

OSTEP

Persistence:

File System API

Questions answered in this lecture:

How to **name** files?

What are **inode numbers**?

How to **lookup** a file based on pathname?

What is a **file descriptor**?

What is the difference between **hard and soft links**?

How can **special requirements** be communicated to file system (fsync)?

Persistent Storage

- Keep a data **intact** even if there is a power loss.
 - Hard disk drive
 - Solid-state storage device
- Two key abstractions in the virtualization of storage
 - File
 - Directory

What is a File?

Array of persistent bytes that can be read/written

Each file has low-level name as **inode number**

The user is not aware of this name.

File system consists of **many** files

Refers to collection of files

Also refers to part of OS that manages those files

Has a responsibility to store data persistently on disk.

Files need **names** to access correct one

File Names

Three types of names

- Unique id: inode numbers
- Path
- File descriptor

Inode Number

Each file has **exactly one inode number**

Inodes are unique (at a given time) within file system

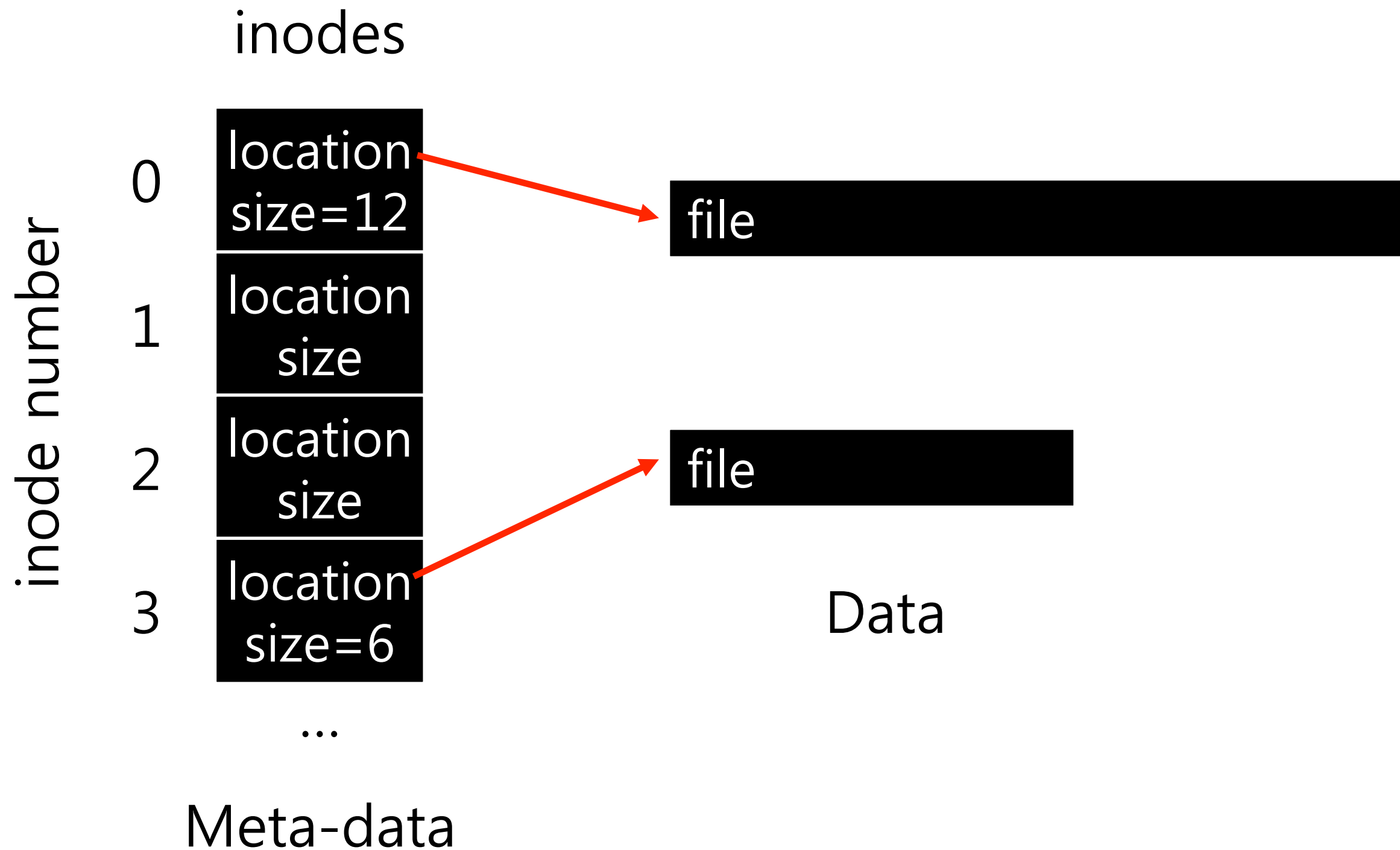
Different file system may use the same number, numbers may be recycled after deletes

See inodes via "ls -li"; see them increment...

What does "i" stand for?

"In truth, I don't know either. It was just a term that we started to use. 'Index' is my best guess, because of the slightly unusual file system structure that stored the access information of files as a flat array on the disk..."

~ Dennis Ritchie



File API (attempt 1)

`read(int inode, void *buf, size_t nbyte)`

`write(int inode, void *buf, size_t nbyte)`

`seek(int inode, off_t offset)`

seek does not cause disk seek until read/write

Disadvantages?

- names hard to remember
- no organization or meaning to inode numbers
- everybody has the same offset
- semantics of offset across multiple processes?

File API (attempt 1)

`pread(int inode, void *buf, off_t offset, size_t nbyte)`
`pwrite(int inode, void *buf, off_t offset, size_t nbyte)`
~~`seek(int inode, off_t offset)`~~

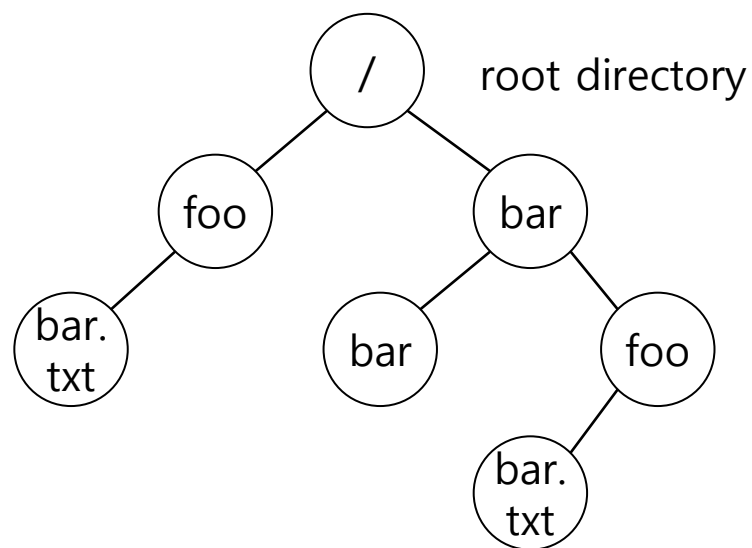
Disadvantages?

- names hard to remember
- no organization or meaning to inode numbers
- ~~— everybody has the same offset~~
- semantics of offset across multiple processes?

Directory

- Directory is like a file, also has a low-level name.
 - It contains a list of (user-readable name, low-level name) pairs.
 - Each entry in a directory refers to either *files* or other *directories*.
- Example)
 - A directory has an entry ("foo", "10")
 - A file "foo" with the low-level name "10"

Directory Tree (Directory Hierarchy)



An Example Directory Tree

Valid files (absolute pathname) :

/foo/bar.txt
/bar/foo/bar.txt

Valid directory :

