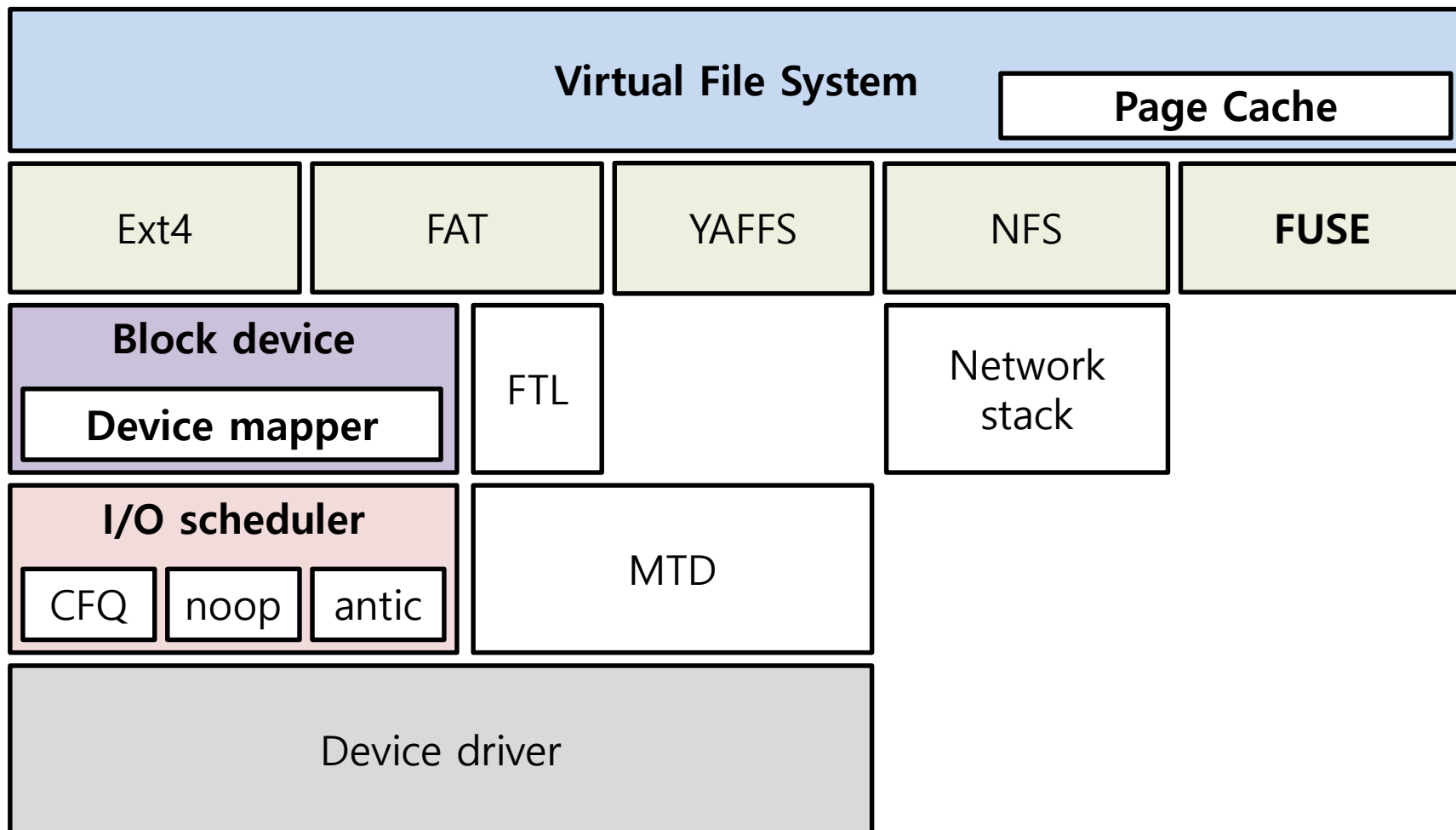


# **Linux virtual file system**



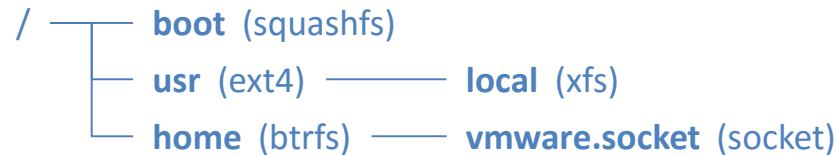
# Storage implementation layers





# Introduction to VFS

- Hordes **different** file system implementations
  - IPC mechanisms (PIPE, FIFO, socket, ...) too



