# Directory

## Youjip Won
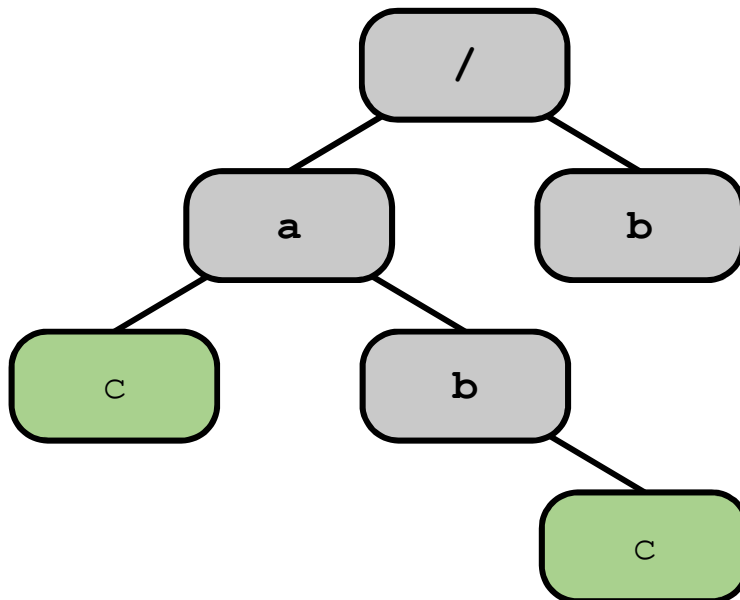
# Directory

- The file that its data is a list of directory entries.

- Directory entry is <user-readable filename, inode number> pair.

Directory (file)

EE209

EE415

EE488

Directory content

<EE209, 10>
<EE415, 24>
<EE488, 25>

# Hierarchical path name

- The directories form a tree, starting at a special directory called the root.

- In xv6, all files and directories appear under the root directory "/".

  - The slash "/" represents the root directory or is used to separate the name of a file or a directory in a path name.
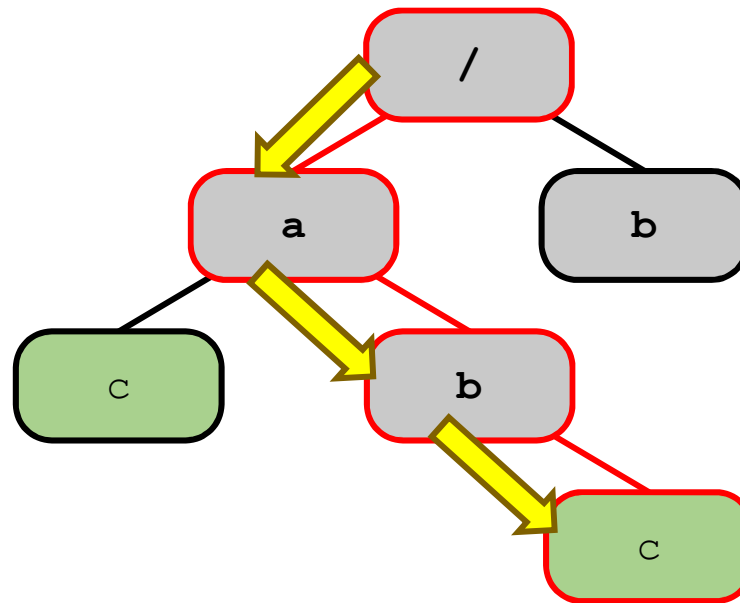
**Examples of directory**

```
/
/b
/a/c
/a/b/c
```

# Path lookup

- A path "`/a/b/c`" refers to the file or directory named `c` inside the directory "`b`" inside the directory "`a`" in the root directory.

- xv6 uses recursive lookup to find the file or directory for a given path.
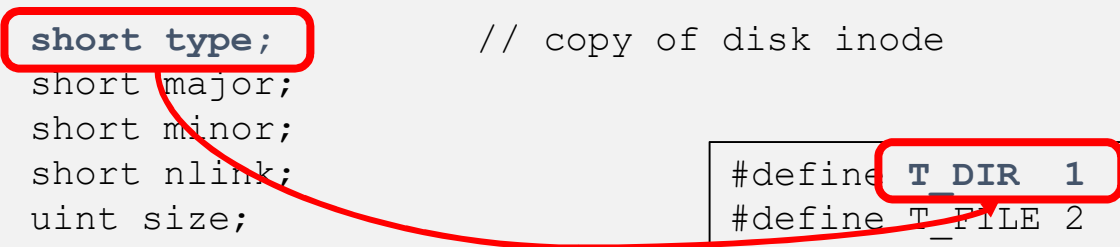
# Directory Inode

