  file_type;
6     char   name[];       /* File name, up to EXT2_NAME_LEN */
7 };

```

4.1.2.4 struct ext2_block_group_desc

This structure stores the description of block groups. It provides the location of *inode bitmap*, *inode table*, *block bitmap*, *free blocks count* and other useful information to manage each block group. Its instances are stored by the *Block Group Descriptor Table*, an array stored immediately after the superblock (see Figure 6 and Table 9). Listing 4.4 presents this structure.

Listing 4.4 – struct ext2_block_group_descriptor

```

1 struct ext2_group_desc {
2     uint32 bg_block_bitmap;    /* Blocks bitmap block */
3     uint32 bg_inode_bitmap;    /* Inodes bitmap block */
4     uint32 bg_inode_table;     /* Inodes table block */
5     uint16 bg_free_blocks_count; /* Free blocks count */
6     uint16 bg_free_inodes_count; /* Free inodes count */
7     uint16 bg_used_dirs_count;  /* Directories count */
8     uint16 bg_pad;
9     uint32 bg_reserved[3];
10 };

```

4.2 Implementation strategy for new filesystems

The task of adding a new filesystem to the kernel is trivial in XV6 VFS. The triviality, however, is not extended to the internal filesystem implementation. It is not a wise idea to implement all filesystem features and only after that start the validation phase. Thus, a good strategy is to create all filesystem-dependent operations as empty operations that call the *panic* function. After this step, configure and register the new filesystem and start coding one operation at time for each system call.

This strategy will let you test and debug your code in parts and locate errors with more precision, thanks to granularity of the code.

Since the EXT2 implementation is not the objective of this work, readers are referred to the Linux kernel implementation (POIRIER, 2011) or to check the file *src/ext2.c* for further details.

4.3 *vfs_operations* for EXT2

As discussed in Section 3.2.1, it is necessary to create an object pointing to the filesystem-dependent general operations. Listing 4.5 shows the *vfs_operations* structure for EXT2 as an example.

Listing 4.5 – *vfs_operations* instance for EXT2.

```
1 struct vfs_operations ext2_ops = {
2     .fs_init = &ext2fs_init,
3     .mount   = &ext2_mount,
4     .unmount = &ext2_unmount,
5     .getroot = &ext2_getroot,
6     .readsb  = &ext2_readsb,
7     .ialloc  = &ext2_ialloc,
8     .balloc  = &ext2_balloc,
9     .bzero   = &ext2_bzero,
10    .bfree   = &ext2_bfree,
11    .brelse  = &brelse,
12    .bwrite  = &bwrite,
13    .bread   = &bread,
14    .namecmp = &ext2_namecmp
15 };
```

Each EXT2-specific function starts with *"ext2_"*. It is an important convention to follow when programming in C because it works similar to a namespace and avoids compilation errors due to multiple definition of identifiers. Operations *brelse*, *bwrite* and

bread are pointing to internal kernel generic operations because they do not need a filesystem-specific implementation.

4.4 *inode_operations* for EXT2

As stated in Section 3.2.2, it is also necessary to create an instance of *inode_operations* pointing to filesystem-dependent inode operations. Listing 4.6 shows the *inode_operations* structure for EXT2.

Listing 4.6 – *inode_operations* instance for EXT2.

```

1 struct inode_operations ext2_iops = {
2     .dirlookup = &ext2_dirlookup,
3     .iupdate   = &ext2_iupdate,
4     .itrunc    = &ext2_itrunc,
5     .cleanup   = &ext2_cleanup,
6     .bmap      = &ext2_bmap,
7     .ilock     = &ext2_ilock,
8     .iunlock   = &generic_iunlock,
9     .stati     = &generic_stati,
10    .readi     = &generic_readi,
11    .writei    = &ext2_writei,
12    .dirlink   = &ext2_dirlink,
13    .unlink    = &ext2_unlink,
14    .isdirempty = &ext2_isdirempty
15 };

```

As may be seen, the structure *inode_operations* also defines some generic operations: *iunlock*, *stati* and *readi*. The same naming convention was used for filesystem-specific functions.

4.5 Configuring and registering the structure *filesystem_type*

One of the most important steps to support a new filesystem in our XV6 VFS is to create the structure *filesystem_type* and populate its variables. It is implemented in *src/ext2.c* because there is no reason to keep a filesystem-specific code global. It basically points to the structures in Sections 4.3 and 4.4 and stores the filesystem name, as be seen in Listing 4.7.

Listing 4.7 – *filesystem_type* instance for EXT2.

```

1 struct filesystem_type ext2fs = {
2     .name = "ext2",

```

```

3     .ops = &ext2_ops,
4     .iops = &ext2_iops
5 };

```

Then it is necessary to inform the kernel that there is a new filesystem to be supported. To implement this, we use the function *register_fs()*, which was introduced in Section 3.1.3. Our EXT2 implementation contains a function named *initext2fs(void)* to initialize all internal data, including the registration of the filesystem, as shown in Listing 4.8.

Listing 4.8 – *initext2fs()* -EXT2 initialization function.

```

1 int initext2fs(void) {
2     initlock(&ext2_sb_pool.lock, "ext2_sb_pool");
3     return register_fs(&ext2fs);
4 }

```

Unlike the Linux kernel, we cannot load modules dynamically on XV6, so the function *initext2fs()* had to be hardcoded into the kernel initialization code. To organize filesystem initializations, the function *initfss()* was created as in Listing 4.9.

Listing 4.9 – Kernel function to initialize filesystems.

```

1 static void initfss(void) {
2     if (inits5fs() != 0) // init s5 fs
3         panic("S5 not registered");
4     if (initext2fs() != 0) // init ext2 fs
5         panic("ext2 not registered");
6 }

```

Following these steps, the mount system call presented in Section 3.3 will be able to support devices formatted with the EXT2 filesystem.

4.6 Final remarks

Our EXT2 implementation was based on Linux's version. However, many changes were necessary because of the internal kernel data manipulation in XV6. The first change was to remove byte endianness compatibility present in Linux. As XV6 is designed to run into a *X86* architecture, we removed all instructions to convert the byte endianness. Also, another change was the usage of *buffer_head* structure to manage blocks of the filesystem in XV6, which diverge from Linux's implementation that uses *page cache*. It is important to know that the structures and functions presented in Sections 4.3 and 4.4 can be implemented based on available implementations of the desired filesystem.

5 XV6 VFS evaluation

5.1 Methodology

In this chapter, we evaluate the VFS operability on XV6. With our EXT2 implementation, it is possible to check the interoperability with other operating systems that support this filesystem. The environment used to perform this experiment was a virtual machine running *Debian 3.2.65-1* with *Linux kernel 3.2.0-4-amd64*, which was used to compile XV6's code and to create EXT2 filesystems with *mkfs.ext2*. The XV6 runs in an i386 machine emulated with *qemu-system-i386* (check [Appendix A](#) to configure an execution environment).

The major goal of this experiment is to show that our XV6 VFS works properly and that EXT2 block devices are totally operational. To this end, operations like create and read directories, read directory entries, write or delete files are going to be performed on XV6, and, after that, the filesystem is going to be mounted on Linux to check if they were performed correctly. Then, the reverse direction is going to be considered. To do that, the filesystem is going to be modified on Linux and then mounted on XV6 to check if it is possible to successfully access the modifications. Considering that Linux's EXT2 implementation is a widely used and stable commercial filesystem, used as base for filesystems like EXT3 and EXT4, its implementation using our XV6 VFS implementation validates our architectural design.

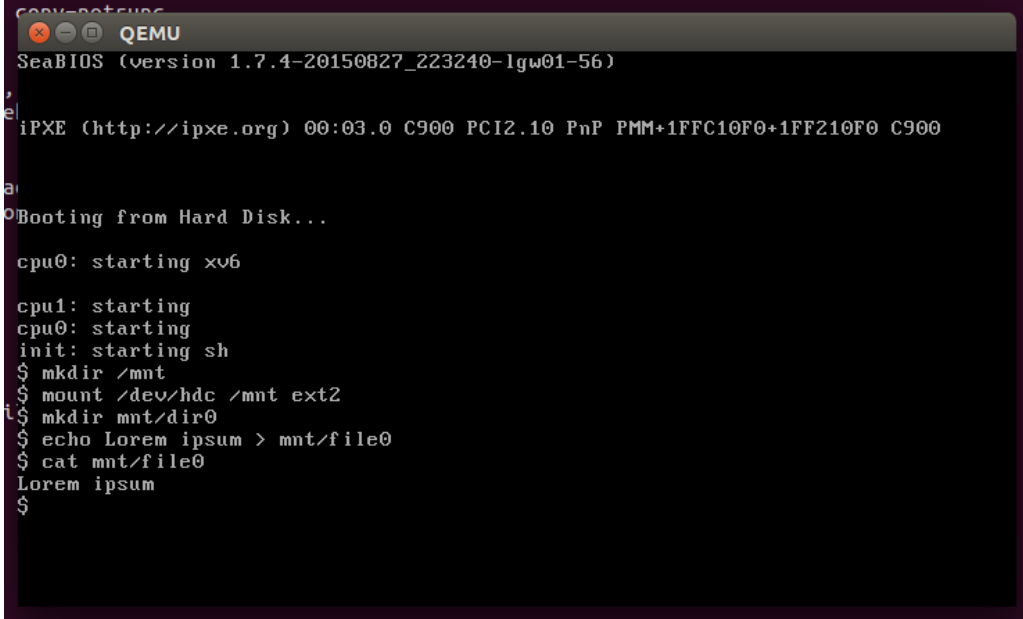
5.2 Experiments and results

XV6's terminal does not support script automation, so our experiments had to be performed manually. In the first part of this evaluation, we ran commands shown in [Listing 5.1](#), and the obtained results are presented in [Figure 7](#).

Listing 5.1 – List of commands that modify the EXT2 filesystem on XV6.

```
1 mkdir /mnt
2 mount /dev/hdc /mnt ext2
3 mkdir mnt/dir0
4 echo Lorem ipsum > mnt/file0
5 cat mnt/file0
```

These commands may look simple, but a lot of work was done by XV6 to process them. The first two commands are necessary to mount the EXT2 filesystem in a directory. The *mount* program, which execute the *mount* system call, requires three parameters: the



```

QEMU
SeaBIOS (version 1.7.4-20150827_223240-lgw01-56)

iPXE (http://ipxe.org) 00:03.0 C900 PCI2.10 PnP PMM+1FFC10F0+1FF210F0 C900

Booting from Hard Disk...

cpu0: starting xv6

cpu1: starting
cpu0: starting
init: starting sh
$ mkdir /mnt
$ mount /dev/hdc /mnt ext2
$ mkdir mnt/dir0
$ echo Lorem ipsum > mnt/file0
$ cat mnt/file0
Lorem ipsum
$

```

Figure 7 – Execution of commands that modify the EXT2 filesystem on XV6.

device to be mounted, the directory where the device will be mounted, and the type of the filesystem contained in this device. After that, a directory named *dir0* is created on this device. Finally, to test if file writing is being done correctly, Line 4 writes "*Lorem ipsum*" into *file0* and Line 5 reads its contents to check if the write operation was successful.

To make sure these modifications are correct in the device, we mounted this filesystem on Linux and ran another sequence of commands, as shown in Listing 5.2 and Figure 8.

Listing 5.2 – List of commands that modify the EXT2 filesystem on Linux.

```

1 sudo losetup /dev/loop0 src/ext2.img
2 sudo mount /dev/loop0 mnt/
3 sudo cat mnt/file0
4 sudo cp -R /usr/include mnt/
5 sudo cat mnt/include/termio.h
6 sudo umount mnt
7 sudo losetup -d /dev/loop0

```

Line 1 uses tool of Linux kernel to copy the filesystem image located at *src/ext2.img* to a *loop device* to allow mounting this image as a virtual block device. Line 2 mounts the filesystem image into *mnt/*. Line 3 prints the content of *file0* in the terminal. Line 4 copies a complex directory structure located in */usr/include/* to *mnt/*. Lines 5 and 6 finish the manipulation of the EXT2 image.

Finally, we check if the manipulations made on Linux can be correctly loaded on XV6. This is done with the commands in Listing 5.3, and the results are shown in Figure 9.


```

vagrant@wheezy-amd64:/vagrant$ sudo losetup /dev/loop0 src/ext2.img
vagrant@wheezy-amd64:/vagrant$ sudo mount /dev/loop0 mnt/
vagrant@wheezy-amd64:/vagrant$ sudo cat mnt/file0
Lorem ipsum
vagrant@wheezy-amd64:/vagrant$ sudo cp -R /usr/include mnt/
vagrant@wheezy-amd64:/vagrant$ sudo cat mnt/include/termio.h
/* Compatible <termio.h> for old 'struct termio' ioctl interface.
   This is obsolete; use the POSIX.1 'struct termios' interface
   defined in <termios.h> instead. */

#include <termios.h>
#include <sys/ioctl.h>
vagrant@wheezy-amd64:/vagrant$ sudo umount mnt
vagrant@wheezy-amd64:/vagrant$ sudo losetup -d /dev/loop0
vagrant@wheezy-amd64:/vagrant$

```

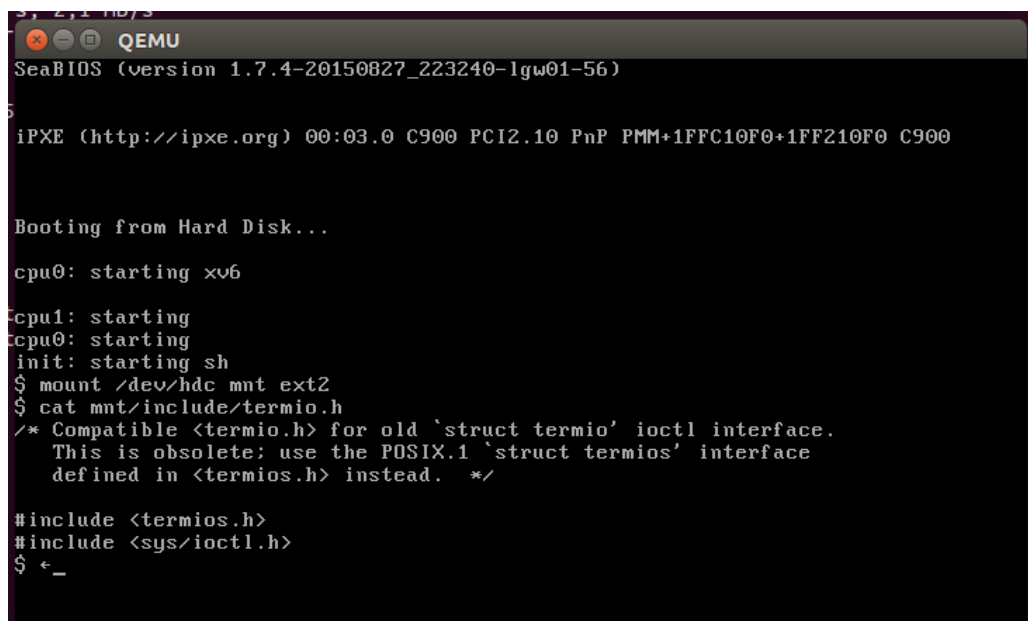
Figure 8 – Execution of commands that modify the EXT2 filesystem on XV6.

Listing 5.3 – List of commands to verify modifications in the EXT2 filesystem on XV6.