- Abstracts **generic** file system implementations
  - Directory traversal
  - Page cache
- Interfaces with **POSIX** system call APIs
  - File descriptor management



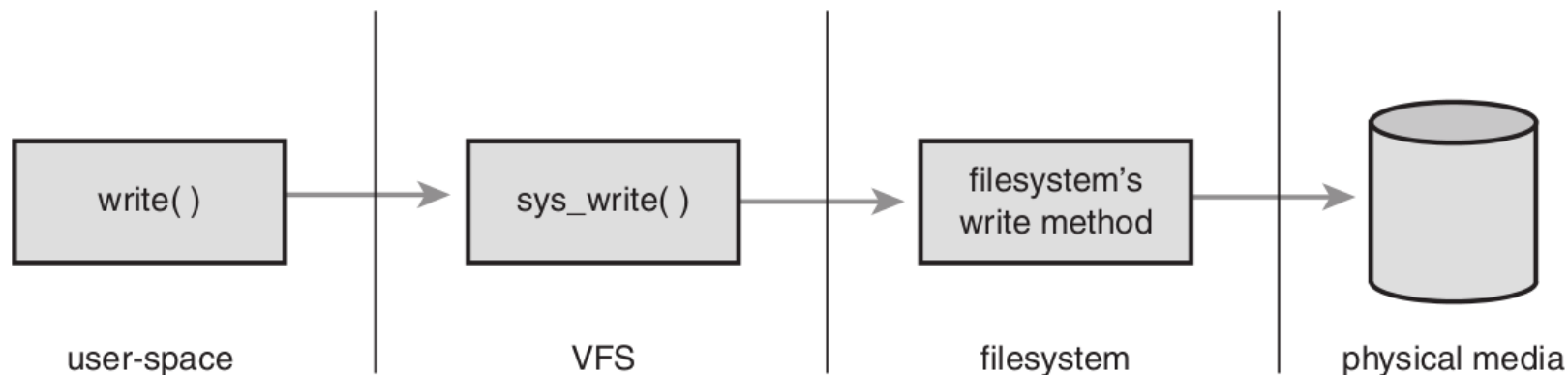
# File system APIs

- **File operations**
  - open, creat, read, write, lseek, close, dup
  - pread, pwrite, readv, writev, preadv, pwritev
- **Directory operations**
  - chdir, mkdir, rmdir
  - link, unlink, rename, chmod, chown, stat
  - opendir, readdir, readdir\_r, closedir, scandir
- **Special commands**
  - fcntl, ioctl, flock



# VFS implementation

- System call to file system's specific methods



- **Generic objects**

- *Superblock*: specific file system
- *Inode*: a file meta-data
- *Dentry*: a directory entry
- *File*: an open file



# VFS operations

- File system specific operations
  - Pseudo object oriented programming model
    - File system specific object: ex. `sb->s_fs_info`
    - File system specific operations: ex. `sb->s_op.sync_fs()`
  - Generic object + FS specific object + operations = VFS
- VFS internal objects
  - To handle file system status
    - struct `file_system_type` file system mounting
    - struct `vfsmount` file system mount point
    - struct `file_struct` file descriptor management
      - struct `file` \*`fd_array`[`NR_OPEN_DEFAULT`]
    - struct `fs_struct` process status (working dir, ...)
  - Dentry cache



# VFS operations example

- `sys_open` (`fs/open.c`)
  - Main routine: `do_sys_open()`
  - Allocate fd: `get_unused_fd_flags()`
  - Walk path and open a file: `do_filp_open()`
    - `lookup_fast()`, `__d_lookup()` : *dentry cache lookup*
    - `i_op->lookup(dir, dentry)` : *repeat to target inode*
    - `d_op->d_hash(dentry, name)`,  
`d_op->d_compare(dentry, name1, name2)`
    - `f_op->open(inode, file)`
    - `i_op->create(dir, dentry, mode)`
    - `s_op->alloc_inode(sb)`



