

Chapter 11: File System Implementation



Operating System Concepts – 8th Edition

Silberschatz, Galvin and Gagne ©2009



Chapter 11: File System Implementation

- File-System Structure
- File-System Implementation
- Directory Implementation
- Allocation Methods
- Free-Space Management
- Efficiency and Performance
- Recovery
- Log-Structured File Systems
- NFS
- Example: WAFL File System

Operating System Concepts – 8th Edition

11.2

Silberschatz, Galvin and Gagne ©2009



Objectives

- To describe the details of implementing local file systems and directory structures
- To describe the implementation of remote file systems
- To discuss block allocation and free-block algorithms and trade-offs

Operating System Concepts – 8th Edition

11.3

Silberschatz, Galvin and Gagne ©2009



File-System Structure

- File structure
 - Logical storage unit
 - Collection of related information
- File system resides on secondary storage (disks)
- File system organized into layers
- **File control block** – storage structure consisting of information about a file

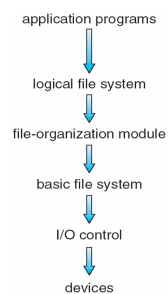
Operating System Concepts – 8th Edition

11.4

Silberschatz, Galvin and Gagne ©2009



Layered File System



Operating System Concepts – 8th Edition

11.5

Silberschatz, Galvin and Gagne ©2009



A Typical File Control Block

file permissions
file dates (create, access, write)
file owner, group, ACL
file size
file data blocks or pointers to file data blocks

Operating System Concepts – 8th Edition

11.6

Silberschatz, Galvin and Gagne ©2009





In-Memory File System Structures

- The following figure illustrates the necessary file system structures provided by the operating systems.
- Figure 12-3(a) refers to opening a file.
- Figure 12-3(b) refers to reading a file.

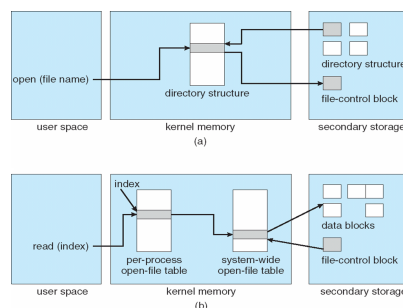
Operating System Concepts – 8th Edition

11.7

Silberschatz, Galvin and Gagne ©2009



In-Memory File System Structures



Operating System Concepts – 8th Edition

11.8

Silberschatz, Galvin and Gagne ©2009



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 to the VFS interface, rather than any specific type of file system.

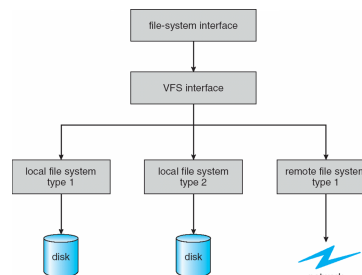
Operating System Concepts – 8th Edition

11.9

Silberschatz, Galvin and Gagne ©2009



Schematic View of Virtual File System



Operating System Concepts – 8th Edition

11.10

Silberschatz, Galvin and Gagne ©2009



Directory Implementation

- **Linear list** of file names with pointer to the data blocks.
 - simple to program
 - time-consuming to execute
- **Hash Table** – linear list with hash data structure.
 - decreases directory search time
 - **collisions** – situations where two file names hash to the same location
 - fixed size

Operating System Concepts – 8th Edition

11.11

Silberschatz, Galvin and Gagne ©2009



Allocation Methods

- An allocation method refers to how disk blocks are allocated for files:
- **Contiguous allocation**
- **Linked allocation**
- **Indexed allocation**

Operating System Concepts – 8th Edition

11.12

Silberschatz, Galvin and Gagne ©2009

Contiguous Allocation

- Each file occupies a set of contiguous blocks on the disk
- Simple – only starting location (block #) and length (number of blocks) are required
- Random access
- Wasteful of space (dynamic storage-allocation problem)
- Files cannot grow

Operating System Concepts – 8th Edition 11.13 Silberschatz, Galvin and Gagne ©2009

Contiguous Allocation

- Mapping from logical to physical

LA/512

Q
R

Block to be accessed = ! + starting address
Displacement into block = R

Operating System Concepts – 8th Edition 11.14 Silberschatz, Galvin and Gagne ©2009

Contiguous Allocation of Disk Space

file	start	length
count	0	2
tr	14	3
mail	19	6
list	28	4
f	6	2

Operating System Concepts – 8th Edition 11.15 Silberschatz, Galvin and Gagne ©2009

Extent-Based Systems

- Many newer file systems (I.e. Veritas File System) use a modified contiguous allocation scheme
- Extent-based file systems allocate disk blocks in **extents**
- An **extent** is a contiguous block of disks
 - Extents are allocated for file allocation
 - A file consists of one or more extents.