/
/foo
/bar
/bar/bar
/bar/foo/

} Sub-directories

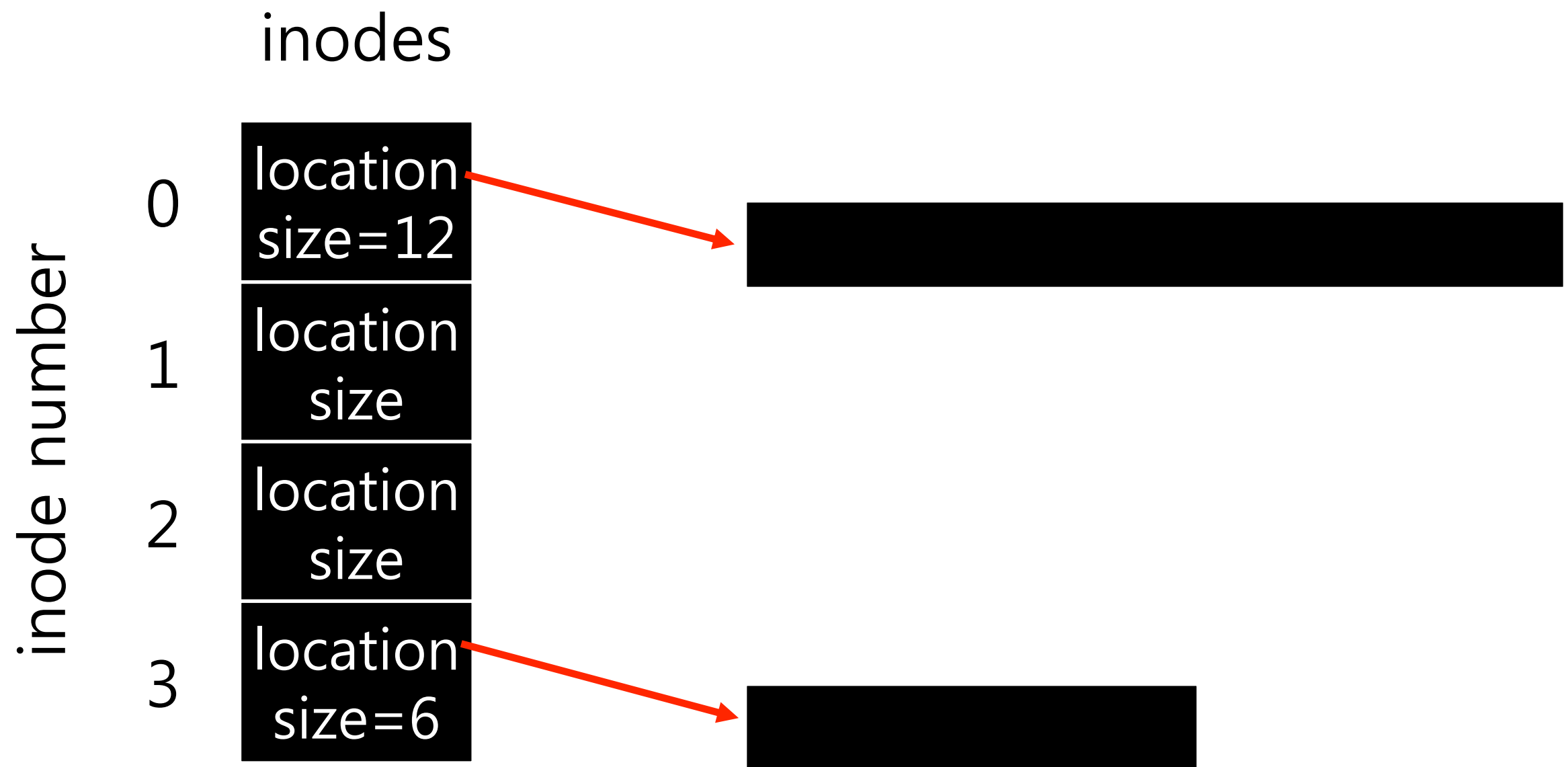
Paths

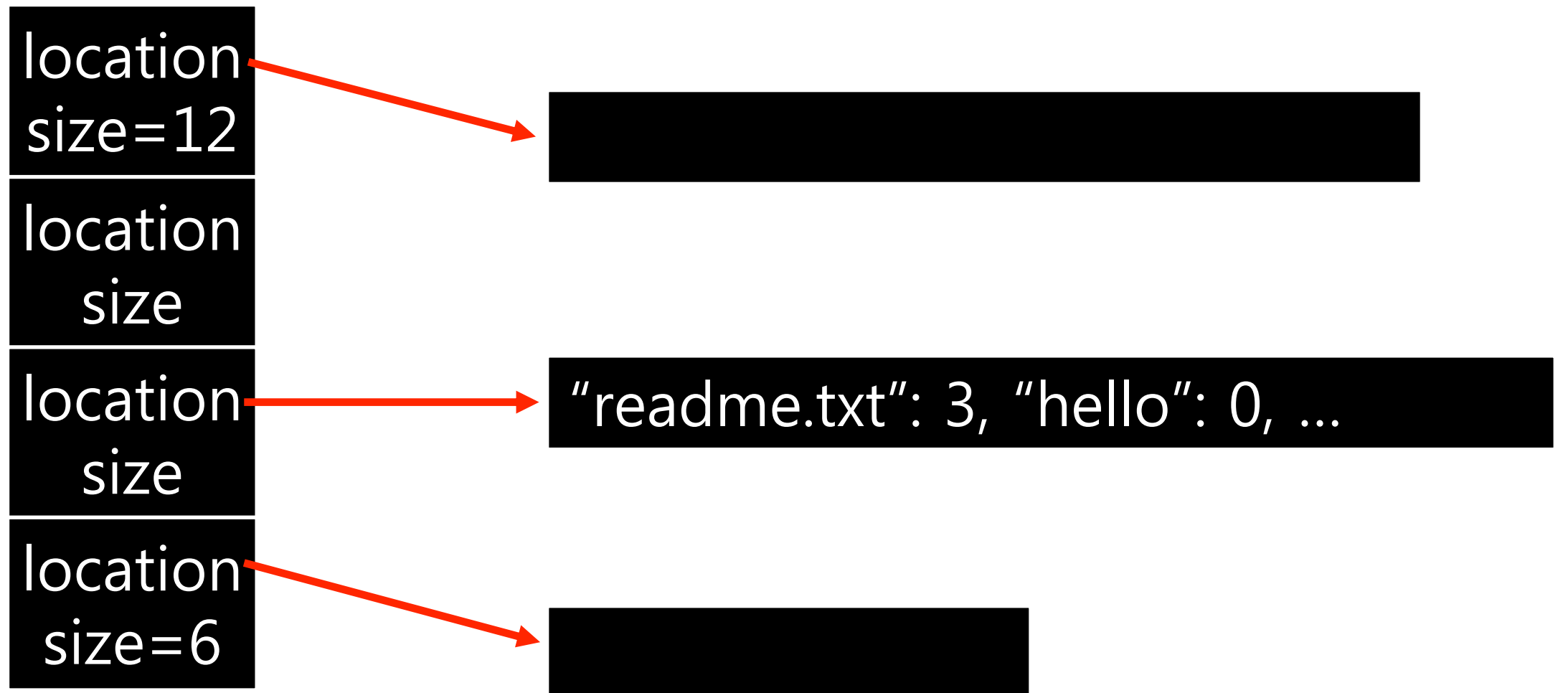
String names are friendlier than **number** names

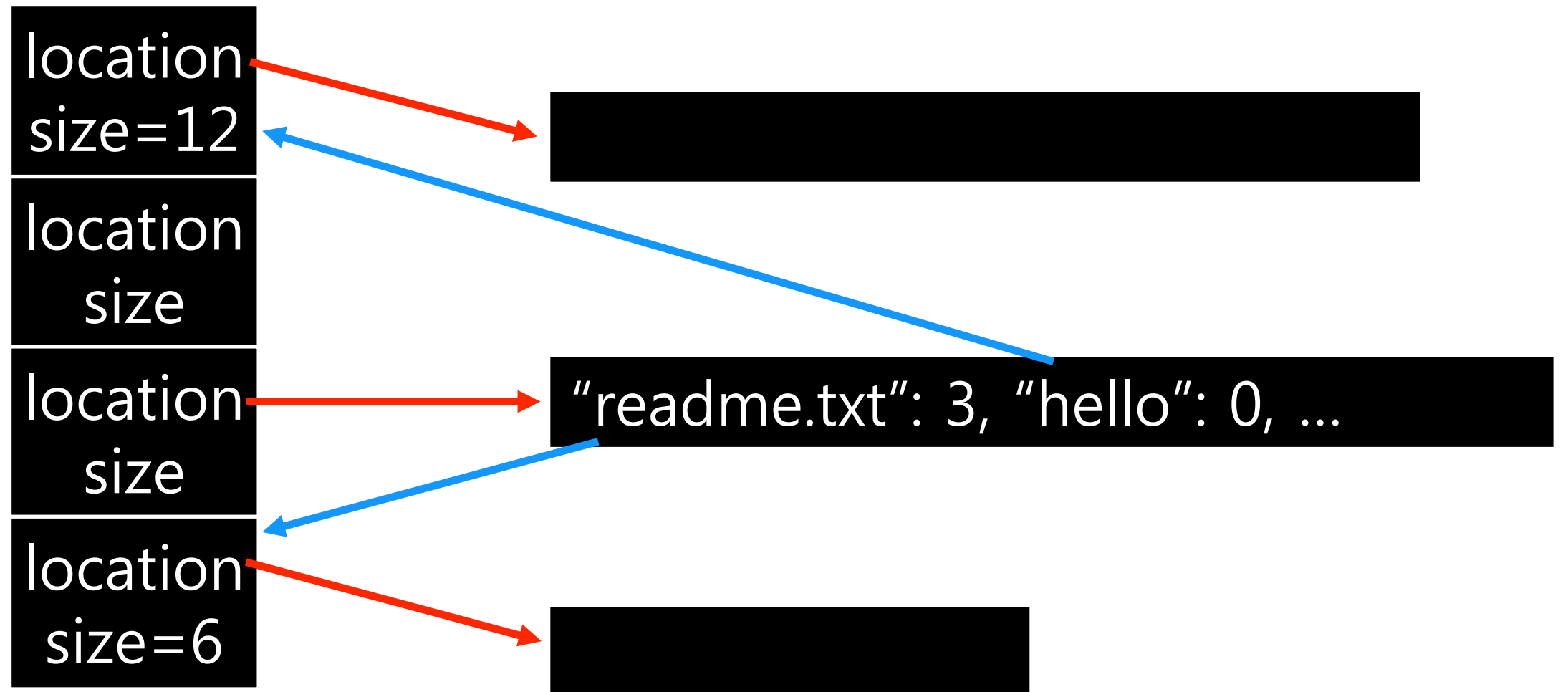
File system still interacts with inode numbers

Store *path-to-inode* mappings in predetermined "root" file (typically inode 2)

Directory!







Paths

Generalize!

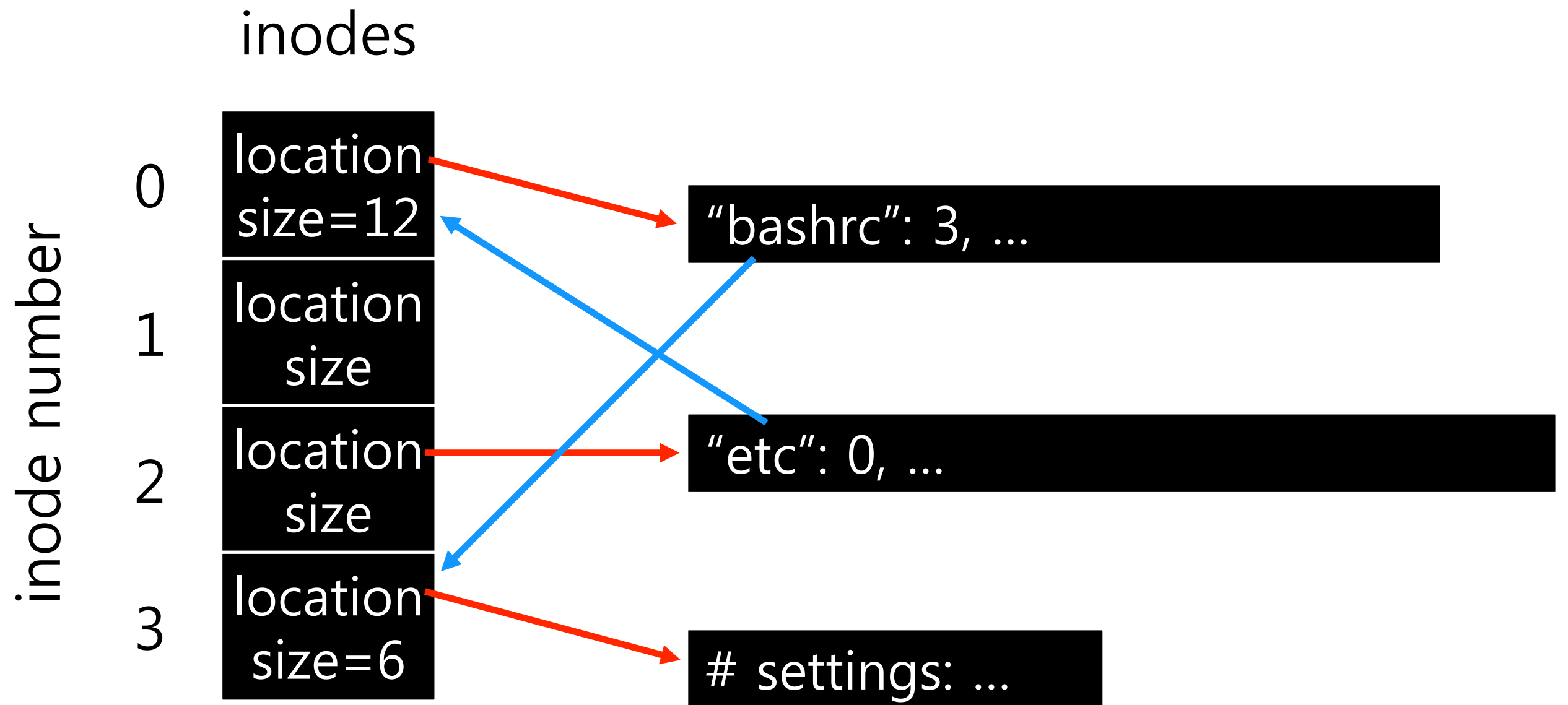
Directory Tree instead of single root directory

