

Chapter 11: File System Implementation

Prof. Li-Pin Chang
CS@NYCU

Chapter 11: File System Implementation

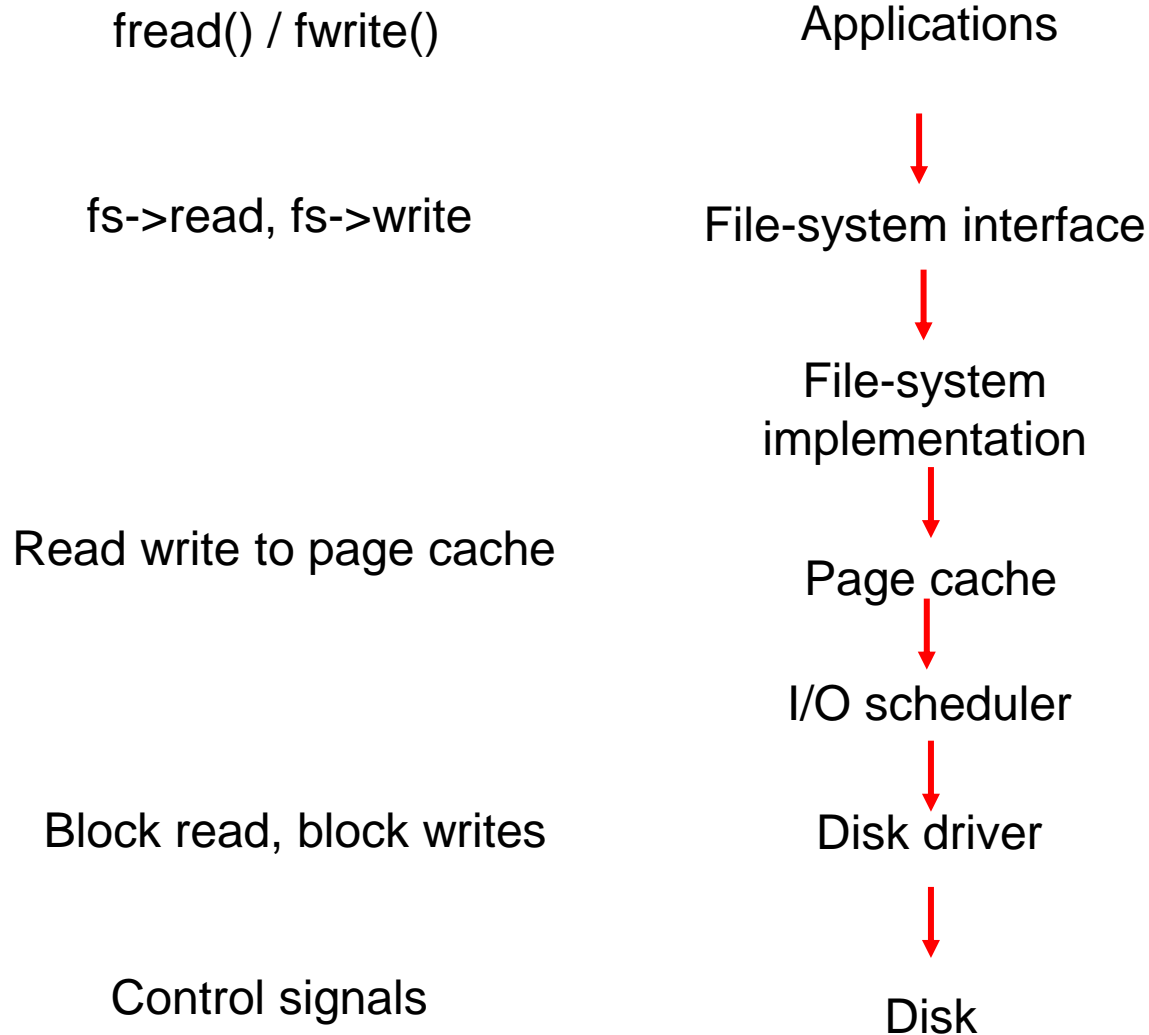
- File-System Structure
- Directory Implementation
- Allocation Methods
- Free-Space Management
- Efficiency and Performance
- Recovery
- Log-Structured File Systems

Objectives

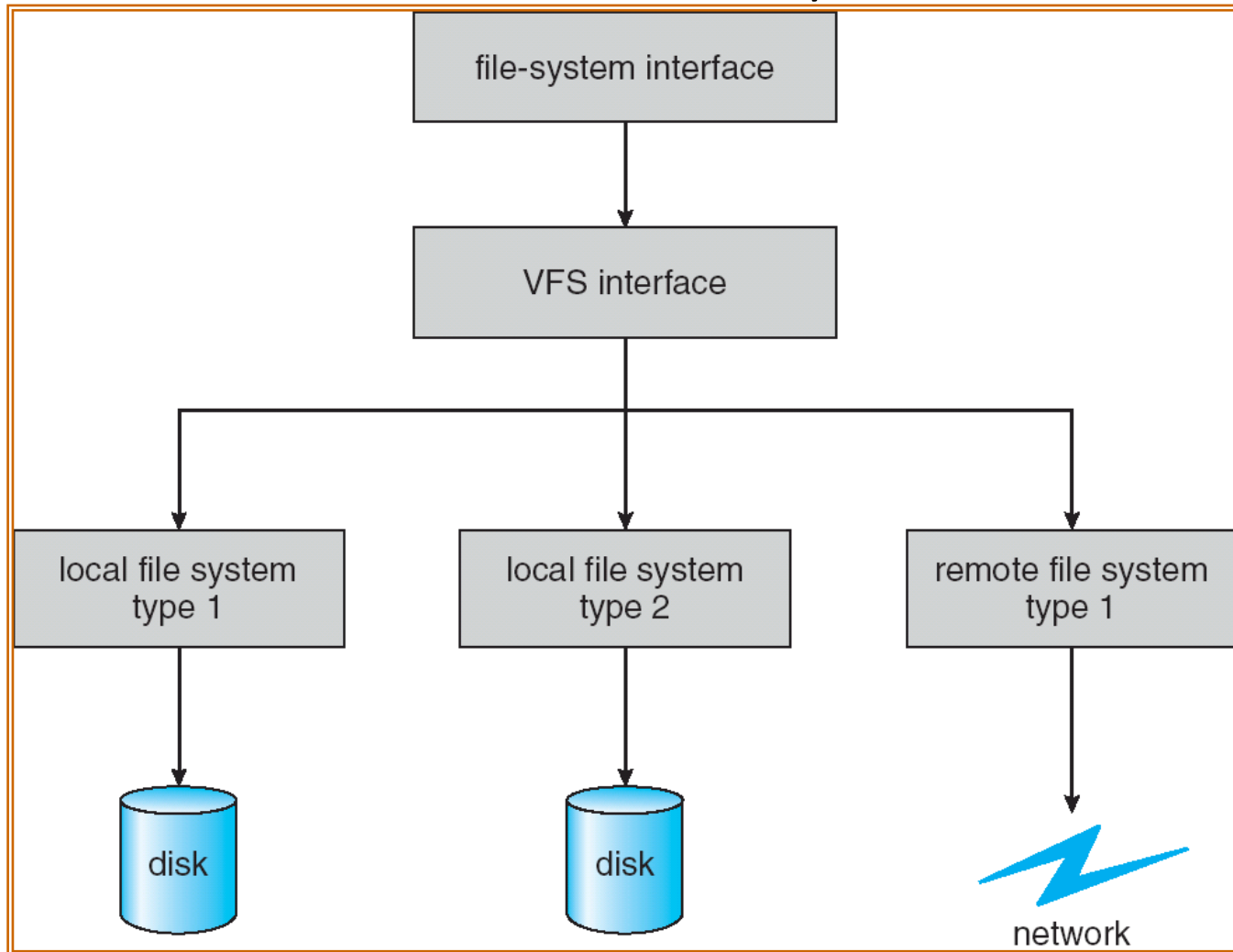
- To describe the details of implementing local file systems and directory structures
- To discuss block allocation and free-block algorithms and trade-offs