Operating System Concepts – 8th Edition 11.16 Silberschatz, Galvin and Gagne ©2009

Linked Allocation

- Each file is a linked list of disk blocks: blocks may be scattered anywhere on the disk.

block =

pointer

Operating System Concepts – 8th Edition 11.17 Silberschatz, Galvin and Gagne ©2009

Linked Allocation (Cont.)

- Simple – need only starting address
- Free-space management system – no waste of space
- No random access
- Mapping

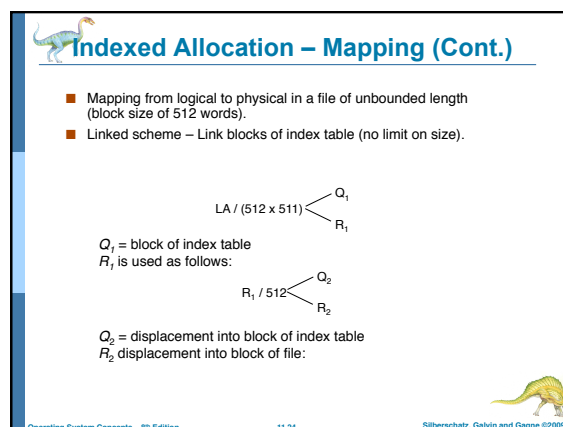
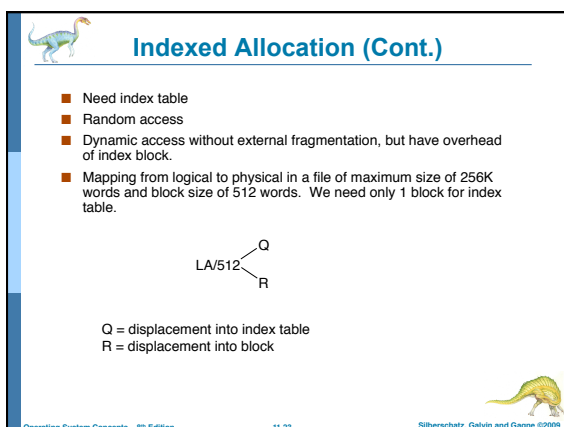
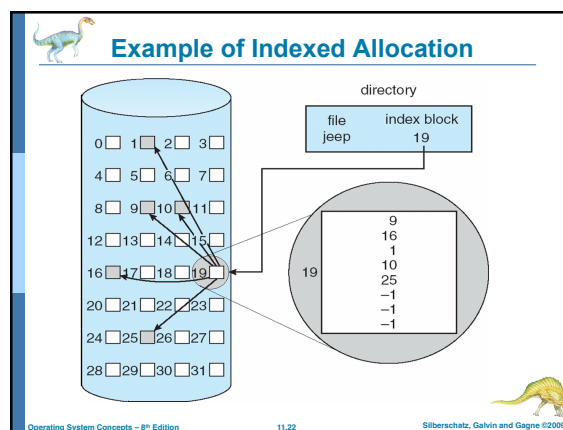
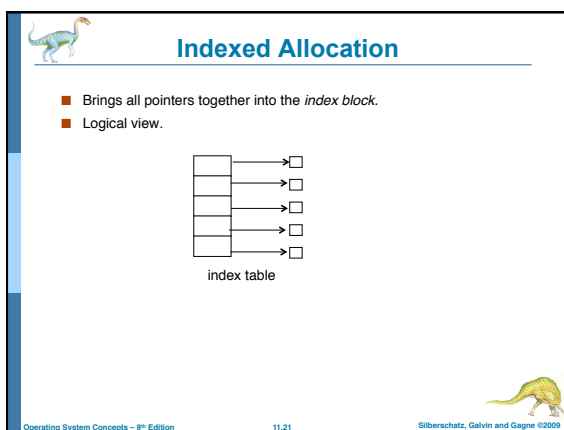
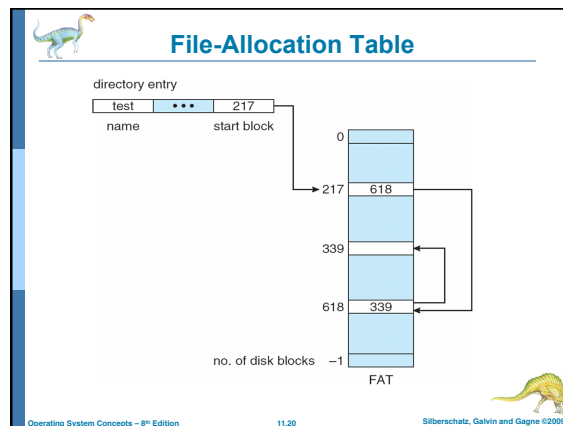
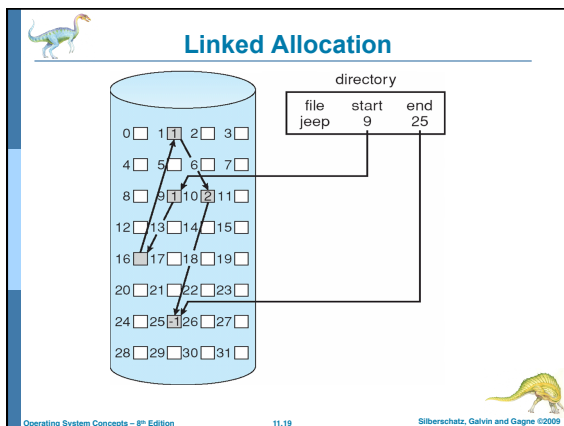
LA/511

