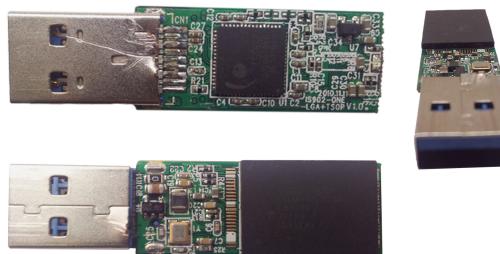
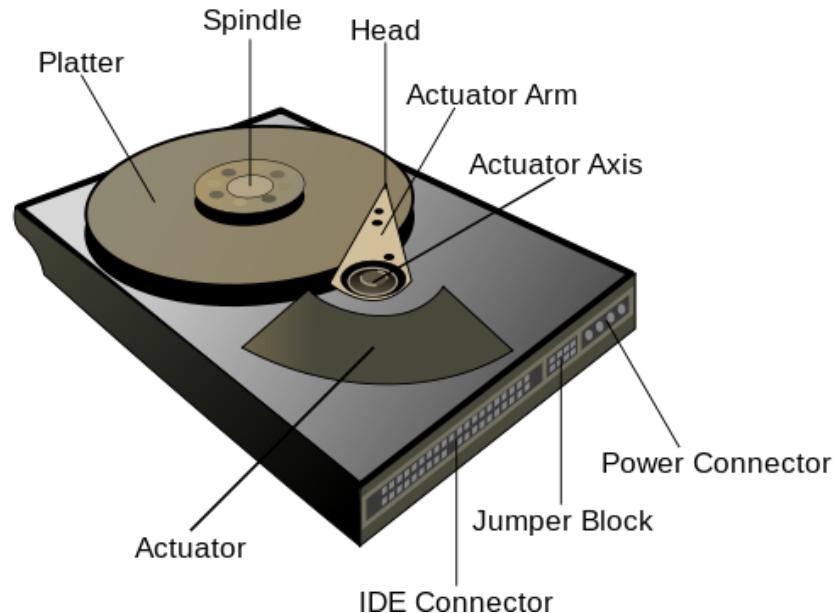


An Intro To File Systems

Pierre Rioux

**ACElab Developer Meeting,
April 2015, Common Era**

A disk device



"SixHardDriveFormFactors" by Paul R. Potts - Provided by Author.

Licensed under CC BY-SA 3.0 us via Wikimedia Commons - <http://commons.wikimedia.org/wiki/File:SixHardDriveFormFactors.jpg#/media/File:SixHardDriveFormFactors.jpg>

"Hard drive-en" by I. Surachit.

Licensed under CC BY-SA 3.0 via Wikimedia Commons - http://commons.wikimedia.org/wiki/File:Hard_drive-en.svg#/media/File:Hard_drive-en.svg

"USB 3.0 Flash Drive PCB" by Ravenperch - camera photograph.

Licensed under CC0 via Wikimedia Commons - http://commons.wikimedia.org/wiki/File:USB_3.0_Flash_Drive_PCB.jpg#/media/File:USB_3.0_Flash_Drive_PCB.jpg

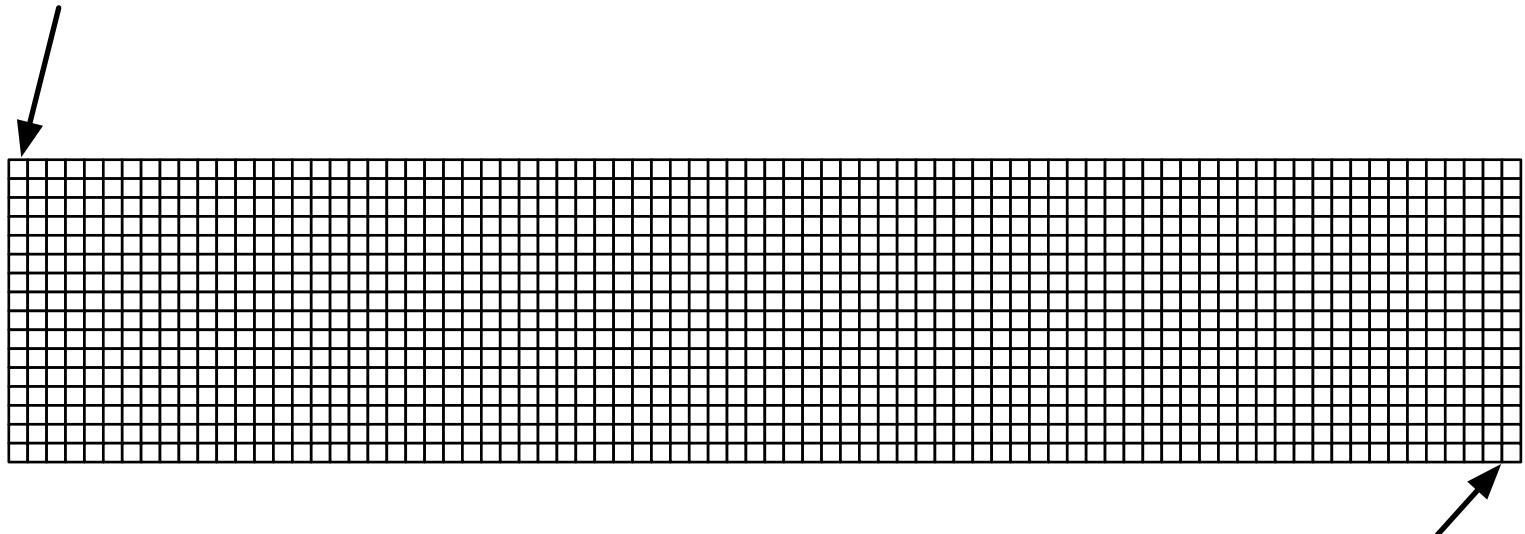
A block device

Basic functionality:

- The device stores a set of *blocks* or *sectors*.
- Blocks are numbered from 0 to n.
- Blocks contain a fixed set of bytes (old: 512, new: 4096).
- Reading and writing is made in units of full blocks.
- The device does not understand anything about files, directories, permissions, etc.
- Once connected, the hardware interface only shuffles blocks back and forth.

All Your Blocks Are Belong To Us

block #0



block #1,000,000

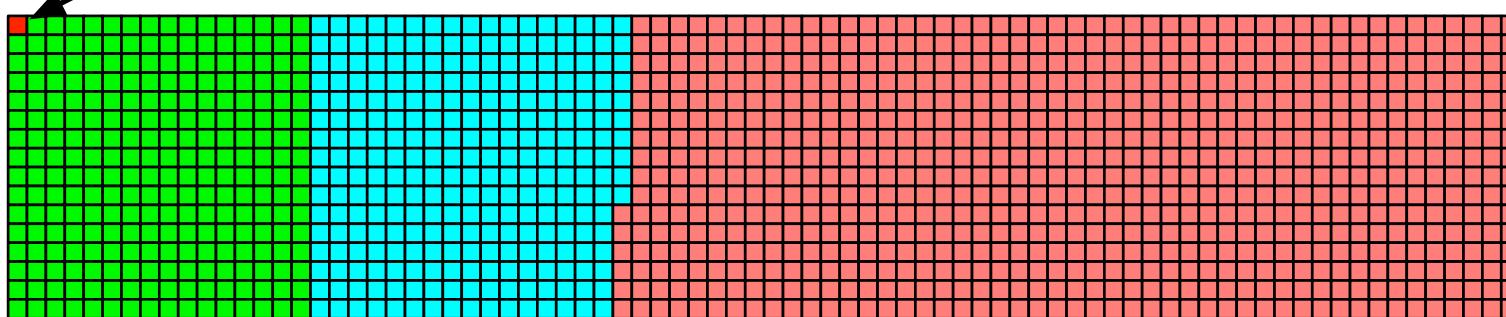
**Each block
contains
4096 bytes
of data**

Partitions

- They are a *convention*.
- There are several ways to record them (dos, gpt, apple...).
- Often recorded within the first few blocks.
- The device does not understand them either.