- The inode that represents a directory has type `T_DIR`.

```
12: // in-memory copy of an inode
13: struct inode {
14:    uint dev;              // Device number
15:    uint inum;             // Inode number
16:    int ref;               // Reference count
17:    struct sleeplock lock; // protects everything below here
18:    int valid;             // inode has been read from disk?
19:
20:    short type;            // copy of disk inode
21:    short major;
22:    short minor;
23:    short nlink;                  #define T_DIR  1   // Directory
24:    uint size;                    #define T_FILE 2   // File
25:    uint addrs[NDIRECT+1];        #define T_DEV  3   // Device
26: };
```
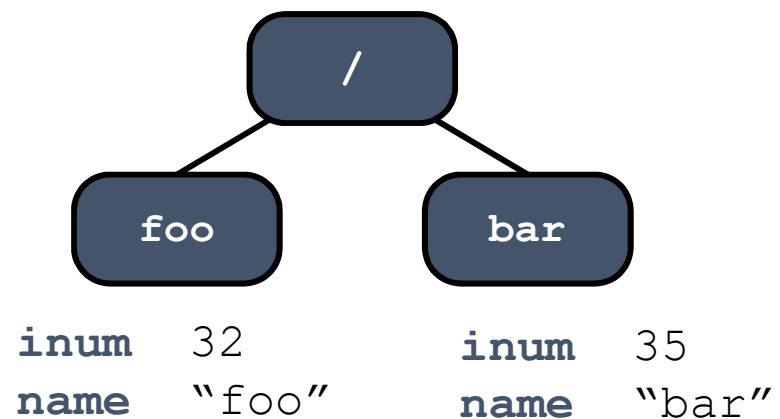
# struct dirent: Directory Entry

- Data structure `dirent` represent a directory entry.

- Each directory entry contains the inode number and the file name.

  - `dirent` wirh zero inode number is *free*.

  - Maximum length of the file name is `DIRSIZ`.

  - If the length of file name is less than `DIRSIZ`, it is terminated by a NUL (0) byte.

```
53: struct dirent {
54:   ushort inum;
55:   char name[DIRSIZ];
56: };
```

14

**16 Byte structure**

/

foo    bar

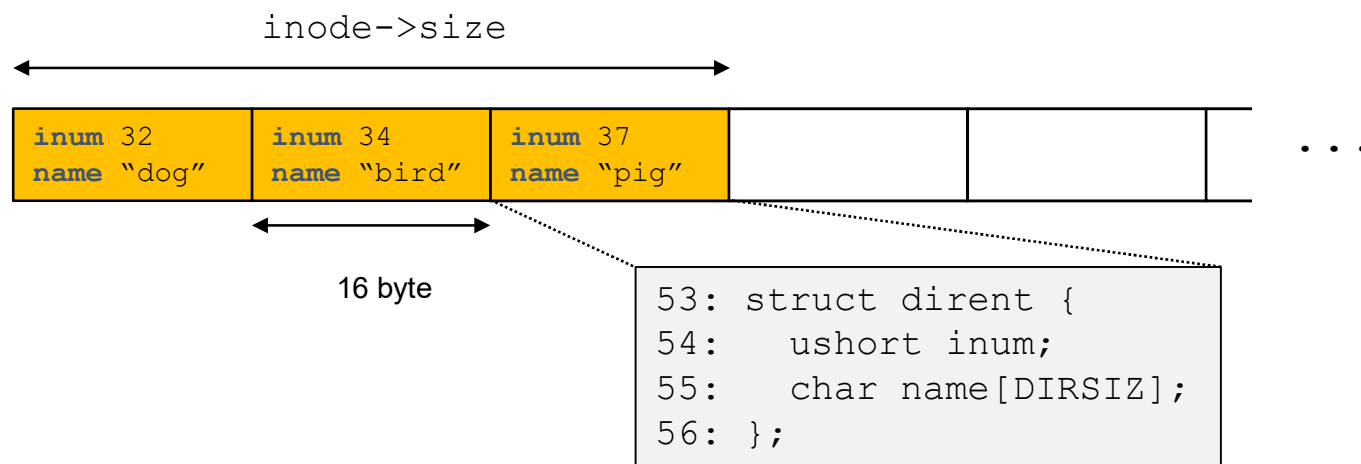inum    32          inum    35
name    "foo"       name    "bar"

# struct dirent: Directory Entry (Cont'd)

- Directory entries are stored in the file block as an array.

```
12: // in-memory copy of an inode
13: struct inode {
14:   uint dev;            // Device number
15:   uint inum;           // Inode number

            ...

24:   uint size;
25:   uint addrs[NDIRECT+1];
26: };
```

inode->size

| inum 32<br>name "dog" | inum 34<br>name "bird" | inum 37<br>name "pig" | | | | ... |

16 byte

```
53: struct dirent {
54:   ushort inum;
55:   char name[DIRSIZ];
56: };
```

# Find or insert an entry in the directory

- `struct inode *dirlookup(inode *dp, char *name,`

  `uint *poff)`

  - Find a file or directory named `name` under the directory that pointed by `dp`.

  - If there is target inode, it returns the pointer of target inode.

  - `poff` is set to the offset of the matched entry in the directory.

- `int dirlink(struct inode *dp, char *name, uint inum)`

  - Add the new directory entry to the directory that pointed by `dp`.

  - The directory entry is a pair of `name` and `inum`.

  - Return 0 on success, -1 on failure.

# `dirlookup(inode *dp, char *name, uint *poff)`

① Check the inode parameter `dp` if it is `T_DIR` typed.

② Read an entry and store into the local variable `dp`.

③ If inode number is zero, it is considered as an empty directory entry.

④ Compare the string `de.name` whether it matches the argument `name`.

⑤ Return the pointer of an inode if found by calling `iget()`.

# dirlookup()

- Search a directory for an entry with the given name `name`.