```

1 mount /dev/hdc /mnt ext2
2 cat mnt/include/termio.h

```



```

SeaBIOS (version 1.7.4-20150827_223240-1gw01-56)

iPXE (http://ipxe.org) 00:03.0 C900 PCI2.10 PnP PMM+1FFC10F0+1FF210F0 C900

Booting from Hard Disk...

cpu0: starting xv6
cpu1: starting
cpu0: starting
init: starting sh
$ mount /dev/hdc mnt ext2
$ cat mnt/include/termio.h
/* Compatible <termio.h> for old 'struct termio' ioctl interface.
   This is obsolete; use the POSIX.1 'struct termios' interface
   defined in <termios.h> instead. */

#include <termios.h>
#include <sys/ioctl.h>
$

```

Figure 9 – Execution of commands to verify modifications in the EXT2 filesystem on XV6.

After Line 2, it is possible to verify that the content of the file *"mnt/include/termio.h"* was correctly printed in the terminal because we see the same output on Linux (see Figure 8). It is important to say that the command *cat* uses the system call *read*, which is implemented using two major VFS operations: *readi()* and *dirlookup()* (see section 3.2.2).

With these experiments, we were able to validate the major features of the EXT2 filesystem running on XV6. Unmodified native XV6 commands were used to manipulate the filesystem, commands that were implemented using XV6 system calls. This behavior could be achieved only because our XV6 VFS implementation does not change the system call interface, just their internal functions, achieving the desired behavior of a VFS

implementation in terms of abstraction. In addition, this experiment shows that our VFS allows adding new filesystems to XV6 without compromising its operation.

6 Conclusion

This work revealed to us one of the major advantages of a VFS, the interoperability of filesystems among operating systems. The compatibility with more than one filesystem is a very important feature in a modern operating system, what gave us the intuition that VFS design and implementation is a topic that should be taught to operating system engineers.

The major contribution of this work is the implementation and documentation of a simple, but powerful VFS layer in an operating system that is designed for academic purposes, making it a great start point for operating system developers. Our EXT2 implementation indicates that the VFS design achieved the desired abstraction. Another important contribution of our work is porting XV6's filesystem to run in the VFS layer, which is also part of the validation, since all operations performed in the root filesystem are using the VFS layer. In addition, it shows that our XV6 VFS allows using more than one type of filesystem at the same time.

6.1 Limitations and future works

As a future work, it would be interesting to implement a non Unix-like filesystem to validate the power of the VFS architecture presented in this work. The same thing is applied to diskless filesystems like *procfs*, *sysfs* or *NFS*.

In addition, there are some limitations in this XV6 VFS implementation that can improve its abstraction level. The first improvement that can be done is to make the VFS compatible with extent-based filesystem. In our implementation, we assume that all filesystems are block-based, which is not true for modern filesystems like BTRFS, EXT4 or ZFS.

Second, it is necessary to implement a new system call to read directory entries. Without this system call, it is not possible implement programs like *ls* in an elegant way.

Third, XV6 should be modified to enable global system time access in order to update the last time an *inode* is modified, and to improve the memory usage of Block I/O subsystem, since it is currently consuming 4KB per block on buffer cache even if the block is smaller than that.

Finally, it is important to create operations that support Access Control List patterns on the VFS layer. It was not implemented in this version because there is only one user on XV6, the root user.

Bibliography

BACH, M. J. *The Design of the UNIX Operating System*. Upper Saddle River, NJ, USA: Prentice-Hall, Inc., 1986. ISBN 0-13-201799-7. Referenced 2 times in pages 27 and 31.

CARD, R.; TS'O, T.; TWEEDIE, S. Design and implementation of the second extended filesystem. *Proceedings of the First Dutch International Symposium on Linux*, 2010. Available from Internet: <<http://web.mit.edu/tytso/www/linux/ext2intro.html>>. Referenced in page 35.

KLEIMAN, S. R. Vnodes: An architecture for multiple file system types in sun unix. In: . [S.l.: s.n.], 1986. p. 238–247. Referenced 2 times in pages 10 and 17.

LOVE, R. *Linux Kernel Development*. 3rd. ed. [S.l.]: Addison-Wesley Professional, 2010. ISBN 0672329468, 9780672329463. Referenced in page 14.

MCKUSICK, M. K. et al. A fast file system for unix. *ACM Trans. Comput. Syst.*, ACM, New York, NY, USA, v. 2, n. 3, p. 181–197, ago. 1984. ISSN 0734-2071. Available from Internet: <<http://doi.acm.org/10.1145/989.990>>. Referenced in page 35.

PATE, S.; BOSCH, F. V. D. *UNIX Filesystems: Evolution, Design and Impemenation*. New York, NY, USA: John Wiley & Sons, Inc., 2003. ISBN 0471164836. Referenced 4 times in pages 9, 11, 14, and 15.

POIRIER, D. The second extended file system: Internal layout. 2011. Available from Internet: <<http://www.nongnu.org/ext2-doc/ext2.pdf>>. Referenced 4 times in pages 35, 36, 37, and 41.

RODEH, O.; BACIK, J.; MASON, C. Btrfs: The linux b-tree filesystem. *Trans. Storage*, ACM, New York, NY, USA, v. 9, n. 3, p. 9:1–9:32, ago. 2013. ISSN 1553-3077. Available from Internet: <<http://doi.acm.org/10.1145/2501620.2501623>>. Referenced in page 10.

TANENBAUM, A. S. *Modern Operating Systems*. 3rd. ed. Upper Saddle River, NJ, USA: Prentice Hall Press, 2007. ISBN 9780136006633. Referenced 2 times in pages 7 and 9.

TECHNOLOGY, I. S. Ata interface reference manual. 1993. Available from Internet: <<ftp://ftp.seagate.com/acrobat/reference/111-1c.pdf>>. Referenced in page 30.

Appendix

APPENDIX A – Configuring an environment to build XV6

There are two ways to build and run the xv6 OS. The first one is using a pre-configured vagrant machine and we recommend follow this way to avoid headaches. The second solution we recommend follow the official documentation that can be found in <https://pdos.csail.mit.edu/6.828/2014/tools.html>

Required software:

- Vagrant: <https://www.vagrantup.com/downloads.html>
- git: <https://git-scm.com/book/en/v2/Getting-Started-Installing-Git>
- QEMU : https://en.wikibooks.org/wiki/QEMU/Installing_QEMU

To configure an environment using Vagrant, the first thing to do is download and install the Vagrant software. You can follow the instructions available in Vagrant's link given above. If you are not use to use vagrant, don't worry. All configurations are already done on this Vagrantfile (<https://gist.github.com/caiolima/fdc6974c1fec0e57caac>). The workflow using the Vagrant is simple: share the folder where the XV6's source code is between Vagrant machine and host operating system. This way, you can decide your own development environment (IDE, Editor, etc.).

After install and configure the vagrant environment, you need to download our XV6 code and place it in the shared folder between Vagrant machine and your host operating system. It is recommended clone the source code on the same folder where the Vagrantfile is placed. You can find this repository on <https://github.com/caiolima/xv6-public> or download the code using the following command:

```
1 git clone https://github.com/caiolima/xv6-public.git
```

After the clone command, your folder tree should have *Vagrantfile* and *xv6-public* folder. Our example we placed these files in */xv6-dev* folder.

Now, we are almost in the end of the build step. We need to run the commands:

```
1 cd ~/xv6-dev
2 vagrant up
3 vagrant ssh
```

The first time you run *vagrant up* you will need to wait for a while because the Vagrant need to download the virtual machine image. After the download, the load time will be faster.

When you successfully ssh into the vagrant virtual machine, run the following commands to build the xv6 OS:

```
1 cd /vagrant/xv6-public/src
2 make
```

To create an valid EXT2 filesystem, you need to run the following commands on /vagrant/xv6-public/src folder:

```
1 dd if=/dev/zero of=ext2.img bs=30M count=1
2 sudo mkfs.ext2 -b 1024 -T "EXT2_TEST" ext2.img
```

To run and test the XV6, you need to open the terminal application of your host operating system and run:

```
1 make qemu
```

Finally, the QEMU will start and boot the XV6.

APPENDIX B – src/s5.c

```

1  // It is the s5 filesystem implementation
2
3  #include "types.h"
4  #include "defs.h"
5  #include "param.h"
6  #include "stat.h"
7  #include "mmu.h"
8  #include "proc.h"
9  #include "spinlock.h"
10 #include "vfs.h"
11 #include "buf.h"
12 #include "file.h"
13 #include "vfsmount.h"
14 #include "s5.h"
15
16 /*
17  * Its is a pool to allocate s5 inodes structs.
18  * We use it because we don't have a kmalloc function.
19  * With an kmalloc implementatios, it need to be removed.
20  */
21 static struct {
22     struct spinlock lock;
23     struct s5_inode s5_i_entry[NINODE];
24 } s5_inode_pool;
25
26 struct s5_inode*
27 alloc_s5_inode()
28 {
29     struct s5_inode *ip;
30
31     acquire(&s5_inode_pool.lock);
32     for (ip = &s5_inode_pool.s5_i_entry[0]; ip < &s5_inode_pool.s5_i_entry[NINODE]; ip++) {
33         if (ip->flag == S5_INODE_FREE) {
34             ip->flag |= S5_INODE_USED;
35             release(&s5_inode_pool.lock);
36
37             return ip;
38         }
39     }
40     release(&s5_inode_pool.lock);
41
42     return 0;
43 }
44
45 static struct {
46     struct spinlock lock;
47     struct s5_superblock sb[MAXVFSSIZE];
48 } s5_sb_pool; // It is a Pool of S5 Superblock Filesystems
49
50 struct s5_superblock*
51 alloc_s5_sb()
52 {
53     struct s5_superblock *sb;
54

```

```

55     acquire(&s5_sb_pool.lock);
56     for (sb = &s5_sb_pool.sb[0]; sb < &s5_sb_pool.sb[MAXVFSSIZE]; sb++) {
57         if (sb->flags == S5_SB_FREE) {
58             sb->flags |= S5_SB_USED;
59             release(&s5_sb_pool.lock);
60
61             return sb;
62         }
63     }
64     release(&s5_sb_pool.lock);
65
66     return 0;
67 }
68
69 struct vfs_operations s5_ops = {
70     .fs_init = &s5fs_init,
71     .mount = &s5_mount,
72     .unmount = &s5_unmount,
73     .getroot = &s5_getroot,
74     .readsb = &s5_readsb,
75     .ialloc = &s5_ialloc,
76     .balloc = &s5_balloc,
77     .bzero = &s5_bzero,
78     .bfree = &s5_bfree,
79     .brelse = &brelse,
80     .bwrite = &bwrite,
81     .bread = &bread,
82     .namecmp = &s5_namecmp
83 };
84
85 struct inode_operations s5_iops = {
86     .dirlookup = &s5_dirlookup,
87     .iupdate = &s5_iupdate,
88     .itrunc = &s5_itrunc,
89     .cleanup = &s5_cleanup,
90     .bmap = &s5_bmap,
91     .ilock = &s5_ilock,
92     .iunlock = &generic_iunlock,
93     .stati = &generic_stati,
94     .readi = &s5_readi,
95     .writei = &s5_writei,
96     .dirlink = &generic_dirlink,
97     .unlink = &s5_unlink,
98     .isdirempty = &s5_isdirempty
99 };
100
101 struct filesystem_type s5fs = {
102     .name = "s5",
103     .ops = &s5_ops,
104     .iops = &s5_iops
105 };
106
107 int
108 inits5fs(void)
109 {
110     initlock(&s5_sb_pool.lock, "s5_sb_pool");
111     initlock(&s5_inode_pool.lock, "s5_inode_pool");
112     return register_fs(&s5fs);
113 }
114

```

```

115 int
116 s5fs_init(void)
117 {
118     return 0;
119 }
120
121 int
122 s5_mount(struct inode *devi, struct inode *ip)
123 {
124     struct mntentry *mp;
125
126     // Read the Superblock
127     s5_ops.readsb(devi->minor, &sb[devi->minor]);
128
129     // Read the root device
130     struct inode *devrtip = s5_ops.getroot(devi->major, devi->minor);
131
132     acquire(&mtable.lock);
133     for (mp = &mtable.mpoint[0]; mp < &mtable.mpoint[MOUNTSIZE]; mp++) {
134         // This slot is available
135         if (mp->flag == 0) {
136             found_slot:
137             mp->dev = devi->minor;
138             mp->m_inode = ip;
139             mp->pdata = &sb[devi->minor];
140             mp->flag |= M_USED;
141             mp->m_rtinode = devrtip;
142
143             release(&mtable.lock);
144
145             initlog(devi->minor);
146             return 0;
147         } else {
148             // The disk is already mounted
149             if (mp->dev == devi->minor) {
150                 release(&mtable.lock);
151                 return -1;
152             }
153
154             if (ip->dev == mp->m_inode->dev && ip->inum == mp->m_inode->inum)
155                 goto found_slot;
156         }
157     }
158     release(&mtable.lock);
159
160     return -1;
161 }
162
163 int
164 s5_unmount(struct inode *devi)
165 {
166     return 0;
167 }
168
169 struct inode *
170 s5_getroot(int major, int minor)
171 {
172     return s5_iget(minor, ROOTINO);
173 }
174

```