# Superblock API

- Superblock **object**
  - Per mounted file system instance
- Superblock **operations**
  - **Superblock** management
    - write\_super(), put\_super()
  - **Inode** management
    - alloc\_inode(), write\_inode()
  - **File system** management
    - sync\_fs(), free\_fs()
- **Initialization:** get\_sb() function

Type	Name
list_head	s_list
dev_t	s_dev
list_head	s_inodes
list_head	s_files
super_operations	<b>s_op</b>
void *	<b>s_fs_info</b>
...	...





# Inode API

- Inode **object**
- Inode **operations**
  - **Inode** management
    - create, truncate, setattr, fallocate,
  - **Directory** management
    - lookup, link, unlink, symlink, mkdir, ...

Type	Name
super_block	i_sb
list_head	<b>i_dentry</b>
unsigned long	<b>i_ino</b>
atomic_t	i_count
uid_t	i_uid
struct timespec	i_atime
loff_t	<b>i_size</b>
<b>address_space</b>	<b>i_mapping</b>
inode_operations	i_op
file_operations	i_fop
void *	i_private



# File API

- File object
- File operations
  - llseek(), read(), write(), mmap()
  - open(), release(), flush()

Type	Name
struct path	f_path
int	f_flags
loff_t	<b>f_pos</b>
address_space	f_mapping
file_operations	<b>f_op</b>
void *	private_data