  ① Check the inode parameter `dp` if it is `T_DIR` typed.

```
524: struct inode*
525: dirlookup(struct inode *dp, char *name, uint *poff)
526: {
527:   uint off, inum;
528:   struct dirent de;
529:
530:   if(dp->type != T_DIR)
531:     panic("dirlookup not DIR");

                      ...
```
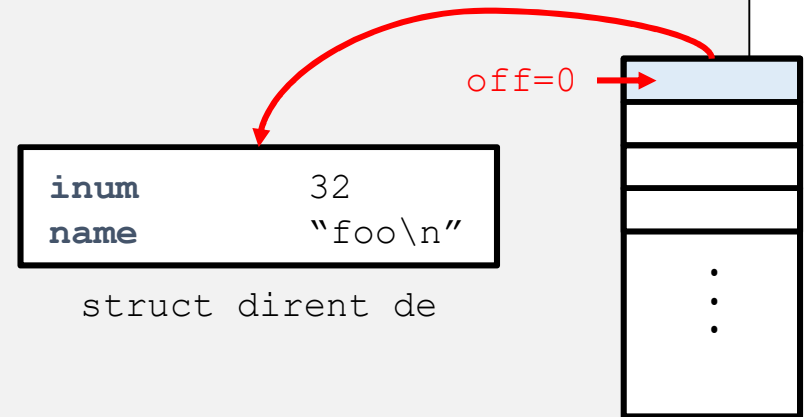
# `dirlookup()` (Cont'd)

- Search a directory for an entry with given name `name`.

  ② Read an entry and store into the local variable `de`.

```
                            ...
533:   for(off = 0; off < dp->size; off += sizeof(de)){
534:     if(readi(dp, (char*)&de, off, sizeof(de)) != sizeof(de))
535:       panic("dirlookup read");
536:     if(de.inum == 0)
537:       continue;
538:     if(namecmp(name, de.name) == 0){
539:       // entry matches path element
540:       if(poff)
541:         *poff = off;
542:       inum = de.inum;
543:       return iget(dp->dev, inum);
544:     }
545:   }
546:
547:   return 0;
548: }
```

Read the directory entry at `off`.

`off=0`

| inum | 32 |
| name | "foo\n" |

struct dirent de

# `dirlookup()` (Cont'd)

- Search a directory for an entry with given name `name`.

  ③ If inode number is zero, it is considered as an empty directory entry.
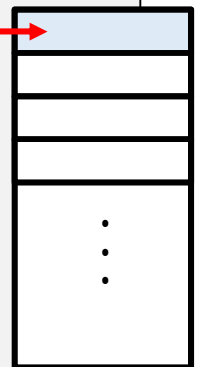
```
                          ...
533:    for(off = 0; off < dp->size; off += sizeof(de)){
534:      if(readi(dp, (char*)&de, off, sizeof(de)) != sizeof(de))
535:        panic("dirlookup read");
536:      if(de.inum == 0)
537:        continue;
538:      if(namecmp(name, de.name) == 0){
539:        // entry matches path element
540:        if(poff)
541:          *poff = off;
542:        inum = de.inum;
543:        return iget(dp->dev, inum);
544:      }
545:    }
546:
547:    return 0;
548: }
```

Check whether its
`inum` is 0 on not?          `off=0` →

| inum | 32 |
|------|-----|
| **name** | "foo\n" |

struct dirent de

KAIST OSLab
Operating Systems Laboratory

# `dirlookup()` (Cont'd)

- Search a directory for an entry with given name `name`.

  ④ Compare the string `de.name` whether it matches the argument `name`.
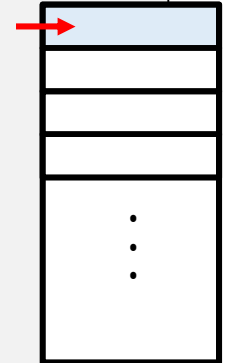
```
                              ...
533:   for(off = 0; off < dp->size; off += sizeof(de)){
534:     if(readi(dp, (char*)&de, off, sizeof(de)) != sizeof(de))
535:       panic("dirlookup read");
536:     if(de.inum == 0)
537:       continue;
538:     if(namecmp(name, de.name) == 0){
539:       // entry matches path element
540:       if(poff)
541:         *poff = off;
542:       inum = de.inum;
543:       return iget(dp->dev, inum);
544:     }
545:   }
546:
547:   return 0;
548: }
```

Check the entry to see if it matches the name it is looking for.

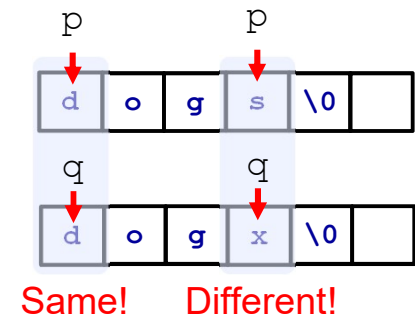off=0

| inum | 32 |
| name | "foo\n" |

struct dirent de

# dirlookup() (Cont'd)

- Search a directory for an entry with given name `name`.

  ④ Compare the string `de.name` whether it matches the argument `name`.

  `strncmp()` : Compare the given string character by character.

```
516: int
517: namecmp(const char *s, const char *t)
518: {
519:    return strncmp(s, t, DIRSIZ);
520: }
```

```
58: int
59: strncmp(const char *p, const char *q, uint n)
60: {
61:    while(n > 0 && *p && *p == *q)     If characters are same,
62:       n--, p++, q++;                   move to the next character
63:    if(n == 0)
64:       return 0;
65:    return (uchar)*p - (uchar)*q;
66: }
```



Same!    Different!