Q
R

Block to be accessed is the Qth block in the linked chain of blocks representing the file.
Displacement into block = R + 1

File-allocation table (FAT) – disk-space allocation used by MS-DOS and OS/2.

Operating System Concepts – 8th Edition 11.18 Silberschatz, Galvin and Gagne ©2009



Indexed Allocation – Mapping (Cont.)

- Two-level index (maximum file size is 512^3)

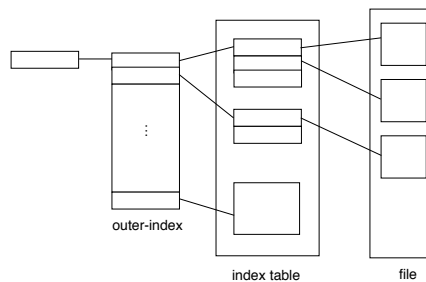
$$LA / (512 \times 512) \begin{cases} Q_1 \\ R_1 \end{cases}$$

Q_1 = displacement into outer-index
 R_1 is used as follows:

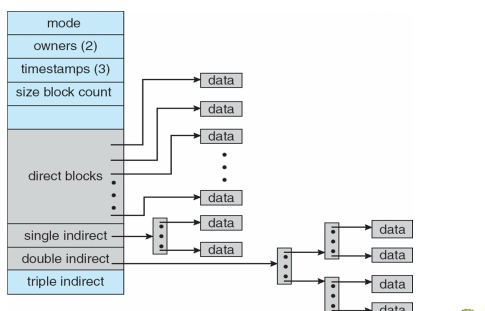
$$R_1 / 512 \begin{cases} Q_2 \\ R_2 \end{cases}$$

Q_2 = displacement into block of index table
 R_2 = displacement into block of file:

Indexed Allocation – Mapping (Cont.)



Combined Scheme: UNIX (4K bytes per block)



Free-Space Management

- Bit vector (n blocks)



$$\text{bit}[i] = \begin{cases} 0 \Rightarrow \text{block}[i] \text{ free} \\ 1 \Rightarrow \text{block}[i] \text{ occupied} \end{cases}$$

Block number calculation

$$(\text{number of bits per word}) * (\text{number of 0-value words}) + \text{offset of first 1 bit}$$

Free-Space Management (Cont.)