Only **file name** needs to be unique

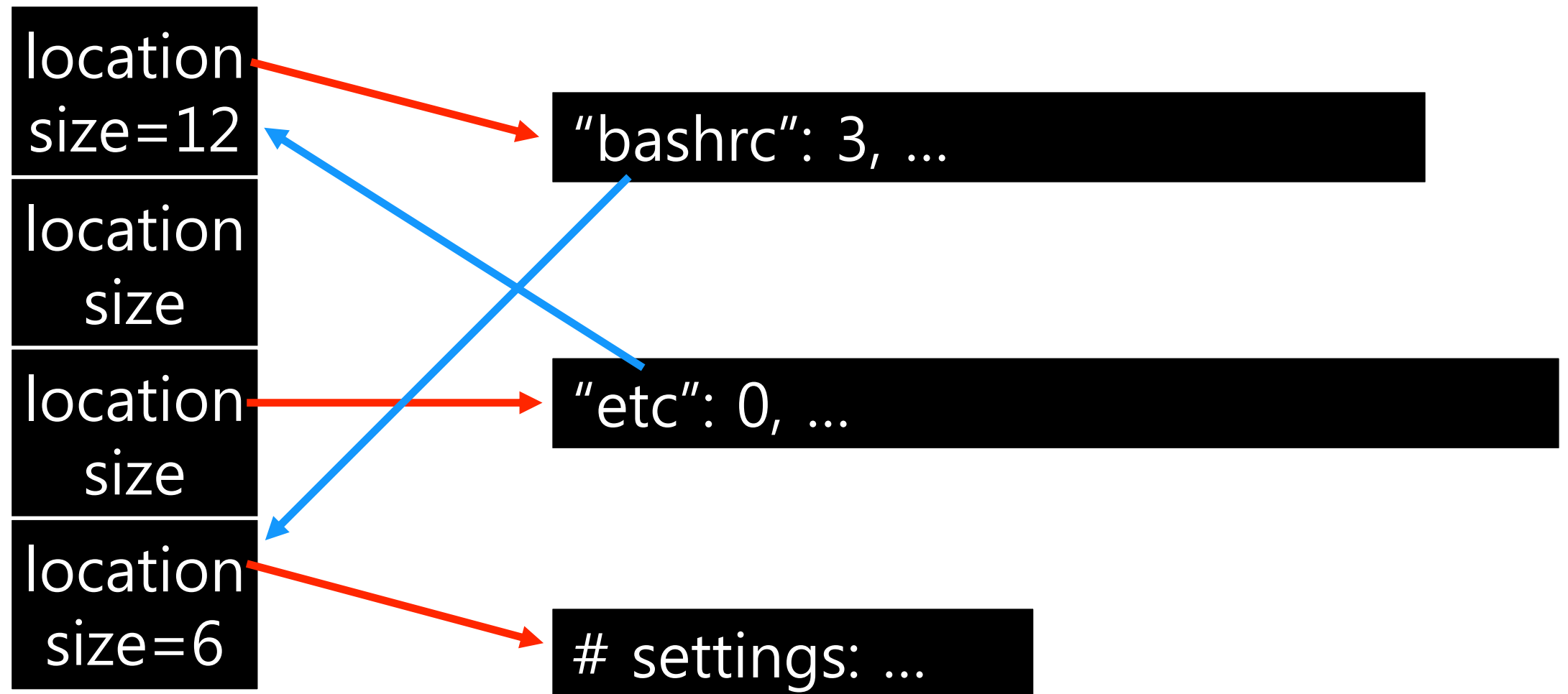
/usr/dusseau/file.txt

/tmp/file.txt

Store file-to-inode mapping for each directory

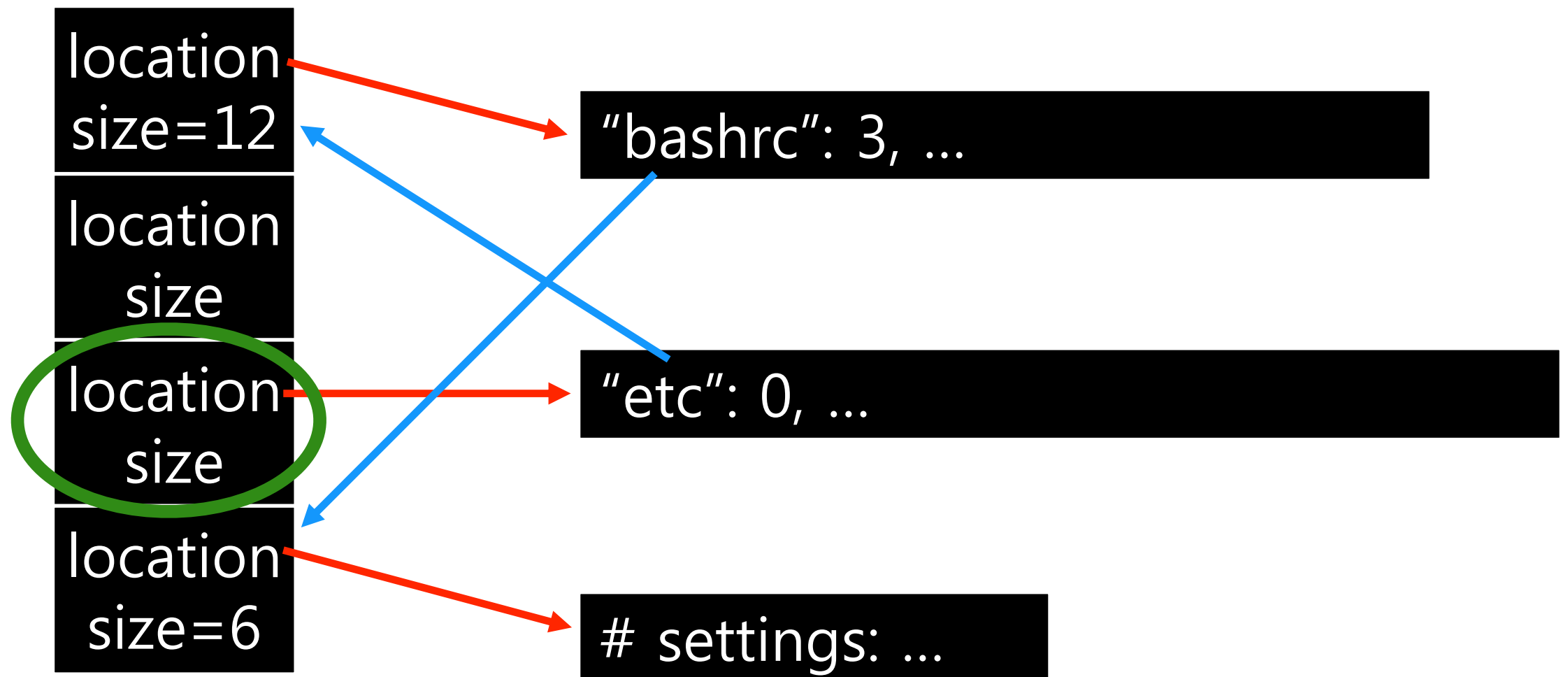


read /etc/bashrc



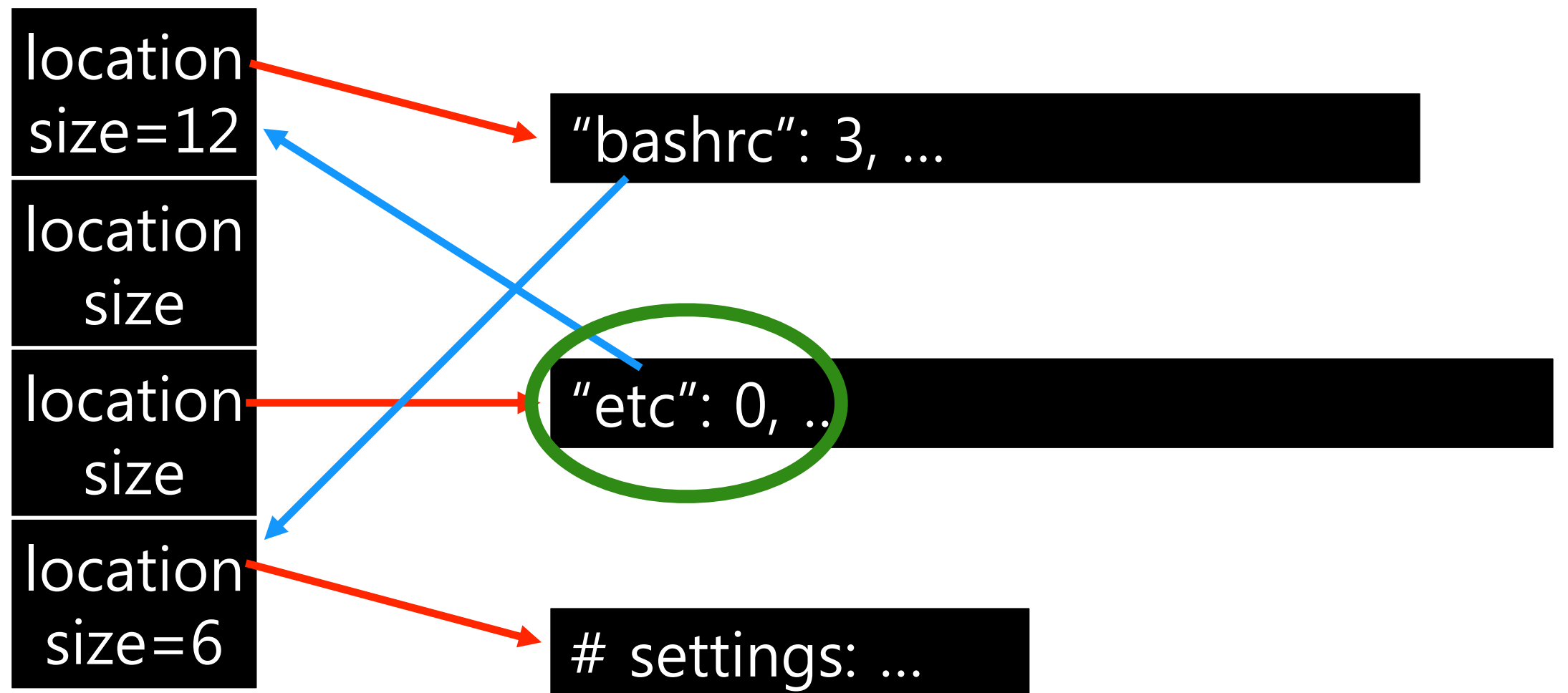
reads: 0

read /etc/bashrc



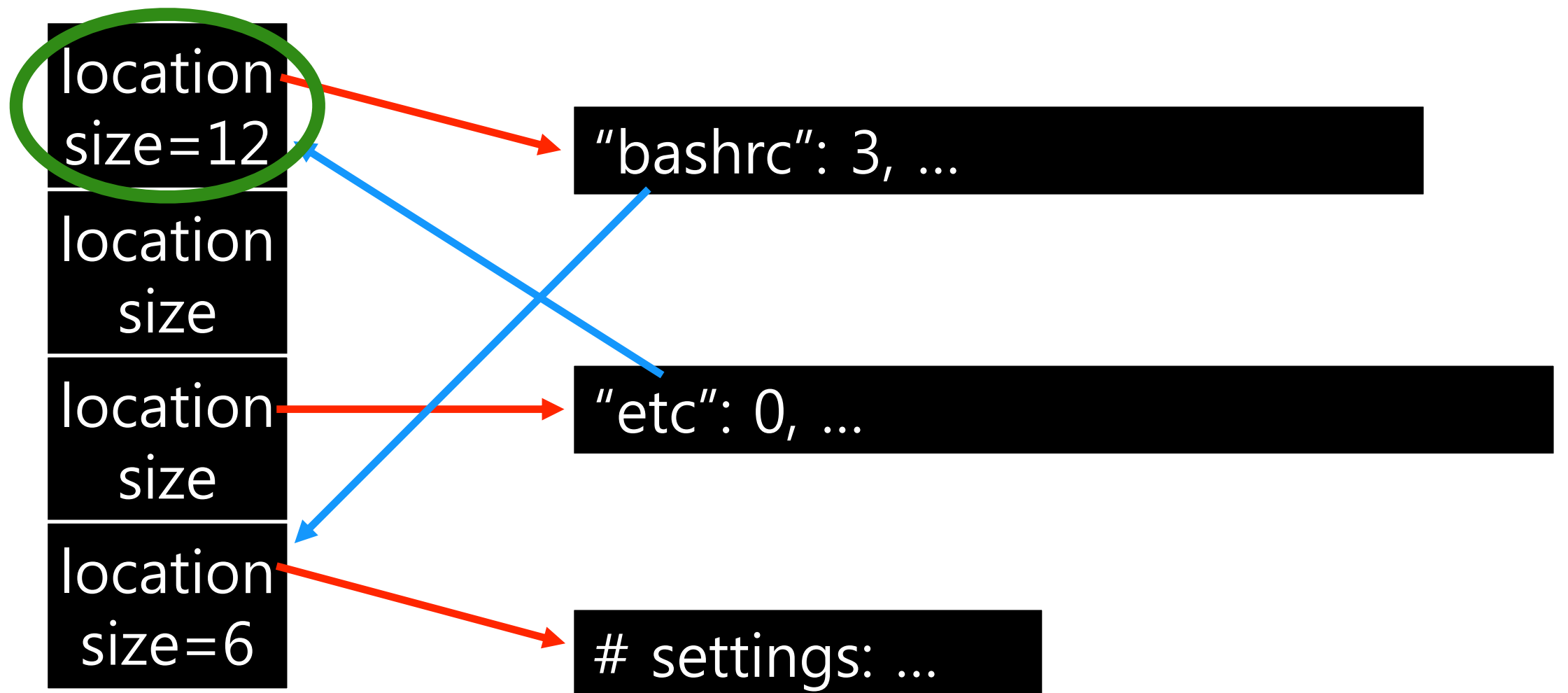
reads: 1

read /etc/bashrc



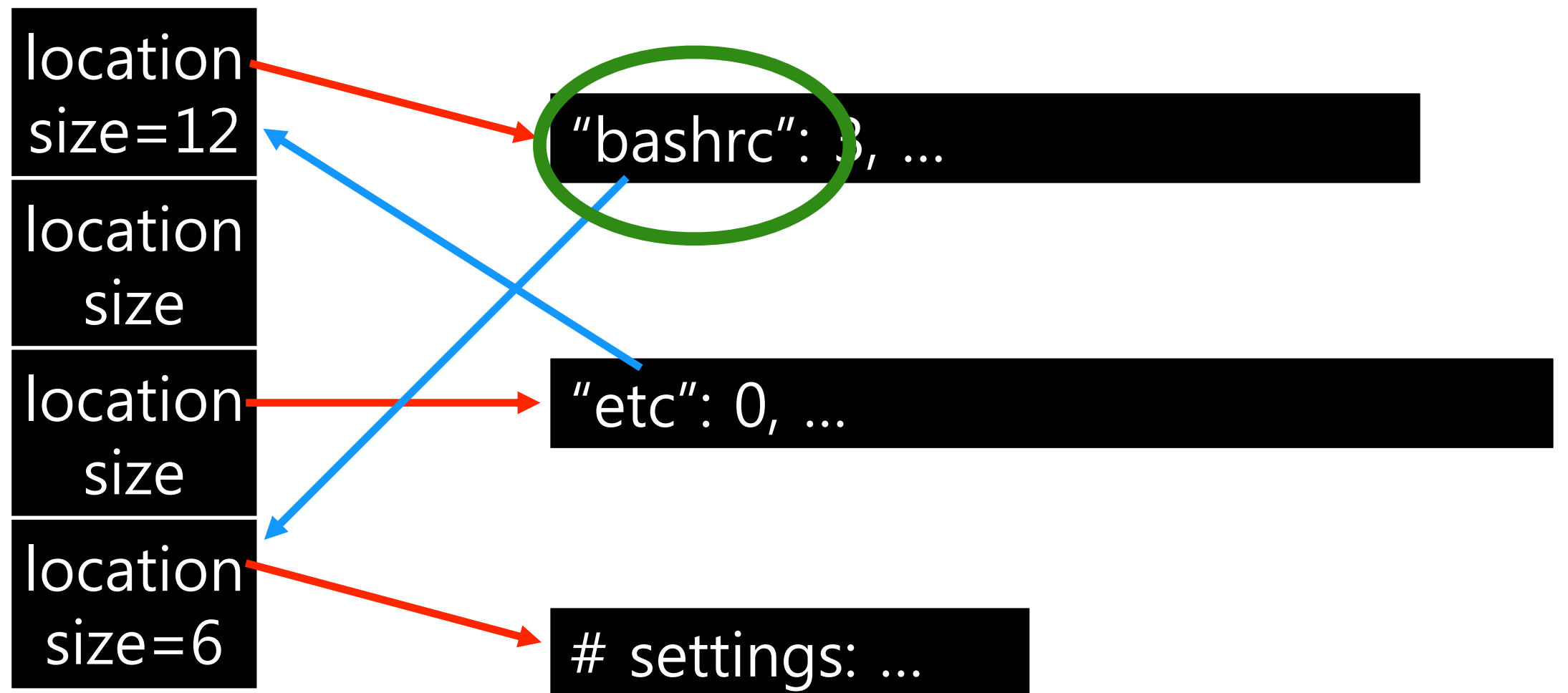
reads: 2

read /etc/bashrc



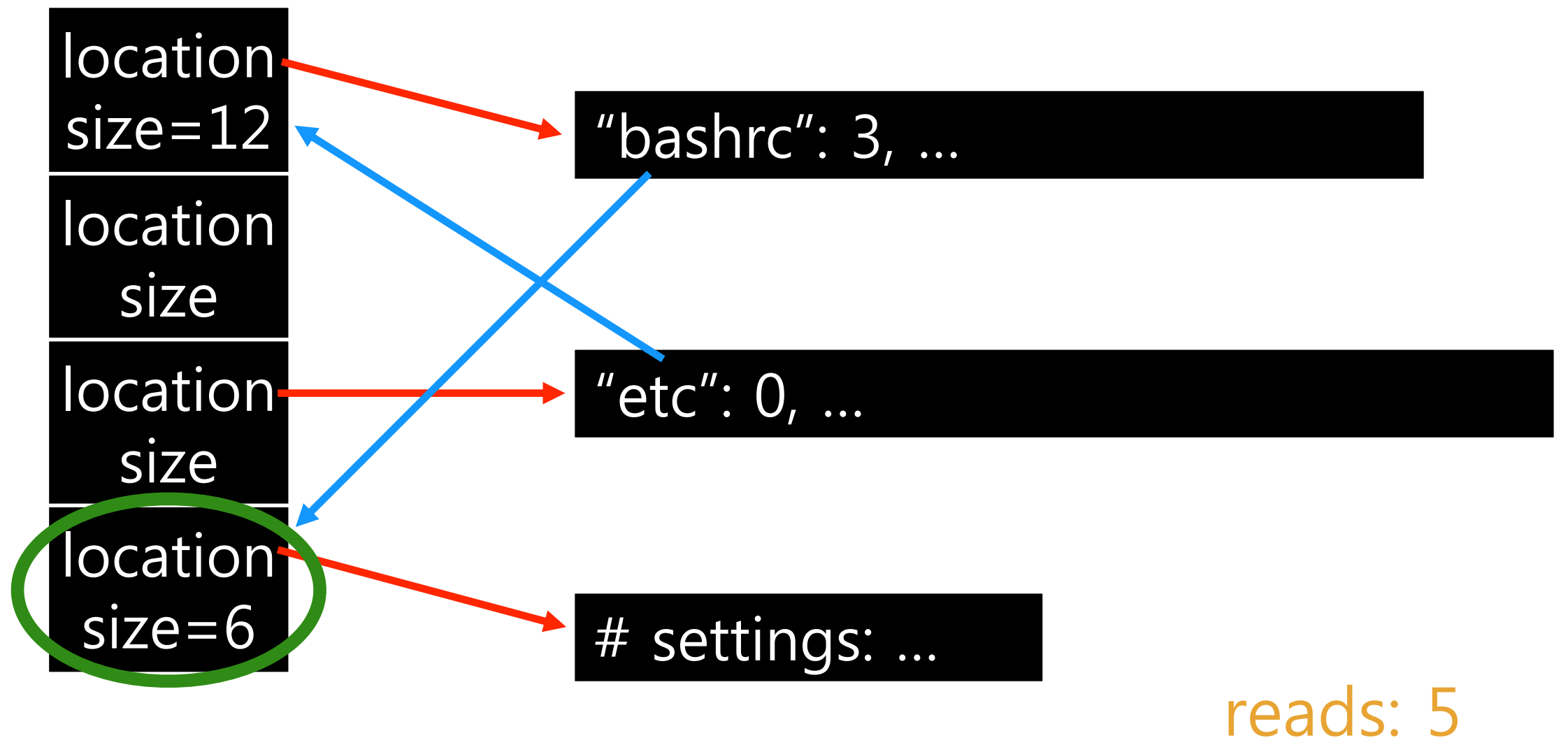
reads: 3

read /etc/bashrc

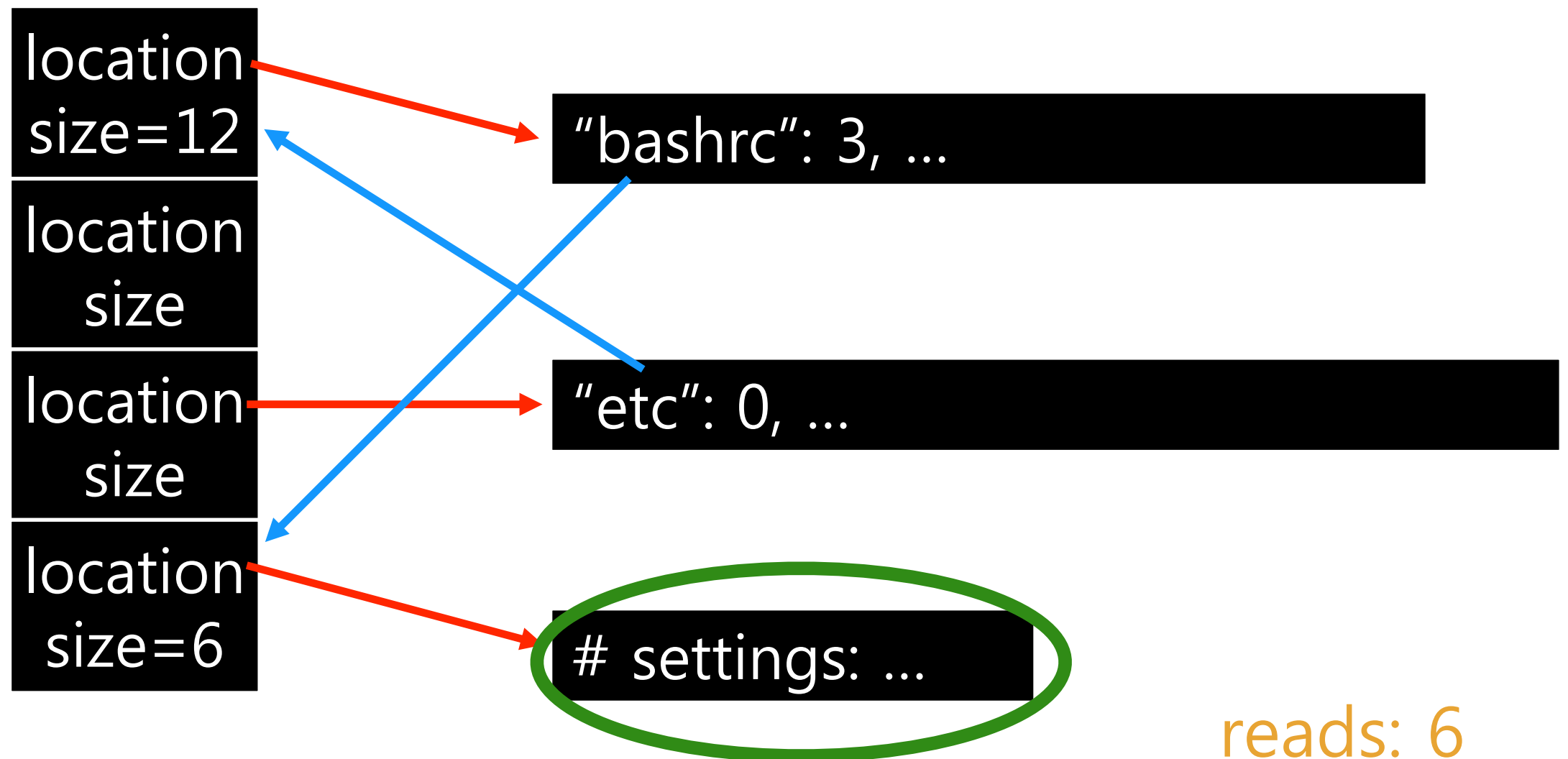


reads: 4

read /etc/bashrc



read /etc/bashrc



Reads for getting final inode called "traversal"

Read root dir (inode and data);
read etc dir (inode and data);
read bashrc file (indode and data)

Directory Calls

`mkdir`: create new directory

`readdir`: read/parse directory entries

Why no `writedir`?

Making Directories

- `mkdir()`: Make a directory

```
prompt> strace mkdir foo
...
mkdir("foo", 0777)          = 0
prompt>
```

- When a directory is created, it is **empty**.
- Empty directory have two entries: `.` (itself), `..` (parent)

```
prompt> ls -a
./      ../
prompt> ls -al
total 8
drwxr-x---  2 remzi remzi    6 Apr 30 16:17 ./
drwxr-x--- 26 remzi remzi 4096 Apr 30 16:17 ../
```

Reading Directories

- A sample code to read directory entries (like `ls`).

```
int main(int argc, char *argv[]) {
    DIR *dp = opendir(".");           // open current directory
    assert(dp != NULL);
    struct dirent *d;
    while ((d = readdir(dp)) != NULL) // read one directory entry
    {
        // print out the name and inode number of each file
        printf("%d %s\n", (int) d->d_ino, d->d_name);
    }
    closedir(dp);                     // close current directory
    return 0;
}
```

- The information available within `struct dirent`

```
struct dirent {
    char          d_name[256]; /* filename */
    ino_t         d_ino;       /* inode number */
    off_t         d_off;       /* offset to the next direct */
    unsigned short d_reclen;    /* length of this record */
    unsigned char  d_type;      /* type of file */
}
```

Deleting Directories

- `rmdir()`: Delete a directory.
 - Require that the directory be **empty**.
 - If you call `rmdir()` to a non-empty directory, it will fail.
 - I.e., Only has "." and ".." entries.