# **dirlookup()** (Cont'd)
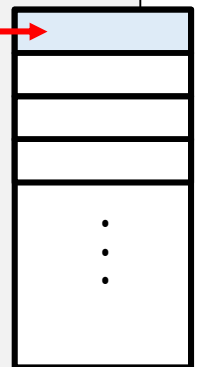
- Search a directory for an entry with given name `name`.

  ⑤ Return the pointer of an inode if found by calling `iget()`.

```
                            . . .
533:    for(off = 0; off < dp->size; off += sizeof(de)){
534:      if(readi(dp, (char*)&de, off, sizeof(de)) != sizeof(de))
535:        panic("dirlookup read");
536:      if(de.inum == 0)
537:        continue;
538:      if(namecmp(name, de.name) == 0){
539:        // entry matches path element
540:        if(poff)
541:          *poff = off;
542:        inum = de.inum;
543:        return iget(dp->dev, inum);
544:      }
545:    }
546:
547:    return 0;
548: }
```

off=0 →

| inum | 32 |
|------|----|
| name | "foo\n" |

struct dirent de
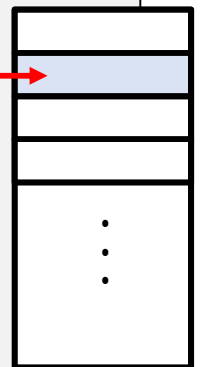
# **`dirlookup()` (Cont'd)**

- Search a directory for an entry with given name `name`.

  ⑤ - 2. Otherwise, move to the next entry and repeat ①~⑤.

```
                              . . .

533:    for(off = 0; off < dp->size; off += sizeof(de)){
534:       if(readi(dp, (char*)&de, off, sizeof(de)) != sizeof(de))
535:         panic("dirlookup read");
536:       if(de.inum == 0)
537:         continue;
538:       if(namecmp(name, de.name) == 0){
539:         // entry matches path element
540:         if(poff)
541:            *poff = off;
542:         inum = de.inum;
543:         return iget(dp->dev, inum);
544:       }
545:    }
546:
547:    return 0;
548: }
```
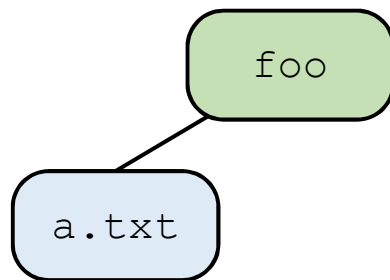
Next offset →

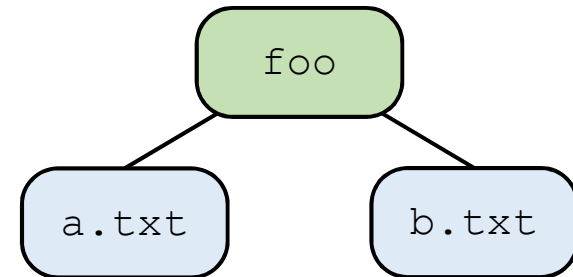# `dirlink(inode *dp, char *name, uint inum)`

- Add the new directory entry with the given `name` and inode number `inum`.

- If the name already exists, `dirlink()` returns an error (-1).

| a.txt | 10 |
|-------|-----|
| "" | 0 |
| ... | |

Data of directory `foo`

**Adding a file**
**b.txt**

| a.txt | 10 |
|-------|-----|
| b.txt | 11 |
| ... | |

Data of directory `foo`

# `dirlink()`

- Call `dirlookup()` to check any directory with the same name exists.

- `dirlookup()` returns zero if a directory entry with name `name` is not found.

```
553 int dirlink(struct inode *dp, char *name, uint inum){
554   int off;
555   struct dirent de;
556   struct inode *ip;
557
558   // Check that name is not present.
559   if((ip = dirlookup(dp, name, 0)) != 0){
560     iput(ip);
561     return -1;
562   }
563
      … // Removed for saving space.
576
577   return 0;
578 }
```

# `dirlink()` (Cont'd)

① Search for an empty directory entry. It is considered empty if inode number is zero.

```
553 int dirlink(struct inode *dp, char *name, uint inum){
      ... // Removed for saving space.
564   // Look for an empty dirent.
565   for(off = 0; off < dp->size; off += sizeof(de)){
566     if(readi(dp, (char*)&de, off, sizeof(de)) != sizeof(de))
567       panic("dirlink read");
568     if(de.inum == 0)
569       break;
570   }
      ... // Removed for saving space.
578 }
```

Read the directory
entry at `off`.

off=0

| inum | 52 |
| name | "foo" |

struct dirent de

If entry is not free, read
the next directory entry.