Partition table:

partition	start	end
1	1	200000
2	200001	421222
3	421223	1000000

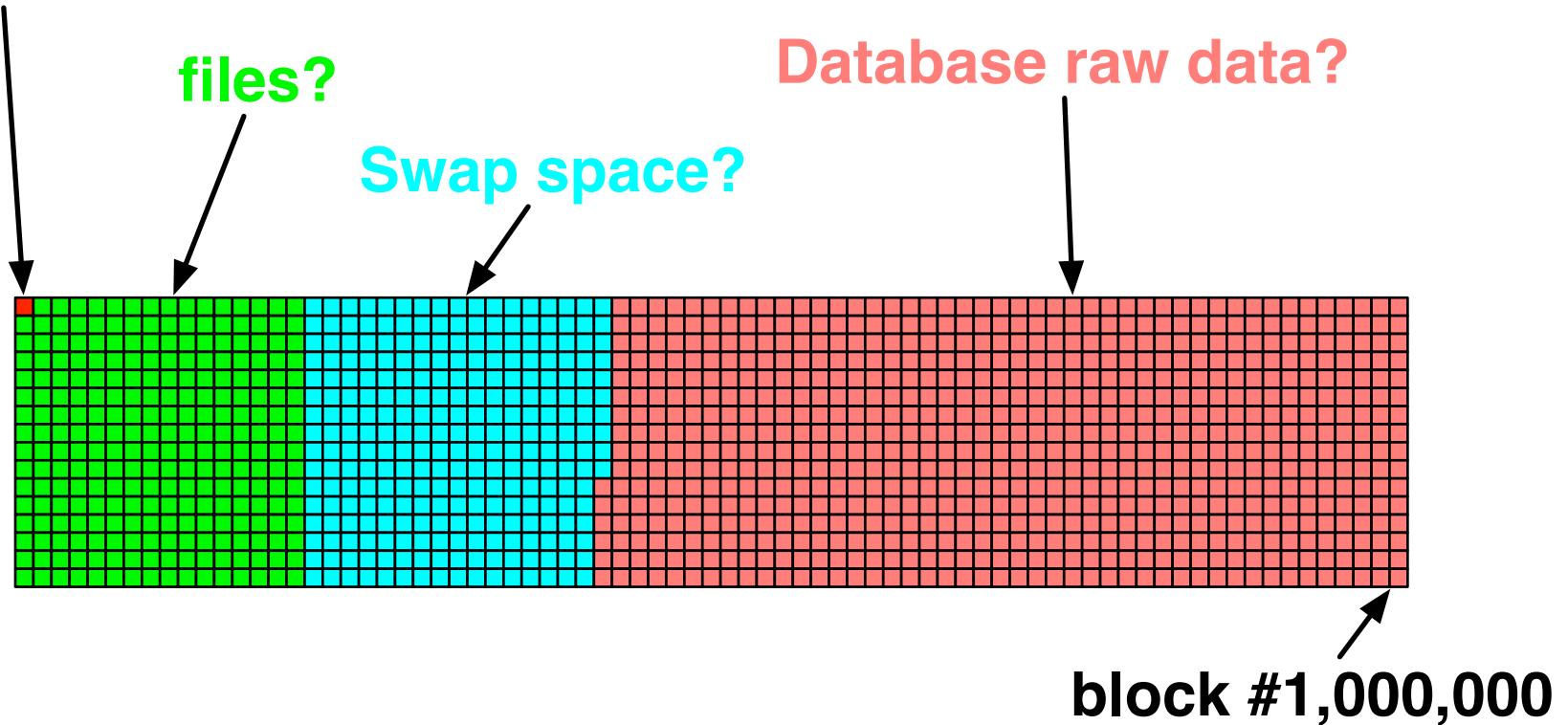


block #1,000,000

Partition roles

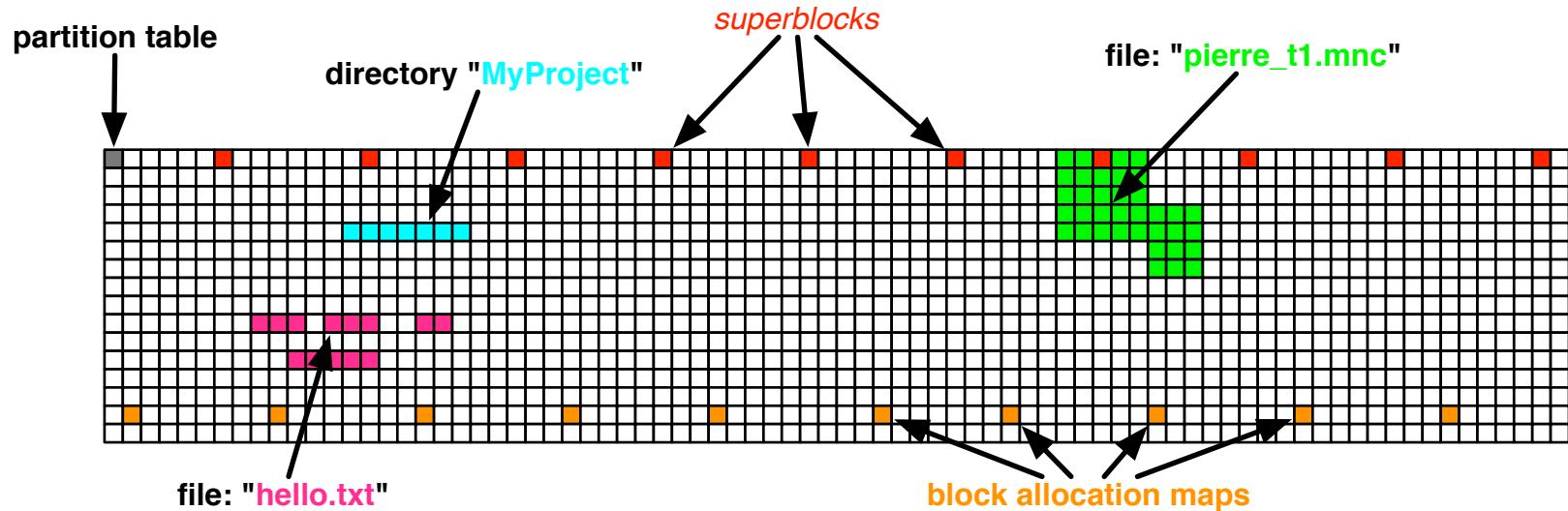
- Roles are not enforced by the device.

Partition table



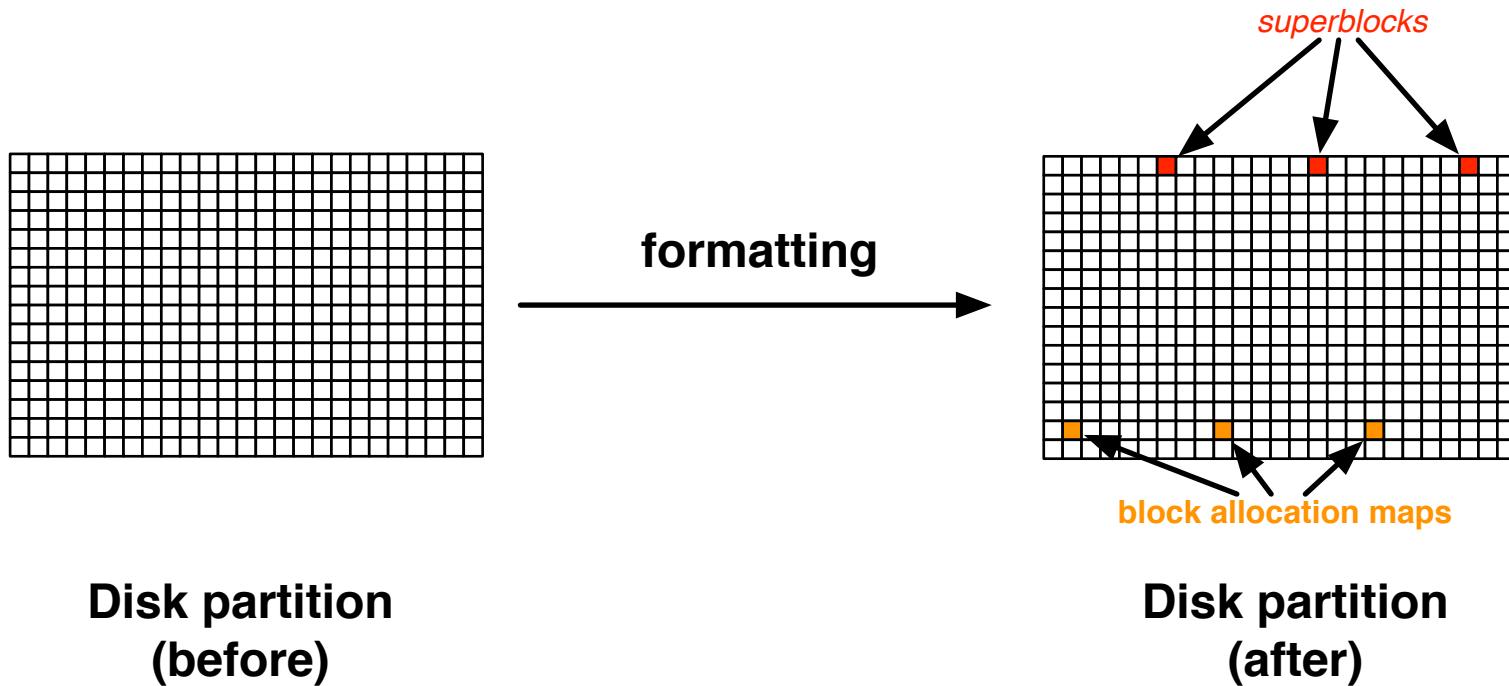
File systems

- They are *software components* of an operating system.
- They manage *files*, *directories*, *etc* within a partition.
- They keep their structures persistent by using special blocks within the partition (think about power failures!).
- By extension, we use the word *file system* to talk about the set of files and special blocks within a partition.



Formatting

- Given a partition with no particular data on it, installing an initial file system is called *formatting*.

