



Stýrkerfi (Operating Systems)

18. Directories & File System Examples

Hans P. Reiser | Spring term 2025

REYKJAVIK UNIVERSITY



Overview

- 1 Implementing Directories
- 2 UNIX File System Structure
- 3 File System Examples

Implementing Directories

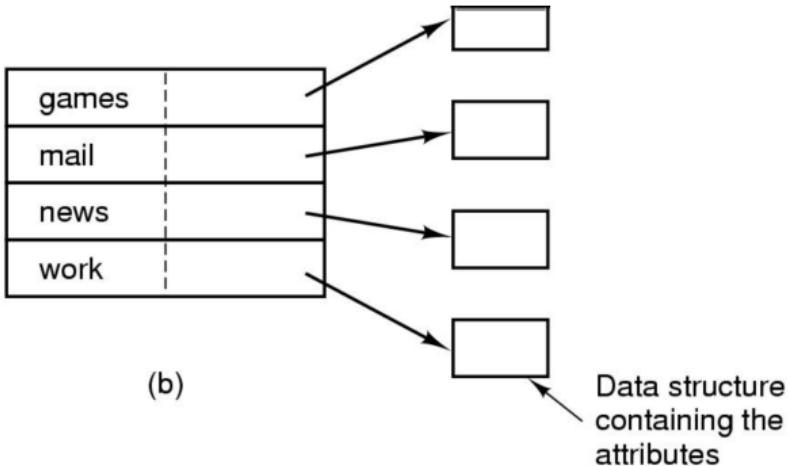
How to find the mapping of a file name to blocks?

- Storing file names
- Deleting entries
- Linear directory lookup
- Hashed directory lookup
- Tree-structured directory lookup

Implementing Directories

| | |
|-------|------------|
| games | attributes |
| mail | attributes |
| news | attributes |
| work | attributes |

(a)

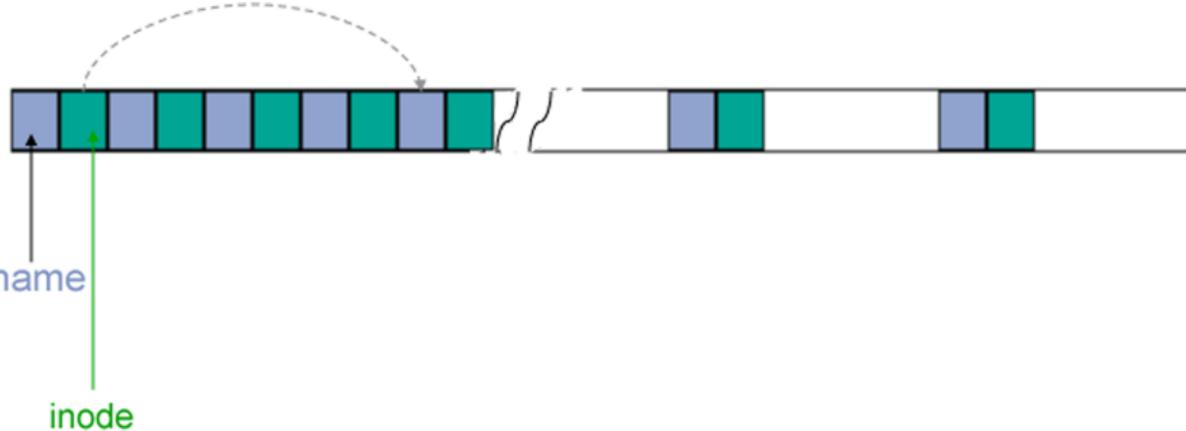


(a) A simple directory (MS-DOS)