File System Structure and Abstraction

Layered File System

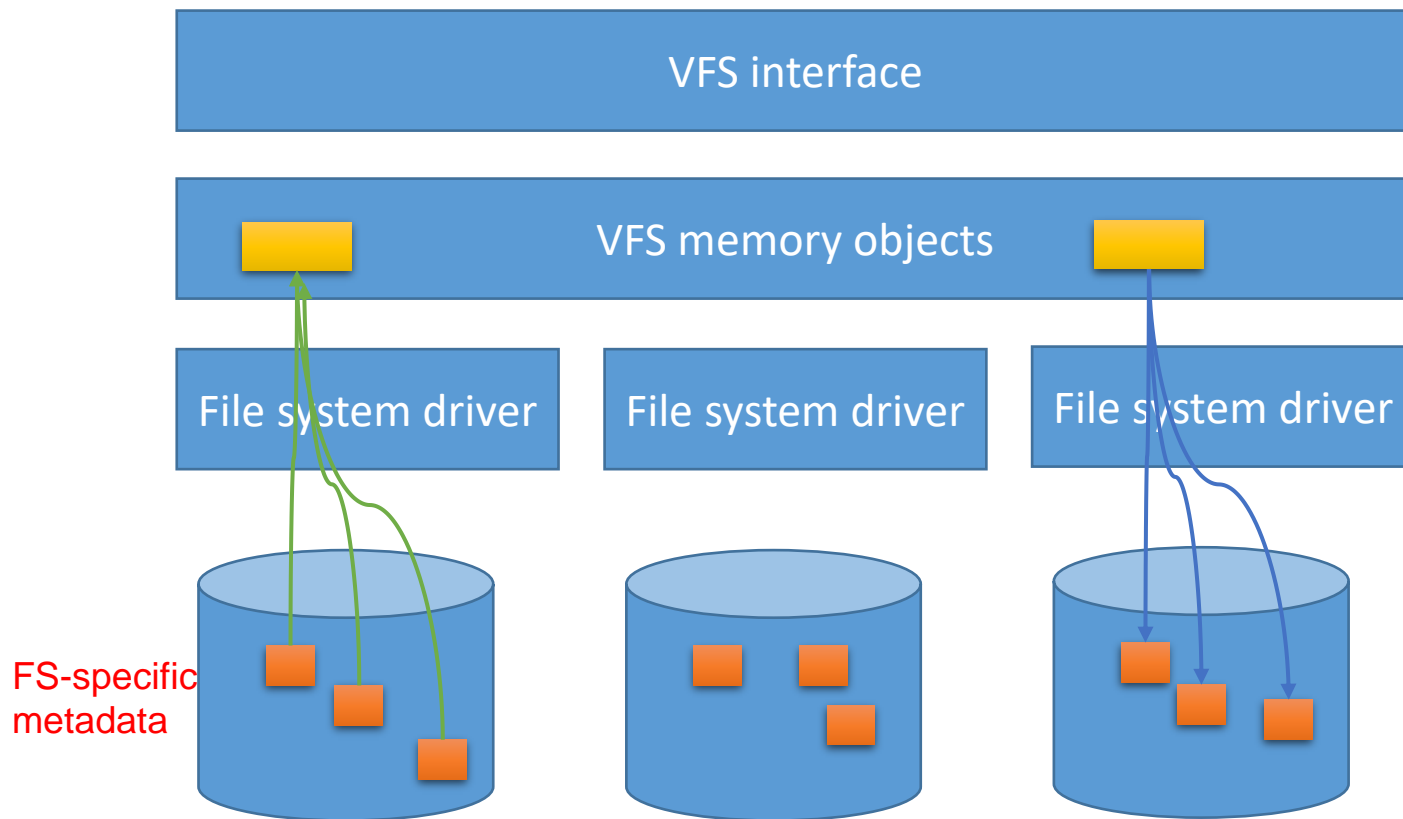


Schematic View of Virtual File System



Linux Virtual File System Architecture

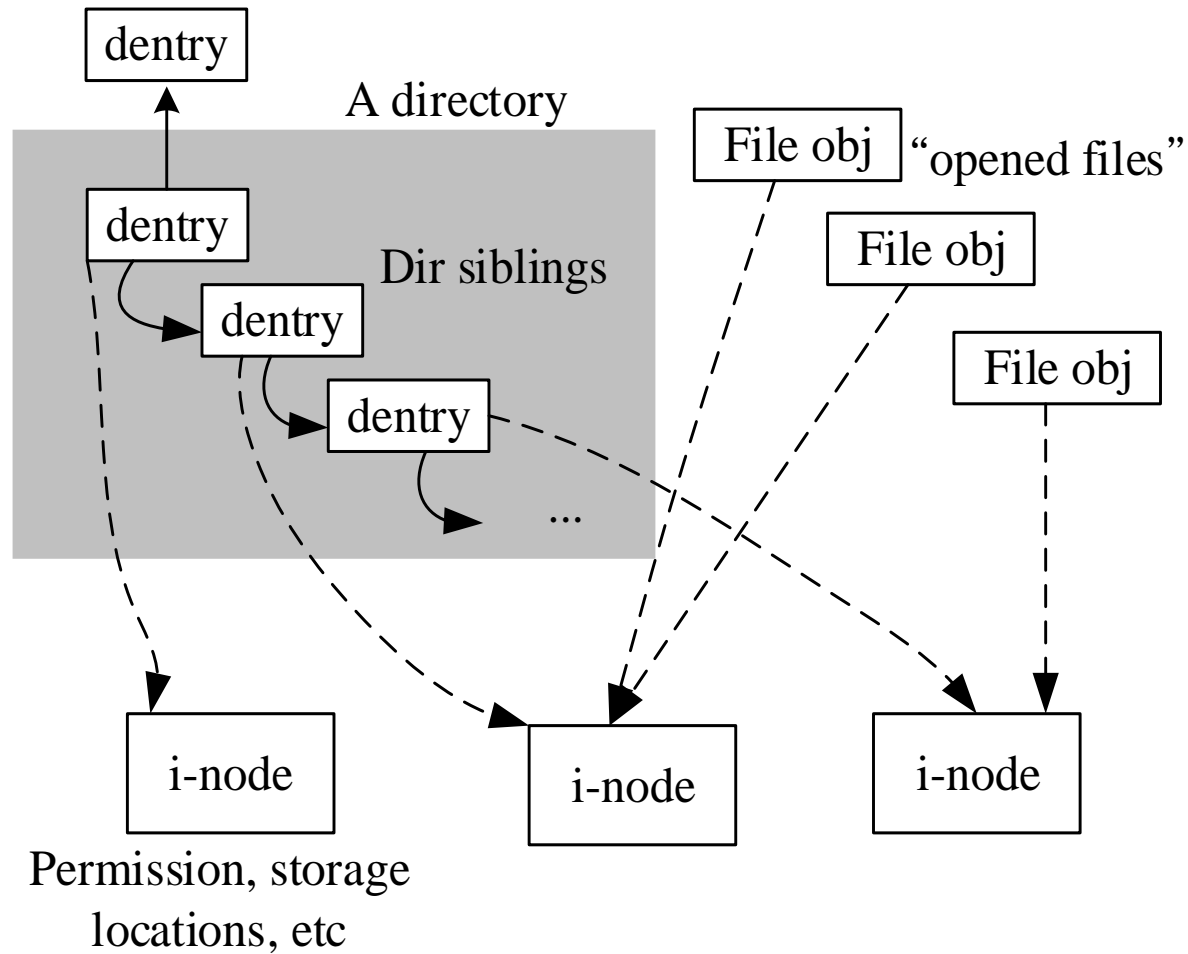
- File system drivers translate between kernel VFS memory objects and disk metadata