- Bit map requires extra space
 - Example:
 - block size = 2^{12} bytes
 - disk size = 2^{30} bytes (1 gigabyte)
 - $n = 2^{30}/2^{12} = 2^{18}$ bits (or 32K bytes)
- Easy to get contiguous files
- Linked list (free list)
 - Cannot get contiguous space easily
 - No waste of space
- Grouping
- Counting

Free-Space Management (Cont.)

- Need to protect:
 - Pointer to free list
 - Bit map
 - Must be kept on disk
 - Copy in memory and disk may differ
 - Cannot allow for block[i] to have a situation where bit[i] = 1 in memory and bit[i] = 0 on disk
- Solution:
 - Set bit[i] = 1 in disk
 - Allocate block[i]
 - Set bit[i] = 1 in memory

Directory Implementation

- Linear list of file names with pointer to the data blocks
 - simple to program
 - time-consuming to execute
- Hash Table – linear list with hash data structure
 - decreases directory search time
 - **collisions** – situations where two file names hash to the same location
 - fixed size

Operating System Concepts – 8th Edition 11.31 Silberschatz, Galvin and Gagne ©2009

Linked Free Space List on Disk

Operating System Concepts – 8th Edition 11.32 Silberschatz, Galvin and Gagne ©2009

Efficiency and Performance

- Efficiency dependent on:
 - disk allocation and directory algorithms
 - types of data kept in file's directory entry
- Performance
 - disk cache – separate section of main memory for frequently used blocks
 - free-behind and read-ahead – techniques to optimize sequential access
 - improve PC performance by dedicating section of memory as virtual disk, or RAM disk

Operating System Concepts – 8th Edition 11.33 Silberschatz, Galvin and Gagne ©2009

Page Cache

- A **page cache** caches pages rather than disk blocks using virtual memory techniques
- Memory-mapped I/O uses a page cache
- Routine I/O through the file system uses the buffer (disk) cache
- This leads to the following figure

Operating System Concepts – 8th Edition 11.34 Silberschatz, Galvin and Gagne ©2009

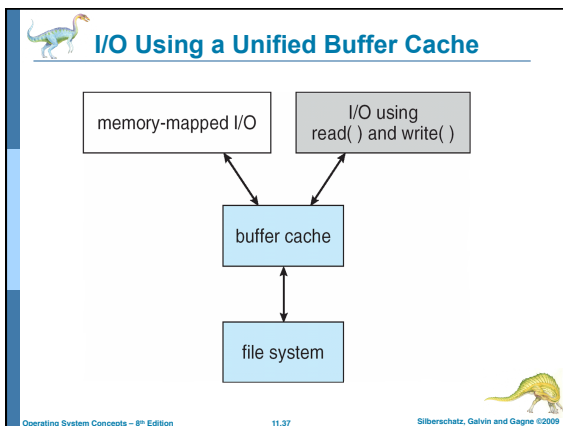
I/O Without a Unified Buffer Cache

Operating System Concepts – 8th Edition 11.35 Silberschatz, Galvin and Gagne ©2009

Unified Buffer Cache

- A unified buffer cache uses the same page cache to cache both memory-mapped pages and ordinary file system I/O

Operating System Concepts – 8th Edition 11.36 Silberschatz, Galvin and Gagne ©2009



Recovery

- Consistency checking – compares data in directory structure with data blocks on disk, and tries to fix inconsistencies
- Use system programs to **back up** data from disk to another storage device (floppy disk, magnetic tape, other magnetic disk, optical)
- Recover lost file or disk by **restoring** data from backup

Operating System Concepts – 8th Edition 11.38 Silberschatz, Galvin and Gagne ©2009

Log Structured File Systems

- Log structured** (or **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

Operating System Concepts – 8th Edition 11.39 Silberschatz, Galvin and Gagne ©2009

The Sun Network File System (NFS)

- An implementation and a specification of a software system for accessing remote files across LANs (or WANs)
- The implementation is part of the Solaris and SunOS operating systems running on Sun workstations using an unreliable datagram protocol (UDP/IP) protocol and Ethernet

Operating System Concepts – 8th Edition 11.40 Silberschatz, Galvin and Gagne ©2009

NFS (Cont.)