- fixed size entries
- disk addresses and attributes in directory entry

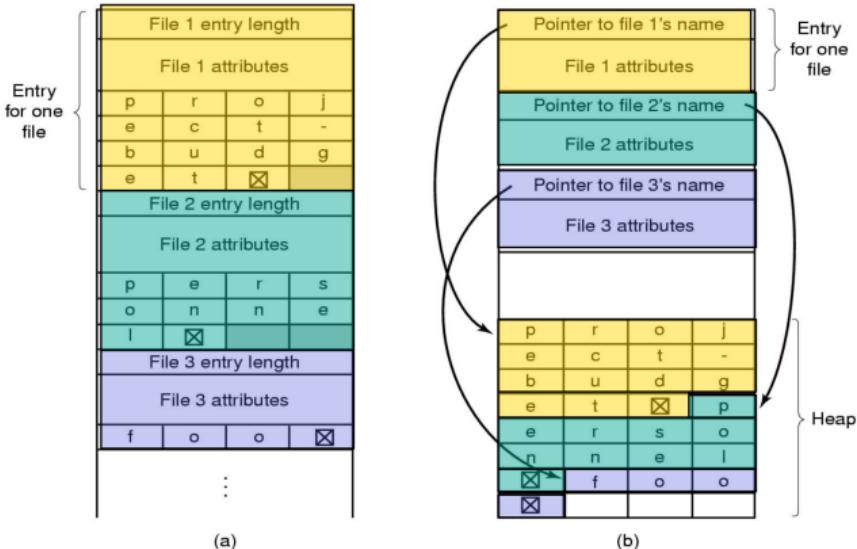
(b) Directory in which each entry just refers to an inode (UNIX)

Implementing Directories



- What to do when some entries are deleted?
 - Never reuse
 - Bridge over the directory holes
 - Compaction, but when?
 - eager or
 - lazy

Directory Entries & Long Filenames



- Two ways of handling long file names in directories

- (a) In-line
 - (b) In a heap

(1) Linear Directory Lookup

- Linear search ⇒ for big directories not efficient
- Space efficient as long as we do compaction
 - Either eagerly after entry deletion or
 - Lazily (but when?)
- With variable file names ⇒ deal with fragmentation
- Alternatives
 - (e.g. extensible) hashing
 - (e.g. B-)tree structures

(1) Linear Directory Lookup

Example: FAT32 file system

- Fixed size (32 byte) directory table entries
- 8 byte filename, 3 byte file extension, attributes
- Special treatment of first byte of filename
 - 00h: End of list
 - E5h: Free entry (used for deleting files)

Example: ext2 file system

- File system object like a file
- Sequence of variable sized directory entries:
inode, entry length, file name (max. 256 characters)
- entry length: can be larger than file name length (4 byte alignment of entries)
- Deleting file: adjust previous length to point to subsequent entry

(2) Hashing a Directory Lookup

- Direct hashing a file name to an inode (+ handling collisions?)
- Extendible hashing (based on Fagin's algorithm):
 - Hash table with 2^k buckets
 - k lowest bits of file name hash used to select bucket
 - File name to be found in block corresponding to that bucket
 - Block full: Split into two (increase k by one if necessary)
- Advantages:
 - Fast lookup and relatively simple
- Disadvantages:
 - Not as efficient as trees for very large directories

Examples: ZFS, GFS2 (and ext3: filename hash + HTree)

(3) Tree Structure for a Directory

- Method:

- Sort files by name
- Store directory entries in a B-tree like structure
- Create/delete/search in that B-tree

- Advantages:

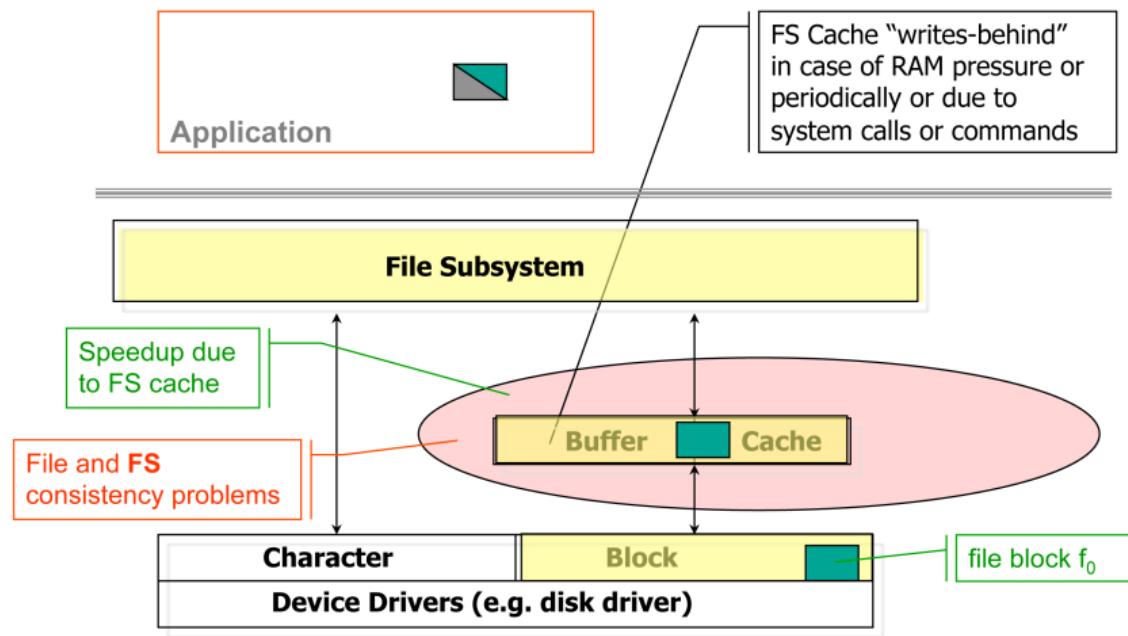
- Efficient for a large number of files per directory