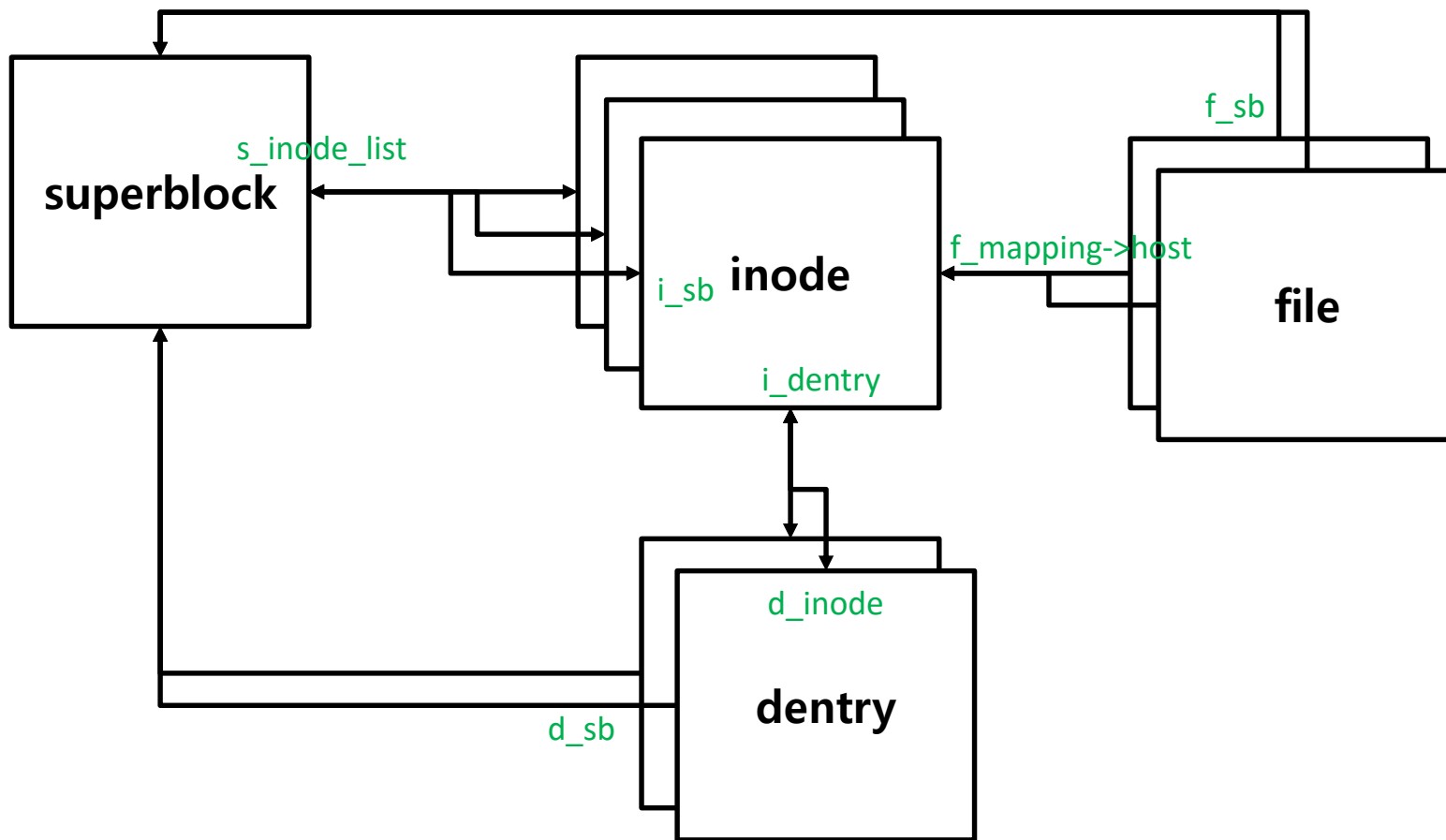
# Directory API

- Dentry object
- Operations
  - d\_revalidate
  - d\_hash
  - d\_compare
  - d\_delete, d\_release
  - d\_iput, d\_dname
- Almost no need to implement
  - Exception: *case insensitive* file name

Type	Name
struct inode	d_inode
hlist_node	d_hash
struct dentry	d_parent
struct qstr	<b>d_name</b>
list_head	d_subdirs
list_head	d_alias
dentry_operations	d_op
void*	d_fsdata
char	<b>d_iname[]</b>




# VFS objects relationships





# Page cache implementation

- Integrated with process virtual memory module
- **Page** associated with inode
  - `page->mapping, index`
- **address\_space** object
  - Target file
  - Page cache tree (radix)
- **address\_space operations**
  - **I/O**: `writepage, readpage, writepages, readpages`
  - **Block allocation and mapping**: `bmap()`  Delayed allocation

Type	Name
struct inode	host
radix_tree_root	page_tree
long	nrpages
address_space_operations	a_ops
...	...



# Page cache implementation

- **Useful API**

- `page = find_get_page(mapping, index)`
- `SetPageDirty(page)`, `ClearPageDirty(page)`