# `dirlink()` (Cont'd)

① Search for an empty directory entry. It is considered empty if inode number
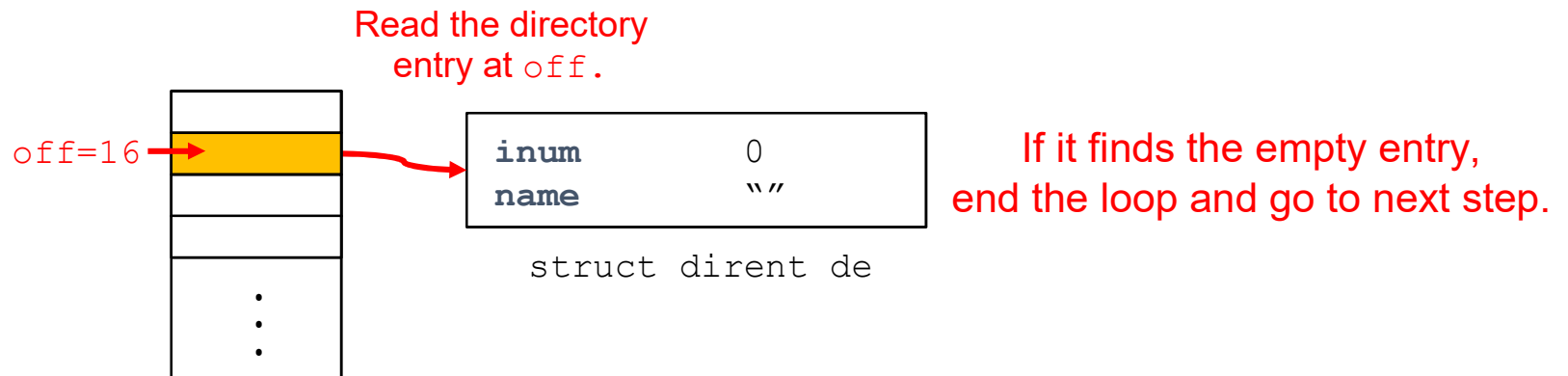
is zero.

```
553 int dirlink(struct inode *dp, char *name, uint inum){
      ... // Removed for saving space.
564   // Look for an empty dirent.
565   for(off = 0; off < dp->size; off += sizeof(de)){
566     if(readi(dp, (char*)&de, off, sizeof(de)) != sizeof(de))
567       panic("dirlink read");
568     if(de.inum == 0)
569       break;
570   }
      ... // Removed for saving space.
578 }
```

Read the directory
entry at `off`.

off=16

| inum | 0 |
| name | "" |

struct dirent de

If it finds the empty entry,
end the loop and go to next step.

# `dirlink()` (Cont'd)

② If found an empty entry, write the new entry to the this by calling `writei()`.

```
553 int dirlink(struct inode *dp, char *name, uint inum){
        … // Removed for saving space.
571
572   strncpy(de.name, name, DIRSIZ);
573   de.inum = inum;
574   if(writei(dp, (char*)&de, off, sizeof(de)) != sizeof(de))
575     panic("dirlink");
576
577   return 0;
578 }
```



Fill in the proper value and write it at the empty entry of the directory

off=16

| inum | inum |
| name | name |

Input parameter.

struct dirent de

# Pathname lookup

- Path: sequence of directories that ends with the filename or directory
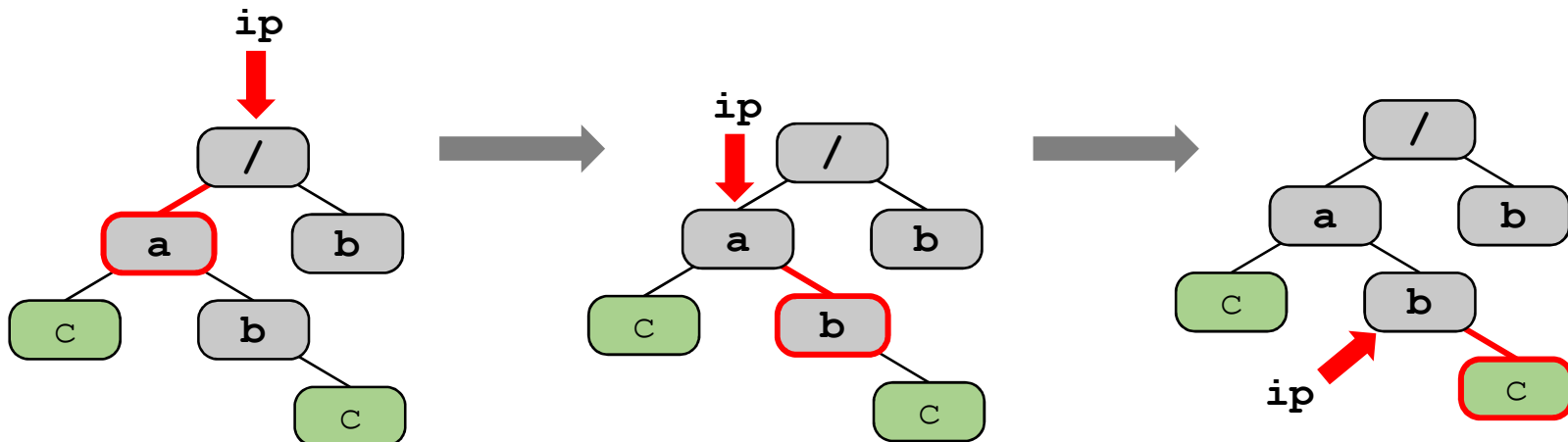
/a/b/c

- Path name lookup involves a succession of `dirlookup()` calls, one for each directory name.

- The lookup, that calls to `dirlookup()`, would start at (1) root directory or (2) process's current directory.

`dirlookup(struct inode *ip, char *name, uint *poff)`

① **dirlookup(ip,"a",)**　② **dirlookup(ip,"b",)**　③ **dirlookup(ip,"c",)**

# Pathname lookup (Cont.)

- If the path begins with a slash, evaluation begins at the root; otherwise, the current directory.