- Disadvantages:

- Complex
- Not that efficient for a small number of files

Examples: BTRFS (B-Trees), XFS (B+Trees for large directories, simple list for small)

UNIX File System Structure



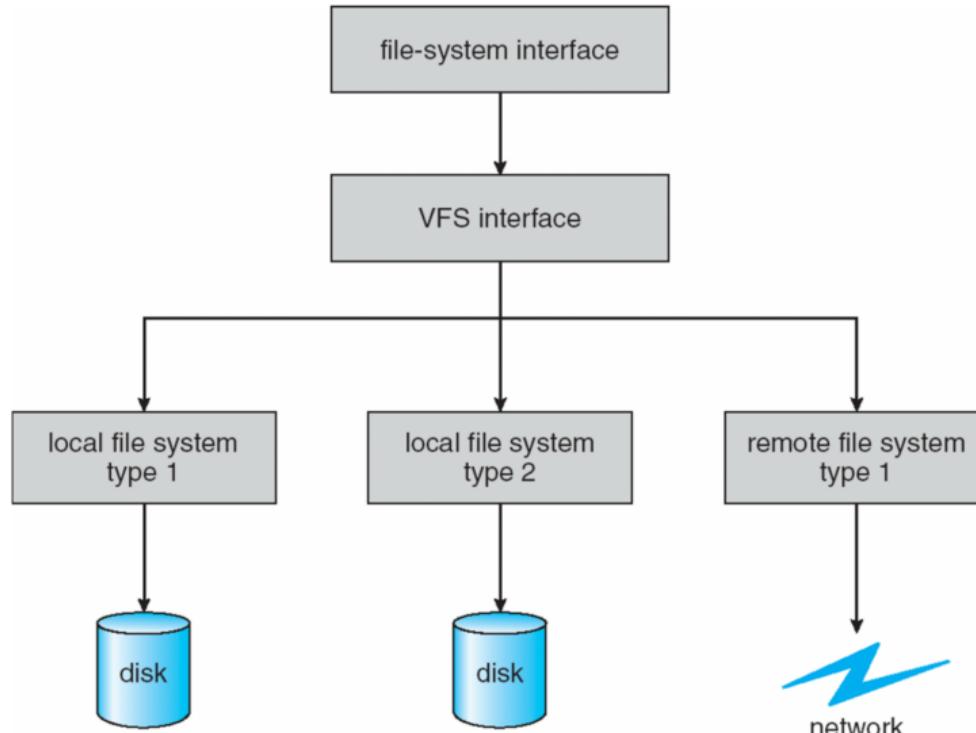
Buffer Cache

- Write-behind policy might lead to
 - data losses in case of system crash or power loss
 - forcing all dirty buffers to disk with `fsync()`
 - `rename()` is implemented atomic
 - inconsistent state of the FS (see OSTEP [ADAD18] chapter 42)
 - requires file system check `fsck`
 - journaling (write-ahead logging)
- Always two copies involved
 - from disk to buffer cache (in kernel space)
 - from buffer to user address space
- FS Cache *wiping* if sequentially reading a very large file from end to end and not accessing it again

Virtual File Systems

- Virtual File Systems (VFS) provide an object-oriented way of implementing file systems.
- VFS allows the same system call interface (the API) to be used for different types of file systems.
- The API is for the VFS interface, rather than any specific type of file system.