In-memory Kernel Objects of Linux VFS

- Superblock
 - Representing the entire filesystem
- Inode
 - Uniquely representing an individual file
- File object
 - Representing an opened file, one for each fopen instance
- Dentry object
 - Representing an individual directory entry

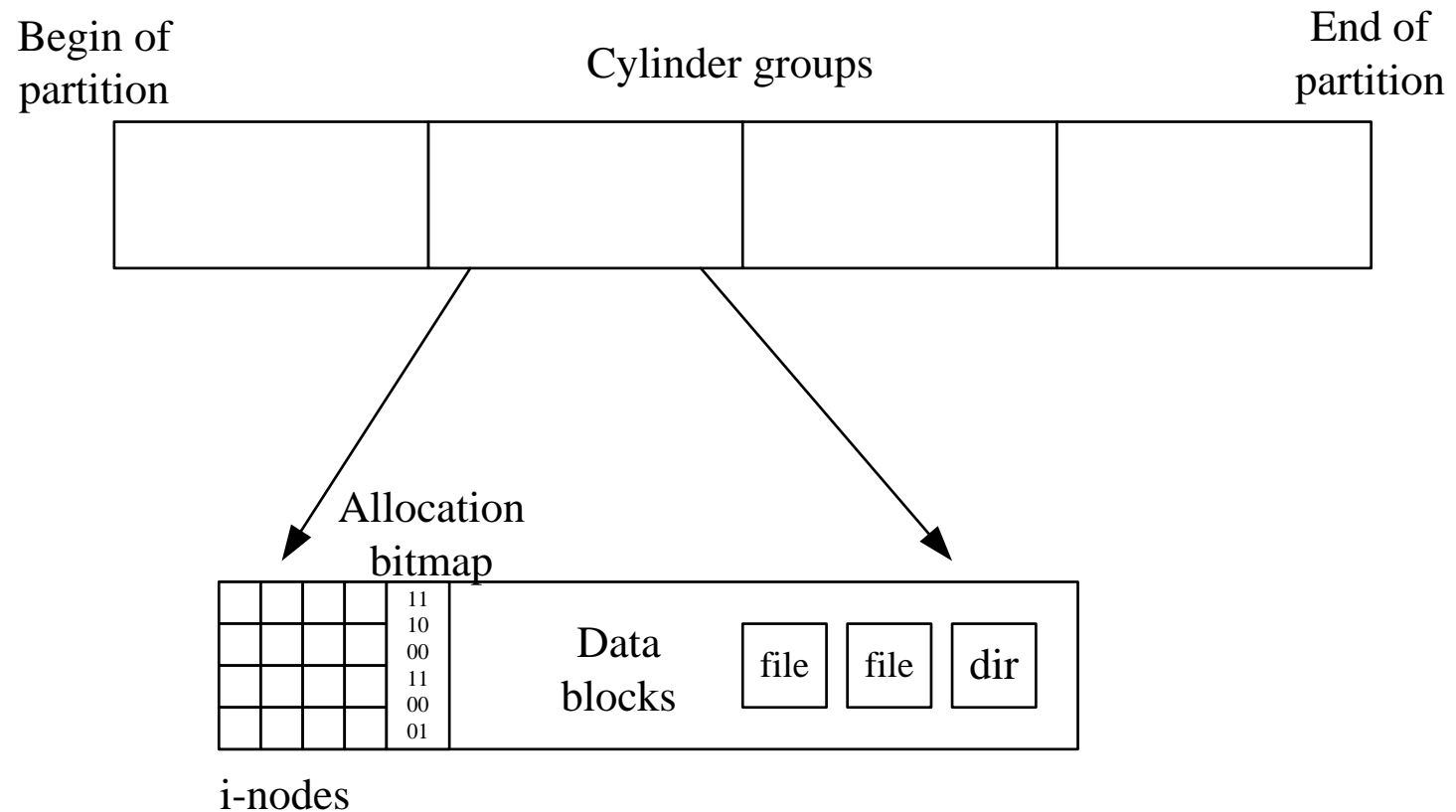
In-memory objects of Linux VFS



Disk Metadata

- File-system-specific; vary from file system to file system
- Linux ext file system
 - Super block, Inodes, Allocation bitmaps
- Microsoft FAT file system
 - File allocation tables, Directories
- File system driver must fill the in-memory objects with the information in disk metadata
 - May not be one-to-one mapped, e.g.,
Ext file system has i-node on disk; FAT file system does not

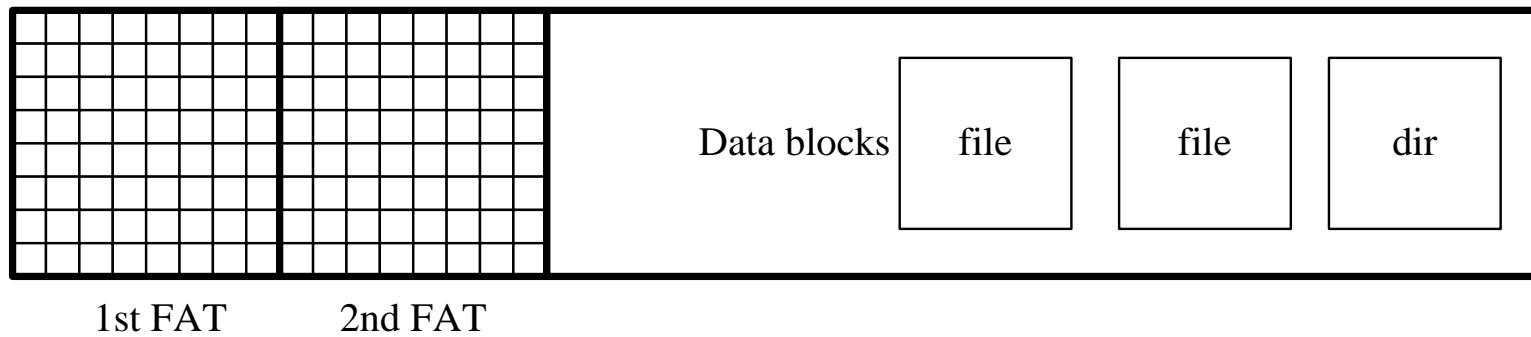
Disk Layout of the Linux ext 2/3/4 file systems



Disk layout of FAT 12/16/32 file systems

Begin of
partition

End of
partition



File System Key Design Issues

Key Design Issues

1. Directory implementation
2. Allocation (index) methods
3. Free-space management

Issue 1: Directory Implementation

- **Linear list** of file names with pointer to the data blocks.
 - simple design
 - time-consuming operations
 - FAT file system
- **B-trees (or variants)**
 - Efficient search
 - XFS, NTFS, ext4 (H-tree, fixed 2 levels) [[ref](#)]
 - Scaling well for large directories