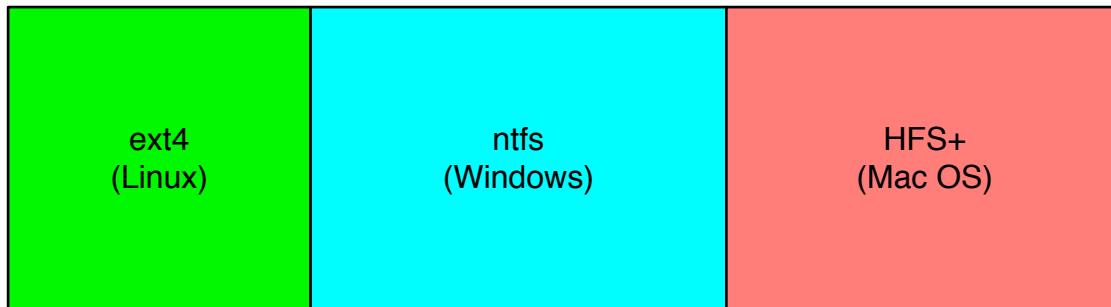


Types of file systems

Names	Years	Typical OS
fat, fat16, fat32	1977, 1984, 1996	Windows, USB keys
ext, ext2, ext3, ext4	1992, 1993, 1999, 2006	Linux
ntfs	1993, 2001	Windows
btrfs	2007	Linux, UNIX
ufs	1994	UNIX
hfs, hfs+	1985, 1998	Mac OS

Hybrid disks

- A single device can have multiple partitions, each with a different type of file system.
- The data in each file system can be read by the same computer, as long as it has the appropriate file system software component in its operating system.



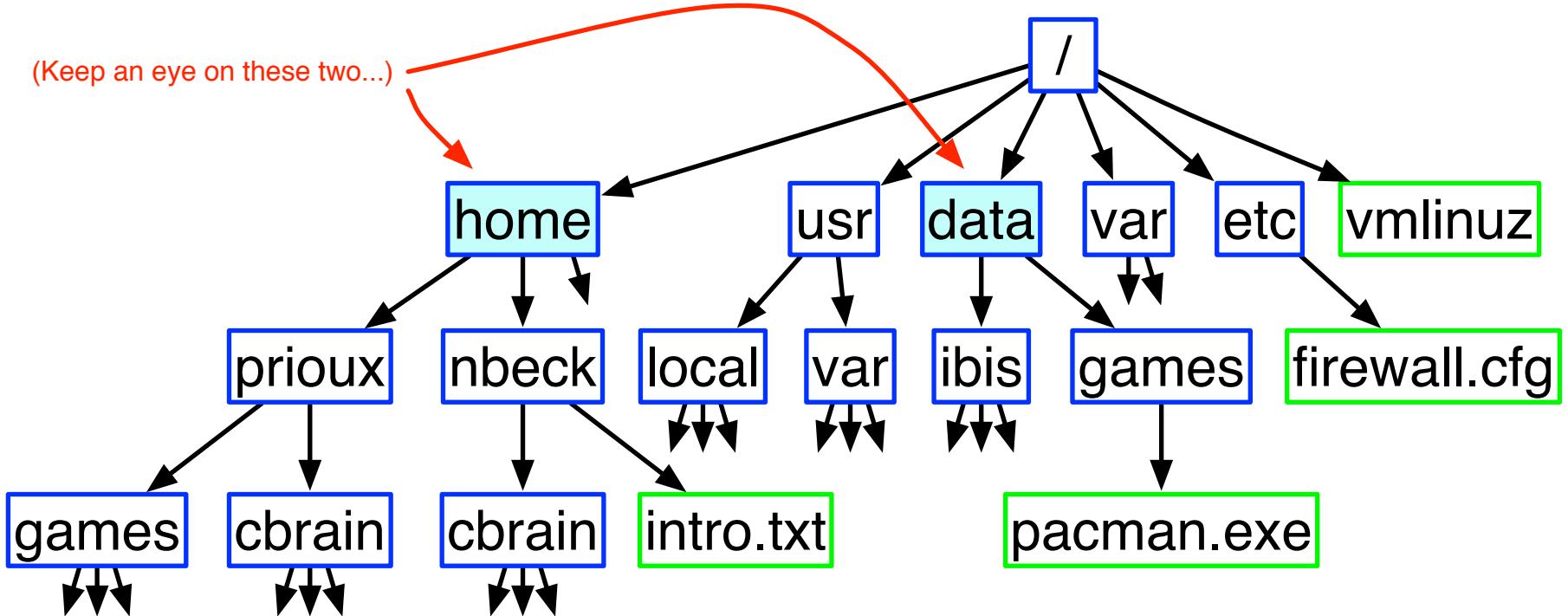
A weird USB key; 3 partitions and 3 file systems types

The UNIX file system 1

- UNIX-like operating systems present all files under a single tree rooted at "/".
- / itself is the main directory of the filesystem of the boot disk.
- The entire set of files visible under / is often called ***THE*** file system, even though it may be comprised of several file systems.
- It is organized in a tree structure.

The UNIX file system 2

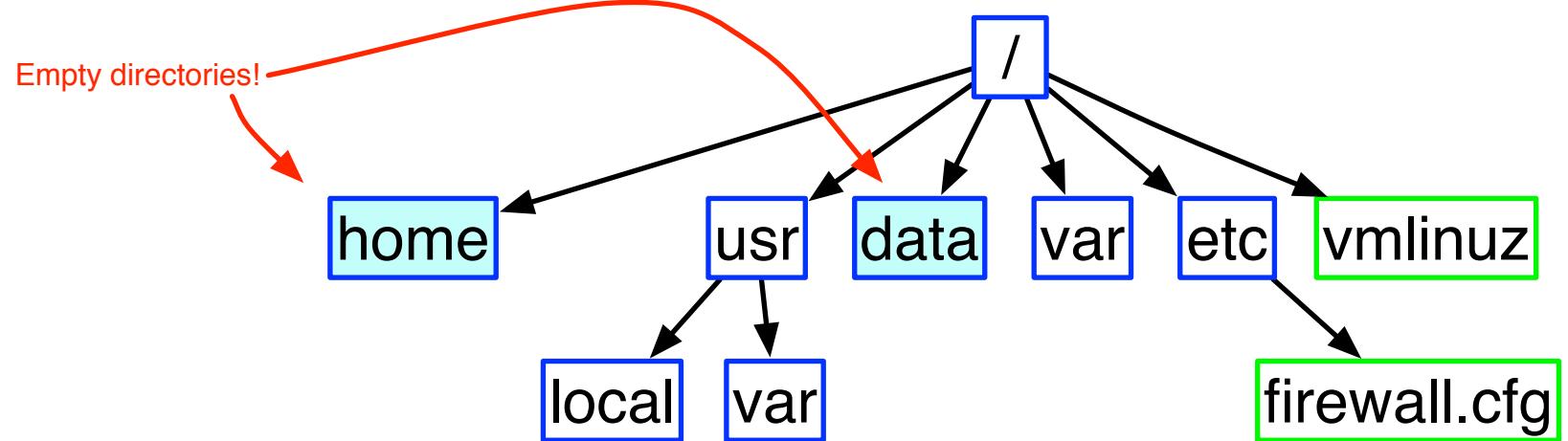
- ... we're all familiar with this view, right?



Directories are in blue, files are in green.

UNIX mountpoints 1

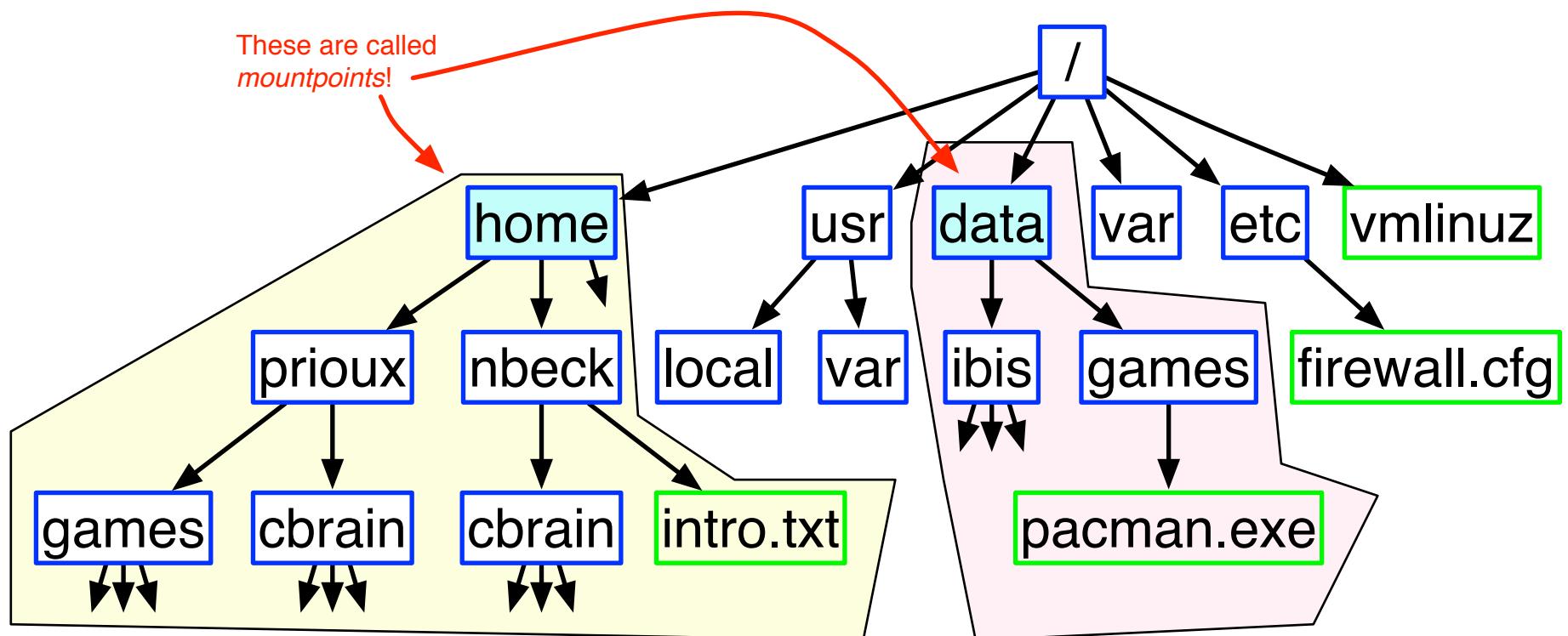
- When the system boots, the file system might be partial...