Schematic View of Virtual File System [SGG12]

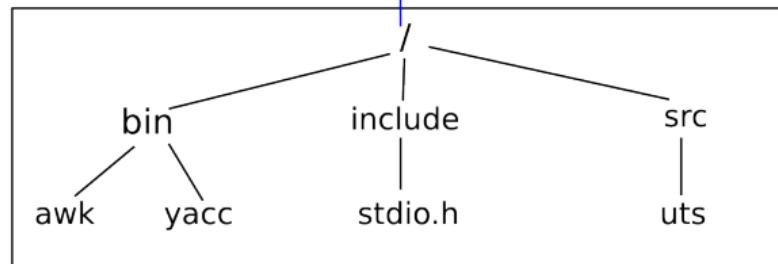
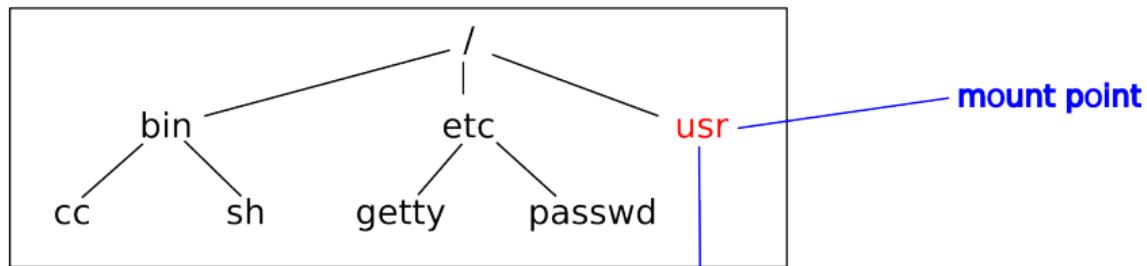


Logical and Physical File System

- A logical file system can consist of different physical file systems
- The only file system that has to be always is the root file system
(in general on a partition of a hard disk)
- A file system can be mounted at any place within another file system
- Using `mount/umount` the OS manages a mount table supporting the resolution of path name crossing file systems

Logical and Physical File System

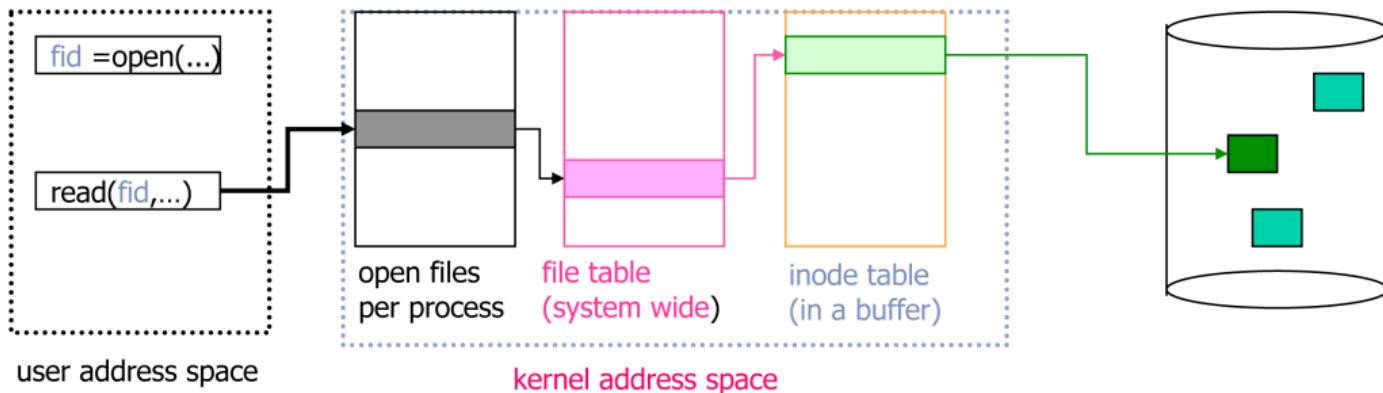
root file system



mountable file system

Using a UNIX File

- Opening a file creates a file descriptor fid
- Used as an index into a process-specific table of open files
- The corresponding table entry points to a system-wide file table
- Via buffered inode table, you finally get the data blocks



Overview

- 1 Implementing Directories
- 2 UNIX File System Structure
- 3 File System Examples