Example: Directory Dump in FAT

Offset	00	01	02	03	04	05	06	07	08	09	0A	0B	0C	0D	0E	0F	0123456789ABCDEF
000167936	41	6D	00	79	00	64	00	69	00	72	00	0F	00	E6	32	00	Am.y.d.i.r....2.
000167952	00	00	FF	FF	FF	FF	FF	FF	FF	FF	00	00	FF	FF	FF	FF
000167968	4D	59	44	49	52	32	20	20	20	20	10	00	00	90	B1		MYDIR2
000167984	A6	42	A6	42	00	00	90	B1	A6	42	04	00	00	00	00	00	.B.B.....B.....
000168000	41	6D	00	79	00	64	00	69	00	72	00	0F	00	DE	31	00	Am.y.d.i.r....1.
000168016	00	00	FF	FF	FF	FF	FF	FF	FF	FF	00	00	FF	FF	FF	FF
000168032	4D	59	44	49	52	31	20	20	20	20	10	00	64	6A	B1		MYDIR1 ..dj.
000168048	A6	42	A6	42	00	00	6A	B1	A6	42	03	00	00	00	00	00	.B.B..j..B.....
000168064	41	6D	00	79	00	66	00	69	00	6C	00	0F	00	8B	65	00	Am.y.f.i.l....e.
000168080	31	00	2E	00	74	00	78	00	74	00	00	00	00	00	FF	FF	1...t.x.t.....
000168096	4D	59	46	49	4C	45	31	20	54	58	54	20	00	64	99	B1	MYFILE1 TXT .d..
000168112	A6	42	A6	42	00	00	99	B1	A6	42	05	00	0F	00	00	00	.B.B.....B.....
000168128	E5	6D	00	79	00	66	00	69	00	6C	00	0F	00	5B	65	00	.m.y.f.i.l...[e.
000168144	32	00	2E	00	74	00	78	00	74	00	00	00	00	00	FF	FF	2...t.x.t.....
000168160	E5	59	46	49	4C	45	32	20	54	58	54	20	00	64	77	8B	.YFILE2 TXT .dw.
000168176	A7	42	A6	42	00	00	77	8B	A7	42	07	00	22	20	09	00	.B.B..w..B.." ..
000168192	41	6C	00	64	00	65	00	5F	00	32	00	0F	00	5D	36	00	Al.d.e._.2...]6.
000168208	31	00	2E	00	74	00	67	00	7A	00	00	00	00	00	FF	FF	1...t.g.z.....
000168224	4C	44	45	5F	32	36	31	20	54	47	5A	20	00	64	77	8B	LDE_261 TGZ .dw.
000168240	A7	42	A6	42	00	00	77	8B	A7	42	07	00	22	20	09	00	.B.B..w..B.." ..

Issue 2: Allocation/Index Methods

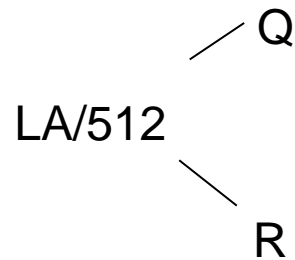
- An allocation method refers to how disk blocks are allocated for files:
 - Contiguous allocation
 - Linked allocation
 - Indexed allocation
 - Extent-based allocation

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
- Files cannot grow beyond the allocated space, unless files are migrated to larger spaces
- Efficient access; perfect for I/O overhead reduction
 - Less I/Os involved
 - Sequential disk operations
- Wasteful of space (dynamic storage-allocation problem)
 - File deletion leaves free holes (external fragmentation)
 - Needs compaction, maybe done in background or downtime

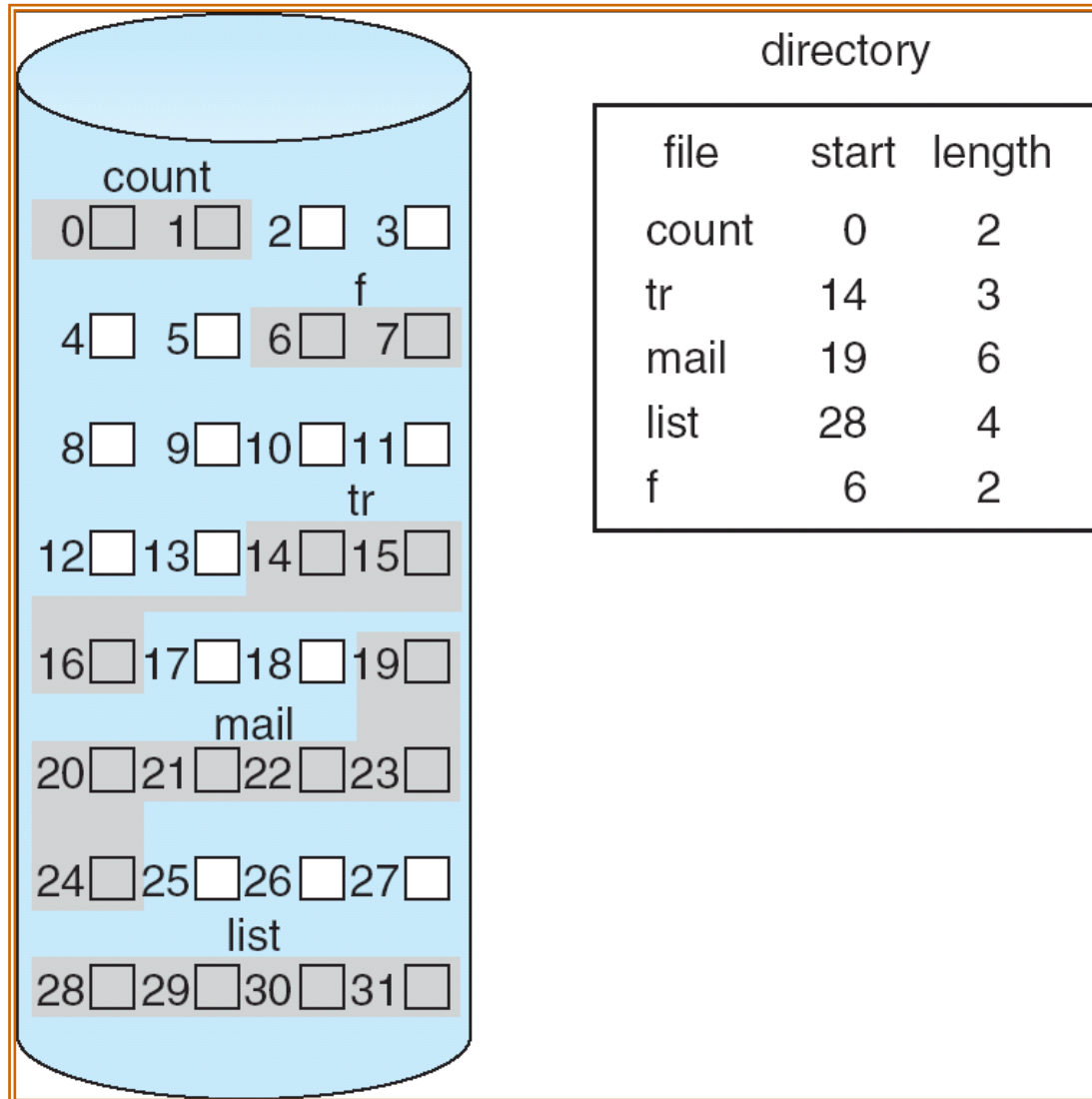
Contiguous Allocation

- Mapping from logical to physical
- LA = file offset (bytes); 1 disk block = 512 bytes



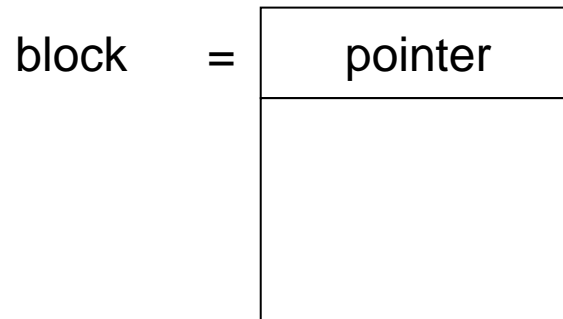
- Block to be accessed = Q + starting address (block)
- Displacement into block = R

Contiguous Allocation of Disk Space



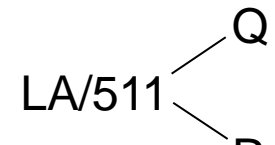
Linked Allocation

- Each file is a linked list of disk blocks
- Physical contiguity of the disk blocks is not absolutely necessary because file data are copied to sequential memory before use