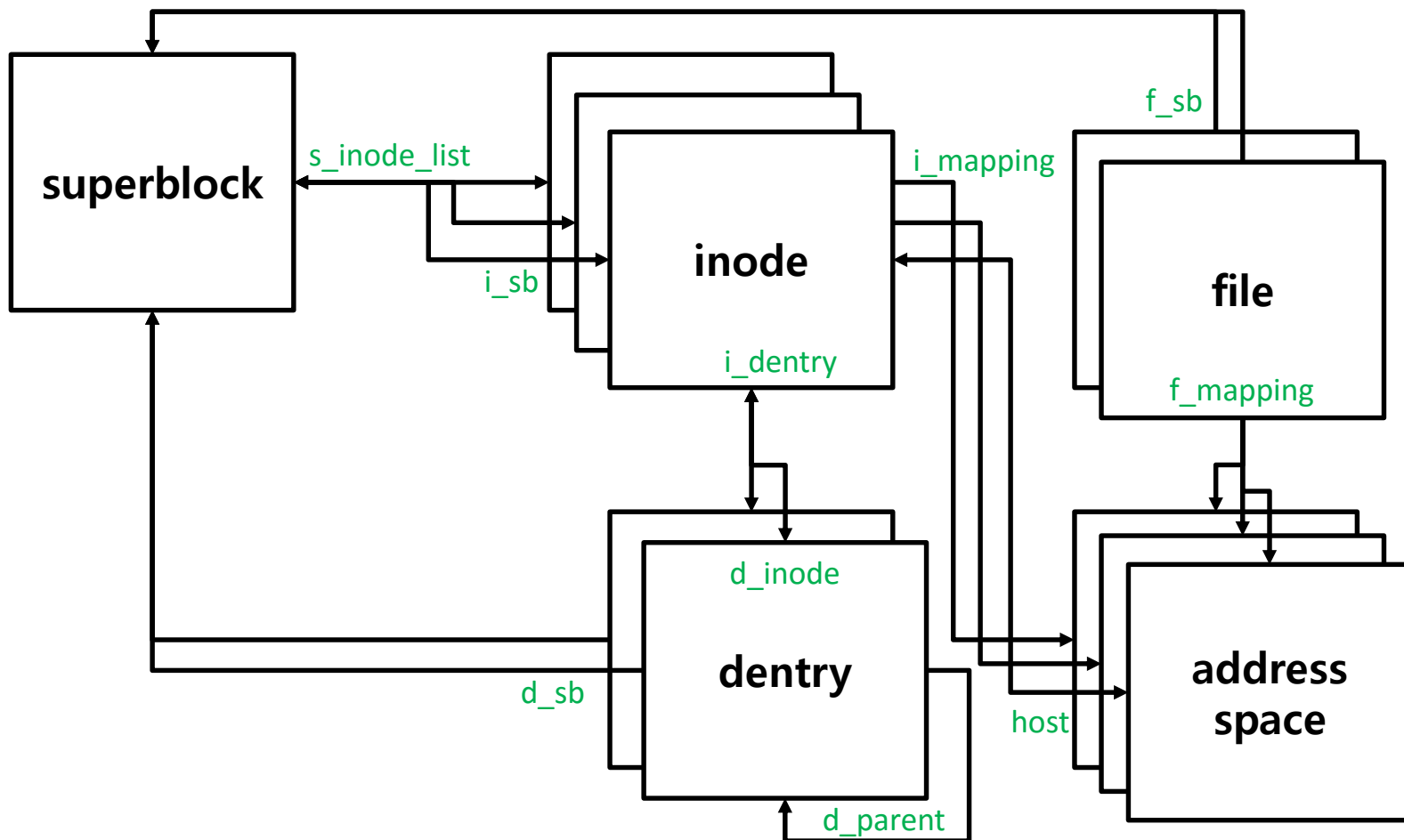
- **Flush thread**

- Write-back dirty pages
  - On-demand, dirty threshold, free threshold
- Per allocation-group
  - Per disk if device mapper is not used
  - Per top-level virtual device if device mapper is applied