- These are the files and directories visible on the *root* file system.

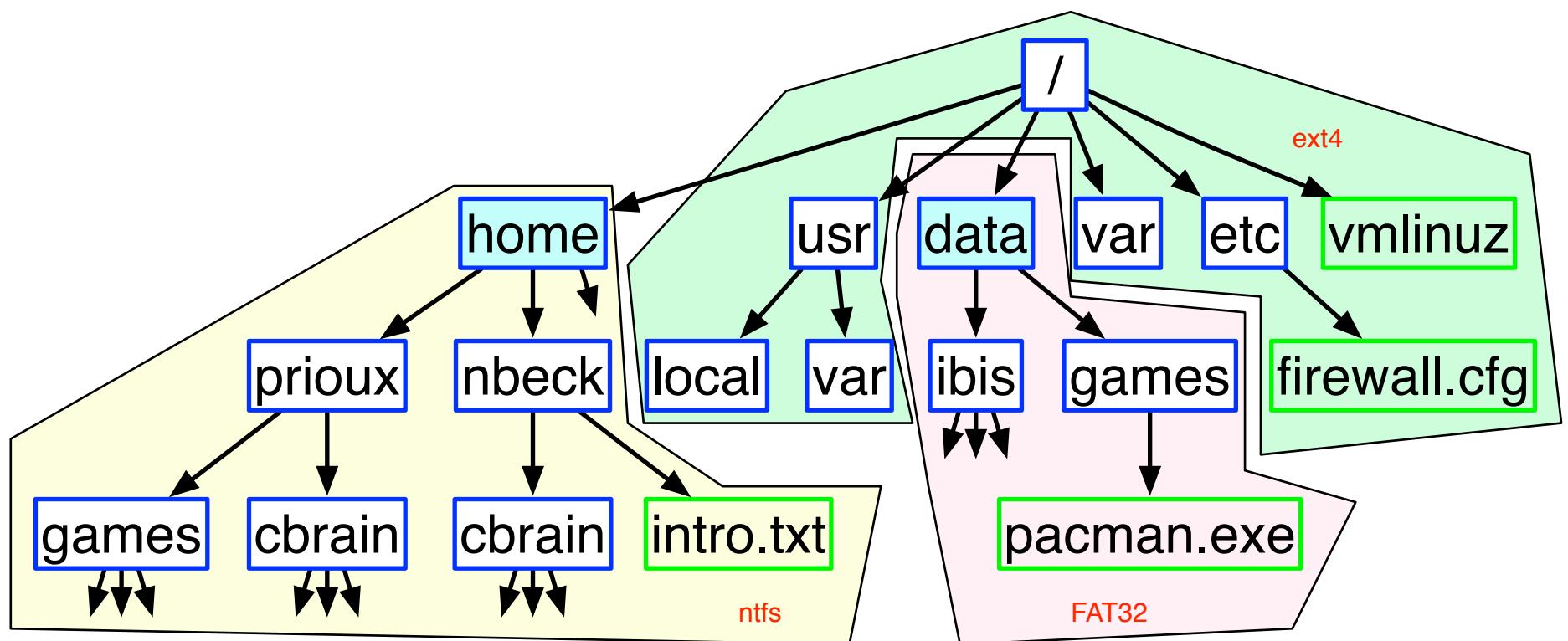
UNIX mountpoints 2

- Other file systems are attached to *mountpoints*.



UNIX mountpoints 3

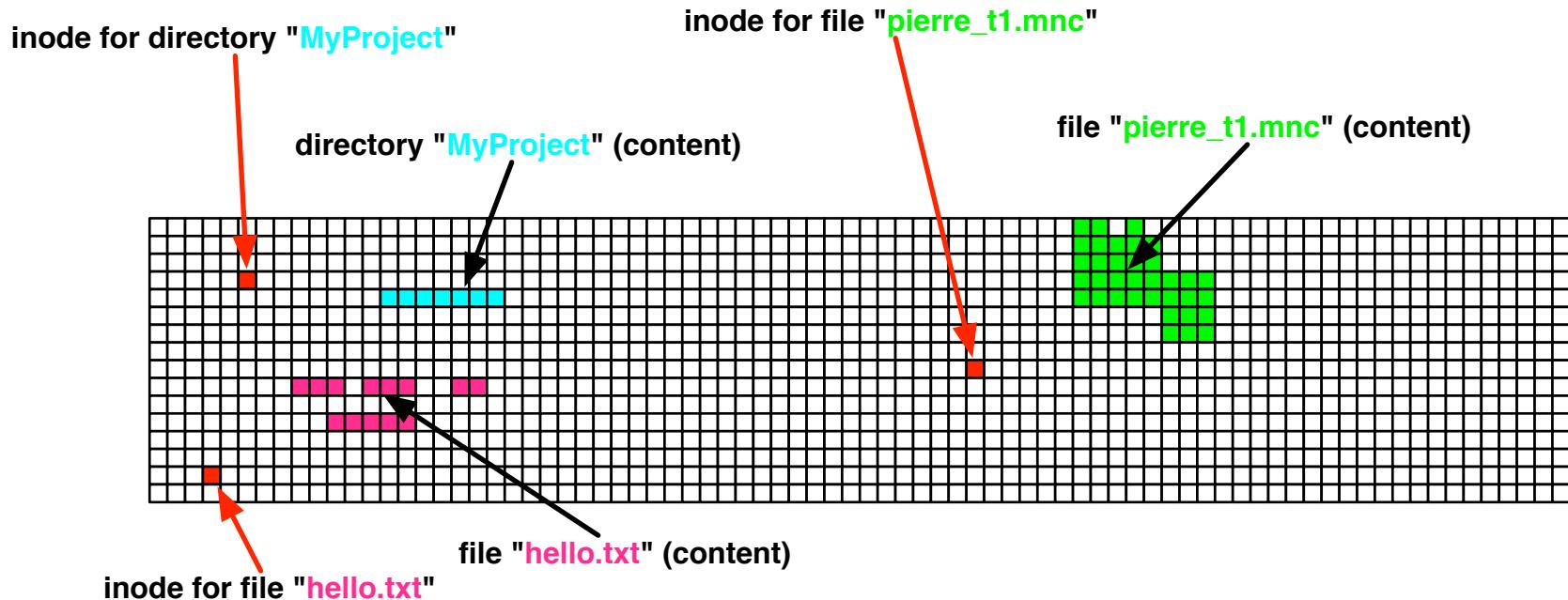
- File systems *mounted* can be of different types.



inodes

(not an Apple product)

- Given a UNIX-like file system, each file or directory on it has an *inode*.



Content of an inode

- An inode contains many pieces of the information familiar to all programmers:
 - user ID
 - group ID
 - size
 - type
 - file permissions
 - several timestamps
 - list of blocks with contents
 - link count (more on that later)

Quick Quiz: something important seems to be missing!

Directory content 1

- The bytes in the blocks storing the *content* of a directory are organized into a flat list of entries.
- Each entry is a simple pair: (name, inode number).
- They are not stored in any particular order.
- You can't tell which entry is which type!

names	inode #
.	25372
..	26471
hello.txt	11242
myproject	4426272
pierre_t1.mnc	2342
games	440224
filesystems.pdf	90210

Directory content 2

- New files are added at the bottom, making the list longer and longer.
- File that are deleted leave a hole in the table.
- There is often NO optimization to compact this list if there are too many holes.

names	inode #
.	25372
..	26471
hello.txt	
myproject	4426272
pierre_t1.mnc	2342
games	440224
filesystems.pdf	90210
.bashrc	227491