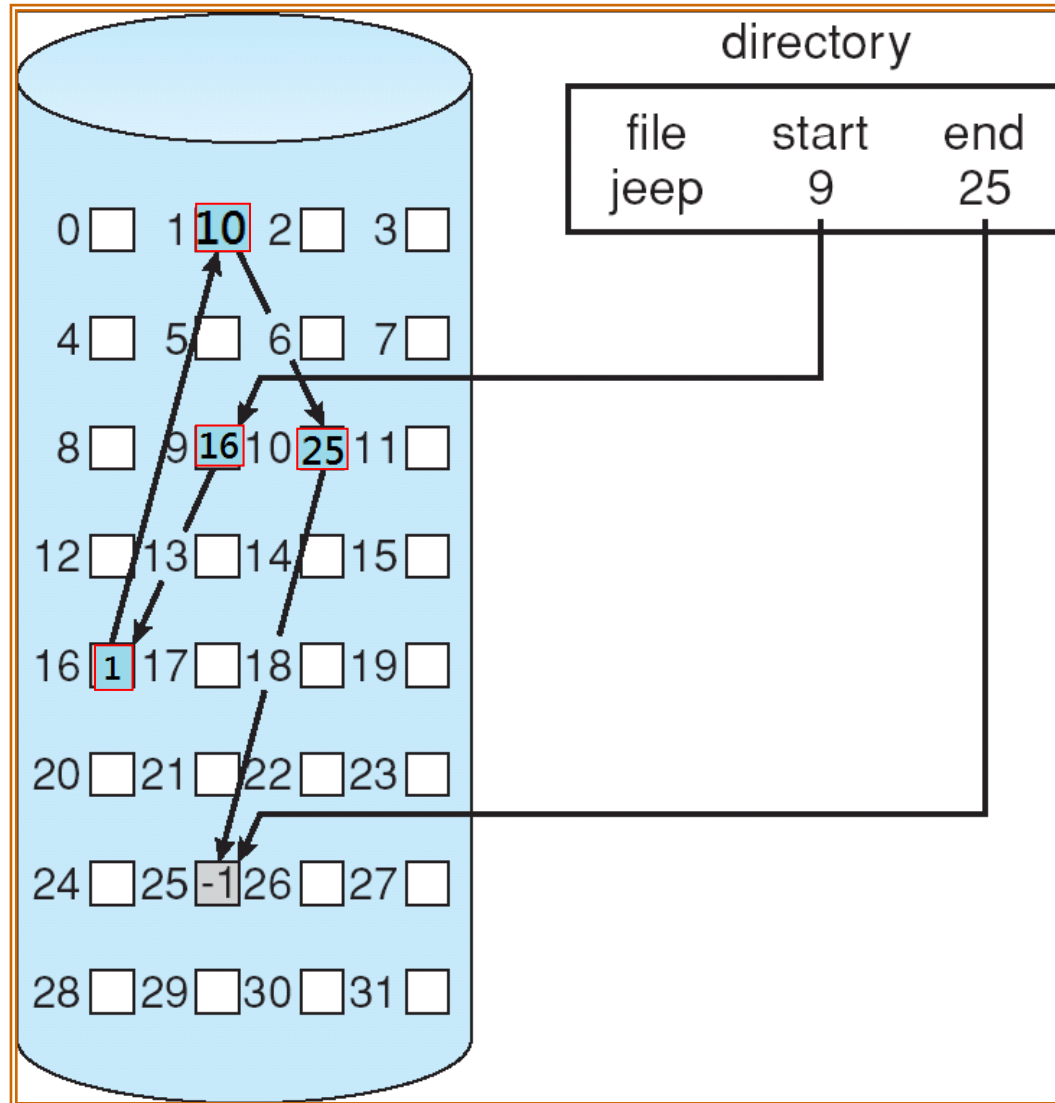


Linked Allocation (Cont.)

- Simple – need only starting address
- Free-space management system
 - no waste of space (no external fragmentation)
 - However, no random access (need to traverse the linked blocks)
- Mapping
 - Embedded pointers, e.g., 511B of data and 1B for ptr
 - Block to be accessed = the Qth block in the file's linked list



Linked Allocation



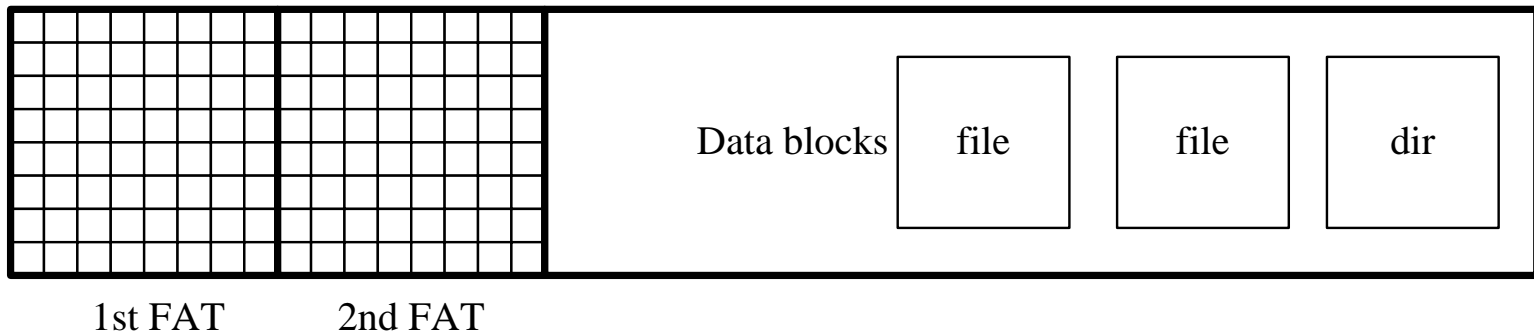
Linked Allocation

- Separating pointers from data blocks
 - Making data size a power of 2; easier to manage
- Example: FAT file system