# Implementing file system

- superblock, inode, file, dentry, address space



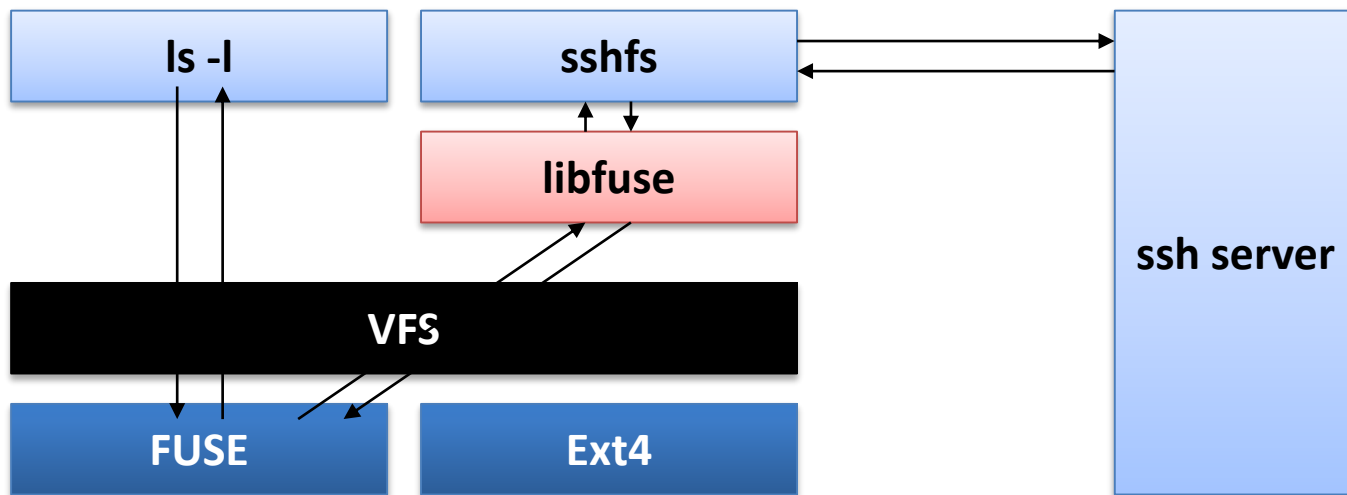
**FUSE: File system in user space**





# Introduction

- Implement simple file system in userspace



- **Compatibility**

- Linux, FreeBSD, NetBSD, Mac OS X, OpenSolaris, ...

- **Example file systems**

- Gnome VFS2, unionfs, ftpfs, sshfs, Android sdcard, ...



- **Installation**

- apt-get install fuse fuse-utils libfuse2 libfuse-dev
- Source: <http://fuse.sourceforge.net/>

- **Mount / unmount**

- *sshfs* user@host /tmp/fuse
- *fusermount* -u /tmp/fuse



# Using SSHFS

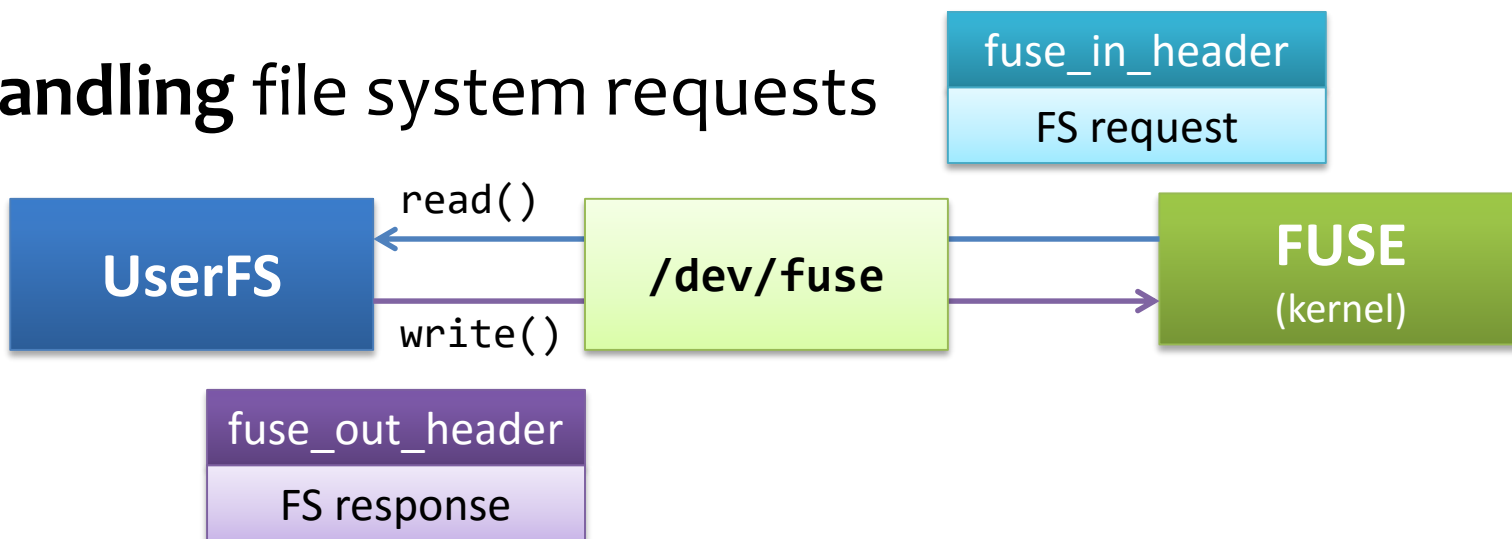
- **Setup ssh server**
  - apt-get install ssh
- **Install sshfs**
  - apt-get install sshfs
- **Mount sshfs**
  - mkdir /home/user/fusemnt
  - sshfs user@localhost:/home/user /home/user/fusemnt
- **Unmount sshfs**
  - fusermount -u /home/user/fusemnt