```

175 void
176 s5_readsb(int dev, struct superblock *sb)
177 {
178     struct buf *bp;
179     struct s5_superblock *s5sb;
180
181     if((sb->flags & SB_NOT_LOADED) == 0) {
182         s5sb = alloc_s5_sb(); // Allocate a new S5 sb struct to the superblock.
183     } else{
184         s5sb = sb->fs_info;
185     }
186
187     // These sets are needed because of bread
188     sb->major = IDEMAJOR;
189     sb->minor = dev;
190     sb->blocksize = BSIZE;
191
192     bp = s5_ops.bread(dev, 1);
193     memmove(s5sb, bp->data, sizeof(*s5sb) - sizeof(s5sb->flags));
194     s5_ops.brelse(bp);
195
196     sb->fs_info = s5sb;
197 }
198
199 struct inode*
200 s5_ialloc(uint dev, short type)
201 {
202     int inum;
203     struct buf *bp;
204     struct dinode *dip;
205     struct s5_superblock *s5sb;
206
207     s5sb = sb[dev].fs_info;
208
209     for(inum = 1; inum < s5sb->ninodes; inum++){
210         bp = s5_ops.bread(dev, IBLOCK(inum, (*s5sb)));
211         dip = (struct dinode*)bp->data + inum%IPB;
212         if(dip->type == 0){ // a free inode
213             memset(dip, 0, sizeof(*dip));
214             dip->type = type;
215             log_write(bp); // mark it allocated on the disk
216             s5_ops.brelse(bp);
217             return s5_iget(dev, inum);
218         }
219         s5_ops.brelse(bp);
220     }
221     panic("ialloc: no inodes");
222 }
223
224 uint
225 s5_balloc(uint dev)
226 {
227     int b, bi, m;
228     struct buf *bp;
229     struct s5_superblock *s5sb;
230
231     s5sb = sb[dev].fs_info;
232     bp = 0;
233     for (b = 0; b < s5sb->size; b += BPB) {
234         bp = s5_ops.bread(dev, BBLOCK(b, (*s5sb)));

```

```

235     for (bi = 0; bi < BPB && b + bi < s5sb->size; bi++) {
236         m = 1 << (bi % 8);
237         if ((bp->data[bi/8] & m) == 0) { // Is block free?
238             bp->data[bi/8] |= m; // Mark block in use.
239             log_write(bp);
240             s5_ops.brelse(bp);
241             s5_ops.bzero(dev, b + bi);
242             return b + bi;
243         }
244     }
245     s5_ops.brelse(bp);
246 }
247 panic("balloc: out of blocks");
248 }
249
250 void
251 s5_bzero(int dev, int bno)
252 {
253     struct buf *bp;
254
255     bp = s5_ops.bread(dev, bno);
256     memset(bp->data, 0, BSIZE);
257     log_write(bp);
258     s5_ops.brelse(bp);
259 }
260
261 void
262 s5_bfree(int dev, uint b)
263 {
264     struct buf *bp;
265     int bi, m;
266     struct s5_superblock *s5sb;
267
268     s5sb = sb[dev].fs_info;
269     s5_ops.readsb(dev, &sb[dev]);
270     bp = s5_ops.bread(dev, BBLOCK(b, (*s5sb)));
271     bi = b % BPB;
272     m = 1 << (bi % 8);
273     if ((bp->data[bi/8] & m) == 0)
274         panic("freeing free block");
275     bp->data[bi/8] &= ~m;
276     log_write(bp);
277     s5_ops.brelse(bp);
278 }
279
280 struct inode*
281 s5_dirlookup(struct inode *dp, char *name, uint *poff)
282 {
283     uint off, inum;
284     struct dirent de;
285
286     if(dp->type == T_FILE || dp->type == T_DEV)
287         panic("dirlookup not DIR");
288
289     for(off = 0; off < dp->size; off += sizeof(de)){
290         if(s5_iops.readi(dp, (char*)&de, off, sizeof(de)) != sizeof(de))
291             panic("dirlink read");
292         if(de.inum == 0)
293             continue;
294         if(s5_ops.namecmp(name, de.name) == 0){

```

```

295     // entry matches path element
296     if (poff)
297         *poff = off;
298     inum = de.inum;
299     return s5_iget(dp->dev, inum);
300 }
301 }
302
303 return 0;
304 }
305
306 void
307 s5_iupdate(struct inode *ip)
308 {
309     struct buf *bp;
310     struct dinode *dip;
311     struct s5_superblock *s5sb;
312     struct s5_inode *s5ip;
313
314     s5ip = ip->i_private;
315     s5sb = sb[ip->dev].fs_info;
316     bp = s5_ops.bread(ip->dev, IBLOCK(ip->inum, (*s5sb)));
317     dip = (struct dinode*)bp->data + ip->inum/IPB;
318     dip->type = ip->type;
319     dip->major = ip->major;
320     dip->minor = ip->minor;
321     dip->nlink = ip->nlink;
322     dip->size = ip->size;
323     memmove(dip->addrs, s5ip->addrs, sizeof(s5ip->addrs));
324     log_write(bp);
325     s5_ops.brelse(bp);
326 }
327
328 void
329 s5_itrunc(struct inode *ip)
330 {
331     int i, j;
332     struct buf *bp;
333     uint *a;
334     struct s5_inode *s5ip;
335
336     s5ip = ip->i_private;
337
338     for (i = 0; i < NDIRECT; i++){
339         if (s5ip->addrs[i]){
340             s5_ops.bfree(ip->dev, s5ip->addrs[i]);
341             s5ip->addrs[i] = 0;
342         }
343     }
344
345     if (s5ip->addrs[NDIRECT]){
346         bp = s5_ops.bread(ip->dev, s5ip->addrs[NDIRECT]);
347         a = (uint*)bp->data;
348         for (j = 0; j < NINDIRECT; j++) {
349             if (a[j])
350                 s5_ops.bfree(ip->dev, a[j]);
351         }
352         s5_ops.brelse(bp);
353         s5_ops.bfree(ip->dev, s5ip->addrs[NDIRECT]);
354         s5ip->addrs[NDIRECT] = 0;

```

```

355     }
356
357     ip->size = 0;
358     s5_iops.iupdate(ip);
359 }
360
361 void
362 s5_cleanup(struct inode *ip)
363 {
364     memset(ip->i_private, 0, sizeof(struct s5_inode));
365 }
366
367 uint
368 s5_bmap(struct inode *ip, uint bn)
369 {
370     uint addr, *a;
371     struct buf *bp;
372     struct s5_inode *s5ip;
373
374     s5ip = ip->i_private;
375
376     if(bn < NDIRECT){
377         if((addr = s5ip->addrs[bn]) == 0)
378             s5ip->addrs[bn] = addr = s5_ops.balloc(ip->dev);
379         return addr;
380     }
381     bn -= NDIRECT;
382
383     if(bn < NINDIRECT){
384         // Load indirect block, allocating if necessary.
385         if((addr = s5ip->addrs[NDIRECT]) == 0)
386             s5ip->addrs[NDIRECT] = addr = s5_ops.balloc(ip->dev);
387         bp = s5_ops.bread(ip->dev, addr);
388         a = (uint*)bp->data;
389         if((addr = a[bn]) == 0){
390             a[bn] = addr = s5_ops.balloc(ip->dev);
391             log_write(bp);
392         }
393         s5_ops.brelse(bp);
394         return addr;
395     }
396
397     panic("bmap: out of range");
398 }
399
400 void
401 s5_iloc(struct inode *ip)
402 {
403     struct buf *bp;
404     struct dinode *dip;
405     struct s5_superblock *s5sb;
406     struct s5_inode *s5ip;
407
408     s5ip = ip->i_private;
409
410     s5sb = sb[ip->dev].fs_info;
411
412     if(ip == 0 || ip->ref < 1)
413         panic("iloc");
414

```

```

415     acquire(&icache.lock);
416     while (ip->flags & I_BUSY)
417         sleep(ip, &icache.lock);
418     ip->flags |= I_BUSY;
419     release(&icache.lock);
420
421     if (!(ip->flags & I_VALID)) {
422         bp = s5_ops.bread(ip->dev, IBLOCK(ip->inum, (*s5sb)));
423         dip = (struct dinode*)bp->data + ip->inum%IPB;
424         ip->type = dip->type;
425         ip->major = dip->major;
426         ip->minor = dip->minor;
427         ip->nlink = dip->nlink;
428         ip->size = dip->size;
429         memmove(s5ip->addrs, dip->addrs, sizeof(s5ip->addrs));
430         s5_ops.brelse(bp);
431         ip->flags |= I_VALID;
432         if (ip->type == 0)
433             panic("ilock: no type");
434     }
435 }
436
437 int
438 s5_readi(struct inode *ip, char *dst, uint off, uint n)
439 {
440     uint tot, m;
441     struct buf *bp;
442
443     if (ip->type == T_DEV){
444         if (ip->major < 0 || ip->major >= NDEV || !devsw[ip->major].read)
445             return -1;
446         return devsw[ip->major].read(ip, dst, n);
447     }
448
449     if (off > ip->size || off + n < off)
450         return -1;
451     if (off + n > ip->size)
452         n = ip->size - off;
453
454     for (tot=0; tot<n; tot+=m, off+=m, dst+=m){
455         bp = ip->fs_t->ops->bread(ip->dev, ip->iops->bmap(ip, off/BSIZE));
456         m = min(n - tot, BSIZE - off%BSIZE);
457         memmove(dst, bp->data + off%BSIZE, m);
458         ip->fs_t->ops->brelse(bp);
459     }
460     return n;
461 }
462
463 int
464 s5_writei(struct inode *ip, char *src, uint off, uint n)
465 {
466     uint tot, m;
467     struct buf *bp;
468
469     if (ip->type == T_DEV){
470         if (ip->major < 0 || ip->major >= NDEV || !devsw[ip->major].write)
471             return -1;
472         return devsw[ip->major].write(ip, src, n);
473     }
474

```

```

475     if(off > ip->size || off + n < off)
476         return -1;
477     if(off + n > MAXFILE*BSIZE)
478         return -1;
479
480     for(tot=0; tot<n; tot+=m, off+=m, src+=m){
481         bp = s5_ops.bread(ip->dev, s5_iops.bmap(ip, off/BSIZE));
482         m = min(n - tot, BSIZE - off%BSIZE);
483         memmove(bp->data + off%BSIZE, src, m);
484         log_write(bp);
485         s5_ops.brelse(bp);
486     }
487
488     if(n > 0 && off > ip->size){
489         ip->size = off;
490         s5_iops.iupdate(ip);
491     }
492     return n;
493 }
494
495 int
496 s5_isdirempty(struct inode *dp)
497 {
498     int off;
499     struct dirent de;
500
501     for(off=2*sizeof(de); off<dp->size; off+=sizeof(de)){
502         if(s5_iops.readi(dp, (char*)&de, off, sizeof(de)) != sizeof(de))
503             panic("isdirempty: readi");
504         if(de.inum != 0)
505             return 0;
506     }
507     return 1;
508 }
509
510 int
511 s5_unlink(struct inode *dp, uint off)
512 {
513     struct dirent de;
514
515     memset(&de, 0, sizeof(de));
516     if(dp->iops->writei(dp, (char*)&de, off, sizeof(de)) != sizeof(de))
517         return -1;
518
519     return 0;
520 }
521
522 int
523 s5_namecmp(const char *s, const char *t)
524 {
525     return strncmp(s, t, DIRSIZ);
526 }
527
528 int
529 s5_fill_inode(struct inode *ip) {
530     struct s5_inode *s5ip;
531
532     s5ip = alloc_s5_inode();
533     if (!s5ip) {
534         panic("No s5 inode available");

```

```
535     }
536
537     ip->i_private = s5ip;
538
539     return 1;
540 }
541
542 struct inode*
543 s5_iget(uint dev, uint inum)
544 {
545     return iget(dev, inum, &s5_fill_inode);
546 }
```

APPENDIX C – src/ext2.c

```

1  #include "types.h"
2  #include "defs.h"
3  #include "param.h"
4  #include "stat.h"
5  #include "mmu.h"
6  #include "proc.h"
7  #include "spinlock.h"
8  #include "vfs.h"
9  #include "buf.h"
10 #include "file.h"
11 #include "vfsmount.h"
12 #include "ext2.h"
13 #include "find_bits.h"
14
15 #define in_range(b, first, len) ((b) >= (first) && (b) <= (first) + (len) - 1)
16 #define ext2_find_next_zero_bit find_next_zero_bit
17 #define ext2_test_bit test_bit
18 #define ext2_set_bit_atomic test_and_set_bit
19 #define ext2_clear_bit_atomic test_and_clear_bit
20
21 static int ext2_block_to_path(struct inode *inode,
22                             long i_block, int offsets[4], int *boundary);
23
24 static struct ext2_inode * ext2_get_inode(struct superblock *sb,
25                                           uint ino, struct buf **bh);
26
27 static struct buf * read_block_bitmap(struct superblock *sb,
28                                       unsigned int block_group);
29
30 static void group_adjust_blocks(struct superblock *sb, int group_no,
31                                struct ext2_group_desc *desc, struct buf *bh,
32                                int count);
33
34 typedef struct {
35     uint32 *p;
36     uint32 key;
37     struct buf *bh;
38 } Indirect;
39
40 static inline void
41 add_chain(Indirect *p, struct buf *bh, uint32 *v)
42 {
43     p->key = *(p->p = v);
44     p->bh = bh;
45 }
46
47 static inline int verify_chain(Indirect *from, Indirect *to)
48 {
49     while (from <= to && from->key == *from->p)
50         from++;
51     return (from > to);
52 }
53
54

```