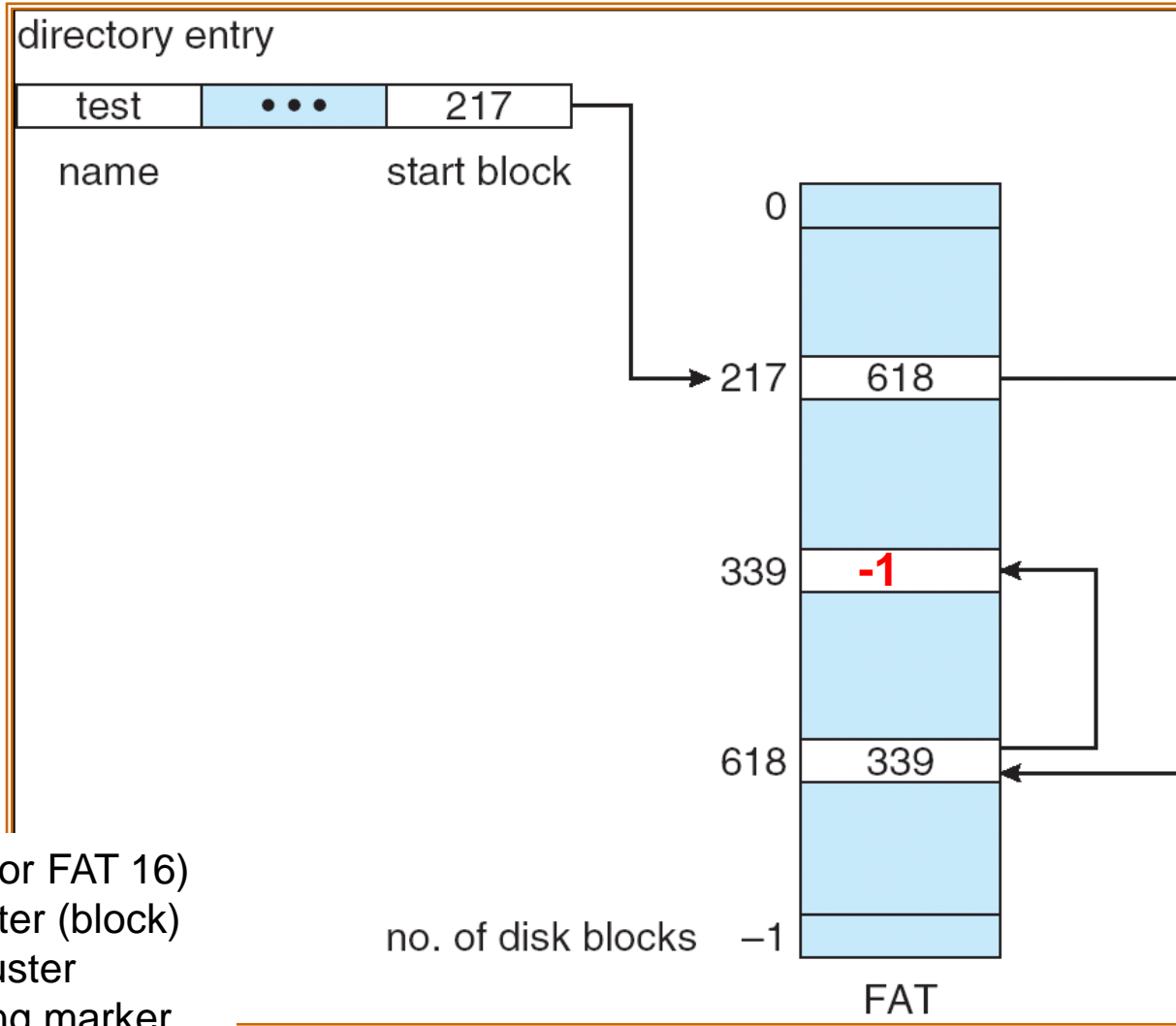
The layout of FAT 12/16/32 file system

Begin of
partition

End of
partition

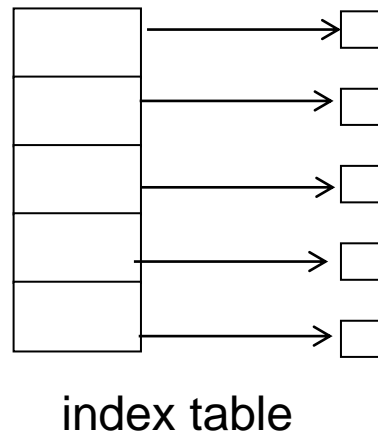


File-Allocation Table

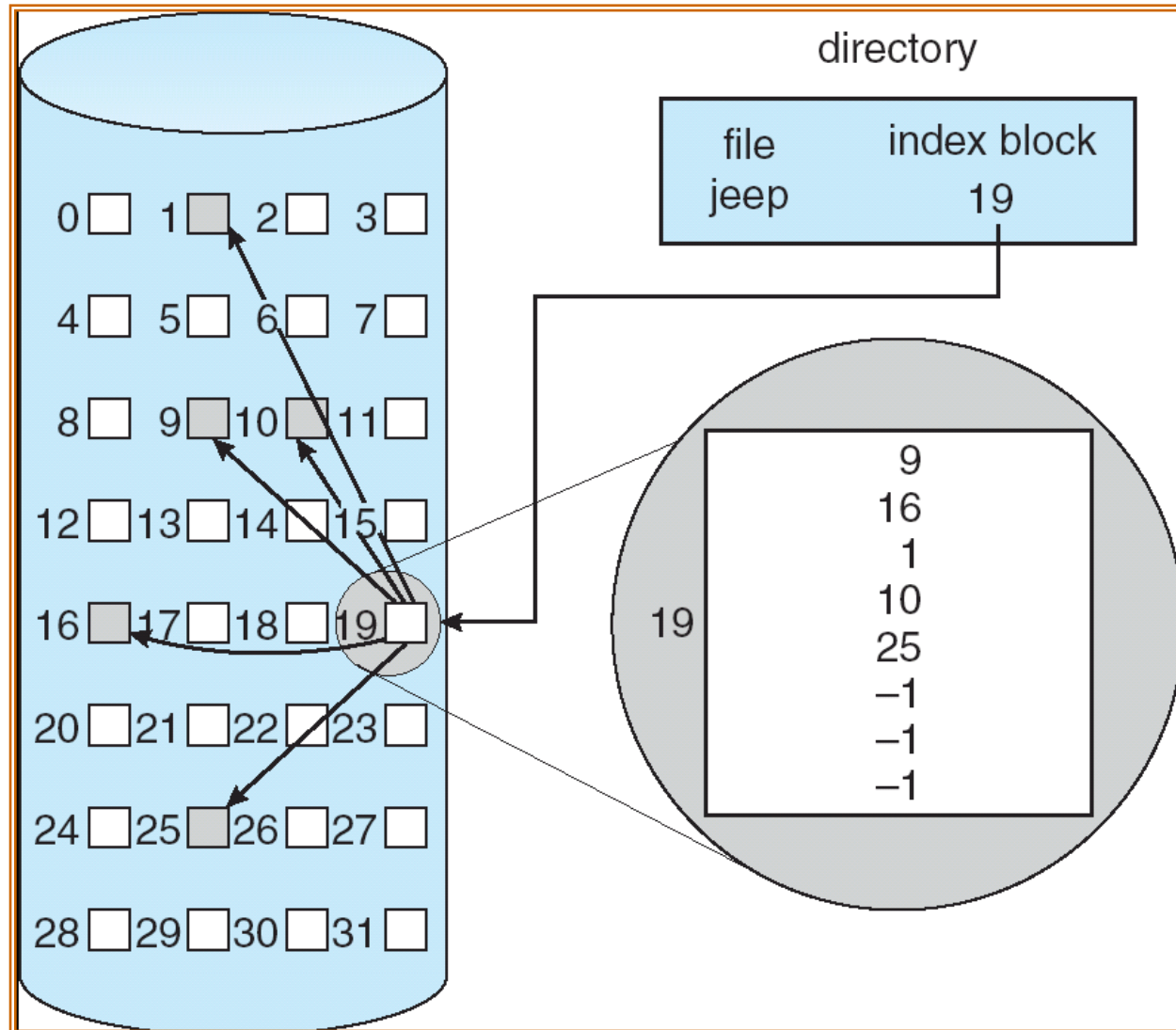


Indexed Allocation

- Brings all pointers together into the index block.
- Logical view.

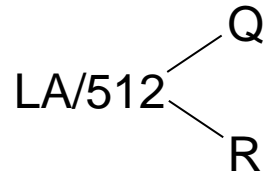


Example of Indexed Allocation



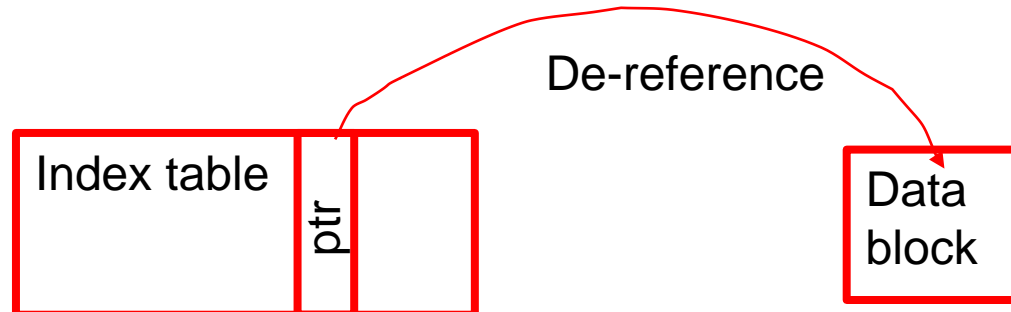
Indexed Allocation (Cont.)