    - The current directory is the per-process attribute.

    - The system call `chdir()` change the current directory.

- Path element or component

    - For the case of path "/a/b/c", there are three elements; a, b, and c.

```
38 struct proc {
       … // Removed for saving space.
49   struct file *ofile[NOFILE];
50   struct inode *cwd;
51   char name[16];
52 };
```

# `namei()` and `nameiparent()`

- `namei()`

  - evaluates `path` and returns the corresponding `inode` of the last element.

  - calls `namex()` with 0 `nameiparent` parameter.

- `nameiparent():` evaluates `path` and returns the inode of the parent of the last element. It copies the last element to `name`.

  - calls `namex()` with 1 `nameiparent` parameter.

```
659 struct inode*
660 namei(char *path)
661 {
662   char name[DIRSIZ];
663   return namex(path, 0, name);
664 }
665
666 struct inode*
667 nameiparent(char *path, char *name)
668 {
669   return namex(path, 1, name);
670 }
```

# `namex()`: path lookup function

- `struct inode *namex(char *path, int nameiparent, char *name)`

- If `nameiparent` is 0,

    - Return the inode pointer for `name` if it is found.

- If not,

    - Copy the final component in the `path` to the `name`.

    - Return the inode pointer of the parent directory for a file `name`.

    - It is usually used when the caller should modify the directory content of a file, such as `link()` or `unlink()`.

# char *skipelem(char *path, char *name)

- char *skipelem(char *path, char *name)

  - Copy the first component of the path into name.

  - Return the pointer to the element following the copied one.

  - Examples:

    - skipelem("a/bb/c", name);

        name is set to "a" and return "bb/c".

    - skipelem("a", name);

        name is set to "a" and return "".

    - skipelem("", name);

        name is set to "" and return NULL.

# `namex()`: Get the start inode pointer.

- If the path begins with a slash, lookup begins at the root directory.

- Otherwise, it begins at the current directory.

- The inode pointer assigned to variable `ip`.

```
626 static struct inode* namex(char *path, int nameiparent, char *name) {
627    struct inode *ip, *next;
628
629    if(*path == '/')
630      ip = iget(ROOTDEV, ROOTINO);
631    else
632      ip = idup(myproc()->cwd);
633
       … // Removed for saving space.
656    return ip;
657 }
```

# `namex()`: Loop for each element in the path.

- `char *skipelem(char *path, char *name)`

  - Copy the first path element from path into `name`.

  - Return a pointer to the element following the copied one.

```
626 static struct inode* namex(char *path, int nameiparent, char *name){
        … // Removed for saving space.
633
634    while((path = skipelem(path, name)) != 0){
635        ilock(ip);


            We explain the detailed implementation, later.


649        iunlockput(ip);
650        ip = next;
651    }
        … // Removed for saving space.
656    return ip;
657 }
```

# `namex()`: Loop for each element in the path. (Cont.)

- namex ("a/b", 1, ...);

  - 1st loop: path = skipelem("a/b", name); ➔ name = "a", path = "b";

  - 2nd loop: path = skipelem("b", name); ➔ name = "b", path = "";

  - 3rd loop: path = skipelem("", name); ➔ name = "", path = NULL; ➔ Stop

```
626 static struct inode* namex(char *path, int nameiparent, char *name){
        ... // Removed for saving space.
633                                              Loop for each element in a path!
634  while((path = skipelem(path, name)) != 0){
635     ilock(ip);



        We explain the detailed implementation, later.



649     iunlockput(ip);
650     ip = next;
651   }
        ... // Removed for saving space.
656   return ip;
657 }
```

# `namex()`: Check whether the `ip` is directory or not.

- For each loop (each element), there are three things to do.

- xv6 finds the element named `name` in the directory `ip` at the third step.

- Before doing the third step, xv6 checks whether the `ip` is directory or not.

```
626 static struct inode* namex(char *path, int nameiparent, char *name){
        … // Removed for saving space.
633
634    while((path = skipelem(path, name)) != 0){
635      ilock(ip);
```

> ① Check whether the `ip` is directory or not.

> ② If `nameiparent` is not 0, stop the lookup one step earlier

> ③ By calling `dirlookup()`, find the inode named `name` in directory `ip`.

```
649      iunlockput(ip);
650      ip = next;
651    }
        … // Removed for saving space.
657 }
```

# `namex()`: Check whether the `ip` is directory or not.

- Before checking it, xv6 acquire the lock for inode `ip`.

- The type of `ip` should be `T_DIR`. Otherwise, release the lock and return the NULL.

```
626 static struct inode* namex(char *path, int nameiparent, char *name){
        … // Removed for saving space.
633
634    while((path = skipelem(path, name)) != 0){
635      ilock(ip);
636      if(ip->type != T_DIR){
637        iunlockput(ip);
638        return 0;
639      }
640

        … // Removed for saving space.
649      iunlockput(ip);
650      ip = next;
651    }
        … // Removed for saving space.
657 }
```

# `namex()`: `nameiparent` is not 0

- If `nameiparent` is not 0,
  - Return **the inode pointer of the parent directory** for the last component in the path.

```
626 static struct inode* namex(char *path, int nameiparent, char *name){
        … // Removed for saving space.
633
634   while((path = skipelem(path, name)) != 0){
635     ilock(ip);
```

① Check whether the `ip` is directory or not.