```

55 static struct {
56     struct spinlock lock;
57     struct ext2_inode_info ei[NINODE];
58 } ext2_ei_pool; // It is a Pool of S5 Superblock Filesystems
59
60 struct ext2_inode_info*
61 alloc_ext2_inode_info()
62 {
63     struct ext2_inode_info *ei;
64
65     acquire(&ext2_ei_pool.lock);
66     for (ei = &ext2_ei_pool.ei[0]; ei < &ext2_ei_pool.ei[NINODE]; ei++) {
67         if (ei->flags == INODE_FREE) {
68             ei->flags |= INODE_USED;
69             release(&ext2_ei_pool.lock);
70
71             return ei;
72         }
73     }
74     release(&ext2_ei_pool.lock);
75
76     return 0;
77 }
78
79 static struct {
80     struct spinlock lock;
81     struct ext2_sb_info sb[MAXVFSSIZE];
82 } ext2_sb_pool; // It is a Pool of S5 Superblock Filesystems
83
84 struct ext2_sb_info*
85 alloc_ext2_sb()
86 {
87     struct ext2_sb_info *sb;
88
89     acquire(&ext2_sb_pool.lock);
90     for (sb = &ext2_sb_pool.sb[0]; sb < &ext2_sb_pool.sb[MAXVFSSIZE]; sb++) {
91         if (sb->flags == SB_FREE) {
92             sb->flags |= SB_USED;
93             release(&ext2_sb_pool.lock);
94
95             return sb;
96         }
97     }
98     release(&ext2_sb_pool.lock);
99
100    return 0;
101 }
102
103 struct vfs_operations ext2_ops = {
104     .fs_init = &ext2fs_init,
105     .mount    = &ext2_mount,
106     .unmount  = &ext2_unmount,
107     .getroot  = &ext2_getroot,
108     .readsb   = &ext2_readsb,
109     .ialloc   = &ext2_ialloc,
110     .balloc   = &ext2_balloc,
111     .bzero    = &ext2_bzero,
112     .bfree    = &ext2_bfree,
113     .brelse   = &brelse,
114     .bwrite   = &bwrite,

```

```

115     .bread    = &bread ,
116     .namecmp  = &ext2_namecmp
117 };
118
119 struct inode_operations ext2_iops = {
120     .dirlookup = &ext2_dirlookup ,
121     .iupdate   = &ext2_iupdate ,
122     .itrunc    = &ext2_itrunc ,
123     .cleanup   = &ext2_cleanup ,
124     .bmap      = &ext2_bmap ,
125     .ilock     = &ext2_ilock ,
126     .iunlock   = &generic_iunlock ,
127     .stati     = &generic_stati ,
128     .readi     = &generic_readi ,
129     .writei    = &ext2_writei ,
130     .dirlink   = &ext2_dirlink ,
131     .unlink    = &ext2_unlink ,
132     .isdirempty = &ext2_isdirempty
133 };
134
135 struct filesystem_type ext2fs = {
136     .name = "ext2" ,
137     .ops = &ext2_ops ,
138     .iops = &ext2_iops
139 };
140
141 int
142 initext2fs(void)
143 {
144     initlock(&ext2_sb_pool.lock , "ext2_sb_pool");
145     /* initlock(&ext2_inode_pool.lock , "ext2_inode_pool"); */
146     return register_fs(&ext2fs);
147 }
148
149 int
150 ext2fs_init(void)
151 {
152     return 0;
153 }
154
155 int
156 ext2_mount(struct inode *devi , struct inode *ip)
157 {
158     struct mntentry *mp;
159
160     // Read the Superblock
161     ext2_ops.readsb(devi->minor , &sb[devi->minor]);
162
163     // Read the root device
164     struct inode *devrtip = ext2_ops.getroot(devi->major , devi->minor);
165
166     acquire(&mtable.lock);
167     for (mp = &mtable.mpoint[0]; mp < &mtable.mpoint[MOUNTSIZE]; mp++) {
168         // This slot is available
169         if (mp->flag == 0) {
170 found_slot:
171             mp->dev = devi->minor;
172             mp->m_inode = ip;
173             mp->pdata = &sb[devi->minor];
174             mp->flag |= M_USED;

```

```

175     mp->m_rtinode = devrtip;
176
177     release(&mtable.lock);
178
179     return 0;
180 } else {
181     // The disk is already mounted
182     if (mp->dev == devi->minor) {
183         release(&mtable.lock);
184         return -1;
185     }
186
187     if (ip->dev == mp->m_inode->dev && ip->inum == mp->m_inode->inum)
188         goto found_slot;
189 }
190 }
191 release(&mtable.lock);
192
193 return -1;
194 }
195
196 int
197 ext2_unmount(struct inode *devi)
198 {
199     panic("ext2_unmount_op_not_defined");
200     return 0;
201 }
202
203 struct inode *
204 ext2_getroot(int major, int minor)
205 {
206     return ext2_iget(minor, EXT2_ROOT_INO);
207 }
208
209 static inline int
210 test_root(int a, int b)
211 {
212     int num = b;
213
214     while (a > num)
215         num *= b;
216     return num == a;
217 }
218
219 static int
220 ext2_group_sparse(int group)
221 {
222     if (group <= 1)
223         return 1;
224     return (test_root(group, 3) || test_root(group, 5) ||
225             test_root(group, 7));
226 }
227
228 /**
229  * ext2_bg_has_super - number of blocks used by the superblock in group
230  * @sb: superblock for filesystem
231  * @group: group number to check
232  *
233  * Return the number of blocks used by the superblock (primary or backup)
234  * in this group. Currently this will be only 0 or 1.

```

```

235  */
236  int
237  ext2_bg_has_super(struct superblock *sb, int group)
238  {
239      if (EXT2_HAS_RO_COMPAT_FEATURE(sb, EXT2_FEATURE_RO_COMPAT_SPARSE_SUPER)&&
240          !ext2_group_sparse(group))
241          return 0;
242      return 1;
243  }
244
245  struct ext2_group_desc *
246  ext2_get_group_desc(struct superblock * sb,
247                      unsigned int block_group,
248                      struct buf ** bh)
249  {
250      unsigned long group_desc;
251      unsigned long offset;
252      struct ext2_group_desc * desc;
253      struct ext2_sb_info *sbi = EXT2_SB(sb);
254
255      if (block_group >= sbi->s_groups_count) {
256          panic("Block_group_# is too large");
257      }
258
259      group_desc = block_group >> EXT2_DESC_PER_BLOCK_BITS(sb);
260      offset = block_group & (EXT2_DESC_PER_BLOCK(sb) - 1);
261      if (!sbi->s_group_desc[group_desc]) {
262          panic("Accessing a group descriptor not loaded");
263      }
264
265      desc = (struct ext2_group_desc *) sbi->s_group_desc[group_desc]->data;
266      if (bh) {
267          *bh = sbi->s_group_desc[group_desc];
268      }
269      return desc + offset;
270  }
271
272  static unsigned long
273  descriptor_loc(struct superblock *sb,
274                unsigned long logic_sb_block,
275                int nr)
276  {
277      unsigned long bg, first_meta_bg;
278      int has_super = 0;
279
280      first_meta_bg = EXT2_SB(sb)->s_es->s_first_meta_bg;
281
282      if (!EXT2_HAS_INCOMPAT_FEATURE(sb, EXT2_FEATURE_INCOMPAT_META_BG) ||
283          nr < first_meta_bg)
284          return (logic_sb_block + nr + 1);
285      bg = EXT2_SB(sb)->s_desc_per_block * nr;
286      if (ext2_bg_has_super(sb, bg))
287          has_super = 1;
288
289      return ext2_group_first_block_no(sb, bg) + has_super;
290  }
291
292  void
293  ext2_readsb(int dev, struct superblock *sb)
294  {

```

```

295  struct buf *bp;
296  struct ext2_sb_info *sbi;
297  struct ext2_superblock *es;
298  uint32 blocksize = EXT2_MIN_BLKSIZE;
299  int db_count, i;
300  unsigned long block;
301  unsigned long logic_sb_block = 1;
302  unsigned long offset = 0;
303
304  if ((sb->flags & SB_NOT_LOADED) == 0) {
305      sbi = alloc_ext2_sb(); // Allocate a new S5 sb struct to the superblock.
306  } else {
307      sbi = sb->fs_info;
308  }
309
310  // These sets are needed because of bread
311  sb->major = IDEMAJOR;
312  sb->minor = dev;
313  sb_set_blocksize(sb, blocksize);
314  sb->fs_info = sbi;
315
316  bp = ext2_ops.bread(dev, logic_sb_block); // Read the 1024 bytes starting from the byte 1024
317  es = (struct ext2_superblock *)bp->data;
318
319  sbi->s_es = es;
320  sbi->s_sbh = bp;
321  if (es->s_magic != EXT2_SUPER_MAGIC) {
322      ext2_ops.brelse(bp);
323      panic("Try to mount a non ext2 fs as an ext2 fs");
324  }
325
326  blocksize = EXT2_MIN_BLKSIZE << es->s_log_block_size;
327
328  /* If the blocksize doesn't match, re-read the thing.. */
329  if (sb->blocksize != blocksize) {
330      ext2_ops.brelse(bp);
331
332      sb_set_blocksize(sb, blocksize);
333
334      logic_sb_block = EXT2_MIN_BLKSIZE / blocksize;
335      offset = EXT2_MIN_BLKSIZE % blocksize;
336      bp = ext2_ops.bread(dev, logic_sb_block);
337
338      if (!bp) {
339          panic("Error on second ext2 superblock read");
340      }
341
342      es = (struct ext2_superblock *) (((char *)bp->data) + offset);
343      sbi->s_es = es;
344
345      if (es->s_magic != EXT2_SUPER_MAGIC) {
346          panic("error: ext2 magic mismatch");
347      }
348  }
349
350  if (es->s_rev_level == EXT2_GOOD_OLD_REV) {
351      sbi->s_inode_size = EXT2_GOOD_OLD_INODE_SIZE;
352      sbi->s_first_ino = EXT2_GOOD_OLD_FIRST_INO;
353  } else {
354      sbi->s_inode_size = es->s_inode_size;

```

```

355     sbi->s_first_ino = es->s_first_ino;
356 }
357
358 sbi->s_blocks_per_group = es->s_blocks_per_group;
359 sbi->s_inodes_per_group = es->s_inodes_per_group;
360
361 sbi->s_inodes_per_block = sb->blocksize / sbi->s_inode_size;
362 sbi->s_itb_per_group = sbi->s_inodes_per_group / sbi->s_inodes_per_block;
363 sbi->s_desc_per_block = sb->blocksize / sizeof(struct ext2_group_desc);
364
365 sbi->s_addr_per_block_bits = ilog2(EXT2_ADDR_PER_BLOCK(sb));
366 sbi->s_desc_per_block_bits = ilog2(EXT2_DESC_PER_BLOCK(sb));
367
368 if (sbi->s_blocks_per_group > sb->blocksize * 8) {
369     panic("error: #blocks_per_group too big");
370 }
371
372 if (sbi->s_inodes_per_group > sb->blocksize * 8) {
373     panic("error: #inodes_per_group too big");
374 }
375
376 sbi->s_groups_count = ((es->s_blocks_count -
377                      es->s_first_data_block - 1)
378                      / sbi->s_blocks_per_group) + 1;
379 db_count = (sbi->s_groups_count + sbi->s_desc_per_block - 1) /
380            sbi->s_desc_per_block;
381
382 if (db_count > EXT2_MAX_BGC) {
383     panic("error: not enough memory to store s_group_desc. Consider change the EXT2_MAX_BGC constant");
384 }
385
386 /* bgl_lock_init(sbi->s_blockgroup_lock); */
387
388 for (i = 0; i < db_count; i++) {
389     block = descriptor_loc(sb, logic_sb_block, i);
390     sbi->s_group_desc[i] = ext2_ops.bread(dev, block);
391     if (!sbi->s_group_desc[i]) {
392         panic("Error on read ext2 group descriptor");
393     }
394 }
395
396 sbi->s_gdb_count = db_count;
397 }
398
399 /*
400  * Read the inode allocation bitmap for a given block_group, reading
401  * into the specified slot in the superblock's bitmap cache.
402  *
403  * Return buffer_head of bitmap on success or NULL.
404  */
405 static struct buf *
406 read_inode_bitmap(struct superblock * sb, unsigned long block_group)
407 {
408     struct ext2_group_desc *desc;
409     struct buf *bh = 0;
410
411     desc = ext2_get_group_desc(sb, block_group, 0);
412     if (!desc)
413         panic("error on read ext2 inode bitmap");
414

```

```

415     bh = ext2_ops.bread(sb->minor, desc->bg_inode_bitmap);
416     if (!bh)
417         panic("error on read ext2 inode bitmap");
418     return bh;
419 }
420
421 /**
422  * It is a dummy implementation of ialloc.
423  * Current Linux implementation uses an heuristic to alloc inodes
424  * in the best place.
425  * Our implementation will take an linear search over the inode bitmap
426  * and get the first free inode.
427  */
428 struct inode*
429 ext2_ialloc(uint dev, short type)
430 {
431     int i, group;
432     unsigned long ino;
433     struct ext2_sb_info *sbi;
434     struct buf *bitmap_bh = 0;
435     struct buf *bh2;
436     struct buf *ibh;
437     struct ext2_group_desc *gdp;
438     struct ext2_inode *raw_inode;
439
440     sbi = EXT2_SB(&sb[dev]);
441
442     group = 0;
443     for(i = 0; i < sbi->s_groups_count; i++) {
444         gdp = ext2_get_group_desc(&sb[dev], group, &bh2);
445
446         if (bitmap_bh)
447             ext2_ops.brelse(bitmap_bh);
448
449         bitmap_bh = read_inode_bitmap(&sb[dev], group);
450         ino = 0;
451
452     repeat_in_this_group:
453         ino = ext2_find_next_zero_bit((unsigned long *)bitmap_bh->data,
454                                     EXT2_INODES_PER_GROUP(&sb[dev]), ino);
455         if (ino >= EXT2_INODES_PER_GROUP(&sb[dev])) {
456             if (++group == sbi->s_groups_count)
457                 group = 0;
458             continue;
459         }
460         if (ext2_set_bit_atomic(ino, (unsigned long *)bitmap_bh->data)) {
461             /* we lost this inode */
462             if (++ino >= EXT2_INODES_PER_GROUP(&sb[dev])) {
463                 /* this group is exhausted, try next group */
464                 if (++group == sbi->s_groups_count)
465                     group = 0;
466                 continue;
467             }
468             /* try to find free inode in the same group */
469             goto repeat_in_this_group;
470         }
471         goto got;
472     }
473
474     /*

```