Special Directory Entries

```
$ ls -la
```

```
total 728
```

```
drwxr-xr-x 34 trh staff 1156 Oct 19 11:41 .
```

```
drwxr-xr-x+ 59 trh staff 2006 Oct 8 15:49 ..
```

```
-rw-r--r--@ 1 trh staff 6148 Oct 19 11:42 .DS_Store
```

```
-rw-r--r-- 1 trh staff 553 Oct 2 14:29 asdf.txt
```

```
-rw-r--r-- 1 trh staff 553 Oct 2 14:05 asdf.txt~
```

```
drwxr-xr-x 4 trh staff 136 Jun 18 15:37 backup
```

```
...
```

```
cd /; ls -lia
```

Creating Files

- Use `open()` system call with `O_CREAT` flag.

```
int fd = open("foo", O_CREAT | O_WRONLY | O_TRUNC);
```

- `O_CREAT` : create file.
 - `O_WRONLY` : only write to that file while opened.
 - `O_TRUNC` : make the file size zero (remove any existing content).
- `open()` system call returns **file descriptor**.
 - *File descriptor* is an integer, and is used to access files.

Reading and Writing Files

- An Example of reading and writing 'foo' file

```
prompt> echo hello > foo  
prompt> cat foo  
hello  
prompt>
```

- `echo` : redirect the output of echo to the file foo
- `cat` : dump the contents of a file to the screen

How does the `cat` program access the file foo ?

We can use `strace` to trace the system calls made by a program.

Reading and Writing Files (Cont.)

```
prompt> strace cat foo
...
open("foo", O_RDONLY|O_LARGEFILE) = 3
read(3, "hello\n", 4096)           = 6
write(1, "hello\n", 6)              = 6 // file descriptor 1: standard out
hello
read(3, "", 4096)                   = 0 // 0: no bytes left in the file
close(3)                            = 0
...
prompt>
```

- `open` (file descriptor, flags)
 - Return file descriptor (3 in example)
 - File descriptor 0, 1, 2, is for standard input/ output/ error.
- `read` (file descriptor, buffer pointer, the size of the buffer)
 - Return the number of bytes it read
- `write` (file descriptor, buffer pointer, the size of the buffer)
 - Return the number of bytes it write

Reading and Writing Files (Cont.)

- Writing a file (A similar set of read steps)
 - A file is opened for writing (`open()`).
 - The `write()` system call is called.
 - Repeatedly called for larger files
 - `close()`

Reading And Writing, But Not Sequentially

- An open file has a **current offset**.
 - Determine **where** the next read or write will begin reading from or writing to within the file.
- Update the current offset
 - **Implicitly**: A read or write of N bytes takes place, N is added to the current offset.
 - **Explicitly**: `lseek()`

Reading And Writing, But Not Sequentially (Cont.)

```
off_t lseek(int fildes, off_t offset, int whence);
```

- `fildes` : File descriptor
- `offset` : Position the file offset to a particular location within the file
- `whence` : Determine how the seek is performed

From the man page:

```
If whence is SEEK_SET, the offset is set to offset bytes.  
If whence is SEEK_CUR, the offset is set to its current  
location plus offset bytes.  
If whence is SEEK_END, the offset is set to the size of the  
file plus offset bytes.
```

File API (attempt 2)

```
pread(char *path, void *buf,  
      off_t offset, size_t nbyte)
```

```
pwrite(char *path, void *buf,  
       off_t offset, size_t nbyte)
```

Disadvantages

Expensive traversal!

Goal: traverse once

File Names

Three types of names:

- inode
- path
- file descriptor

File Descriptor (fd)

Idea:

- Do expensive traversal once (open file)
- store inode in **descriptor** object (kept in memory).
- Do reads/writes via descriptor, which tracks offset

Each process:

- File-descriptor **table** contains pointers to open file descriptors

Integers used for file I/O are indexes into this table
stdin: 0, stdout: 1, stderr: 2

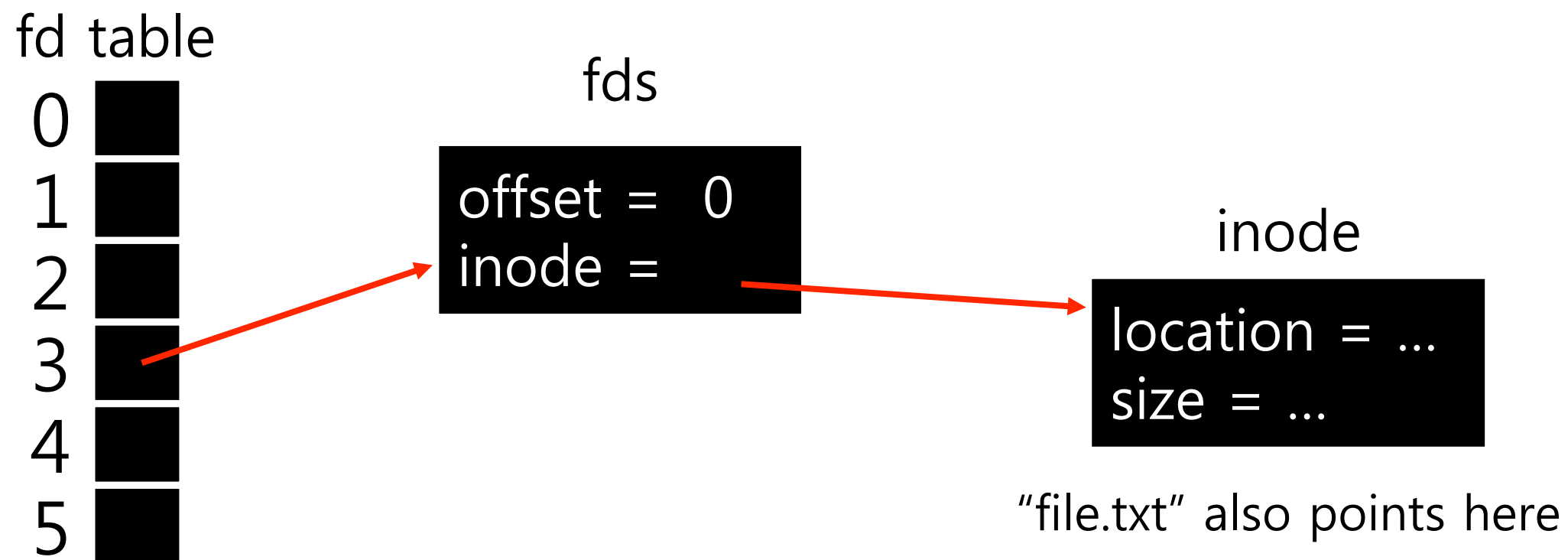
FD Table (xv6)

```
struct file {  
    ...  
    struct inode *ip;  
    uint off;  
};  
  
// Per-process state  
struct proc {  
    ...  
    struct file *ofile[NOFILE]; // Open files  
    ...  
}
```

Code Snippet

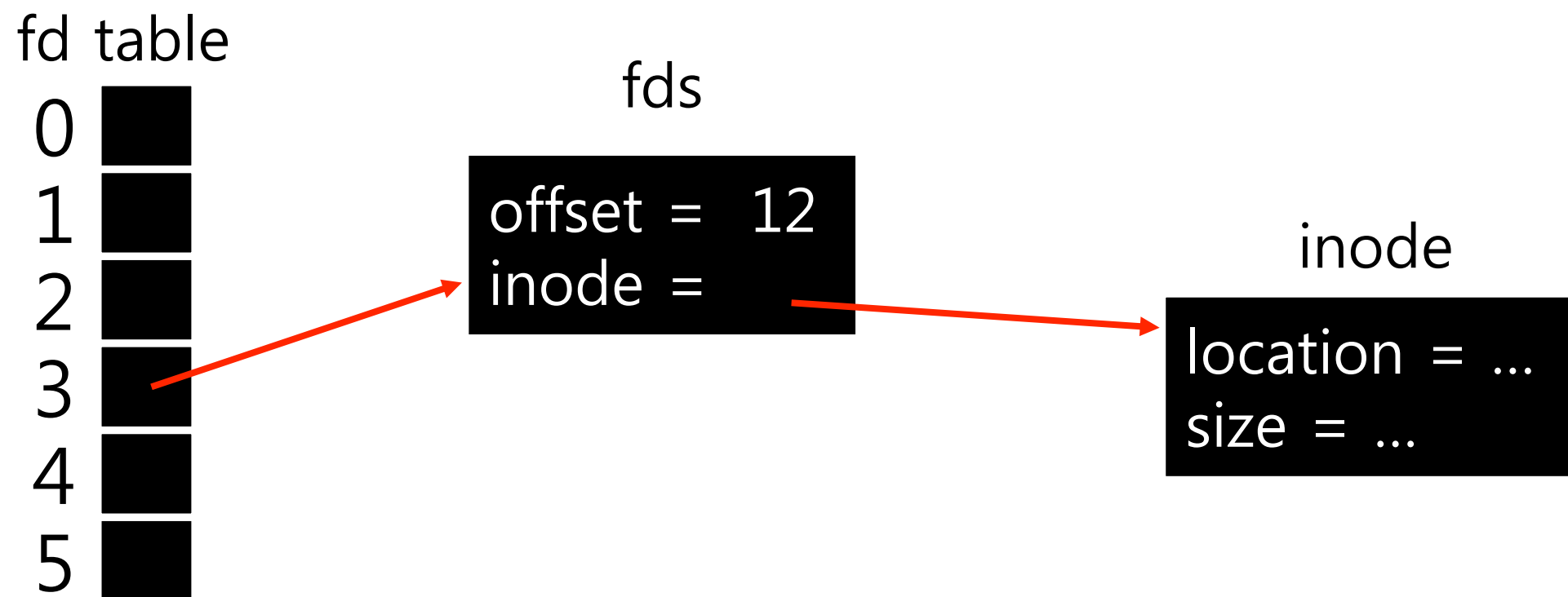
```
int fd1 = open("file.txt"); // returns 3
read(fd1, buf, 12);
int fd2 = open("file.txt"); // returns 4
int fd3 = dup(fd2);          // returns 5
```