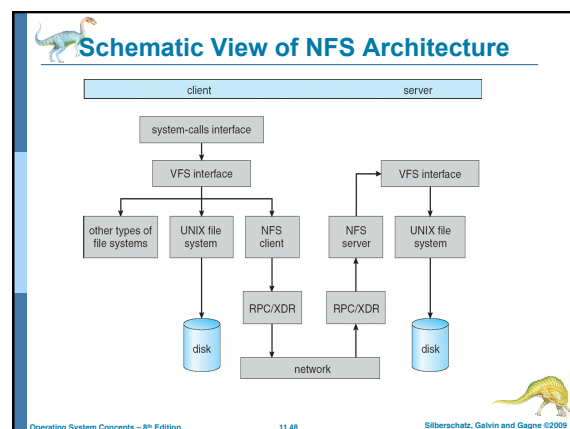
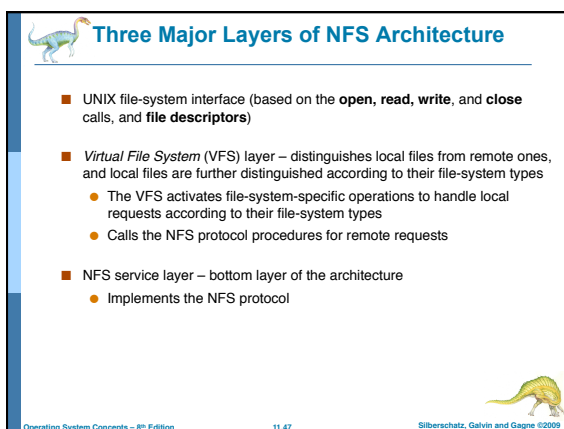
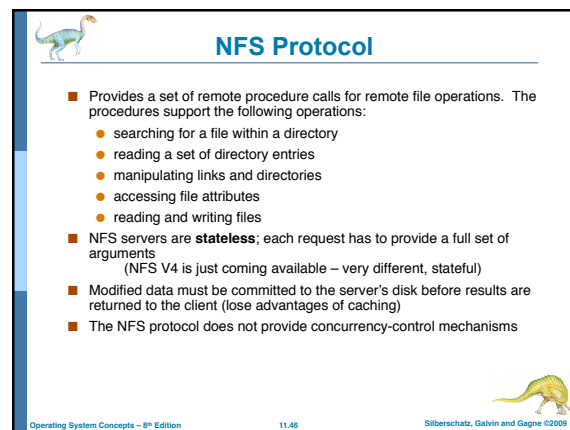
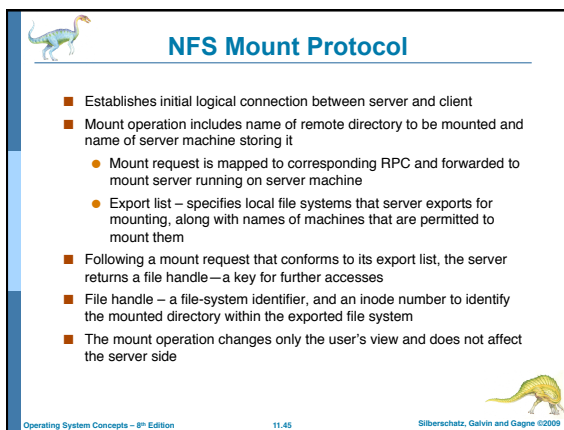
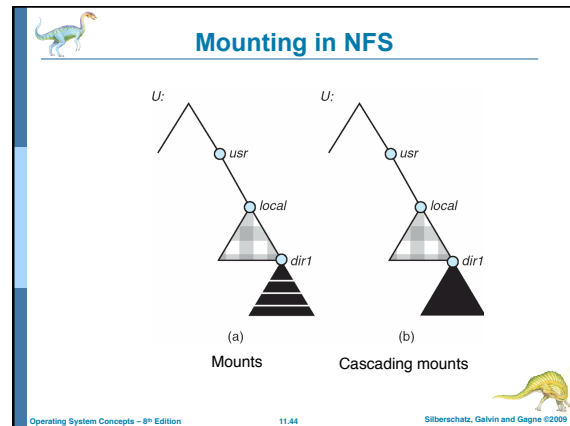
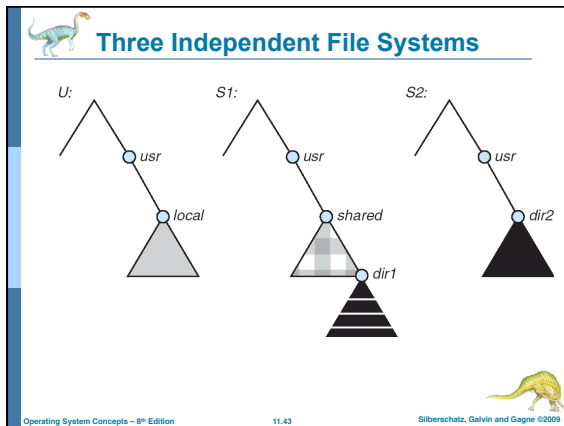
- Interconnected workstations viewed as a set of independent machines with independent file systems, which allows sharing among these file systems in a transparent manner
 - A remote directory is mounted over a local file system directory
 - The mounted directory looks like an integral subtree of the local file system, replacing the subtree descending from the local directory
 - Specification of the remote directory for the mount operation is nontransparent; the host name of the remote directory has to be provided
 - Files in the remote directory can then be accessed in a transparent manner
 - Subject to access-rights accreditation, potentially any file system (or directory within a file system), can be mounted remotely on top of any local directory