- Need a index table
- Capable of “random” access; no list traversing
- Per-file overhead of an index table (block)



Q = displacement into index table (entry #)

R = displacement into the referred block



Indexed Allocation – Mapping

- Assuming two-level index

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

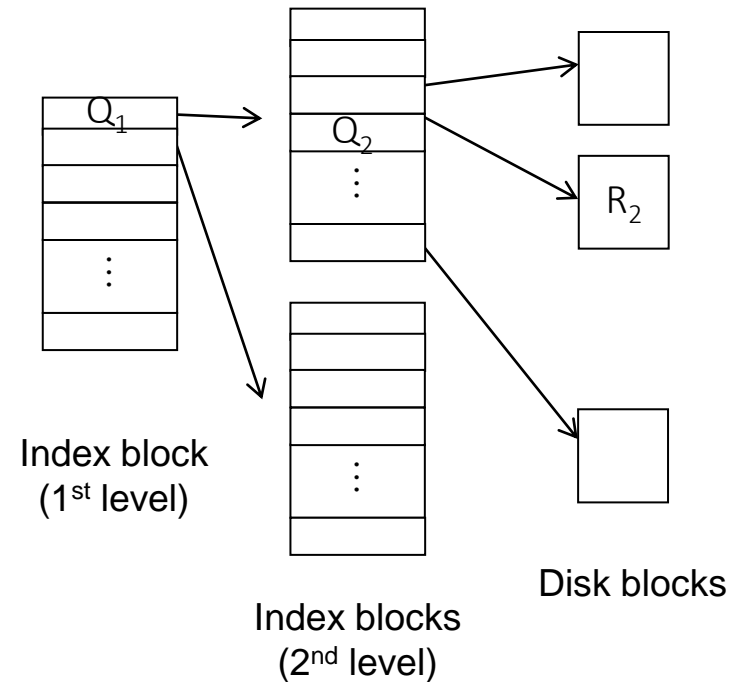
Q_1 = displacement into outer-index

R_1 is used as follows:

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

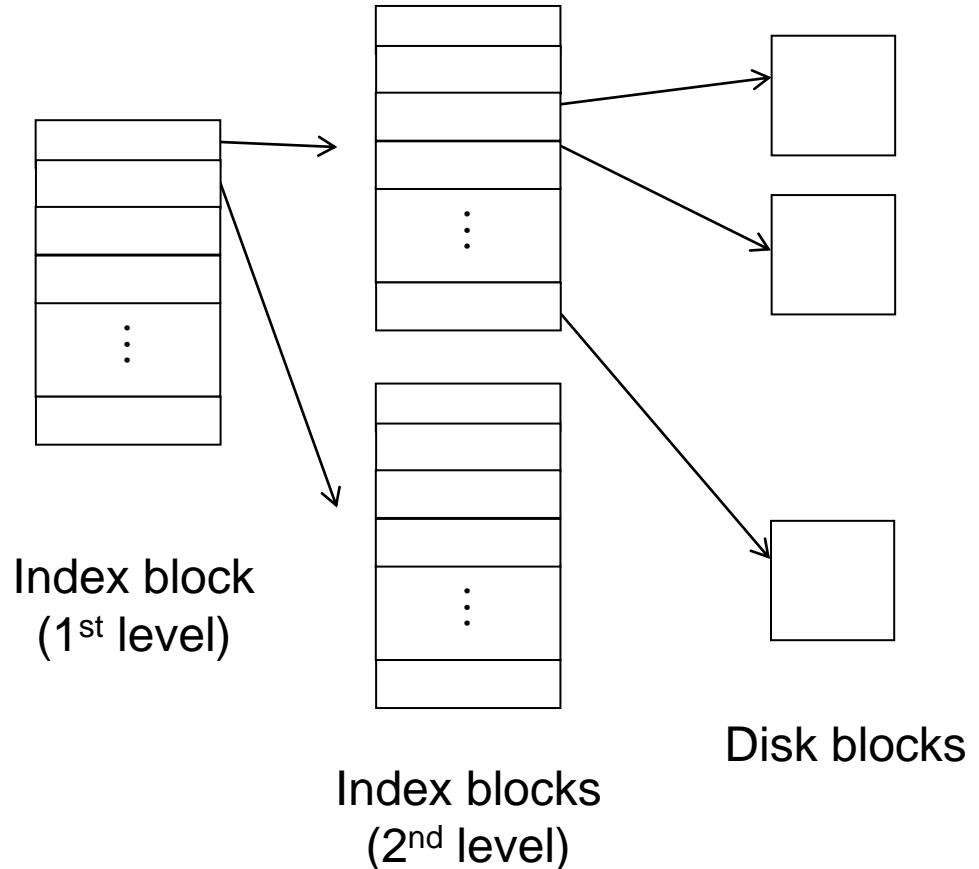
Q_2 = displacement into block of index table

R_2 displacement into block of file

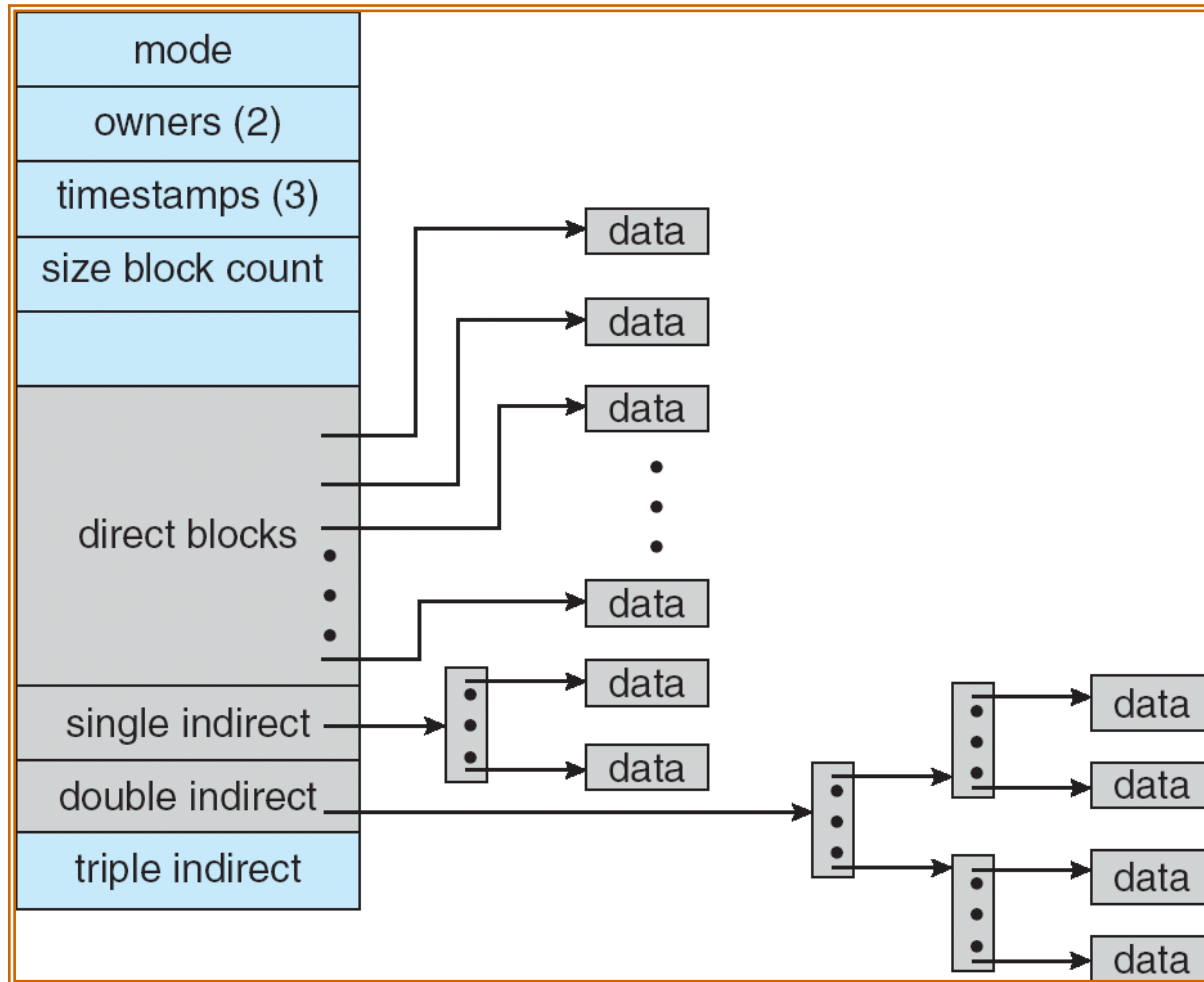


Indexed Allocation – Mapping (Cont.)