Code Snippet



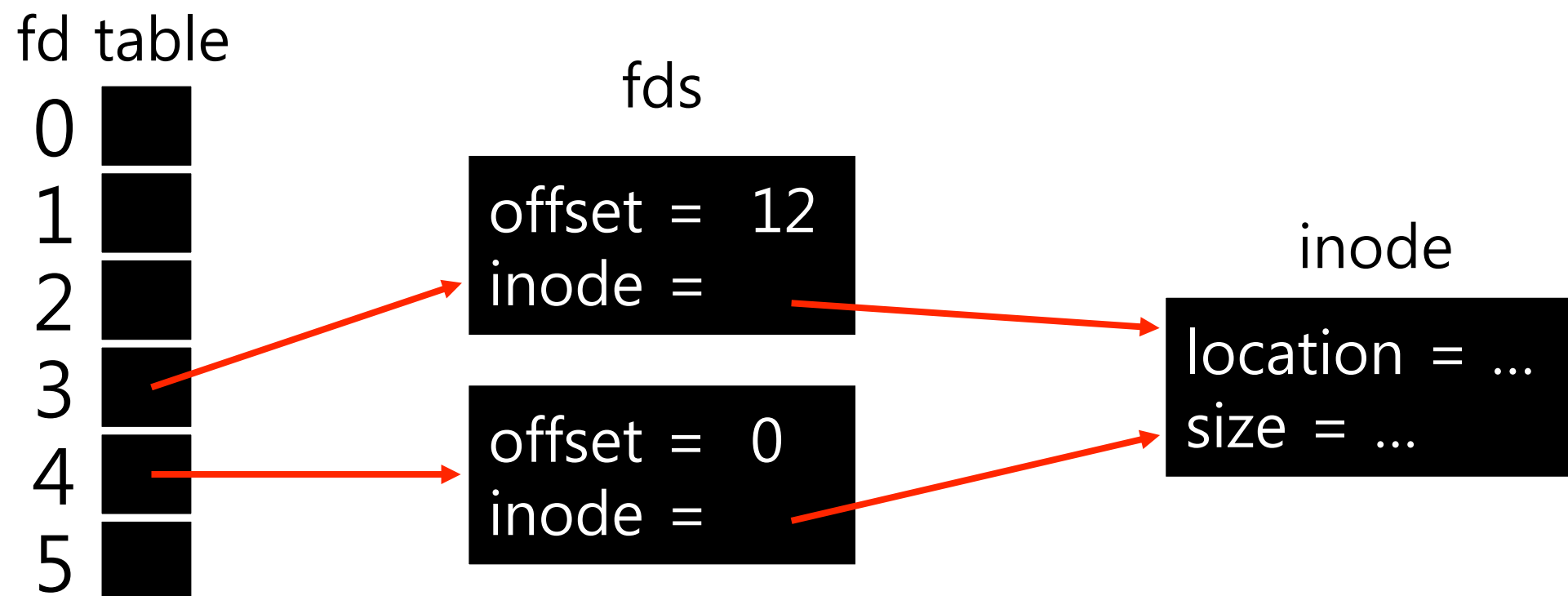
```
int fd1 = open("file.txt"); // returns 3
```

Code Snippet



```
int fd1 = open("file.txt"); // returns 3
read(fd1, buf, 12);
```

Code Snippet

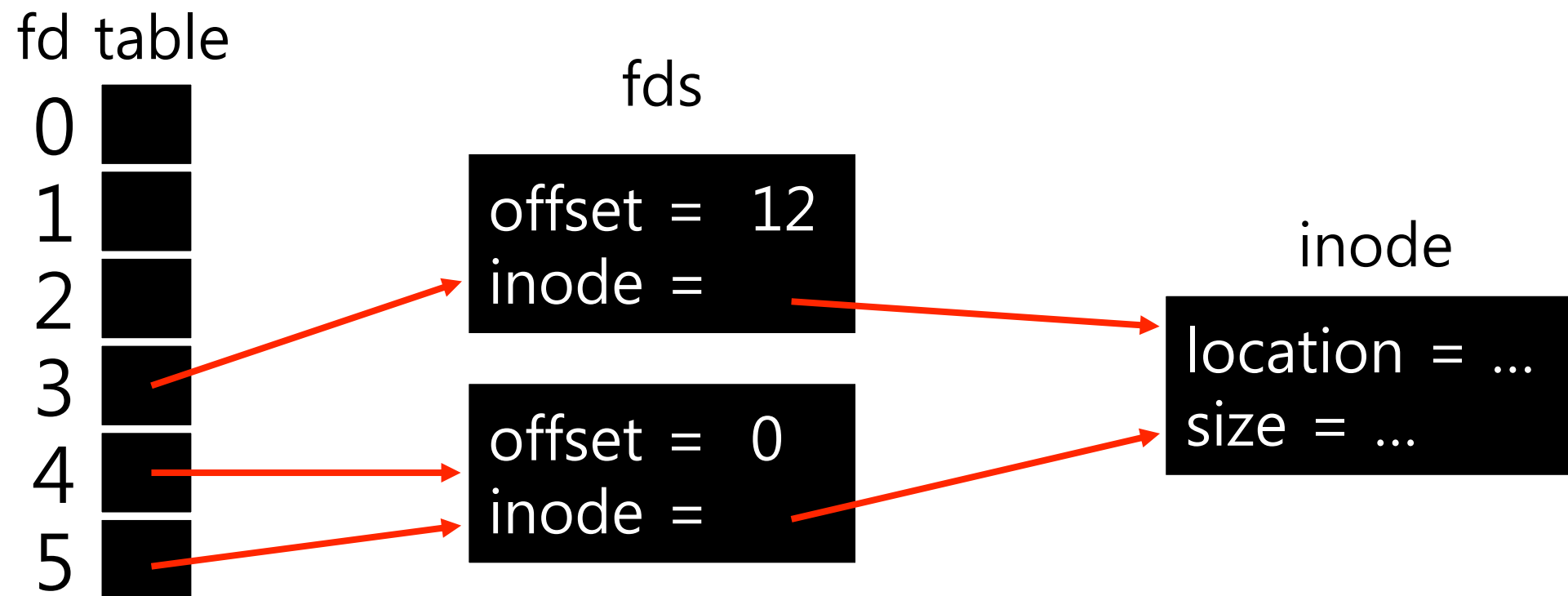


```
int fd1 = open("file.txt"); // returns 3
```

```
read(fd1, buf, 12);
```

```
int fd2 = open("file.txt"); // returns 4
```

Code Snippet



```
int fd1 = open("file.txt"); // returns 3
read(fd1, buf, 12);
int fd2 = open("file.txt"); // returns 4
int fd3 = dup(fd2);          // returns 5
```

File API (attempt 3)

```
int fd = open(char *path, int flag, mode_t mode)  
read(int fd, void *buf, size_t nbyte)  
write(int fd, void *buf, size_t nbyte)  
close(int fd)
```

Deleting Files

There is no system call for deleting files!

Inode (and associated file) is **garbage collected** when there are no references (from **paths** or **fds**)

Paths are deleted when: `unlink()` is called

FDs are deleted when: `close()` or process quits

Network File System Designers

A process can open a file, then remove the directory entry for the file so that it has no name anywhere in the file system, and still read and write the file. This is a disgusting bit of UNIX trivia and at first we were just not going to support it, but it turns out that all of the programs we didn't want to have to fix (csh, send mail, etc.) use this for temporary files.

~ Sandberg *etal.*

Deleting Directories

Directories can also be unlinked with `unlink()`. But only if empty!

How does “`rm -rf`” work? Let’s find out with **strace**!


```
void recursiveDelete(char* dirname) {  
    char filename[FILENAME_MAX];  
    DIR *dp = opendir (dirname);  
    struct dirent *ep;  
    while(ep = readdir (dp)) {  
        snprintf(filename, FILENAME_MAX,  
                 "%s/%s", dirname, ep->d_name);  
        if(is_dir(ep))  
            recursiveDelete(filename);  
        else  
            unlink(filename);  
    }  
    unlink(dirname);  
}
```

my worst bug ever

Hard Links

- `link(old pathname, new one)`
 - **Link** a new file name to an old one
 - Create another way to refer to *the same file*
 - The command-line link program : `ln`

```
prompt> echo hello > file
prompt> cat file
hello
prompt> ln file file2 // create a hard link, link file to file2
prompt> cat file2
hello
```

Hard Links (Cont.)

- The way `link` works:
 - **Create** another name in the directory.
 - **Refer** it to the same inode number of the original file.
 - The file is not copied in any way.
 - Then, we now just have two human names (`file` and `file2`) that both refer to the same file.

Hard Links (Cont.)

- The result of `link()`

```
prompt> ls -i file file2
67158084 file /* inode value is 67158084 */
67158084 file2 /* inode value is 67158084 */
prompt>
```

- Two files have **same inode** number, but two human name (file, file2).
- There is **no difference** between file and file2.
 - Both just links to the underlying metadata about the file.

Hard Links (Cont.)

- Thus, to remove a file, we call `unlink()`.

```
prompt> rm file
removed 'file'
prompt> cat file2           // Still access the file
hello
```

- ***reference count***

- Track how many different file names have been linked to this inode.
- When `unlink()` is called, the reference count decrements.
- If the reference count reaches zero, the filesystem free the inode and related data blocks. → truly "delete" the file

Hard Links (Cont.)

- The result of `unlink()`
 - `stat()` shows the reference count of a file.

```
prompt> echo hello > file          /* create file*/
prompt> stat file
... Inode: 67158084 Links: 1 ...    /* Link count is 1 */
prompt> ln file file2              /* hard link file2 */
prompt> stat file
... Inode: 67158084 Links: 2 ...    /* Link count is 2 */
prompt> stat file2
... Inode: 67158084 Links: 2 ...    /* Link count is 2 */
prompt> ln file2 file3             /* hard link file3 */
prompt> stat file
... Inode: 67158084 Links: 3 ...    /* Link count is 3 */
prompt> rm file                    /* remove file */
prompt> stat file2
... Inode: 67158084 Links: 2 ...    /* Link count is 2 */
prompt> rm file2                   /* remove file2 */
prompt> stat file3
... Inode: 67158084 Links: 1 ...    /* Link count is 1 */
prompt> rm file3
```

Symbolic Links (Soft Link)

- Symbolic link is more **useful** than Hard link.
 - Hard Link cannot create to a directory.
 - Hard Link cannot create to a file to other partition.
 - Because inode numbers are only unique within a file system.
- Create a symbolic link: `ln -s`

```
prompt> echo hello > file
prompt> ln -s file file2 /* option -s : create a symbolic link, */
prompt> cat file2
hello
```

Symbolic Links (Cont.)

- What is different between *Symbolic link* and *Hard Link*?
 - Symbolic links are a **third type** the file system knows about.

```
prompt> stat file
... regular file ...
prompt> stat file2
... symbolic link ...      // Actually a file it self of a different type
```

- The size of symbolic link (file2) is **4 bytes**.

```
prompt> ls -al
drwxr-x---  2 remzi remzi   29 May 3 19:10 ./
drwxr-x--- 27 remzi remzi 4096 May 3 15:14 ../          // directory
-rw-r----- 1 remzi remzi    6 May 3 19:10 file         // regular file
lrwxrwxrwx  1 remzi remzi    4 May 3 19:10 file2 -> file  // symbolic link
```

- A symbolic link holds the pathname of the linked-to file as the data of the link file.

Symbolic Links (Cont.)

- If we link to a longer pathname, our link file would be bigger.

```
prompt> echo hello > alongerfilename
prompt> ln -s alongerfilename file3
prompt> ls -al alongerfilename file3
-rw-r----- 1 remzi remzi  6 May 3 19:17 alongerfilename
lrwxrwxrwx 1 remzi remzi 15 May 3 19:17 file3 -> alongerfilename
```

Symbolic Links (Cont.)

- **Dangling reference**

- When remove a original file, symbolic link points nothing.

```
prompt> echo hello > file
prompt> ln -s file file2
prompt> cat file2
hello
prompt> rm file           // remove the original file
prompt> cat file2
cat: file2: No such file or directory
```

Links: Demonstrate

Show hard links: Both path names use same inode number

File does not disappear until all removed; cannot link directories

```
Echo "Beginning.." > file1
```

```
"ln file1 link"
```

```
"cat link"
```

```
"ls -li" to see reference count
```

```
Echo "More info.." >> file1
```

```
"mv file1 file2"
```

```
"rm file2" decreases reference count
```

Soft or symbolic links: Point to second path name; can softlink to dirs

```
"ln -s oldfile softlink"
```

Confusing behavior: "file does not exist"!

Confusing behavior: "cd linked_dir; cd ..; in different parent!"

Making and Mounting a File System

- `mkfs` tool : Make a file system
 - Write an empty file system, starting with *a root directory*, on to a disk partition.
 - Input:
 - A device (such as a disk partition, e.g., `/dev/sda1`)
 - A file system type (e.g., `ext3`)

Making and Mounting a File System (Cont.)

- `mount ()`
 - Take an existing directory as a target **mount point**.
 - Essentially paste a new file system onto the directory tree at that point.

- **Example)**

```
prompt> mount -t ext3 /dev/sda1 /home/users  
prompt> ls /home/users  
a b
```

- The pathname `/home/users/` now refers to the root of the newly-mounted directory.

Making and Mounting a File System (Cont.)