```

475     * Scanned all blockgroups.
476     */
477     panic("no_space_to_alloc_inode");
478
479 got:
480     ext2_ops.bwrite(bitmap_bh);
481     ext2_ops.brelse(bitmap_bh);
482
483     ino += group * EXT2_INODES_PER_GROUP(&sb[dev]) + 1;
484     if (ino < EXT2_FIRST_INO(&sb[dev]) || ino > sbi->s_es->s_inodes_count) {
485         panic("ext2_invalid_inode_number_allocated");
486     }
487
488     /* spin_lock(sb_bgl_lock(sbi, group)); */
489     gdp->bg_free_inodes_count -= 1;
490     /* spin_unlock(sb_bgl_lock(sbi, group)); */
491
492     ext2_ops.bwrite(bh2);
493
494     raw_inode = ext2_get_inode(&sb[dev], ino, &ibh);
495
496     // Erase the current inode
497     memset(raw_inode, 0, sbi->s_inode_size);
498     // Translate the xv6 to inode type type
499     if (type == T_DIR) {
500         raw_inode->i_mode = S_IFDIR;
501     } else if (type == T_FILE) {
502         raw_inode->i_mode = S_IFREG;
503     } else {
504         // We did not treat char and block devices with difference.
505         panic("ext2: invalid inode mode");
506     }
507
508     ext2_ops.bwrite(ibh);
509     ext2_ops.brelse(ibh);
510
511     return ext2_iget(dev, ino);
512 }
513
514 uint
515 ext2_balloc(uint dev)
516 {
517     panic("ext2_balloc_op_not_defined");
518 }
519
520 void
521 ext2_bzero(int dev, int bno)
522 {
523     panic("ext2_bzero_op_not_defined");
524 }
525
526 void
527 ext2_bfree(int dev, uint b)
528 {
529     panic("ext2_bfree_op_not_defined");
530 }
531
532 struct inode*
533 ext2_dirlookup(struct inode *dp, char *name, uint *poff)
534 {

```

```

535     uint off, inum, currbk;
536     struct ext2_dir_entry_2 *de;
537     struct buf *bh;
538     int namelen = strlen(name);
539
540     for (off = 0; off < dp->size;) {
541         currbk = off / sb[dp->dev].blocksize;
542
543         bh = ext2_ops.bread(dp->dev, ext2_iops.bmap(dp, currbk));
544
545         de = (struct ext2_dir_entry_2 *) (bh->data + (off % sb[dp->dev].blocksize));
546
547         if (de->inode == 0 || de->name_len != namelen) {
548             off += de->rec_len;
549             ext2_ops.brelse(bh);
550             continue;
551         }
552
553         if (strncmp(name, de->name, de->name_len) == 0) {
554             // entry matches path element
555             if (poff)
556                 *poff = off;
557             inum = de->inode;
558             ext2_ops.brelse(bh);
559             return ext2_iget(dp->dev, inum);
560         }
561         off += de->rec_len;
562         ext2_ops.brelse(bh);
563     }
564
565     return 0;
566 }
567
568 void
569 ext2_iupdate(struct inode *ip)
570 {
571     struct buf *bp;
572     struct ext2_inode_info *ei;
573     struct ext2_inode *raw_inode;
574
575     ei = ip->i_private;
576     raw_inode = ext2_get_inode(&sb[ip->dev], ip->inum, &bp);
577
578     raw_inode->i_mode = ei->i_ei.i_mode;
579     raw_inode->i_blocks = ei->i_ei.i_blocks;
580     raw_inode->i_links_count = ip->nlink;
581     memmove(raw_inode->i_block, ei->i_ei.i_block, sizeof(ei->i_ei.i_block));
582     raw_inode->i_size = ip->size;
583
584     ext2_ops.bwrite(bp);
585     ext2_ops.brelse(bp);
586 }
587
588 /**
589  * ext2_free_blocks() — Free given blocks and update quota and i_blocks
590  * @inode:      inode
591  * @block:      start physical block to free
592  * @count:      number of blocks to free
593  */
594 void

```

```

595 ext2_free_blocks(struct inode * inode, unsigned long block,
596                 unsigned long count)
597 {
598     struct buf *bitmap_bh = 0;
599     struct buf * bh2;
600     unsigned long block_group;
601     unsigned long bit;
602     unsigned long i;
603     unsigned long overflow;
604     struct superblock * superb = &sb[inode->dev];
605     struct ext2_sb_info * sbi = EXT2_SB(&sb[inode->dev]);
606     struct ext2_group_desc * desc;
607     struct ext2_superblock * es = sbi->s_es;
608     unsigned freed = 0, group_freed;
609
610     if (block < es->s_first_data_block ||
611         block + count < block ||
612         block + count > es->s_blocks_count) {
613         panic("ext2_free_blocks in not datazone");
614     }
615
616     do_more:
617     overflow = 0;
618     block_group = (block - es->s_first_data_block) / EXT2_BLOCKS_PER_GROUP(superb);
619     bit = (block - es->s_first_data_block) % EXT2_BLOCKS_PER_GROUP(superb);
620     /*
621      * Check to see if we are freeing blocks across a group
622      * boundary.
623      */
624     if (bit + count > EXT2_BLOCKS_PER_GROUP(superb)) {
625         overflow = bit + count - EXT2_BLOCKS_PER_GROUP(superb);
626         count -= overflow;
627     }
628     if (bitmap_bh)
629         brelse(bitmap_bh);
630
631     bitmap_bh = read_block_bitmap(superb, block_group);
632     if (!bitmap_bh)
633         goto error_return;
634
635     desc = ext2_get_group_desc(superb, block_group, &bh2);
636     if (!desc)
637         goto error_return;
638
639     if (in_range(desc->bg_block_bitmap, block, count) ||
640         in_range(desc->bg_inode_bitmap, block, count) ||
641         in_range(block, desc->bg_inode_table,
642                 sbi->s_itb_per_group) ||
643         in_range(block + count - 1, desc->bg_inode_table,
644                 sbi->s_itb_per_group)) {
645         panic("Freeing blocks on system zone");
646         goto error_return;
647     }
648
649     for (i = 0, group_freed = 0; i < count; i++) {
650         if (!ext2_clear_bit_atomic(bit + i, (unsigned long *)bitmap_bh->data)) {
651             panic("ext2 bit already cleared for block");
652         } else {
653             group_freed++;
654         }

```

```

655     }
656
657     ext2_ops.bwrite(bitmap_bh);
658     group_adjust_blocks(superb, block_group, desc, bh2, group_freed);
659     freed += group_freed;
660
661     if (overflow) {
662         block += count;
663         count = overflow;
664         goto do_more;
665     }
666 error_return:
667     ext2_ops.brelse(bitmap_bh);
668 }
669
670 /**
671  * ext2_free_data - free a list of data blocks
672  * @inode: inode we are dealing with
673  * @p: array of block numbers
674  * @q: points immediately past the end of array
675  *
676  * We are freeing all blocks referred from that array (numbers are
677  * stored as little-endian 32-bit) and updating @inode->i_blocks
678  * appropriately.
679  */
680 static inline void
681 ext2_free_data(struct inode *inode, uint32 *p, uint32 *q)
682 {
683     unsigned long block_to_free = 0, count = 0;
684     unsigned long nr;
685
686     for ( ; p < q ; p++) {
687         nr = *p;
688         if (nr) {
689             *p = 0;
690             /* accumulate blocks to free if they're contiguous */
691             if (count == 0)
692                 goto free_this;
693             else if (block_to_free == nr - count)
694                 count++;
695             else {
696                 ext2_free_blocks(inode, block_to_free, count);
697                 /* mark_inode_dirty(inode); */
698 free_this:
699                 block_to_free = nr;
700                 count = 1;
701             }
702         }
703     }
704     if (count > 0) {
705         ext2_free_blocks(inode, block_to_free, count);
706         /* mark_inode_dirty(inode); */
707     }
708 }
709
710 /**
711  * ext2_free_branches - free an array of branches
712  * @inode: inode we are dealing with
713  * @p: array of block numbers
714  * @q: pointer immediately past the end of array

```

```

715  * @depth: depth of the branches to free
716  */
717  static void
718  ext2_free_branches(struct inode *inode, uint32 *p, uint32 *q, int depth)
719  {
720      struct buf *bh;
721      unsigned long nr;
722
723      if (depth--) {
724          int addr_per_block = EXT2_ADDR_PER_BLOCK(&sb[inode->dev]);
725          for ( ; p < q ; p++) {
726              nr = *p;
727              if (!nr)
728                  continue;
729              *p = 0;
730              bh = ext2_ops.bread(inode->dev, nr);
731              /*
732               * A read failure? Report error and clear slot
733               * (should be rare).
734               */
735              if (!bh) {
736                  panic("ext2_block_read_failure");
737                  continue;
738              }
739              ext2_free_branches(inode,
740                               (uint32*)bh->data,
741                               (uint32*)bh->data + addr_per_block,
742                               depth);
743              ext2_ops.brelse(bh);
744              ext2_free_blocks(inode, nr, 1);
745              /* mark_inode_dirty(inode); */
746          }
747      } else {
748          ext2_free_data(inode, p, q);
749      }
750  }
751
752  static void
753  ext2_release_inode(struct superblock *sb, int group, int dir)
754  {
755      struct ext2_group_desc *desc;
756      struct buf *bh;
757
758      desc = ext2_get_group_desc(sb, group, &bh);
759      if (!desc) {
760          panic("Error on get_group_descriptor");
761          return;
762      }
763
764      /* spin_lock(sb_bgl_lock(EXT2_SB(sb), group)); */
765      desc->bg_free_inodes_count += 1;
766      if (dir)
767          desc->bg_used_dirs_count -= 1;
768      /* spin_unlock(sb_bgl_lock(EXT2_SB(sb), group)); */
769      ext2_ops.bwrite(bh);
770  }
771
772  /*
773   * NOTE! When we get the inode, we're the only people
774   * that have access to it, and as such there are no

```

```

775  * race conditions we have to worry about. The inode
776  * is not on the hash-lists, and it cannot be reached
777  * through the filesystem because the directory entry
778  * has been deleted earlier.
779  *
780  * HOWEVER: we must make sure that we get no aliases,
781  * which means that we have to call "clear_inode()"
782  * _before_ we mark the inode not in use in the inode
783  * bitmaps. Otherwise a newly created file might use
784  * the same inode number (not actually the same pointer
785  * though), and then we'd have two inodes sharing the
786  * same inode number and space on the harddisk.
787  */
788  void
789  ext2_free_inode (struct inode * inode)
790  {
791      struct superblock *superb = &sb[inode->dev];
792      int is_directory;
793      unsigned long ino;
794      struct buf *bitmap_bh;
795      unsigned long block_group;
796      unsigned long bit;
797      struct ext2_superblock * es;
798      struct ext2_inode_info *ei;
799
800      ino = inode->inum;
801      ei = inode->i_private;
802
803      es = EXT2_SB(superb)->s_es;
804      is_directory = S_ISDIR(ei->i_ei.i_mode);
805
806      if (ino < EXT2_FIRST_INO(superb) ||
807          ino > es->s_inodes_count) {
808          panic("ext2_reserved_or_non_existent_inode");
809          return;
810      }
811
812      block_group = (ino - 1) / EXT2_INODES_PER_GROUP(superb);
813      bit = (ino - 1) % EXT2_INODES_PER_GROUP(superb);
814      bitmap_bh = read_inode_bitmap(superb, block_group);
815      if (!bitmap_bh)
816          return;
817
818      /* Ok, now we can actually update the inode bitmaps.. */
819      if (!ext2_clear_bit_atomic(bit, (void *) bitmap_bh->data))
820          panic("ext2_bit_already_cleared");
821      else
822          ext2_release_inode(superb, block_group, is_directory);
823
824      ext2_ops.bwrite(bitmap_bh);
825      ext2_ops.brelse(bitmap_bh);
826  }
827
828  void
829  ext2_itrunc(struct inode *ip)
830  {
831      uint32 *i_data;
832      int offsets[4];
833      uint32 nr = 0;
834      int n;

```

```

835     long iblock;
836     unsigned blocksize;
837     blocksize = sb[ip->dev].blocksize;
838     iblock = (blocksize - 1) >> EXT2_BLOCK_SIZE_BITS(&sb[ip->dev]);
839     n = ext2_block_to_path(ip, iblock, offsets, 0);
840
841     struct ext2_inode_info *ei = ip->i_private;
842
843     i_data = ei->i_ei.i_block;
844
845     if (n == 0)
846         return;
847
848     /* lock block here */
849
850     if (n == 1) {
851         ext2_free_data(ip, i_data+offsets[0],
852                        i_data + EXT2_NDIR_BLOCKS);
853     }
854
855     /* Kill the remaining (whole) subtrees */
856     switch (offsets[0]) {
857     default:
858         nr = i_data[EXT2_IND_BLOCK];
859         if (nr) {
860             i_data[EXT2_IND_BLOCK] = 0;
861             /* mark_inode_dirty(inode); */
862             ext2_free_branches(ip, &nr, &nr+1, 1);
863         }
864     case EXT2_IND_BLOCK:
865         nr = i_data[EXT2_DIND_BLOCK];
866         if (nr) {
867             i_data[EXT2_DIND_BLOCK] = 0;
868             /* mark_inode_dirty(inode); */
869             ext2_free_branches(ip, &nr, &nr+1, 2);
870         }
871     case EXT2_DIND_BLOCK:
872         nr = i_data[EXT2_TIND_BLOCK];
873         if (nr) {
874             i_data[EXT2_TIND_BLOCK] = 0;
875             /* mark_inode_dirty(inode); */
876             ext2_free_branches(ip, &nr, &nr+1, 3);
877         }
878     case EXT2_TIND_BLOCK:
879         ;
880     }
881
882     // unlock the inode here
883     ext2_free_inode(ip);
884
885     ext2_iops.iupdate(ip);
886 }
887
888 void
889 ext2_cleanup(struct inode *ip)
890 {
891     memset(ip->i_private, 0, sizeof(struct ext2_inode_info));
892 }
893
894 /**

```