② If `nameiparent` is not 0, stop the lookup one step earlier.

③ By calling `dirlookup()`, find the inode named `name` in directory `ip`.

```
649     iunlockput(ip);
650     ip = next;
651   }
        … // Removed for saving space.
657 }
```

# `namex(): nameiparent` is not 0

- If the first character of `path` is '\0', there is no more component in `path`.

- Since the next `skipelem()` call will return `NULL`, the loop stops at the next step.

- So it returns current `ip`, which is the parent directory of the last component in the path.

```
626 static struct inode* namex(char *path, int nameiparent, char *name){
        … // Removed for saving space.
633
634   while((path = skipelem(path, name)) != 0){
635     ilock(ip);
          … // Removed for saving space.
641     if(nameiparent && *path == '\0'){
642       iunlock(ip);
643       return ip;
644     }
          … // Removed for saving space.
649     iunlockput(ip);
650     ip = next;
651   }
        … // Removed for saving space.
657 }
```

# `namex()`: find the inode named `name` in directory `ip`.

- `namex()` calls the `dirlookup(ip, name, 0)` for finding the inode for `name`.

```
626 static struct inode* namex(char *path, int nameiparent, char *name){
       … // Removed for saving space.
633
634    while((path = skipelem(path, name)) != 0){
635       ilock(ip);
```

> ① Check whether the `ip` is directory or not.

> ② If `nameiparent` is not 0, stop the lookup one step earlier.

> ③ By calling `dirlookup()`, find the inode named `name` in directory `ip`.

```
649       iunlockput(ip);
650       ip = next;
651    }
       … // Removed for saving space.
657 }
```

# Example: `namex("a/b", 1, …)`

- 1st loop: `name = "a", path = "b"; ➔ next = dirlookup(cwd, "a", 0)`
  - ➔ `ip = next = 0xdeadbeef` // inode pointer of "a".

- 2nd loop: `name = "b", path = ""; ➔ next = dirloopup(0xdeadbeef, "b", 0)`
  - ➔ `ip = next = 0x8badf00d` // inode pointer of "a/b".

```
626 static struct inode* namex(char *path, int nameiparent, char *name){
        … // Removed for saving space.
633
634   while((path = skipelem(path, name)) != 0){
635     ilock(ip);
        … // Removed for saving space.
645     if((next = dirlookup(ip, name, 0)) == 0){
646       iunlockput(ip);
647       return 0;
648     }
649     iunlockput(ip);
650     ip = next;
651   }
        … // Removed for saving space.
657 }
```

# namex("",1,…)

- If input parameter `path` of `namex()` is not empty string, `namex()` calls `return` within the loop.

- Otherwise, it does not go into the loop and return NULL.

```
626 static struct inode* namex(char *path, int nameiparent, char *name) {
       … // Removed for saving space.
633
634    while((path = skipelem(path, name)) != 0){
         … // Removed for saving space.
651    }
652    if(nameiparent){
653      iput(ip);
654      return 0;
655    }
656    return ip;
657 }
```

# `namex()`: Acquire and release per-inode lock.

- Each iteration of the loop begins by locking `ip` and find the inode named `name` in `ip`.

- Then, release the lock of `ip` before the end of the iteration.

- `namex()` locks each directory in the path separately.