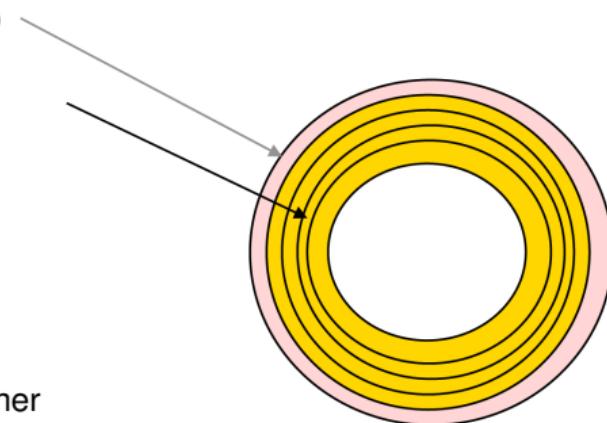
Original UNIX File System

- Simple disk layout

- Block size = sector size (512 bytes)
- Inodes on outermost cylinders
- Data blocks on the inner cylinders
- Freelist as a linked list

- Issues

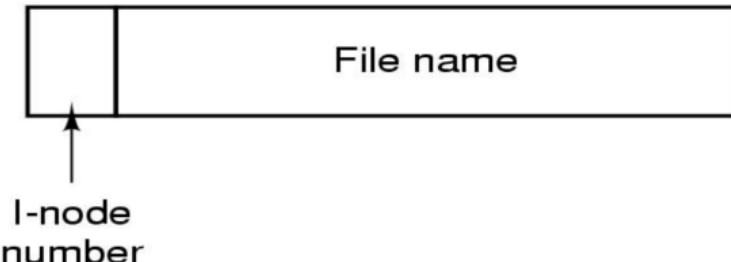
- Index is large
- Fixed number of files
- Inodes far away from data blocks
- Inodes for directory not close together
- Consecutive file blocks can be anywhere
- Poor bandwidth for sequential access



UNIX File Names

- Historically (Version 7) only 14 characters [TB15]

Bytes 2 14



- System V up to 255 ASCII characters

<filename>. <extension>

UNIX Inode

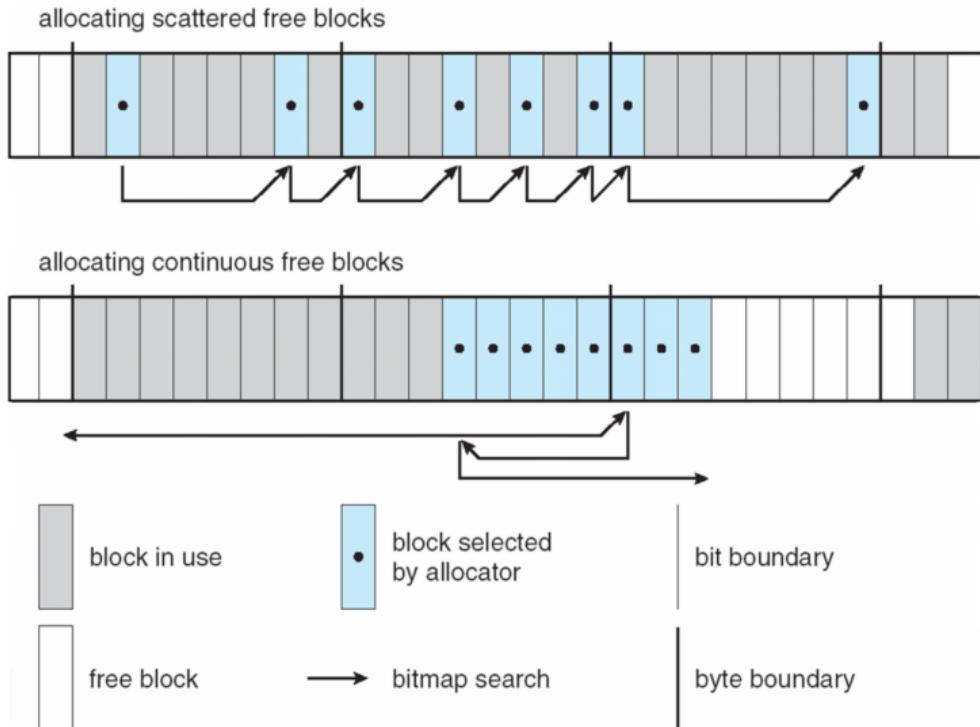
| Field | Bytes | Description |
|--------|-------|---|
| Mode | 2 | File type, protection bits, setuid, setgid bits |
| Nlinks | 2 | Number of directory entries pointing to this i-node |
| UID | 2 | UID of the file owner |
| GID | 2 | GID of the file owner |
| Size | 4 | File size in bytes |
| Addr | 39 | Address of first 10 disk blocks, then 3 indirect blocks |
| Gen | 1 | Generation number (incremented every time i-node is reused) |
| Atime | 4 | Time the file was last accessed |
| Mtime | 4 | Time the file was last modified |
| Ctime | 4 | Time the i-node was last changed (except the other times) |