2) this file was deleted



1) this file was added



Directory content 3

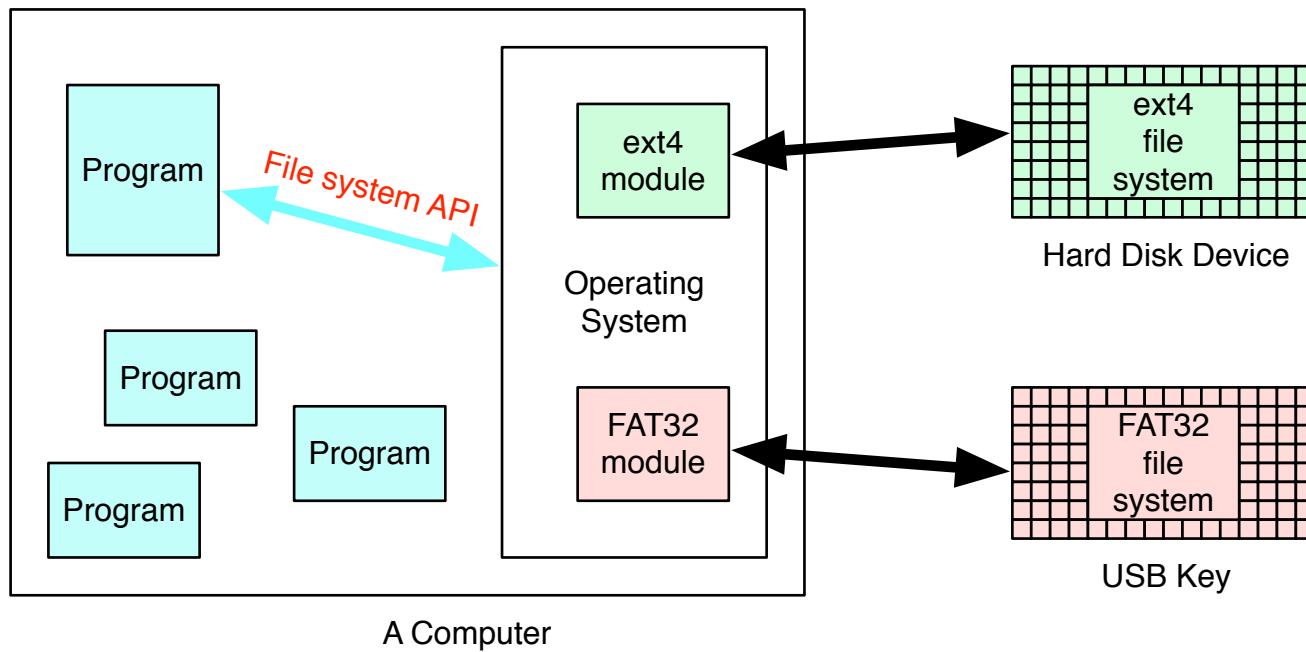
- If millions of entries are put in a single directory, the list becomes enormously inefficient.
- The operating system scans the list linearly from beginning to end whenever an operation is performed.
- "echo abc > def" must scan for file 'def' anywhere in the directory, if it exists, and otherwise create a new entry.

Directory content 4

```
root@brainstorm cbrain # ls -la
total 108
drwx----- 5 cbrain cbrain 4096 Apr  8 12:29 .
drwxr-xr-x 10 root  root  4096 Aug 20  2013 ..
-rw-----  1 cbrain cbrain 3402 Apr 15 17:05 .bash_history
-rw-r--r--  1 cbrain cbrain   33 Apr 16 2012 .bash_logout
-rw-r--r--  1 cbrain cbrain  176 Apr 16 2012 .bash_profile
-rw-r--r--  1 cbrain cbrain  124 Apr 16 2012 .bashrc
-rw-----  1 cbrain cbrain   35 Nov  2 2012 .lessht
drwxr-xr-x  5 cbrain cbrain 4096 Nov  6 2013 .mozilla
drwx-----  2 cbrain cbrain 4096 Apr 17 2012 .ssh
-rw-r--r--  1 cbrain cbrain 5988 Apr 24 2014 SSLCerts-2014-04.tar.gz
drwxr-xr-x  2 cbrain cbrain 4096 Sep 22 2014 .storage_system
-rw-----  1 cbrain cbrain  956 Jul 25 2013 .viminfo
-rw-----  1 cbrain cbrain  292 Apr  8 12:29 .Xauthority
```

```
root@brainstorm cbrain # ls -1fi
28737551 SSLCerts-2014-04.tar.gz
28737555 .storage_system
2 ..
28737545 .ssh
28737547 .lessht
28737549 .viminfo
28737597 .Xauthority
28737543 .bash_logout
28737537 .
28737538 .bash_profile
28737544 .bash_history
28737542 .bashrc
28737539 .mozilla
```

API



API for files 1

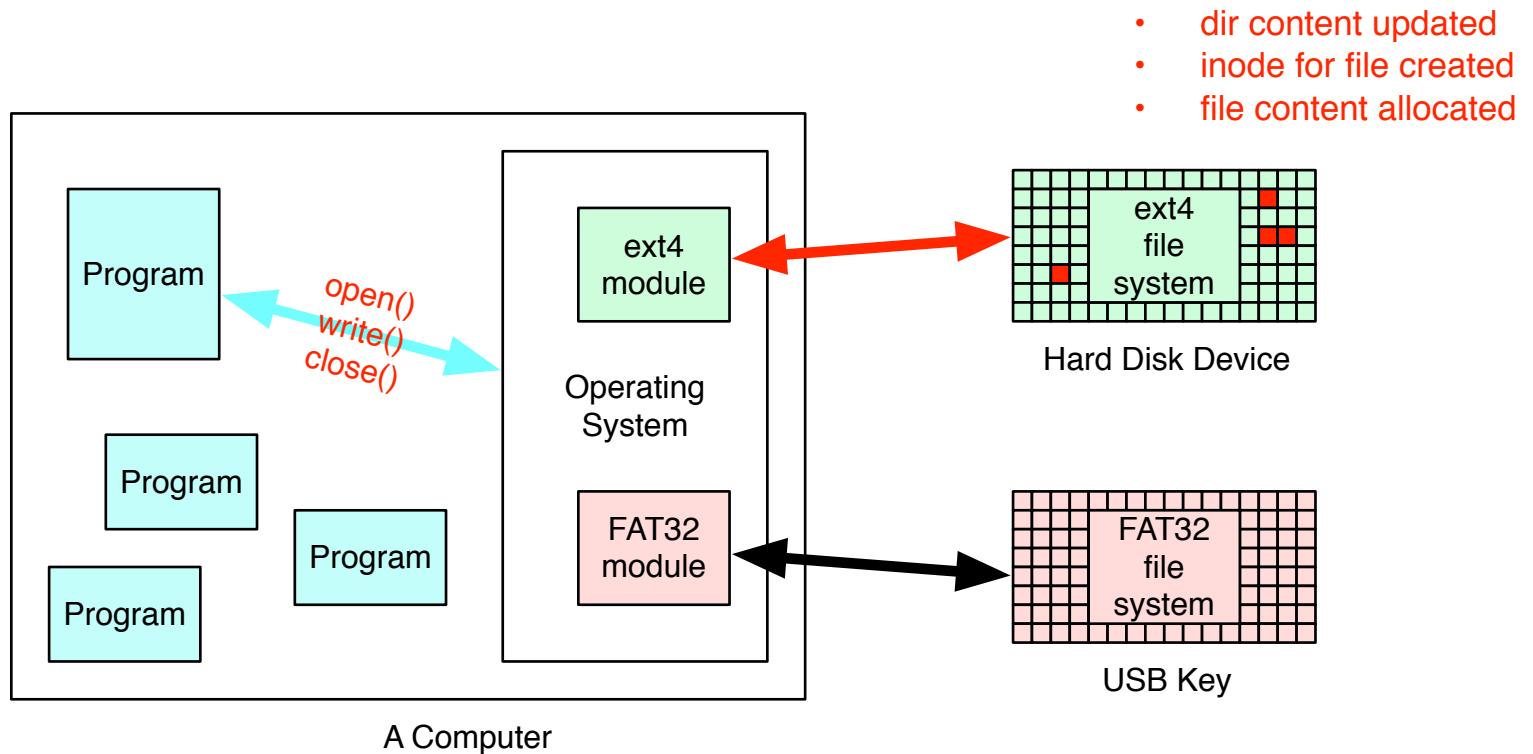
- open()
- read()
- write()
- close()
- unlink()
- seek()

```
#include <unistd.h>
#include <fcntl.h>

int main(int argc, char *argv[]) {
    int fd;

    fd = open("hello.txt", O_CREAT|O_WRONLY, 0700);
    write(fd, "hi!\n", 4);
    close(fd);
}
```

API for files 2



API for directories

- opendir()
- readdir()
- closedir()
- rmdir()

IMPORTANT! No search!



```
#include <stdio.h>
#include <unistd.h>
#include <dirent.h>

int main(int argc, char *argv[]) {
    DIR *dh;
    struct dirent *entry;

    dh=opendir(".");
    while ((entry = readdir(dh)) != NULL) {
        printf("%lld %s\n",entry->d_ino,entry->d_name);
    }
    closedir(dh);
}
```

API for files or directories

- stat()
- rename()
- link()
- symlink()
- lstat()
- readlink()

... and others for changing inode properties

Symbolic links 1

The best way to understand a symbolic link is to imagine that it is...

- a TEXT file...
- ... with a SINGLE line of text...
- ... which describe where to find another file or directory.