# Low-level FUSE interface

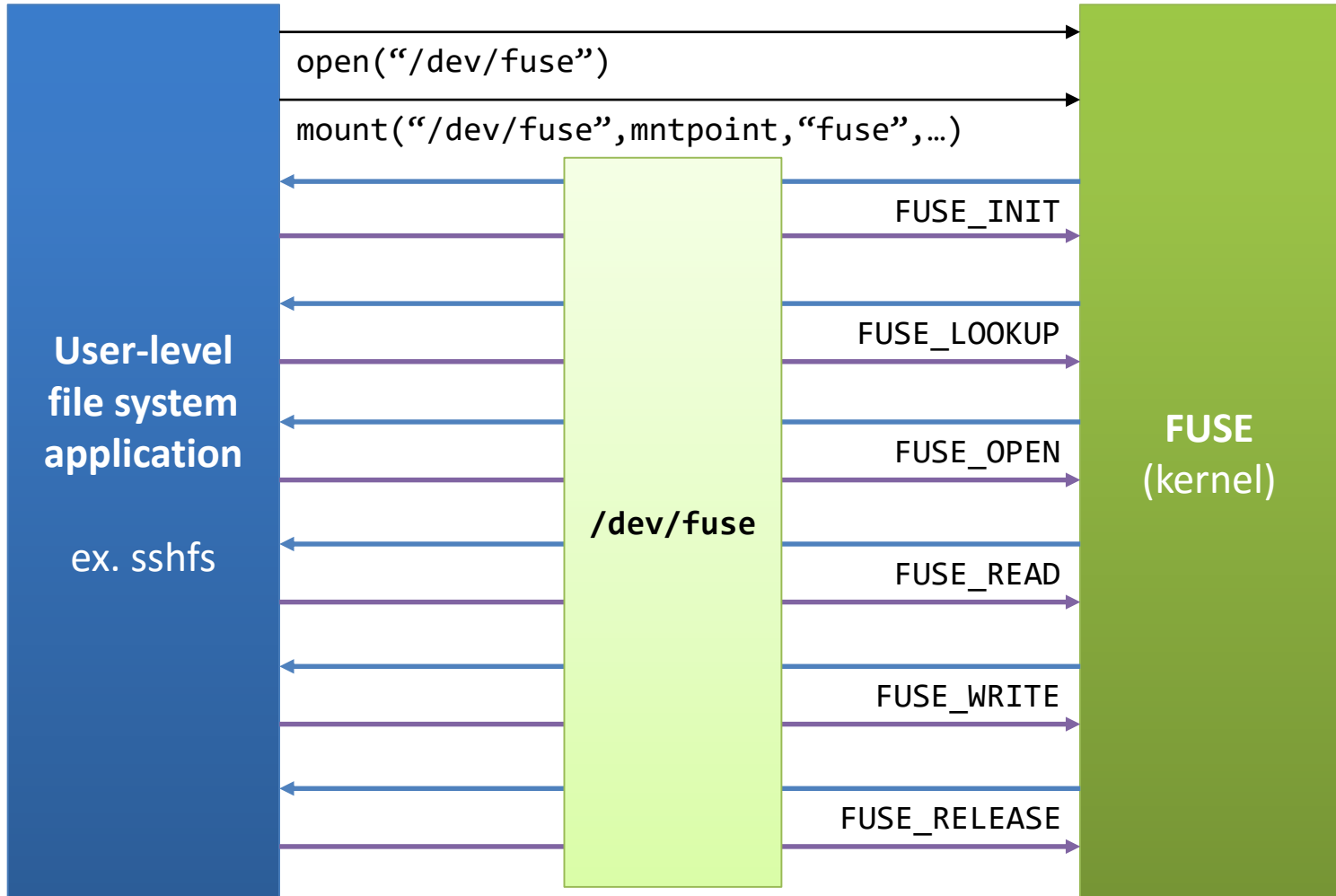
- **FUSE communication channel**
  - `/dev/fuse`      `fd = open("/dev/fuse", O_RDWR)`
- **Mounting a file system**
  - `mount("/dev/fuse", mount_point, "fuse", flags, opts)`
  - `opts: "fd=n,..."`