- `mount` program: show **what is mounted** on a system.

```
/dev/sda1 on / type ext3 (rw)
proc on /proc type proc (rw)
sysfs on /sys type sysfs (rw)
/dev/sda5 on /tmp type ext3 (rw)
/dev/sda7 on /var/vice/cache type ext3 (rw)
tmpfs on /dev/shm type tmpfs (rw)
AFS on /afs type afs (rw)
```

- `ext3`: A standard disk-based file system
- `proc`: A file system for accessing information about current processes
- `tmpfs`: A file system just for temporary files
- `AFS`: A distributed file system

Many File Systems

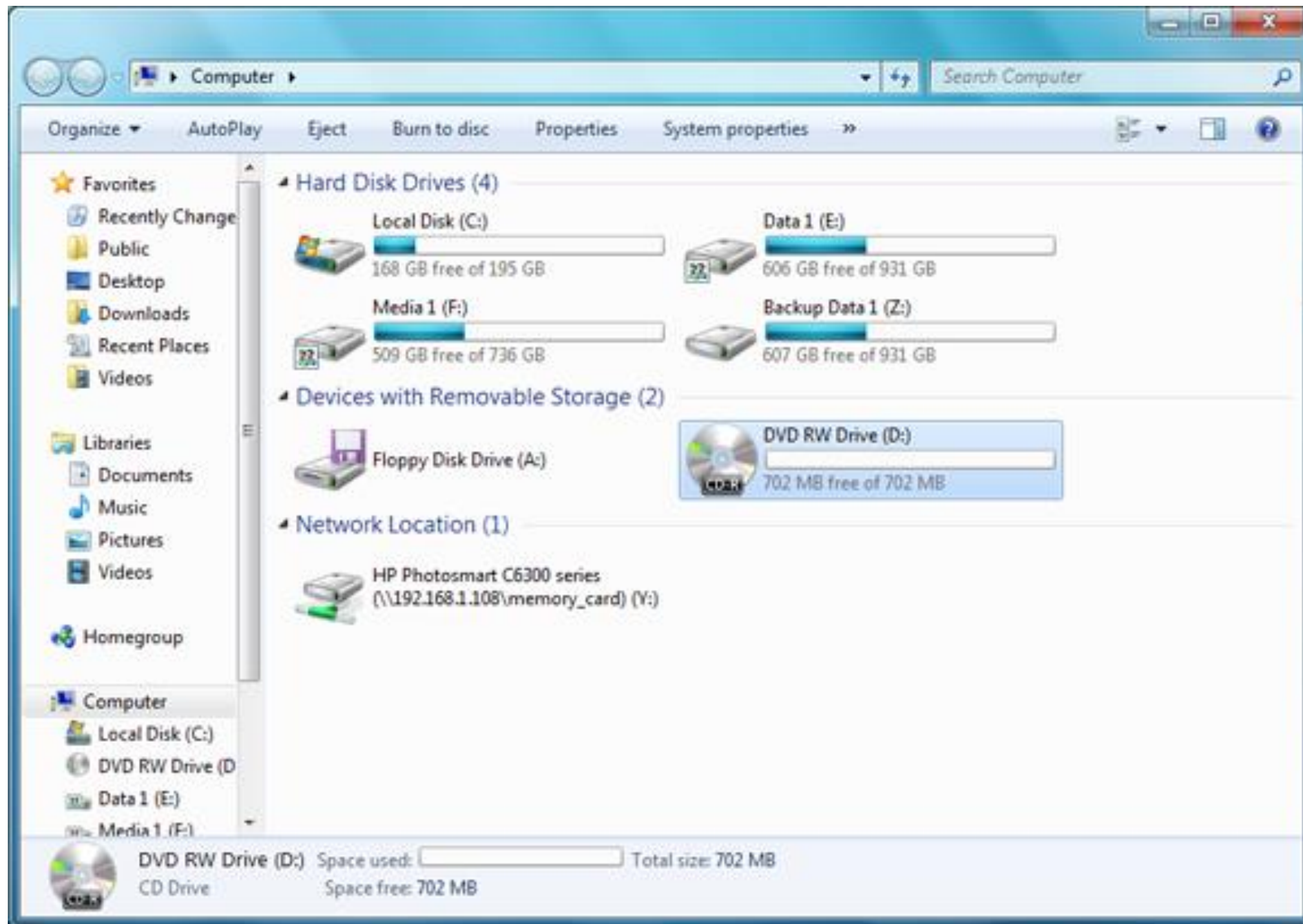
Users often want to use many file systems

For example:

- main disk
- backup disk
- AFS
- thumb drives

What is the most **elegant** way to support this?

Many File Systems: Approach 1



Many File Systems: Approach 2

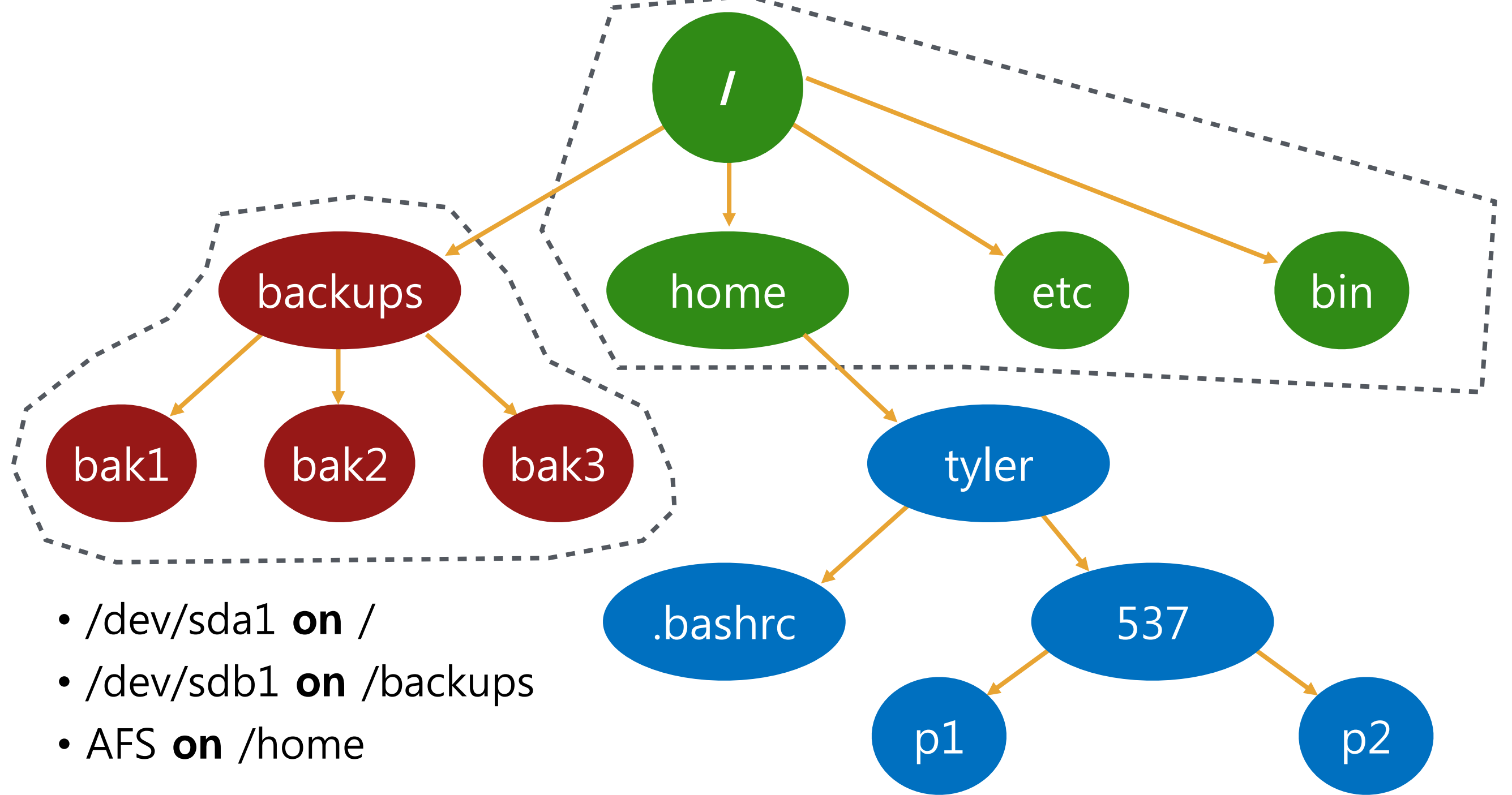
Idea: stitch all the file systems together into a super file system!

```
sh> mount
```

```
/dev/sda1 on / type ext4 (rw)
```

```
/dev/sdb1 on /backups type ext4 (rw)
```

```
AFS on /home type afs (rw)
```



Writing Immediately with `fsync()`

- The file system will **buffer** writes in memory for some time.
 - Ex) 5 seconds, or 30
 - Performance reasons
- At that later point in time, the write(s) will **actually be issued** to the storage device.
 - Write seem to complete quickly.
 - Data can be lost (e.g., the machine crashes).

Writing Immediately with `fsync()`

(Cont.)

- However, some applications require more than eventual guarantee.
 - Ex) DBMS requires force writes to disk from time to time.
- `off_t fsync(int fd)`
 - Filesystem forces all dirty (i.e., not yet written) data to disk for the file referred to by the file description.
 - `fsync()` returns once all of these writes are complete.

Writing Immediately with `fsync()` (Cont.)

- An Example of `fsync()`.

```
int fd = open("foo", O_CREAT | O_WRONLY | O_TRUNC);  
assert (fd > -1)  
int rc = write(fd, buffer, size);  
assert (rc == size);  
rc = fsync(fd);  
assert (rc == 0);
```

- In some cases, this code needs to `fsync()` the directory that contains the file `foo`.

Communicating Requirements: fsync

File system keeps newly written data in memory for awhile

Write buffering improves performance (why?)

But what if system **crashes** before buffers are flushed?

If application cares:

`fsync(int fd)` forces buffers to flush to disk, and (usually) tell
s disk to flush its write cache too

Makes data **durable**

Renaming Files

- `rename(char* old, char *new)`
 - Rename a file to different name.
 - It implemented as an **atomic call**.
 - Ex) Change from `foo` to `bar`:

```
prompt> mv foo bar          // mv uses the system call rename()
```

- Ex) How to update a file atomically:

```
int fint fd = open("foo.txt.tmp", O_WRONLY|O_CREAT|O_TRUNC);  
write(fd, buffer, size); // write out new version of file  
fsync(fd);  
close(fd);  
rename("foo.txt.tmp", "foo.txt");
```