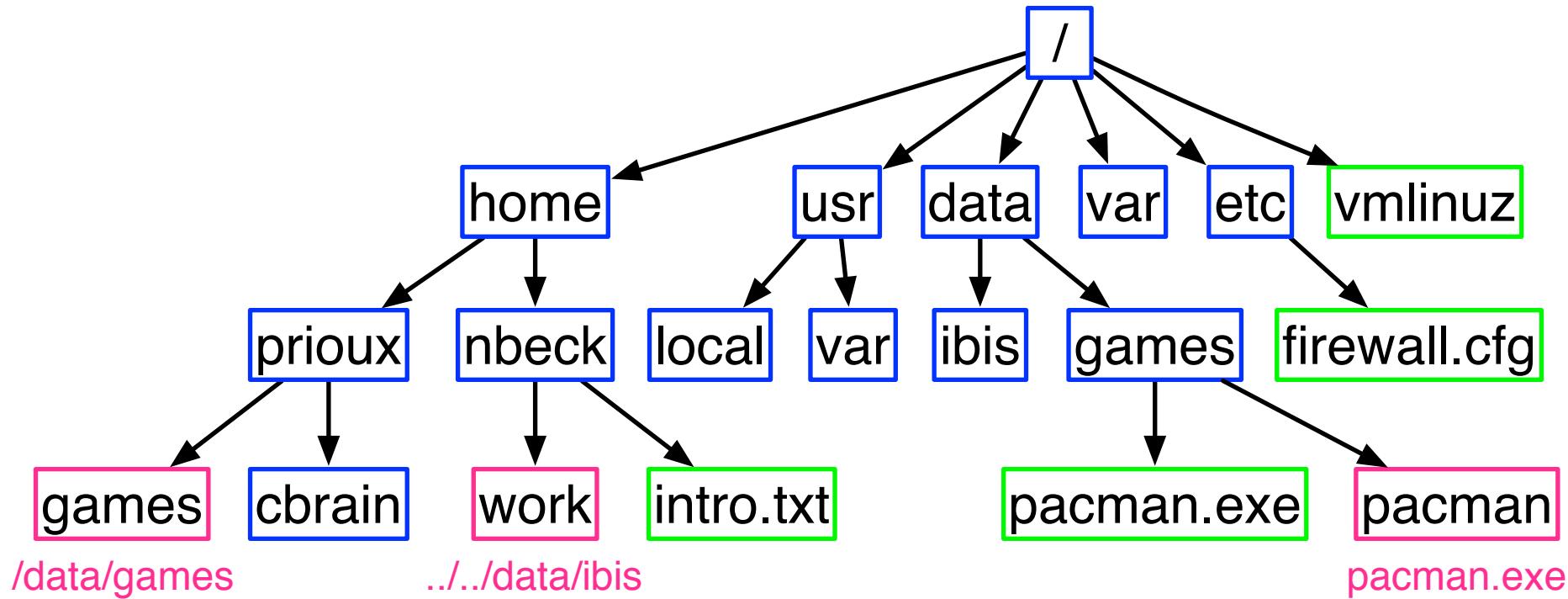
Called the *link value*



The link value can be...

- a relative path: "`../../abc.txt`" or "`hi/how/do/you/do.txt`"
- an absolute path: "`/Volumes/Maps/Earth.pdf`"

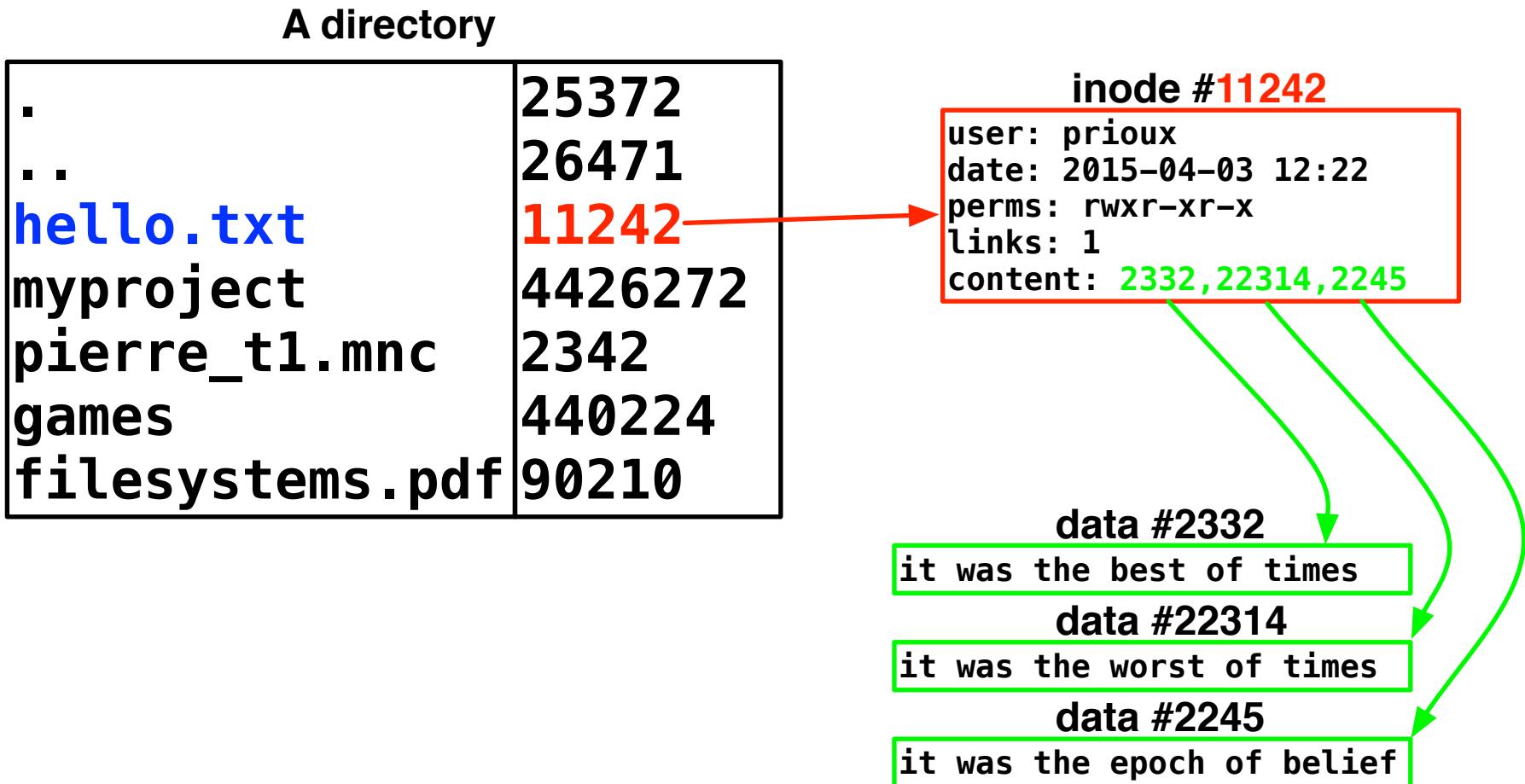
Symbolic links 2



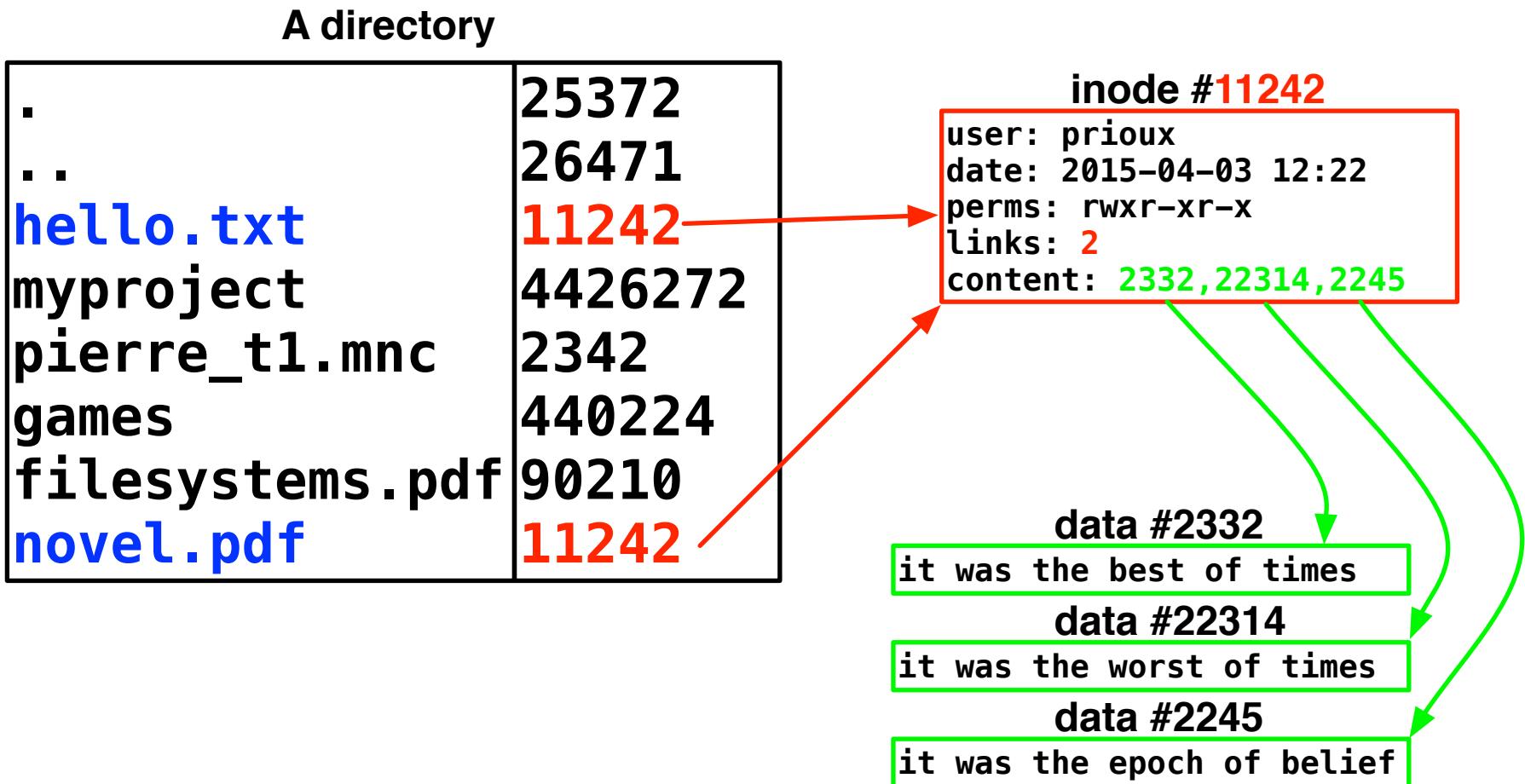
Notes:

- Not all file system types support symbolic links.
- Symbolic links can point to stuff on other filesystems.
- The link value can be anything and might not even point to anything that exists!