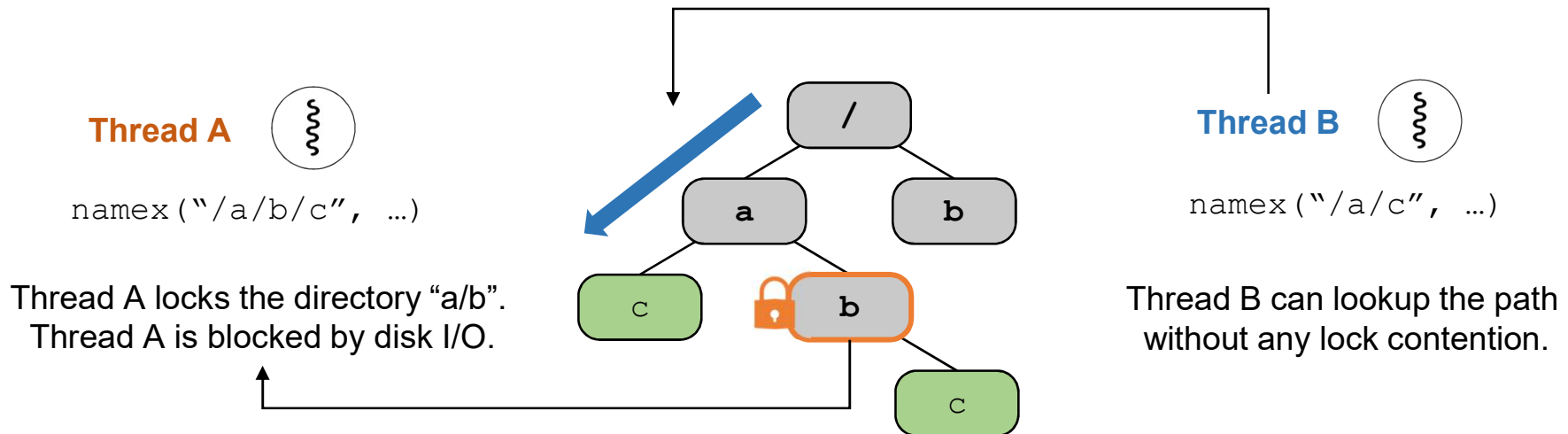  - ➜ Lookups in different directories can proceed in parallel.

```
626 static struct inode* namex(char *path, int nameiparent, char *name){
      … // Removed for saving space.
633
634   while((path = skipelem(path, name)) != 0){
635     ilock(ip);


        By calling dirlookup(), find the inode named name in directory ip.


649     iunlockput(ip);
650     ip = next;
651   }
      … // Removed for saving space.
656   return ip;
657 }
```
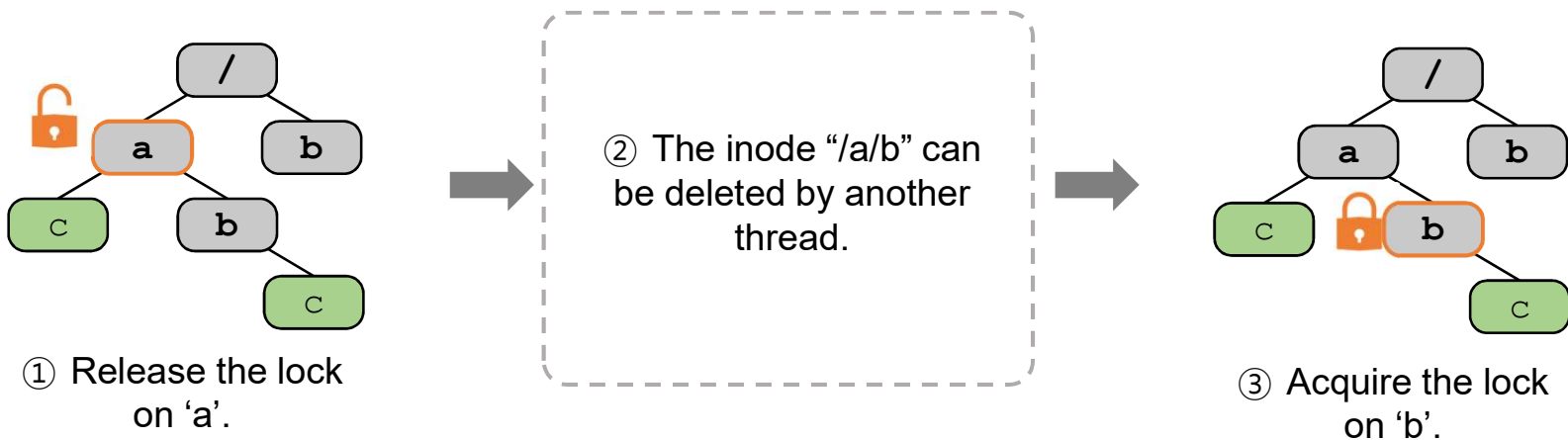
# Concurrency of `namex()`

- The procedure `namex()` may take a long time to complete.

  - It could involve several disk operations.

  - `ilock()` could read on-disk inodes to load the inode structure in memory.

  - `dirlookup()` could read file blocks of directories to traverse its entries.

- `namex()` locks each directory in the path separately.

- If a thread invokes `namex()`, another thread looking up a different pathname can proceed concurrently.

**Thread A**

`namex("/a/b/c", …)`

Thread A locks the directory "a/b".
Thread A is blocked by disk I/O.

**Thread B**

`namex("/a/c", …)`

Thread B can lookup the path
without any lock contention.

# Risk of concurrency: race condition

- In `namex()`, each iteration only locks a single inode.

  - `dirlookup()` returns the pointer of next inode.

  - The returned inode pointer is locked after releasing the lock of parent directory.

- There can be following situation in xv6.



① Release the lock on 'a'.

② The inode "/a/b" can be deleted by another thread.

③ Acquire the lock on 'b'.

**Is the inode pointer returned from `dirlookup()` still valid?**
**Can xv6 invoke `ilock()` for this inode pointer?**

# Risk of concurrency: race condition (Cont.)

- `dirlookup()` returns an inode pointer that was obtained using `iget()`.

  - `iget()` increases the reference count of the inode.

- In xv6, **if reference count is larger than 0,** the inode is not deleted from inode cache and from the file system. (`iput()`)

- By separating the `iget()` and `ilock()`, xv6 avoids the race condition.

```
333 void iput(struct inode *ip){
334    acquiresleep(&ip->lock);
335    if(ip->valid && ip->nlink == 0){
336       acquire(&icache.lock);
337       int r = ip->ref;
338       release(&icache.lock);
339       if(r == 1){

               Remove the in-memory inode as well as on-disk inode.
345       }
346    }
347    releasesleep(&ip->lock);
                              ...
352 }
```

If this process is the last reference, xv6 removes this inode.

# Risk of concurrency: deadlock

- What happen if locking the next inode before releasing the lock on the parent directory?

  - It may result in a deadlock.

  - If `namex("/./a",…)` is invoked, the next inode "/." is same with parent directory "/" in the first iteration.

  - In this case, the thread may try to acquire the lock that already held.

# Summary

- Directory layer

  - `dirlookup()` and `dirlink()`

- Path lookup

  - `namex()`, `namei()`, and `nameiparent()`

  - Concurrency of `namex()`