BSD FFS (Fast File System)

- Use a larger block size: 4 KB
 - Allow large blocks to be chopped into 2,4 or 8 sub-blocks
 - Used for little files and pieces at the end of files
 - As the file grows, continue to allocate sub-blocks.
If 4KB reached, copy sub-blocks to a full 4 KB block.
- Use [bitmap](#) instead of a free list
 - Try to allocate more contiguously
 - 10% free space reserve for system administrator
- Optimizations to better match disk characteristics
 - Cylinder groups for exploiting locality
 - Data and meta data in same cylinder group
 - Items of one directory in same or nearby cylinder groups

See OSTEP [ADAD18] chapter 41

Ext2fs Block-Allocation Policies [SGG12]



FFS/Ext2 Directory (1)

- Directory entry needs three elements:
 - inode number (index to a table of inodes)
 - length of dir-entry (variable length of file names)
 - entry type
 - file name length
 - file name (up to 255 characters)
- Each directory contains at least two entries:
 - . . = link to the parent directory (forming the directory tree)
 - . = link to itself
- FFS offers a “tree-like structure” (like Multics), supporting human preference, ordering hierarchically

Inode Table – closeup view

The Inode Table (Closeup)

| | | | | iblock 0 | iblock 1 | iblock 2 | iblock 3 | iblock 4 | | | | |
|-------|--------|--------|--|-------------|-------------|-------------|-------------|-------------|--|--|--|--|
| Super | i-bmap | d-bmap | | 0 1 2 3 | 16 17 18 19 | 32 33 34 35 | 48 49 50 51 | 64 65 66 67 | | | | |
| | | | | 4 5 6 7 | 20 21 22 23 | 36 37 38 39 | 52 53 54 55 | 68 69 70 71 | | | | |
| | | | | 8 9 10 11 | 24 25 26 27 | 40 41 42 43 | 56 57 58 59 | 72 73 74 75 | | | | |
| | | | | 12 13 14 15 | 28 29 30 31 | 44 45 46 47 | 60 61 62 63 | 76 77 78 79 | | | | |

- Superblock: Generic information about file system (type/version, blocksize, number of inodes, etc.)
- i-bmap: free inode bitmap
- d-bmap: free data block bitmap

[ADAD18]

Journaling File Systems

- Journaling file systems record each update to the file system as a **transaction**
- All transactions are written to a **log**
 - A transaction is considered **committed** once it is written to the log
 - However, the file system may not yet be updated
- The transactions in the log are asynchronously written to the file system
 - When the file system is modified, the transaction is removed from the log
- If the file system crashes, all remaining transactions in the log must still be performed

Log-Structured File Systems

- Log-structured FS: use disk as a circular buffer
- Write all updates, including inodes, meta data and data to end of log
 - have all writes initially buffered in memory
 - periodically write these within 1 segment (1 MB)
 - when file opened, locate i-node, then find blocks
- From the other end, clear all data , no longer used

References I

- [ADAD18] R.H. Arpaci-Dusseau and A.C. Arpaci-Dusseau. *Operating Systems: Three Easy Pieces*. 2018.
retrieved from <http://pages.cs.wisc.edu/~remzi/OSTEP/>.
- [SGG12] Abraham Silberschatz, Peter B. Galvin, and Greg Gagne. *Operating System Concepts*. Wiley Publishing, 9th edition, 2012.
- [TB15] Andrew S Tanenbaum and Herbert Bos. *Modern operating systems*. Pearson, 4th edition, 2015.