```

895  * ext2_block_to_path - parse the block number into array of offsets
896  * @inode: inode in question (we are only interested in its superblock)
897  * @i_block: block number to be parsed
898  * @offsets: array to store the offsets in
899  * @boundary: set this non-zero if the referred-to block is likely to be
900  *             followed (on disk) by an indirect block.
901  * To store the locations of file's data ext2 uses a data structure common
902  * for UNIX filesystems - tree of pointers anchored in the inode, with
903  * data blocks at leaves and indirect blocks in intermediate nodes.
904  * This function translates the block number into path in that tree -
905  * return value is the path length and @offsets[n] is the offset of
906  * pointer to (n+1)th node in the nth one. If @block is out of range
907  * (negative or too large) warning is printed and zero returned.
908  *
909  * Note: function doesn't find node addresses, so no IO is needed. All
910  * we need to know is the capacity of indirect blocks (taken from the
911  * superblock).
912  */
913
914 /*
915  * Portability note: the last comparison (check that we fit into triple
916  * indirect block) is spelled differently, because otherwise on an
917  * architecture with 32-bit longs and 8Kb pages we might get into trouble
918  * if our filesystem had 8Kb blocks. We might use long long, but that would
919  * kill us on x86. Oh, well, at least the sign propagation does not matter -
920  * i_block would have to be negative in the very beginning, so we would not
921  * get there at all.
922  */
923
924 static int
925 ext2_block_to_path(struct inode *inode,
926                   long i_block, int offsets[4], int *boundary)
927 {
928     int ptrs = EXT2_ADDR_PER_BLOCK(&sb[inode->dev]);
929     int ptrs_bits = EXT2_ADDR_PER_BLOCK_BITS(&sb[inode->dev]);
930     const long direct_blocks = EXT2_NDIR_BLOCKS,
931             indirect_blocks = ptrs,
932             double_blocks = (1 << (ptrs_bits * 2));
933     int n = 0;
934     int final = 0;
935
936     if (i_block < 0) {
937         panic("block_to_path: invalid block num");
938     } else if (i_block < direct_blocks) {
939         offsets[n++] = i_block;
940         final = direct_blocks;
941     } else if ((i_block == direct_blocks) < indirect_blocks) {
942         offsets[n++] = EXT2_IND_BLOCK;
943         offsets[n++] = i_block;
944         final = ptrs;
945     } else if ((i_block == indirect_blocks) < double_blocks) {
946         offsets[n++] = EXT2_DIND_BLOCK;
947         offsets[n++] = i_block >> ptrs_bits;
948         offsets[n++] = i_block & (ptrs - 1);
949         final = ptrs;
950     } else if (((i_block == double_blocks) >> (ptrs_bits * 2)) < ptrs) {
951         offsets[n++] = EXT2_TIND_BLOCK;
952         offsets[n++] = i_block >> (ptrs_bits * 2);
953         offsets[n++] = (i_block >> ptrs_bits) & (ptrs - 1);
954         offsets[n++] = i_block & (ptrs - 1);

```

```

955     final = ptrs;
956 } else {
957     panic("This block is out of bounds from this ext2 fs");
958 }
959
960 if (boundary)
961     *boundary = final - 1 - (i_block & (ptrs - 1));
962
963 return n;
964 }
965
966 static void
967 ext2_update_branch(struct inode *inode, uint bn, Indirect *chain)
968 {
969     int ptrs = EXT2_ADDR_PER_BLOCK(&sb[inode->dev]);
970     int ptrs_bits = EXT2_ADDR_PER_BLOCK_BITS(&sb[inode->dev]);
971     const long direct_blocks = EXT2_NDIR_BLOCKS,
972             indirect_blocks = ptrs,
973             double_blocks = (1 << (ptrs_bits * 2));
974     struct ext2_inode_info *ei;
975
976     ei = inode->i_private;
977
978     // Update inode block
979     if (bn < 0) {
980         panic("block_to_path invalid block num");
981     } else if (bn < direct_blocks) {
982         if (ei->i_block[bn] == 0)
983             ei->i_block[bn] = chain[0].key;
984     } else if ((bn -= direct_blocks) < indirect_blocks) {
985         if (ei->i_block[EXT2_IND_BLOCK] == 0)
986             ei->i_block[EXT2_IND_BLOCK] = chain[0].key;
987     } else if ((bn -= indirect_blocks) < double_blocks) {
988         if (ei->i_block[EXT2_DIND_BLOCK] == 0)
989             ei->i_block[EXT2_DIND_BLOCK] = chain[0].key;
990     } else if (((bn -= double_blocks) >> (ptrs_bits * 2)) < ptrs) {
991         if (ei->i_block[EXT2_TIND_BLOCK] == 0)
992             ei->i_block[EXT2_TIND_BLOCK] = chain[0].key;
993     } else {
994         panic("This block is out of bounds from this ext2 fs");
995     }
996
997     return;
998 }
999
1000 /**
1001  * ext2_get_branch - read the chain of indirect blocks leading to data
1002  * @inode: inode in question
1003  * @depth: depth of the chain (1 - direct pointer, etc.)
1004  * @offsets: offsets of pointers in inode/indirect blocks
1005  * @chain: place to store the result
1006  * @err: here we store the error value
1007  *
1008  * Function fills the array of triples <key, p, bh> and returns %NULL
1009  * if everything went OK or the pointer to the last filled triple
1010  * (incomplete one) otherwise. Upon the return chain[i].key contains
1011  * the number of (i+1)-th block in the chain (as it is stored in memory,
1012  * i.e. little-endian 32-bit), chain[i].p contains the address of that
1013  * number (it points into struct inode for i==0 and into the bh->b_data
1014  * for i>0) and chain[i].bh points to the buffer_head of i-th indirect

```

```

1015  * block for  $i > 0$  and NULL for  $i = 0$ . In other words, it holds the block
1016  * numbers of the chain, addresses they were taken from (and where we can
1017  * verify that chain did not change) and buffer_heads hosting these
1018  * numbers.
1019  *
1020  * Function stops when it stumbles upon zero pointer (absent block)
1021  * (pointer to last triple returned,  $*@err == 0$ )
1022  * or when it gets an IO error reading an indirect block
1023  * (ditto,  $*@err == -EIO$ )
1024  * or when it notices that chain had been changed while it was reading
1025  * (ditto,  $*@err == -EAGAIN$ )
1026  * or when it reads all  $@depth - 1$  indirect blocks successfully and finds
1027  * the whole chain, all way to the data (returns %NULL,  $*err == 0$ ).
1028  */
1029 static Indirect *ext2_get_branch(struct inode *inode,
1030                                int depth,
1031                                int *offsets,
1032                                Indirect chain[4])
1033 {
1034     Indirect *p = chain;
1035     struct buf *bh;
1036     struct ext2_inode_info *ei = inode->i_private;
1037
1038     add_chain(chain, 0, ei->i_ei.i_block + *offsets);
1039     if (!p->key)
1040         goto no_block;
1041     while (--depth) {
1042         bh = ext2_ops.bread(inode->dev, p->key);
1043         if (!bh)
1044             panic("error on ext2_get_branch");
1045         if (!verify_chain(chain, p))
1046             panic("ext2_get_branch chain changed");
1047         add_chain(++p, bh, (uint32*)bh->data + *++offsets);
1048         if (!p->key)
1049             goto no_block;
1050     }
1051     return 0;
1052
1053 no_block:
1054     return p;
1055 }
1056
1057 /**
1058  * ext2_find_near - find a place for allocation with sufficient locality
1059  * @inode: owner
1060  * @ind: descriptor of indirect block.
1061  *
1062  * This function returns the preferred place for block allocation.
1063  * It is used when heuristic for sequential allocation fails.
1064  * Rules are:
1065  * + if there is a block to the left of our position - allocate near it.
1066  * + if pointer will live in indirect block - allocate near that block.
1067  * + if pointer will live in inode - allocate in the same cylinder group.
1068  *
1069  * In the latter case we colour the starting block by the callers PID to
1070  * prevent it from clashing with concurrent allocations for a different inode
1071  * in the same block group. The PID is used here so that functionally related
1072  * files will be close-by on-disk.
1073  *
1074  * Caller must make sure that @ind is valid and will stay that way.

```

```

1075  */
1076
1077  static ext2_fsblk_t ext2_find_near(struct inode *inode, Indirect *ind)
1078  {
1079      struct ext2_inode_info *ei = inode->i_private;
1080      uint32 *start = ind->bh ? (uint32 *) ind->bh->data : ei->i_ei.i_block;
1081      uint32 *p;
1082      ext2_fsblk_t bg_start;
1083      ext2_fsblk_t colour;
1084      ext2_grpblk_t i_block_group;
1085
1086      /* Try to find previous block */
1087      for (p = ind->p - 1; p >= start; p--)
1088          if (*p)
1089              return *p;
1090
1091      /* No such thing, so let's try location of indirect block */
1092      if (ind->bh)
1093          return ind->bh->blockno;
1094
1095      /*
1096       * It is going to be referred from inode itself? OK, just put it into
1097       * the same cylinder group then.
1098       */
1099      i_block_group = (inode->inum - 1) / EXT2_INODES_PER_GROUP(&sb[inode->dev]);
1100      bg_start = ext2_group_first_block_no(&sb[inode->dev], i_block_group);
1101      colour = (proc->pid % 16) *
1102          (EXT2_BLOCKS_PER_GROUP(&sb[inode->dev]) / 16);
1103      return bg_start + colour;
1104  }
1105
1106  static inline ext2_fsblk_t ext2_find_goal(struct inode *inode, long block,
1107      Indirect *partial)
1108  {
1109      return ext2_find_near(inode, partial);
1110  }
1111
1112  /**
1113   * ext2_blks_to_allocate: Look up the block map and count the number
1114   * of direct blocks need to be allocated for the given branch.
1115   *
1116   * @branch: chain of indirect blocks
1117   * @k: number of blocks need for indirect blocks
1118   * @blks: number of data blocks to be mapped.
1119   * @blocks_to_boundary: the offset in the indirect block
1120   *
1121   * return the total number of blocks to be allocate, including the
1122   * direct and indirect blocks.
1123   */
1124  static int
1125  ext2_blks_to_allocate(Indirect * branch, int k, unsigned long blks,
1126      int blocks_to_boundary)
1127  {
1128      unsigned long count = 0;
1129
1130      /*
1131       * Simple case, [t,d]Indirect block(s) has not allocated yet
1132       * then it's clear blocks on that path have not allocated
1133       */
1134      if (k > 0) {

```

```

1135     /* right now don't hanel cross boundary allocation */
1136     if (blks < blocks_to_boundary + 1)
1137         count += blks;
1138     else
1139         count += blocks_to_boundary + 1;
1140     return count;
1141 }
1142
1143 count++;
1144 while (count < blks && count <= blocks_to_boundary
1145        && *(branch[0].p + count) == 0) {
1146     count++;
1147 }
1148 return count;
1149 }
1150
1151 /*
1152  * Read the bitmap for a given block_group, and validate the
1153  * bits for block/inode/inode tables are set in the bitmaps
1154  *
1155  * Return buffer_head on success or NULL in case of failure.
1156  */
1157 static struct buf *
1158 read_block_bitmap(struct superblock *sb, unsigned int block_group)
1159 {
1160     struct ext2_group_desc * desc;
1161     struct buf * bh;
1162     ext2_fsblk_t bitmap_blk;
1163
1164     desc = ext2_get_group_desc(sb, block_group, 0);
1165     if (!desc)
1166         return 0;
1167     bitmap_blk = desc->bg_block_bitmap;
1168     bh = ext2_ops.bread(sb->minor, bitmap_blk);
1169     if (!bh) {
1170         return 0;
1171     }
1172
1173     /* ext2_valid_block_bitmap(sb, desc, block_group, bh); */
1174     /*
1175      * file system mounted not to panic on error, continue with corrupt
1176      * bitmap
1177      */
1178     return bh;
1179 }
1180
1181 /**
1182  * bitmap_search_next_usable_block()
1183  * @start: the starting block (group relative) of the search
1184  * @bh:    bufferhead contains the block group bitmap
1185  * @maxblocks: the ending block (group relative) of the reservation
1186  *
1187  * The bitmap search — search forward through the actual bitmap on disk until
1188  * we find a bit free.
1189  */
1190 static ext2_grpblk_t
1191 bitmap_search_next_usable_block(ext2_grpblk_t start, struct buf *bh,
1192                                ext2_grpblk_t maxblocks)
1193 {
1194     ext2_grpblk_t next;

```



```

1195
1196     next = ext2_find_next_zero_bit((unsigned long *)bh->data, maxblocks, start);
1197     if (next >= maxblocks)
1198         return -1;
1199     return next;
1200 }
1201
1202 /**
1203  * find_next_usable_block()
1204  * @start: the starting block (group relative) to find next
1205  *         allocatable block in bitmap.
1206  * @bh:    bufferhead contains the block group bitmap
1207  * @maxblocks: the ending block (group relative) for the search
1208  *
1209  * Find an allocatable block in a bitmap. We perform the "most
1210  * appropriate allocation" algorithm of looking for a free block near
1211  * the initial goal; then for a free byte somewhere in the bitmap;
1212  * then for any free bit in the bitmap.
1213  */
1214 static ext2_grpblk_t
1215 find_next_usable_block(int start, struct buf *bh, int maxblocks)
1216 {
1217     ext2_grpblk_t here, next;
1218     char *p, *r;
1219
1220     if (start > 0) {
1221         /*
1222          * The goal was occupied; search forward for a free
1223          * block within the next XX blocks.
1224          *
1225          * end_goal is more or less random, but it has to be
1226          * less than EXT2_BLOCKS_PER_GROUP. Aligning up to the
1227          * next 64-bit boundary is simple..
1228          */
1229         ext2_grpblk_t end_goal = (start + 63) & ~63;
1230         if (end_goal > maxblocks)
1231             end_goal = maxblocks;
1232         here = ext2_find_next_zero_bit((unsigned long *)bh->data, end_goal, start);
1233         if (here < end_goal)
1234             return here;
1235     }
1236
1237     here = start;
1238     if (here < 0)
1239         here = 0;
1240
1241     p = ((char *)bh->data) + (here >> 3);
1242     r = memscan(p, 0, ((maxblocks + 7) >> 3) - (here >> 3));
1243     next = (r - ((char *)bh->data)) << 3;
1244
1245     if (next < maxblocks && next >= here)
1246         return next;
1247
1248     here = bitmap_search_next_usable_block(here, bh, maxblocks);
1249     return here;
1250 }
1251
1252 /**
1253  * ext2_try_to_allocate()
1254  * @sb:    superblock

```