- **Handling file system requests**





# FUSE communication example





# FUSE message headers

- struct ***fuse\_in\_header***
  - Describes a file system request
  - enum ***fuse\_opcode*** opcode:
    - init
    - open, read, write
    - mknod, mkdir, readlink, symlink, unlink
    - rmdir, rename
    - lookup, forget, getattr, setattr
    - statfs, release, flush
    - opendir, readdir, releasedir
- struct ***fuse\_out\_header***
  - Return execution result

## fuse\_in\_header

u32 len

**u32 opcode**

u64 unique

u64 nodeid

u32 uid

u32 gid

u32 pid

u32 padding

## fuse\_out\_header

u32 len

s32 error

u64 unique



# Directory access

- **FUSE\_LOOKUP**

- Find file “*data*” in directory *fuse->nodeid*
- Returns struct **fuse\_entry\_out**

## fuse\_entry\_out

u64	nodeid
u64	generation
u64	entry_valid
u64	attr_valid
	fuse_attr

- struct **fuse\_open\_in**

- u32 flags
- Open file identified by *fuse->nodeid*
- User defines file handle *fh*

## fuse\_open\_out

u64	fh
u32	open_flags

- struct **fuse\_release\_in**

- Close file *fh*
- No return

## fuse\_release\_in

u64	fh
u32	flags
u32	release_flags
u64	lock_owner



# Read and write

- struct ***fuse\_read\_in***
  - Next to `fuse_in_header`
  - Read file (offset, size)
  - Send result to `/dev/fuse`
- struct ***fuse\_write\_in***
  - Identical to `fuse_read_in`
  - File data comes next to `fuse_write_in`
  - Return `fuse_write_out` to `/dev/fuse`

## ***fuse\_read\_in***

u64	fh
u64	offset
u32	size
u32	read_flags
u64	lock_owner
u32	flags
u32	padding

## ***fuse\_write\_out***

u32	size
u32	padding





# Other primitives

- Refer to *include/fuse\_kernel.h*
  - struct fuse\_(opcode)\_in/out
- **FUSE for Python**
  - <http://sourceforge.net/apps/mediawiki/fuse/index.php?title=SimpleFilesystemHowto>
- **Other language bindings**
  - C++, Java, C#, Go, Haskell, TCL, Perl, Ocaml, Ruby, Lua, ...
  - <http://sourceforge.net/apps/mediawiki/fuse/index.php?title=LanguageBindings>



# libfuse: C binding

- **High-level API**
  - Implements common fuse  $\Leftrightarrow$  kernel interactions
    - Mount command-line parsing
    - Interpret fuse\_in\_header and build fuse\_out\_header
  - Simplifies file system implementation
- **Developing with libfuse**
  - **Include:** <fuse.h>
  - **Compile:** `gcc -Wall `pkg-config fuse --cflags --libs``
  - **Example:** `/usr/share/doc/libfuse-dev/examples/hello.c`
  - **Document:** <http://fuse.sourceforge.net/doxygen/index.html>



# Implementing FS with libfuse

- **struct fuse\_operations**
  - Define and override file system operations
  - Example
    - `.open(char *path, struct fuse_file_info *fi)`
    - `.read(char *path, char *buf, size_t size, off_t offset, struct fuse_file_info *fi)`
- **fuse\_main(argc, argv, &op, user\_data)**
  - Mounts the file system
  - Fetches file system requests and calls defined functions
  - Returns when unmounted: `fusermount -u`

**Project 3**  
**FUSE file system**



- **Project goal**
  - Implement a simple FUSE file system using libfuse
- **Due: 17<sup>th</sup> June, 23:59**
- **Baseline**
  - Implement root directory (no subdirectory)
  - **File read/write:** run IOzone to prove it working
  - **File allocation:** sequential (no fragmented) block allocation
- **You need to implement**
  - getattr, readdir, open, read, write



# File system layout

- **Use a 1 GB file as a storage media**

```
$ truncate -s 1G disk.img
```

- **Layout**

- 4 KB blocks

- Block bitmap size: 32 KB (8 blocks)

- Inode bitmap size: 1 block ← you define # of inodes

- Array-based directory (max length: 100 bytes)