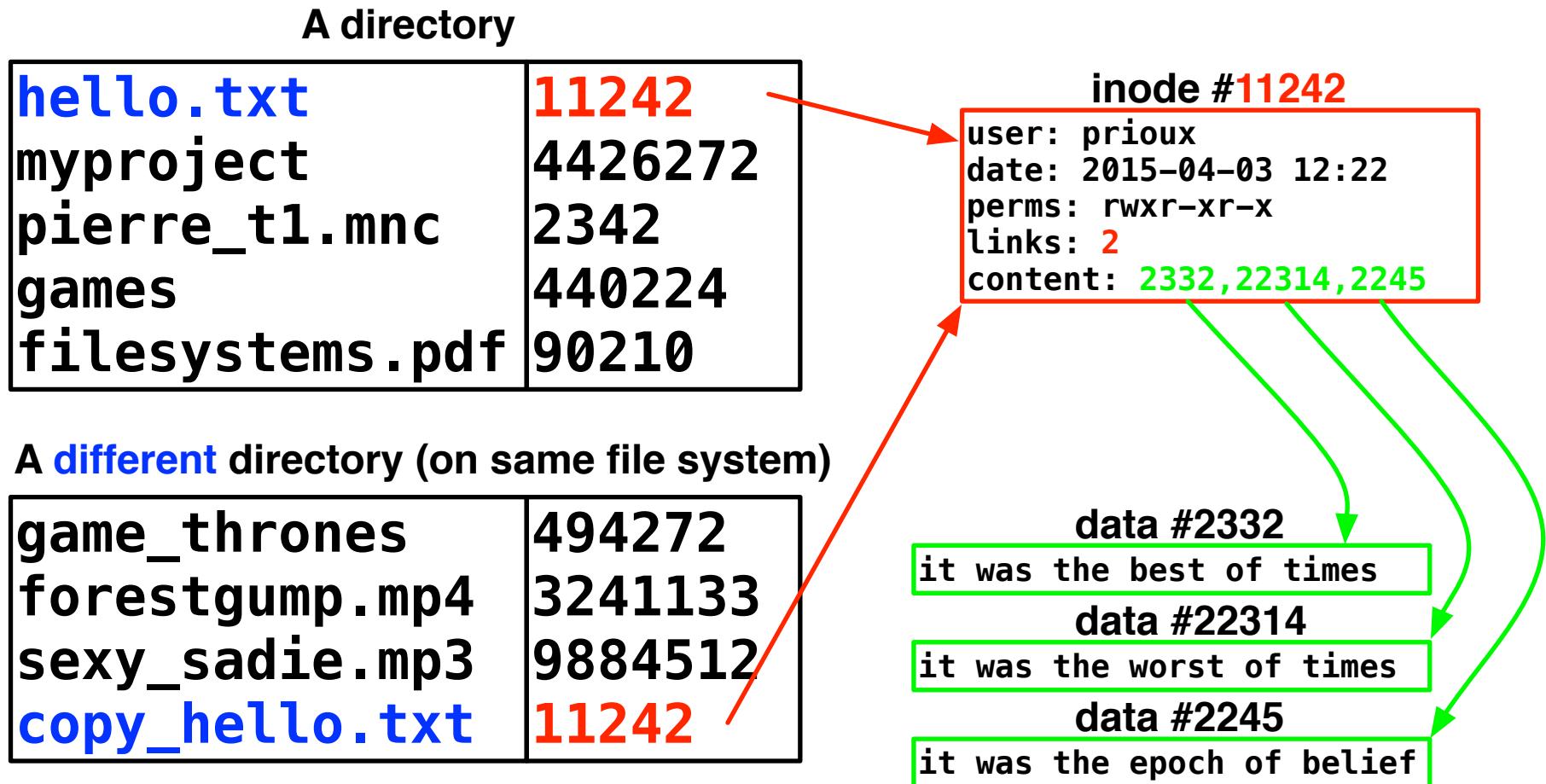
Hard links 1



Hard links 2



Hard links 3



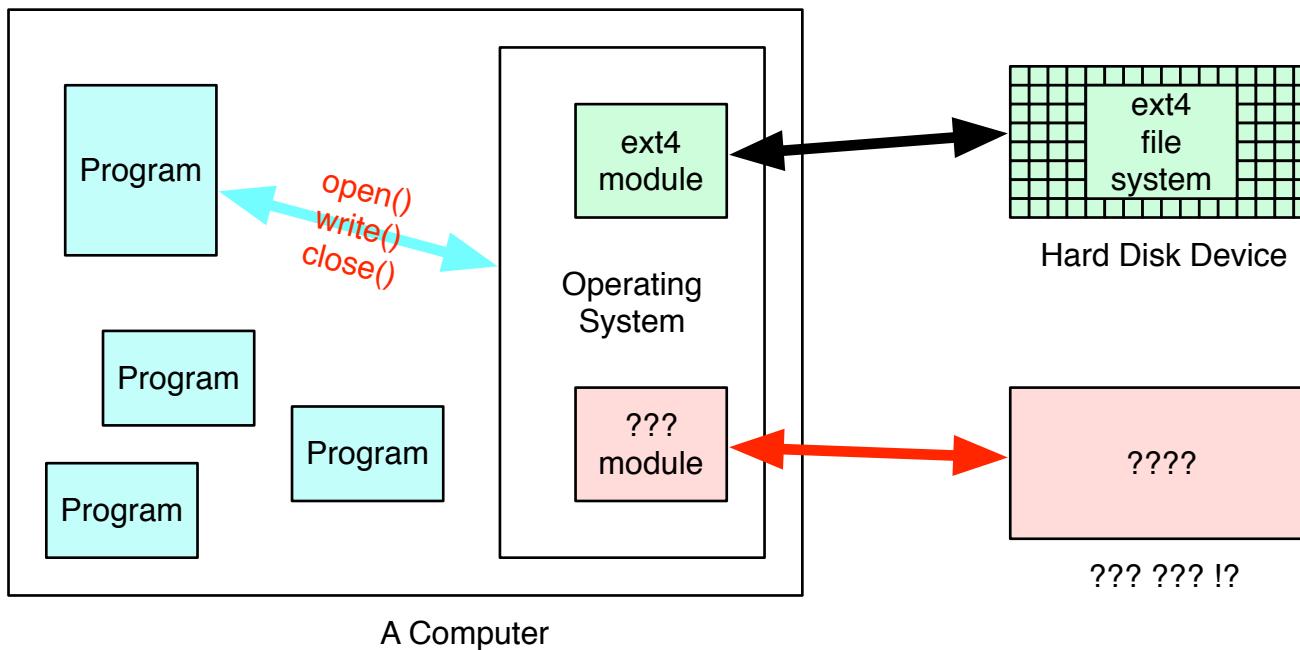
Hard links 4

- Hard links can only be created between contents on the SAME filesystem.
- None of the directory entries is more special than the others.
- The content and the inode is removed only when the last link is deleted.
- Not all file system types allow hard links.
- Modifying the content modifies it for all links!
- Hard links are NOT permitted for directories; the exceptions are "." and ".." (but done by OS).