```

1255  * @group:  given allocation block group
1256  * @bitmap_bh:  bufferhead holds the block bitmap
1257  * @grp_goal:  given target block within the group
1258  * @count:  target number of blocks to allocate
1259  * @my_rsv:  reservation window
1260  *
1261  * Attempt to allocate blocks within a give range. Set the range of allocation
1262  * first, then find the first free bit(s) from the bitmap (within the range),
1263  * and at last, allocate the blocks by claiming the found free bit as allocated.
1264  *
1265  * To set the range of this allocation:
1266  *   if there is a reservation window, only try to allocate block(s)
1267  *   from the file's own reservation window;
1268  *   Otherwise, the allocation range starts from the give goal block,
1269  *   ends at the block group's last block.
1270  *
1271  * If we failed to allocate the desired block then we may end up crossing to a
1272  * new bitmap.
1273  */
1274  static int
1275  ext2_try_to_allocate(struct superblock *sb, int group,
1276                     struct buf *bitmap_bh, ext2_grpblk_t grp_goal,
1277                     unsigned long *count)
1278  {
1279      ext2_grpblk_t start, end;
1280      unsigned long num = 0;
1281
1282      if (grp_goal > 0)
1283          start = grp_goal;
1284      else
1285          start = 0;
1286      end = EXT2_BLOCKS_PER_GROUP(sb);
1287
1288      repeat:
1289      if (grp_goal < 0) {
1290          grp_goal = find_next_usable_block(start, bitmap_bh, end);
1291          if (grp_goal < 0)
1292              goto fail_access;
1293
1294          int i;
1295
1296          for (i = 0; i < 7 && grp_goal > start &&
1297               !ext2_test_bit(grp_goal - 1, (unsigned long *)bitmap_bh->data);
1298               i++, grp_goal--)
1299              ;
1300      }
1301      start = grp_goal;
1302
1303      if (ext2_set_bit_atomic(grp_goal,
1304                             (unsigned long *)bitmap_bh->data)) {
1305          /*
1306           * The block was allocated by another thread, or it was
1307           * allocated and then freed by another thread
1308           */
1309          start++;
1310          grp_goal++;
1311          if (start >= end)
1312              goto fail_access;
1313          goto repeat;
1314      }

```

```

1315     num++;
1316     grp_goal++;
1317     while (num < *count && grp_goal < end &&
1318            !ext2_set_bit_atomic(grp_goal, (unsigned long*)bitmap_bh->data)) {
1319         num++;
1320         grp_goal++;
1321     }
1322     *count = num;
1323     return grp_goal - num;
1324 fail_access:
1325     *count = num;
1326     return -1;
1327 }
1328
1329 static void
1330 group_adjust_blocks(struct superblock *sb, int group_no,
1331                    struct ext2_group_desc *desc, struct buf *bh,
1332                    int count)
1333 {
1334     if (count) {
1335         /* struct ext2_sb_info *sbi = EXT2_SB(sb); */
1336         unsigned free_blocks;
1337
1338         /* spin_lock(sb_bgl_lock(sbi, group_no)); */
1339         free_blocks = desc->bg_free_blocks_count;
1340         desc->bg_free_blocks_count = free_blocks + count;
1341         /* spin_unlock(sb_bgl_lock(sbi, group_no)); */
1342         ext2_ops.bwrite(bh);
1343     }
1344 }
1345
1346 /*
1347  * ext2_new_blocks() — core block(s) allocation function
1348  * @inode: file inode
1349  * @goal: given target block(filesystem wide)
1350  * @count: target number of blocks to allocate
1351  * @errp: error code
1352  *
1353  * ext2_new_blocks uses a goal block to assist allocation. If the goal is
1354  * free, or there is a free block within 32 blocks of the goal, that block
1355  * is allocated. Otherwise a forward search is made for a free block; within
1356  * each block group the search first looks for an entire free byte in the block
1357  * bitmap, and then for any free bit if that fails.
1358  * This function also updates quota and i_blocks field.
1359  */
1360 ext2_fsblk_t
1361 ext2_new_blocks(struct inode *inode, ext2_fsblk_t goal,
1362                unsigned long *count, int *errp)
1363 {
1364     struct buf *bitmap_bh = 0;
1365     struct buf *gdp_bh;
1366     int group_no;
1367     ext2_grpblk_t grp_target_blk; /* blockgroup relative goal block */
1368     ext2_grpblk_t grp_alloc_blk; /* blockgroup-relative allocated block */
1369     ext2_fsblk_t ret_block; /* filesystem-wide allocated block */
1370     int bgi; /* blockgroup iteration index */
1371     ext2_grpblk_t free_blocks; /* number of free blocks in a group */
1372     struct superblock *superb;
1373     struct ext2_group_desc *gdp;
1374     struct ext2_superblock *es;

```

```

1375 struct ext2_sb_info *sbi;
1376 unsigned long ngroups;
1377 unsigned long num = *count;
1378
1379 *errp = -1;
1380 superb = &sb[inode->dev];
1381
1382 sbi = EXT2_SB(superb);
1383 es = sbi->s_es;
1384
1385 /* if (!ext2_has_free_blocks(sbi)) { */
1386 /*     *errp = -ENOSPC; */
1387 /*     goto out; */
1388 /* } */
1389
1390 /*
1391  * First, test whether the goal block is free.
1392  */
1393 if (goal < es->s_first_data_block ||
1394     goal >= es->s_blocks_count) {
1395     goal = es->s_first_data_block;
1396 }
1397
1398 group_no = (goal - es->s_first_data_block) / EXT2_BLOCKS_PER_GROUP(superb);
1399 retry_alloc:
1400 gdp = ext2_get_group_desc(superb, group_no, &gdp_bh);
1401 if (!gdp)
1402     goto io_error;
1403
1404 free_blocks = gdp->bg_free_blocks_count;
1405
1406 if (free_blocks > 0) {
1407     grp_target_blk = ((goal - es->s_first_data_block) %
1408                      EXT2_BLOCKS_PER_GROUP(superb));
1409     bitmap_bh = read_block_bitmap(superb, group_no);
1410     if (!bitmap_bh)
1411         goto io_error;
1412     grp_alloc_blk = ext2_try_to_allocate(superb, group_no,
1413                                         bitmap_bh, grp_target_blk, &num);
1414     if (grp_alloc_blk >= 0)
1415         goto allocated;
1416 }
1417
1418 ngroups = EXT2_SB(superb)->s_groups_count;
1419
1420 /*
1421  * Now search the rest of the groups. We assume that
1422  * group_no and gdp correctly point to the last group visited.
1423  */
1424 for (bgi = 0; bgi < ngroups; bgi++) {
1425     group_no++;
1426     if (group_no >= ngroups)
1427         group_no = 0;
1428     gdp = ext2_get_group_desc(superb, group_no, &gdp_bh);
1429     if (!gdp)
1430         goto io_error;
1431
1432     free_blocks = gdp->bg_free_blocks_count;
1433     /*
1434      * skip this group (and avoid loading bitmap) if there

```

```

1435     * are no free blocks
1436     */
1437     if (!free_blocks)
1438         continue;
1439
1440     ext2_ops.brelse(bitmap_bh);
1441     bitmap_bh = read_block_bitmap(superb, group_no);
1442     if (!bitmap_bh)
1443         goto io_error;
1444     /*
1445      * try to allocate block(s) from this group, without a goal(-1).
1446      */
1447     grp_alloc_blk = ext2_try_to_allocate(superb, group_no,
1448                                         bitmap_bh, -1, &num);
1449     if (grp_alloc_blk >= 0)
1450         goto allocated;
1451 }
1452
1453 goto out;
1454
1455 allocated:
1456
1457     ret_block = grp_alloc_blk + ext2_group_first_block_no(superb, group_no);
1458
1459     if (in_range(gdp->bg_block_bitmap, ret_block, num) ||
1460         in_range(gdp->bg_inode_bitmap, ret_block, num) ||
1461         in_range(ret_block, gdp->bg_inode_table,
1462                 EXT2_SB(superb)->s_itb_per_group) ||
1463         in_range(ret_block + num - 1, gdp->bg_inode_table,
1464                 EXT2_SB(superb)->s_itb_per_group)) {
1465         goto retry_alloc;
1466     }
1467
1468     if (ret_block + num - 1 >= es->s_blocks_count) {
1469         panic("Error on ext2 block alloc");
1470     }
1471
1472     group_adjust_blocks(superb, group_no, gdp, gdp_bh, -num);
1473
1474     ext2_ops.bwrite(bitmap_bh);
1475
1476     *errp = 0;
1477     ext2_ops.brelse(bitmap_bh);
1478     /* if (num < *count) { */
1479     /*     dqot_free_block_nodirty(inode, *count-num); */
1480     /*     mark_inode_dirty(inode); */
1481     /*     *count = num; */
1482     /* } */
1483     return ret_block;
1484
1485 io_error:
1486     *errp = -2;
1487 out:
1488     /*
1489      * Undo the block allocation
1490      */
1491     /* if (!performed_allocation) { */
1492     /*     dqot_free_block_nodirty(inode, *count); */
1493     /*     mark_inode_dirty(inode); */
1494     /* } */

```

```

1495     ext2_ops.brelse(bitmap_bh);
1496     return 0;
1497 }
1498
1499 /**
1500  * ext2_alloc_blocks: multiple allocate blocks needed for a branch
1501  * @indirect_blks: the number of blocks need to allocate for indirect
1502  *   blocks
1503  *
1504  * @new_blocks: on return it will store the new block numbers for
1505  * the indirect blocks(if needed) and the first direct block,
1506  * @blks: on return it will store the total number of allocated
1507  *   direct blocks
1508  */
1509 static int
1510 ext2_alloc_blocks(struct inode *inode,
1511                  ext2_fsblk_t goal, int indirect_blks, int blks,
1512                  ext2_fsblk_t new_blocks[4], int *err)
1513 {
1514     int target;
1515     unsigned long count = 0;
1516     int index = 0;
1517     ext2_fsblk_t current_block = 0;
1518     int ret = 0;
1519
1520     /*
1521      * Here we try to allocate the requested multiple blocks at once,
1522      * on a best-effort basis.
1523      * To build a branch, we should allocate blocks for
1524      * the indirect blocks(if not allocated yet), and at least
1525      * the first direct block of this branch. That's the
1526      * minimum number of blocks need to allocate(required)
1527      */
1528     target = blks + indirect_blks;
1529
1530     while (1) {
1531         count = target;
1532         /* allocating blocks for indirect blocks and direct blocks */
1533         current_block = ext2_new_blocks(inode, goal, &count, err);
1534         if (*err)
1535             goto failed_out;
1536
1537         target -= count;
1538         /* allocate blocks for indirect blocks */
1539         while (index < indirect_blks && count) {
1540             new_blocks[index++] = current_block++;
1541             count--;
1542         }
1543
1544         if (count > 0)
1545             break;
1546     }
1547
1548     /* save the new block number for the first direct block */
1549     new_blocks[index] = current_block;
1550
1551     /* total number of blocks allocated for direct blocks */
1552     ret = count;
1553     *err = 0;
1554     return ret;

```

```

1555 failed_out:
1556     panic("ext2_error_on_ext2_alloc_blocks");
1557     return ret;
1558 }
1559
1560 /**
1561  * ext2_alloc_branch - allocate and set up a chain of blocks.
1562  * @inode: owner
1563  * @num: depth of the chain (number of blocks to allocate)
1564  * @offsets: offsets (in the blocks) to store the pointers to next.
1565  * @branch: place to store the chain in.
1566  *
1567  * This function allocates @num blocks, zeroes out all but the last one,
1568  * links them into chain and (if we are synchronous) writes them to disk.
1569  * In other words, it prepares a branch that can be spliced onto the
1570  * inode. It stores the information about that chain in the branch[], in
1571  * the same format as ext2_get_branch() would do. We are calling it after
1572  * we had read the existing part of chain and partial points to the last
1573  * triple of that (one with zero ->key). Upon the exit we have the same
1574  * picture as after the successful ext2_get_block(), except that in one
1575  * place chain is disconnected - *branch->p is still zero (we did not
1576  * set the last link), but branch->key contains the number that should
1577  * be placed into *branch->p to fill that gap.
1578  *
1579  * If allocation fails we free all blocks we've allocated (and forget
1580  * their buffer_heads) and return the error value the from failed
1581  * ext2_alloc_block() (normally -ENOSPC). Otherwise we set the chain
1582  * as described above and return 0.
1583  */
1584
1585 static int
1586 ext2_alloc_branch(struct inode *inode,
1587                  int indirect_blks, int *blks, ext2_fsblk_t goal,
1588                  int *offsets, Indirect *branch)
1589 {
1590     int blocksize = sb[inode->dev].blocksize;
1591     int i, n = 0;
1592     int err = 0;
1593     struct buf *bh;
1594     int num;
1595     ext2_fsblk_t new_blocks[4];
1596     ext2_fsblk_t current_block;
1597
1598     num = ext2_alloc_blocks(inode, goal, indirect_blks,
1599                             *blks, new_blocks, &err);
1600     if (err)
1601         return err;
1602
1603     branch[0].key = new_blocks[0];
1604     /*
1605      * metadata blocks and data blocks are allocated.
1606      */
1607     for (n = 1; n <= indirect_blks; n++) {
1608         /*
1609          * Get buffer_head for parent block, zero it out
1610          * and set the pointer to new one, then send;
1611          * parent to disk.
1612          */
1613         bh = ext2_ops.bread(inode->dev, new_blocks[n-1]);
1614         if (!bh) {

```

```

1615     goto failed;
1616 }
1617 branch[n].bh = bh;
1618 memset(bh->data, 0, blocksize);
1619 branch[n].p = (uint32 *) bh->data + offsets[n];
1620 branch[n].key = new_blocks[n];
1621 *branch[n].p = branch[n].key;
1622 if ( n == indirect_blks ) {
1623     current_block = new_blocks[n];
1624     /*
1625      * End of chain, update the last new metablock of
1626      * the chain to point to the new allocated
1627      * data blocks numbers
1628      */
1629     for ( i=1; i < num; i++)
1630         *(branch[n].p + i) = ++current_block;
1631 }
1632 ext2_ops.bwrite(bh);
1633 }
1634 *blks = num;
1635 return err;
1636
1637 failed:
1638 panic("ext2_error_on_allocate_blocks_branch");
1639 return err;
1640 }
1641
1642 uint
1643 ext2_bmap(struct inode *ip, uint bn)
1644 {
1645     /* struct buf *bp; */
1646     int depth;
1647     Indirect chain[4];
1648     Indirect *partial;
1649     int offsets[4];
1650     int indirect_blks;
1651     uint blkcn;
1652     int blocks_to_boundary;
1653     ext2_fsblk_t goal;
1654     int count;
1655     unsigned long maxblocks;
1656     int err;
1657
1658     depth = ext2_block_to_path(ip, bn, offsets, &blocks_to_boundary);
1659
1660     if (depth == 0)
1661         panic("Wrong_depth_value");
1662
1663     partial = ext2_get_branch(ip, depth, offsets, chain);
1664
1665     if (!partial) {
1666         goto got_it;
1667     }
1668
1669     maxblocks = sb[ip->dev].blocksize >> EXT2_BLOCK_SIZE_BITS(&sb[ip->dev]);
1670
1671     // The requested block is not allocated yet
1672     goal = ext2_find_goal(ip, bn, partial);
1673
1674     /* the number of blocks need to allocate for [d,t]indirect blocks */

```