rename

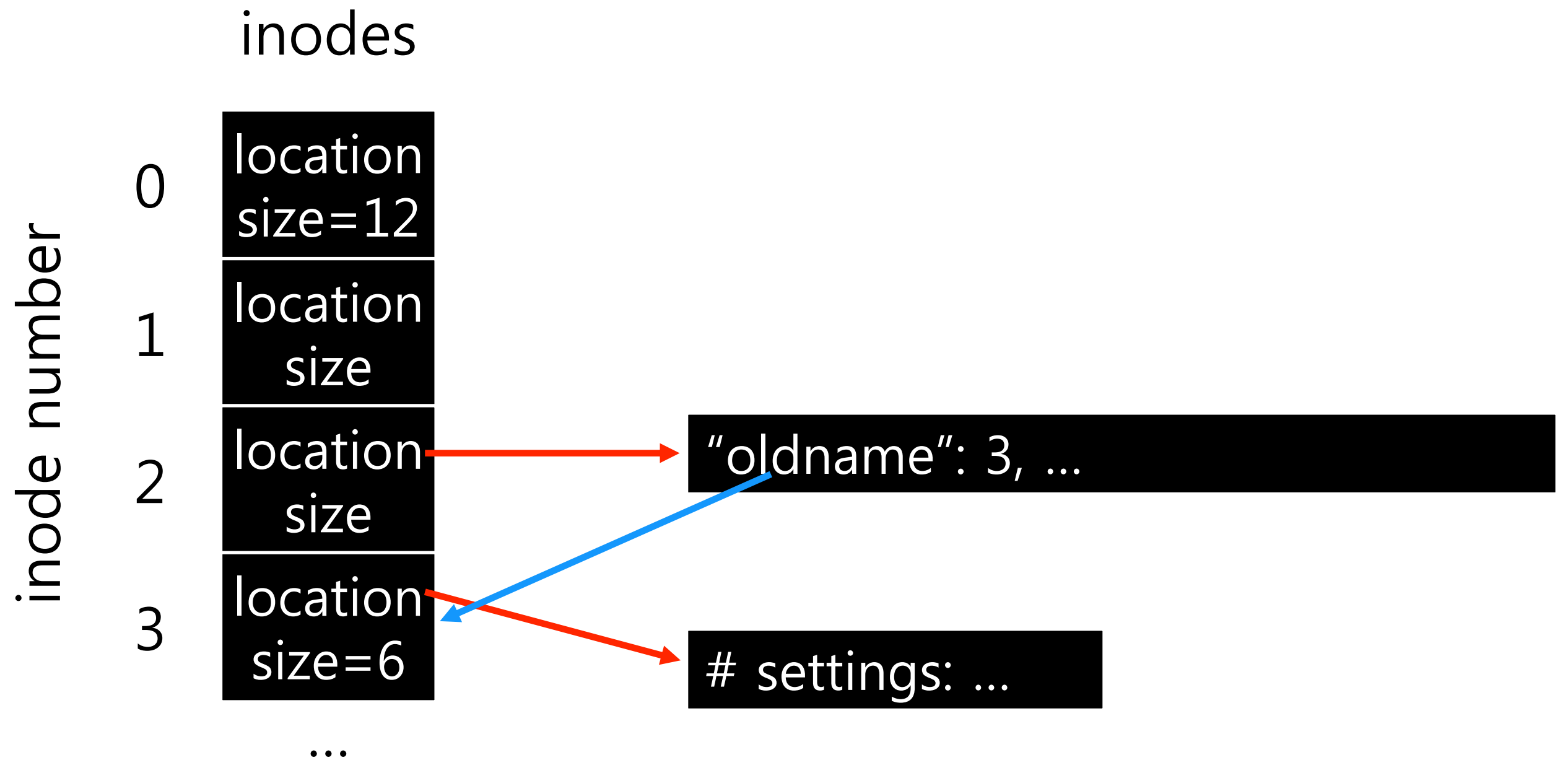
rename(char *old, char *new):

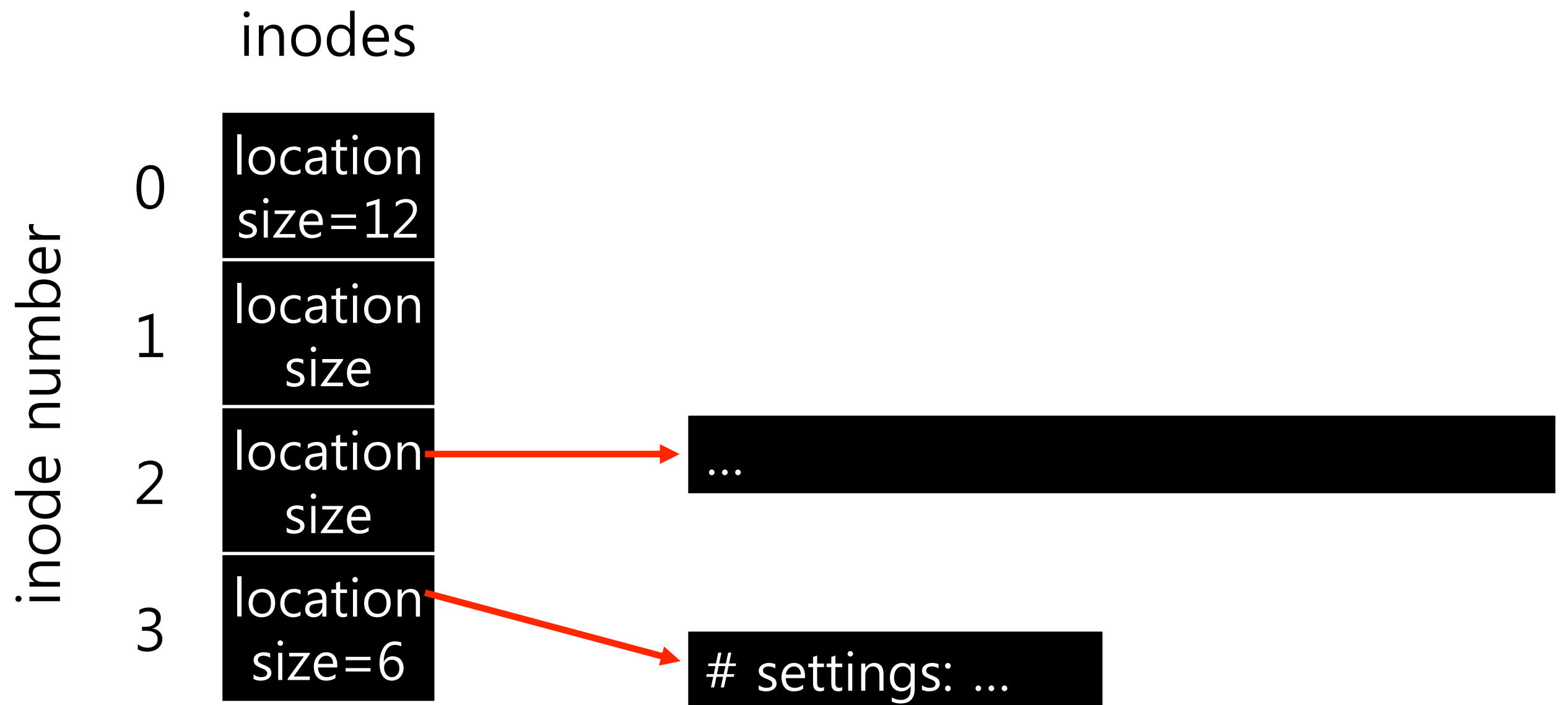
- deletes an old link to a file
- creates a new link to a file

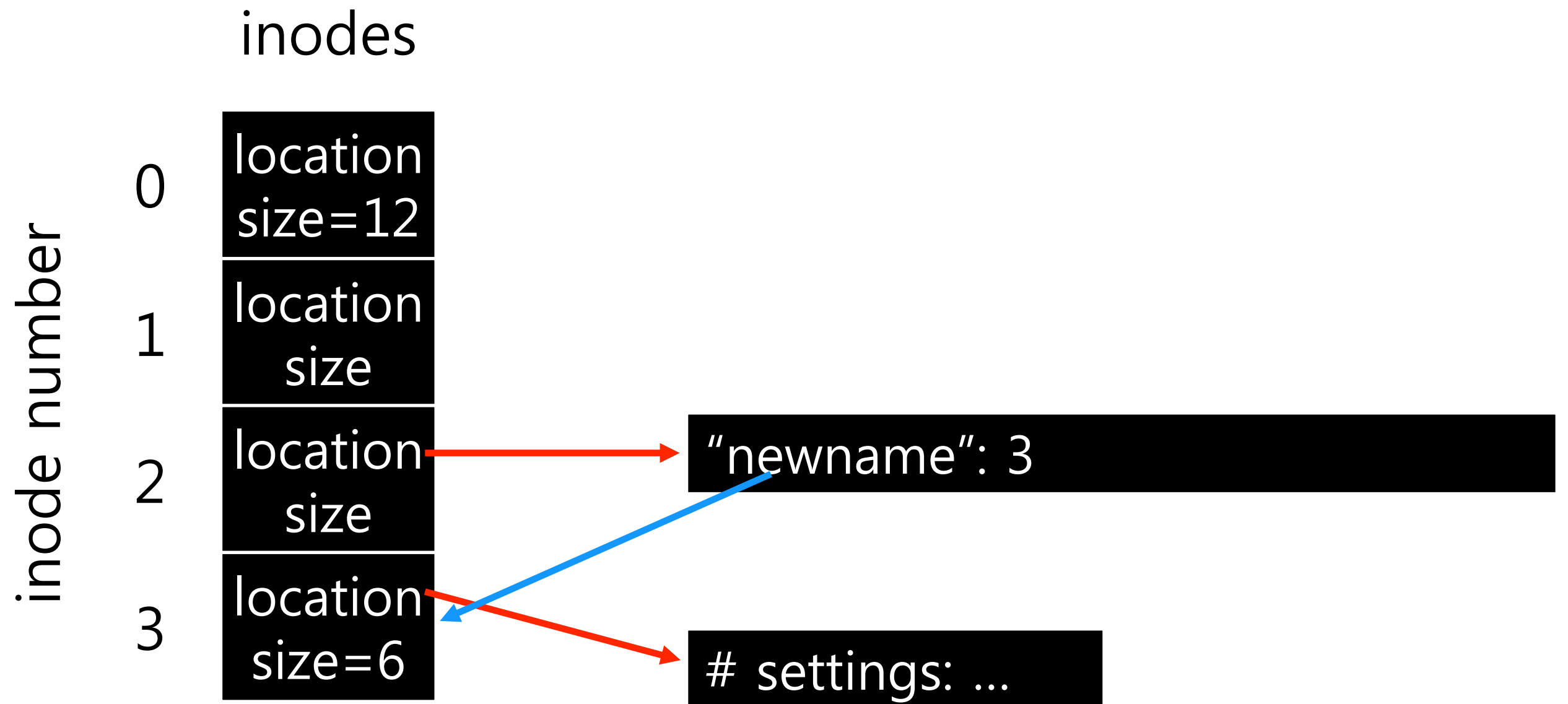
Just changes name of file, does not move data

Even when renaming to new directory (unless...?)

What can go wrong if system crashes at wrong time?







rename

rename(char *old, char *new):

- deletes an old link to a file
- creates a new link to a file

What if we crash?

FS does extra work to guarantee atomicity; return to this issue later...

Getting Information About Files

- `stat()`, `fstat()` : Show the file metadata
 - **Metadata** is information about each file.
 - Ex) Size, Low-level name, Permission, ...
 - `stat` structure is below:

```
struct stat {
    dev_t st_dev; /* ID of device containing file */
    ino_t st_ino; /* inode number */
    mode_t st_mode; /* protection */
    nlink_t st_nlink; /* number of hard links */
    uid_t st_uid; /* user ID of owner */
    gid_t st_gid; /* group ID of owner */
    dev_t st_rdev; /* device ID (if special file) */
    off_t st_size; /* total size, in bytes */
    blksize_t st_blksize; /* blocksize for filesystem I/O */
    blkcnt_t st_blocks; /* number of blocks allocated */
    time_t st_atime; /* time of last access */
    time_t st_mtime; /* time of last modification */
    time_t st_ctime; /* time of last status change */
};
```

Getting Information About Files (Cont.)

- To see stat information, you can use the command line tool `stat`.

```
prompt> echo hello > file
prompt> stat file

File: `file'
Size: 6 Blocks: 8 IO Block: 4096 regular file
Device: 811h/2065d Inode: 67158084 Links: 1
Access: (0640/-rw-r-----)  Uid: (30686/  root)  Gid: (30686/  remzi)
Access: 2011-05-03 15:50:20.157594748 -0500
Modify: 2011-05-03 15:50:20.157594748 -0500
Change: 2011-05-03 15:50:20.157594748 -0500
```

- File system keeps this type of information in a `inode` structure.

Removing Files

- `rm` is Linux command to remove a file
 - `rm` call `unlink()` to remove a file.

```
prompt> strace rm foo
...
unlink("foo")          = 0    // return 0 upon success
...
prompt>
```

Why it calls `unlink()`? not "remove or delete"
We can get the answer later.

Atomic File Update

Say application wants to update file.txt atomically

If crash, should see only old contents or only new contents

1. write new data to file.txt.tmp file
2. fsync file.txt.tmp
3. rename file.txt.tmp over file.txt, replacing it

Concurrency

How can multiple processes avoid updating the same file at the **same time**?

Normal locks don't work, as developers may have developed their programs independently.

Use **flock()**, for example:

- flock(fd, LOCK_EX)
- flock(fd, LOCK_UN)

Summary

Using multiple types of name provides

- convenience
- efficiency

Mount and link features provide flexibility.

Special calls (fsync, rename) let developers communicate special requirements to file system