Operating System Concepts – 8th Edition 11.41 Silberschatz, Galvin and Gagne ©2009

NFS (Cont.)

- NFS is designed to operate in a heterogeneous environment of different machines, operating systems, and network architectures; the NFS specifications independent of these media
- This independence is achieved through the use of RPC primitives built on top of an External Data Representation (XDR) protocol used between two implementation-independent interfaces
- The NFS specification distinguishes between the services provided by a mount mechanism and the actual remote-file-access services

Operating System Concepts – 8th Edition 11.42 Silberschatz, Galvin and Gagne ©2009





NFS Path-Name Translation

- Performed by breaking the path into component names and performing a separate NFS lookup call for every pair of component name and directory vnode
- To make lookup faster, a directory name lookup cache on the client's side holds the vnodes for remote directory names

Operating System Concepts – 8th Edition

11.49

Silberschatz, Galvin and Gagne ©2009



NFS Remote Operations

- Nearly one-to-one correspondence between regular UNIX system calls and the NFS protocol RPCs (except opening and closing files)
- NFS adheres to the remote-service paradigm, but employs buffering and caching techniques for the sake of performance
- File-blocks cache – when a file is opened, the kernel checks with the remote server whether to fetch or revalidate the cached attributes
 - Cached file blocks are used only if the corresponding cached attributes are up to date
- File-attribute cache – the attribute cache is updated whenever new attributes arrive from the server
- Clients do not free delayed-write blocks until the server confirms that the data have been written to disk

Operating System Concepts – 8th Edition

11.50

Silberschatz, Galvin and Gagne ©2009



Example: WAFL File System

- Used on Network Appliance "Filers" – distributed file system appliances
- "Write-anywhere file layout"
- Serves up NFS, CIFS, http, ftp
- Random I/O optimized, write optimized
 - NVRAM for write caching
- Similar to Berkeley Fast File System, with extensive modifications

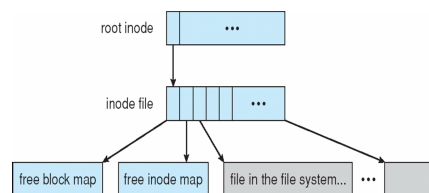
Operating System Concepts – 8th Edition

11.51

Silberschatz, Galvin and Gagne ©2009



The WAFL File Layout



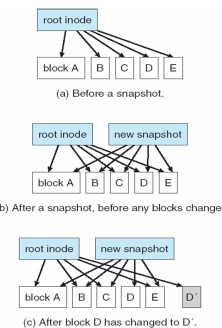
Operating System Concepts – 8th Edition

11.52

Silberschatz, Galvin and Gagne ©2009



Snapshots in WAFL



Operating System Concepts – 8th Edition

11.53

Silberschatz, Galvin and Gagne ©2009



11.02

file permissions
file dates (create, access, write)
file owner, group, ACL
file size
file data blocks or pointers to file data blocks

Operating System Concepts – 8th Edition

11.54

Silberschatz, Galvin and Gagne ©2009



End of Chapter 11