```

1675     indirect_blks = (chain + depth) - partial - 1;
1676
1677     indirect_blks = (chain + depth) - partial - 1;
1678     /*
1679      * Next look up the indirect map to count the totoal number of
1680      * direct blocks to allocate for this branch.
1681      */
1682     count = ext2_blks_to_allocate(partial, indirect_blks,
1683                                   maxblocks, blocks_to_boundary);
1684
1685     err = ext2_alloc_branch(ip, indirect_blks, &count, goal,
1686                            offsets + (partial - chain), partial);
1687
1688     if (err < 0)
1689         panic("error on ext2_alloc_branch");
1690
1691     got_it:
1692         blkn = chain[depth-1].key;
1693         ext2_update_branch(ip, bn, chain);
1694
1695         /* Clean up and exit */
1696         partial = chain + depth - 1; /* the whole chain */
1697         /* cleanup: */
1698         while (partial > chain) {
1699             brelse(partial->bh);
1700             partial--;
1701         }
1702
1703         return blkn;
1704     }
1705
1706     void
1707     ext2_ilock(struct inode *ip)
1708     {
1709         struct buf *bp;
1710         struct ext2_inode *raw_inode;
1711         struct ext2_inode_info *ei;
1712
1713         ei = ip->i_private;
1714
1715         if(ip == 0 || ip->ref < 1)
1716             panic("ilock");
1717
1718         acquire(&icache.lock);
1719         while (ip->flags & I_BUSY)
1720             sleep(ip, &icache.lock);
1721         ip->flags |= I_BUSY;
1722         release(&icache.lock);
1723
1724         if (!(ip->flags & I_VALID)) {
1725             raw_inode = ext2_get_inode(&sb[ip->dev], ip->inum, &bp);
1726             // Translate the inode type to xv6 type
1727             if (S_ISDIR(raw_inode->i_mode)) {
1728                 ip->type = T_DIR;
1729             } else if (S_ISREG(raw_inode->i_mode)) {
1730                 ip->type = T_FILE;
1731             } else if (S_ISCHR(raw_inode->i_mode) || S_ISBLK(raw_inode->i_mode)) {
1732                 ip->type = T_DEV;
1733             } else {
1734                 panic("ext2: invalid file mode");

```

```

1735     }
1736     ip->nlink = raw_inode->i_links_count;
1737     ip->size = raw_inode->i_size;
1738     memmove(&ei->i_ei, raw_inode, sizeof(ei->i_ei));
1739
1740     ext2_ops.brelse(bp);
1741     ip->flags |= I_VALID;
1742     if (ip->type == 0)
1743         panic("ext2 ilock: no type");
1744 }
1745 }
1746
1747 int
1748 ext2_writei(struct inode *ip, char *src, uint off, uint n)
1749 {
1750     uint tot, m;
1751     struct buf *bp;
1752
1753     if (ip->type == T_DEV) {
1754         if (ip->major < 0 || ip->major >= NDEV || !devsw[ip->major].write)
1755             return -1;
1756         return devsw[ip->major].write(ip, src, n);
1757     }
1758
1759     if (off > ip->size || off + n < off)
1760         return -1;
1761
1762     // TODO: Verify the max file size
1763
1764     for (tot = 0; tot < n; tot += m, off += m, src += m) {
1765         bp = ext2_ops.bread(ip->dev, ext2_iops.bmap(ip, off / sb[ip->dev].blocksize));
1766         m = min(n - tot, sb[ip->dev].blocksize - off % sb[ip->dev].blocksize);
1767         memmove(bp->data + off % sb[ip->dev].blocksize, src, m);
1768         ext2_ops.bwrite(bp);
1769         ext2_ops.brelse(bp);
1770     }
1771
1772     if (n > 0 && off > ip->size){
1773         ip->size = off;
1774         ext2_iops.iupdate(ip);
1775     }
1776
1777     return n;
1778 }
1779
1780 /*
1781  * Return the offset into page 'page_nr' of the last valid
1782  * byte in that page, plus one.
1783  */
1784 static unsigned
1785 ext2_last_byte(struct inode *inode, unsigned long page_nr)
1786 {
1787     unsigned last_byte = inode->size;
1788     last_byte -= page_nr * sb[inode->dev].blocksize;
1789     if (last_byte > sb[inode->dev].blocksize)
1790         last_byte = sb[inode->dev].blocksize;
1791     return last_byte;
1792 }
1793
1794

```

```

1795 int
1796 ext2_dirlink(struct inode *dp, char *name, uint inum, uint type)
1797 {
1798     int namelen = strlen(name);
1799     struct buf *bh;
1800     unsigned chunk_size = sb[dp->dev].blocksize;
1801     unsigned reclen = EXT2_DIR_REC_LEN(namelen);
1802     unsigned short rec_len, name_len;
1803     char *dir_end;
1804     struct ext2_dir_entry_2 *de;
1805     int n;
1806     int numblocks = (dp->size + chunk_size - 1) / chunk_size;
1807     char *kaddr;
1808
1809     if (ext2_iops.dirlookup(dp, name, 0) != 0) {
1810         return -1;
1811     }
1812
1813     for (n = 0; n <= numblocks; n++) {
1814         bh = ext2_ops.bread(dp->dev, ext2_iops.bmap(dp, n));
1815         kaddr = (char *) bh->data;
1816         de = (struct ext2_dir_entry_2 *) kaddr;
1817         dir_end = kaddr + ext2_last_byte(dp, n);
1818         kaddr += chunk_size - reclen;
1819
1820         while ((char *)de <= kaddr) {
1821             if ((char *)de == dir_end) {
1822                 /* We hit i_size */
1823                 name_len = 0;
1824                 rec_len = chunk_size;
1825                 de->rec_len = chunk_size;
1826                 de->inode = 0;
1827                 goto got_it;
1828             }
1829
1830             if (de->rec_len == 0) {
1831                 return -1;
1832             }
1833
1834             name_len = EXT2_DIR_REC_LEN(de->name_len);
1835             rec_len = de->rec_len;
1836             if (!de->inode && rec_len >= reclen)
1837                 goto got_it;
1838             if (rec_len >= name_len + reclen)
1839                 goto got_it;
1840             de = (struct ext2_dir_entry_2 *) ((char *) de + rec_len);
1841         }
1842
1843         ext2_ops.brelse(bh);
1844     }
1845
1846     return -1;
1847
1848 got_it:
1849     if (de->inode) {
1850         struct ext2_dir_entry_2 *del = (struct ext2_dir_entry_2 *) ((char *) de + name_len);
1851         del->rec_len = rec_len - name_len;
1852         de->rec_len = name_len;
1853         de = del;
1854     }

```

```

1855     de->name_len = namelen;
1856     strncpy(de->name, name, namelen);
1857     de->inode = inum;
1858
1859     // Translate the xv6 to inode type type
1860     if (type == T_DIR) {
1861         de->file_type = EXT2_FT_DIR;
1862     } else if (type == T_FILE) {
1863         de->file_type = EXT2_FT_REG_FILE;
1864     } else {
1865         // We did not treat char and block devices with difference.
1866         panic("ext2: invalid inode mode");
1867     }
1868
1869     ext2_ops.bwrite(bh);
1870     ext2_ops.brelse(bh);
1871
1872     if ((n + 1) * chunk_size > dp->size) {
1873         dp->size += rec_len;
1874         ext2_iops.iupdate(dp);
1875     }
1876
1877     return 0;
1878 }
1879
1880 int
1881 ext2_isdirempty(struct inode *dp)
1882 {
1883     struct buf *bh;
1884     unsigned long i;
1885     char *kaddr;
1886     struct ext2_dir_entry_2 *de;
1887     int chunk_size = sb[dp->dev].blocksize;
1888     int numblocks = (dp->size + chunk_size - 1) / chunk_size;
1889
1890     for (i = 0; i < numblocks; i++) {
1891         bh = ext2_ops.bread(dp->dev, ext2_iops.bmap(dp, i));
1892
1893         if (!bh) {
1894             panic("ext2_isemptydir error");
1895         }
1896
1897         kaddr = (char *)bh->data;
1898         de = (struct ext2_dir_entry_2 *)kaddr;
1899         kaddr += ext2_last_byte(dp, i) - EXT2_DIR_REC_LEN(1);
1900
1901         while ((char *)de <= kaddr) {
1902             if (de->rec_len == 0) {
1903                 goto not_empty;
1904             }
1905             if (de->inode != 0) {
1906                 /* check for . and .. */
1907                 if (de->name[0] != '.')
1908                     goto not_empty;
1909                 if (de->name_len > 2)
1910                     goto not_empty;
1911                 if (de->name_len < 2) {
1912                     if (de->inode != dp->inum)
1913                         goto not_empty;
1914                 } else if (de->name[1] != '.')

```

```

1915         goto not_empty;
1916     }
1917     de = (struct ext2_dir_entry_2 *)((char *)de + de->rec_len);
1918 }
1919 ext2_ops.brelse(bh);
1920 }
1921 return 1;
1922
1923 not_empty:
1924     ext2_ops.brelse(bh);
1925     return 0;
1926 }
1927
1928 int
1929 ext2_unlink(struct inode *dp, uint off)
1930 {
1931     struct buf *bh;
1932     uint bn, offset;
1933     struct ext2_dir_entry_2 *dir;
1934     int chunk_size;
1935
1936     chunk_size = sb[dp->dev].blocksize;
1937     bn = off / sb[dp->dev].blocksize;
1938     offset = off % sb[dp->dev].blocksize;
1939     bh = ext2_ops.bread(dp->dev, ext2_iops.bmap(dp, bn));
1940
1941     dir = (struct ext2_dir_entry_2 *)(bh->data + offset);
1942     char *kaddr = (char *)bh->data;
1943
1944     unsigned from = ((char *)dir - kaddr) & ~(chunk_size - 1);
1945     unsigned to = ((char *)dir - kaddr) + dir->rec_len;
1946
1947     struct ext2_dir_entry_2 *pde = 0;
1948     struct ext2_dir_entry_2 *de = (struct ext2_dir_entry_2 *) (kaddr + from);
1949
1950     while ((char *)de < (char *)dir) {
1951         if (de->rec_len == 0)
1952             panic("ext2_unlink: invalid dir content");
1953         pde = de;
1954         de = (struct ext2_dir_entry_2 *)((char *)de + de->rec_len);
1955     }
1956
1957     if (pde) {
1958         from = (char *)pde - (char *)bh->data;
1959         pde->rec_len = to - from;
1960     }
1961
1962     dir->inode = 0;
1963
1964     ext2_ops.bwrite(bh);
1965     ext2_ops.brelse(bh);
1966
1967     return 0;
1968 }
1969
1970 int
1971 ext2_namecmp(const char *s, const char *t)
1972 {
1973     unsigned short slen = strlen(s), tlen = strlen(t);
1974     unsigned short size = slen;

```

```

1975
1976     if (slen != tlen)
1977         return -1;
1978
1979     if (tlen > slen)
1980         size = tlen;
1981
1982     return strncmp(s, t, size);
1983 }
1984
1985 static struct ext2_inode *
1986 ext2_get_inode(struct superblock *sb, uint ino, struct buf **bh)
1987 {
1988     struct buf * bp;
1989     unsigned long block_group;
1990     unsigned long block;
1991     unsigned long offset;
1992     struct ext2_group_desc *gdp;
1993     struct ext2_inode *raw_inode;
1994
1995     if ((ino != EXT2_ROOT_INO && ino < EXT2_FIRST_INO(sb)) ||
1996         ino > EXT2_SB(sb)->s_es->s_inodes_count)
1997         panic("Ext2_invalid_inode_number");
1998
1999     block_group = (ino - 1) / EXT2_INODES_PER_GROUP(sb);
2000     gdp = ext2_get_group_desc(sb, block_group, 0);
2001     if (!gdp)
2002         panic("Invalid_group_descriptor_at_ext2_get_inode");
2003
2004     /*
2005      * Figure out the offset within the block group inode table
2006      */
2007     offset = ((ino - 1) % EXT2_INODES_PER_GROUP(sb)) * EXT2_INODE_SIZE(sb);
2008     block = gdp->bg_inode_table +
2009         (offset >> EXT2_BLOCK_SIZE_BITS(sb));
2010
2011     if (!(bp = ext2_ops.bread(sb->minor, block)))
2012         panic("Error_on_read_the_block_inode");
2013
2014     offset &= (EXT2_BLOCK_SIZE(sb) - 1);
2015     raw_inode = (struct ext2_inode *) (bp->data + offset);
2016     if (bh)
2017         *bh = bp;
2018
2019     return raw_inode;
2020 }
2021
2022 /**
2023  * Its is called because the icache lookup failed
2024  */
2025 int
2026 ext2_fill_inode(struct inode *ip) {
2027     struct ext2_inode_info *ei;
2028     struct ext2_inode *raw_inode;
2029     struct buf *bh;
2030
2031     ei = alloc_ext2_inode_info();
2032
2033     if (ei == 0)
2034         panic("No_memory_to_alloc_ext2_inode");

```

```
2035
2036     raw_inode = ext2_get_inode(&sb[ip->dev], ip->inum, &bh);
2037     memmove(&ei->i_ei, raw_inode, sizeof(ei->i_ei));
2038     ip->i_private = ei;
2039
2040     ext2_ops.brelse(bh);
2041
2042     // Translate the inode type to xv6 type
2043     if (S_ISDIR(ei->i_ei.i_mode)) {
2044         ip->type = T_DIR;
2045     } else if (S_ISREG(ei->i_ei.i_mode)) {
2046         ip->type = T_FILE;
2047     } else if (S_ISCHR(ei->i_ei.i_mode) || S_ISBLK(ei->i_ei.i_mode)) {
2048         ip->type = T_DEV;
2049     } else {
2050         panic("ext2: invalid file mode");
2051     }
2052
2053     ip->nlink = ei->i_ei.i_links_count;
2054     ip->size = ei->i_ei.i_size;
2055     return 1;
2056 }
2057
2058 struct inode*
2059 ext2_iget(uint dev, uint inum)
2060 {
2061     return iget(dev, inum, &ext2_fill_inode);
2062 }
```