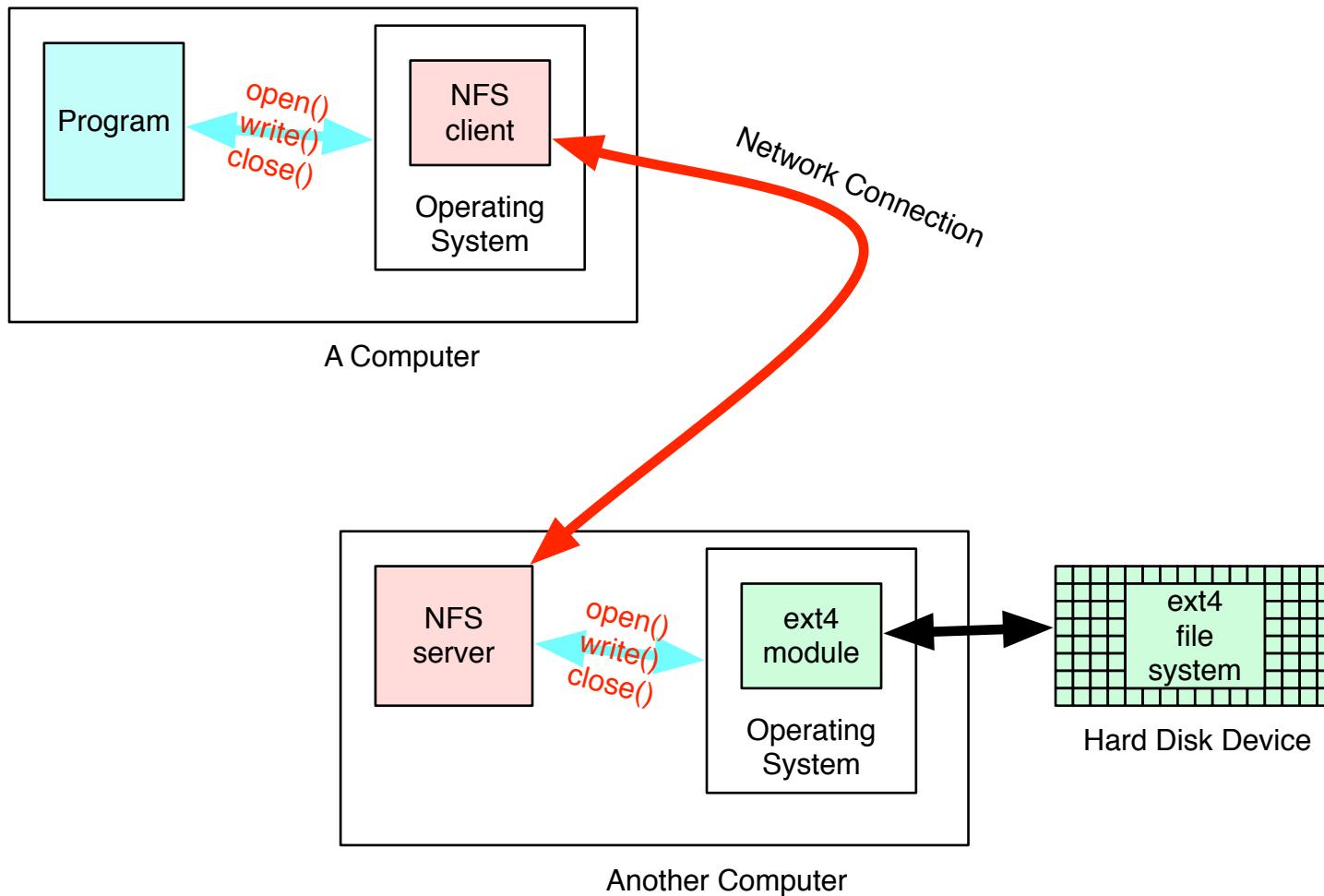
Links Demo

```
Bianca home Demo >
Bianca home Demo > echo "some content" > hello.txt
Bianca home Demo > ls -l
total 8
-rw-r--r-- 1 prioux staff 13 Apr 23 14:14 hello.txt
Bianca home Demo > ln -s hello.txt sym_to_hello.txt
Bianca home Demo > ls -l
total 16
-rw-r--r-- 1 prioux staff 13 Apr 23 14:14 hello.txt
lrwxr-xr-x 1 prioux staff 9 Apr 23 14:14 sym_to_hello.txt@ -> hello.txt
Bianca home Demo > ln hello.txt hardlink_hello.txt
Bianca home Demo > ls -l
total 24
-rw-r--r-- 2 prioux staff 13 Apr 23 14:14 hardlink_hello.txt
-rw-r--r-- 2 prioux staff 13 Apr 23 14:14 hello.txt
lrwxr-xr-x 1 prioux staff 9 Apr 23 14:14 sym_to_hello.txt@ -> hello.txt
Bianca home Demo > echo "more content" >> hello.txt
Bianca home Demo > ls -l
total 24
-rw-r--r-- 2 prioux staff 26 Apr 23 14:15 hardlink_hello.txt
-rw-r--r-- 2 prioux staff 26 Apr 23 14:15 hello.txt
lrwxr-xr-x 1 prioux staff 9 Apr 23 14:14 sym_to_hello.txt@ -> hello.txt
Bianca home Demo > cat sym_to_hello.txt
some content
more content
Bianca home Demo > rm hello.txt
remove hello.txt? y
Bianca home Demo > cat sym_to_hello.txt
cat: sym_to_hello.txt: No such file or directory
Bianca home Demo > ls -l
total 16
-rw-r--r-- 1 prioux staff 26 Apr 23 14:15 hardlink_hello.txt
lrwxr-xr-x 1 prioux staff 9 Apr 23 14:14 sym_to_hello.txt@ -> hello.txt
```

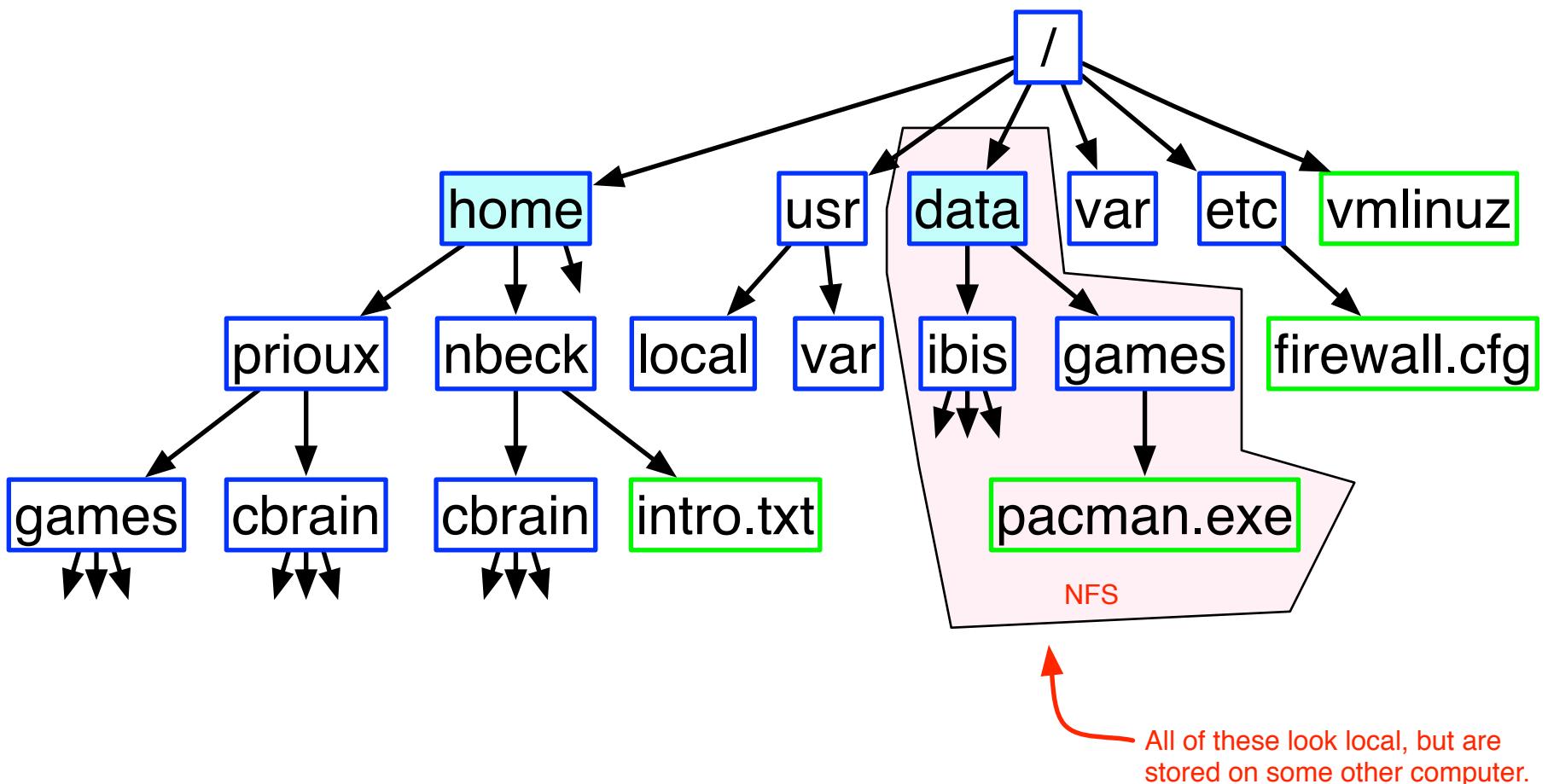
API abstraction



NFS 1



NFS 2



Operations that are FAST within a file system

- Renaming an entry (file or directory) in place.
- Creating a hard link.
- Creating a symbolic link.
- Unlinking a file with more than one link.
- Erasing files, as long as they are small.
- Creating directory entries in directories with few entries.
- Modifying inode information.

Operations that are SLOW within a file system

- Anything that involves manipulating a directory entry in a directory with already a large number of entries.
- Erasing large files.
- Resolving symbolic link spaghetti.

The END