- 1block=512B, 1ptr=1B
- 1 idx. block has 512 ptrs
- 1st level: pointers to index tables
- 2nd level: pointers to data blocks
- Max. file size = $512 * 512 * 512$ bytes
- Isn't it similar to two-level page tables?



Example: UNIX inode



Small files use only embedded pointers to direct blocks

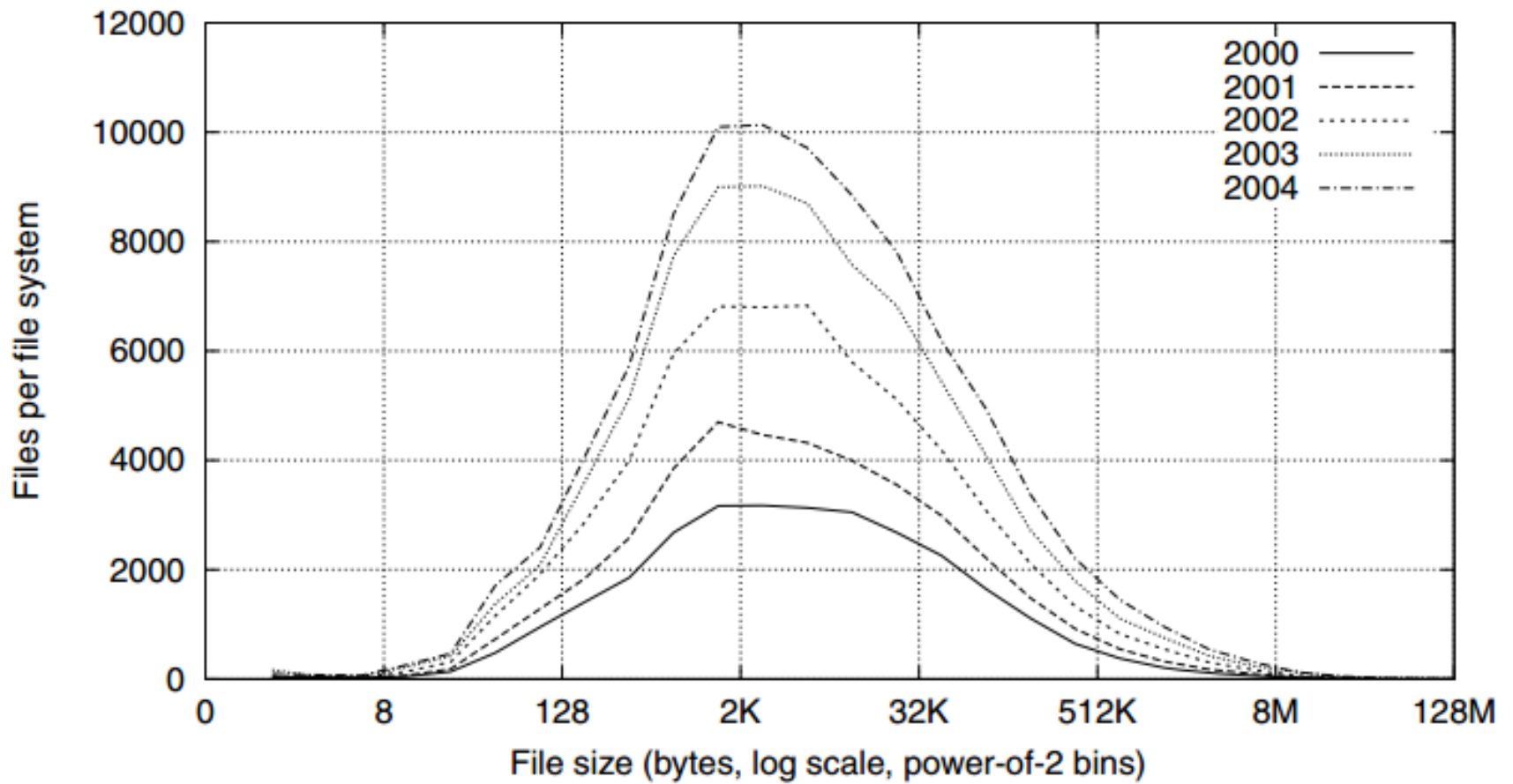
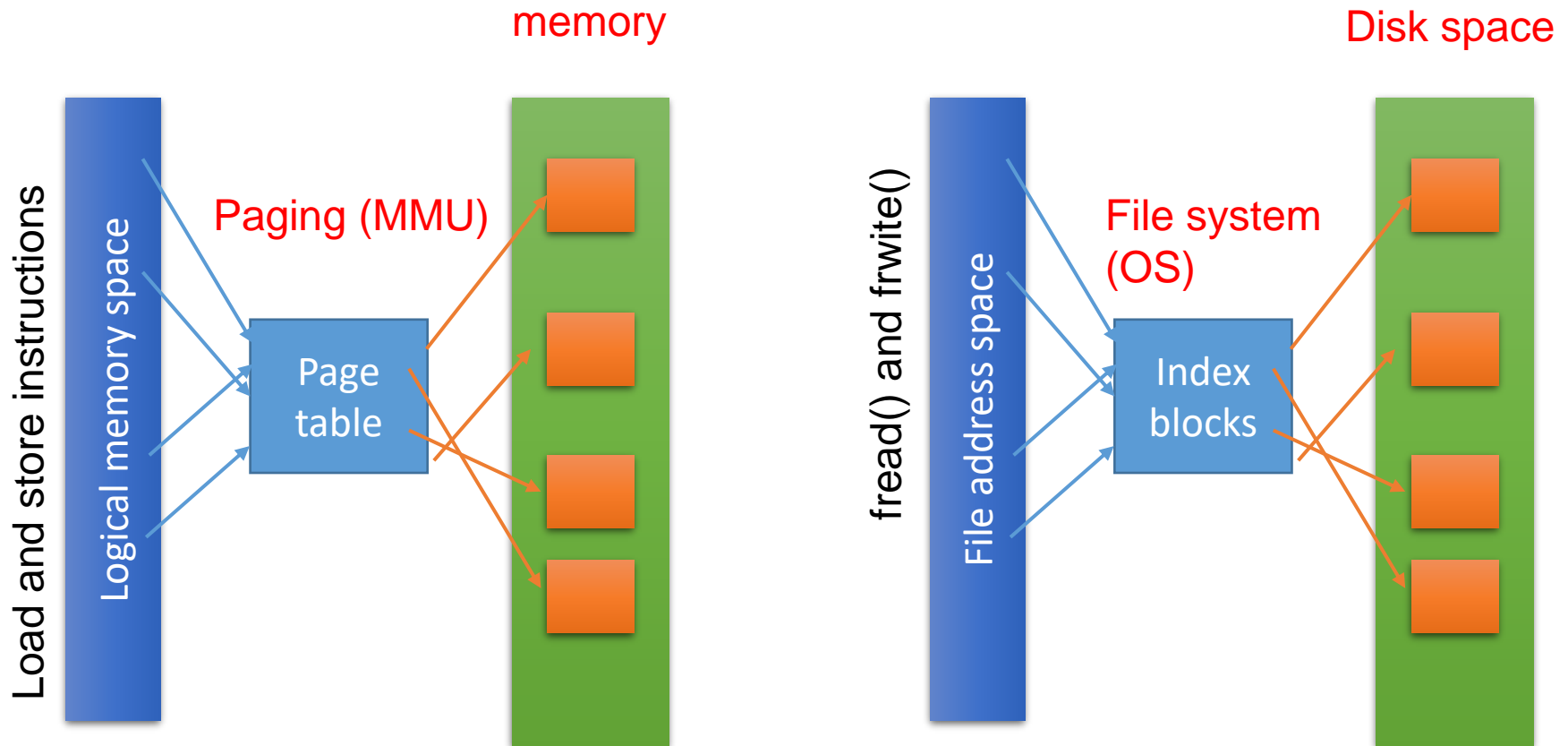


Fig. 2. Histograms of files by size.

Indirection, indirection, indirection ...

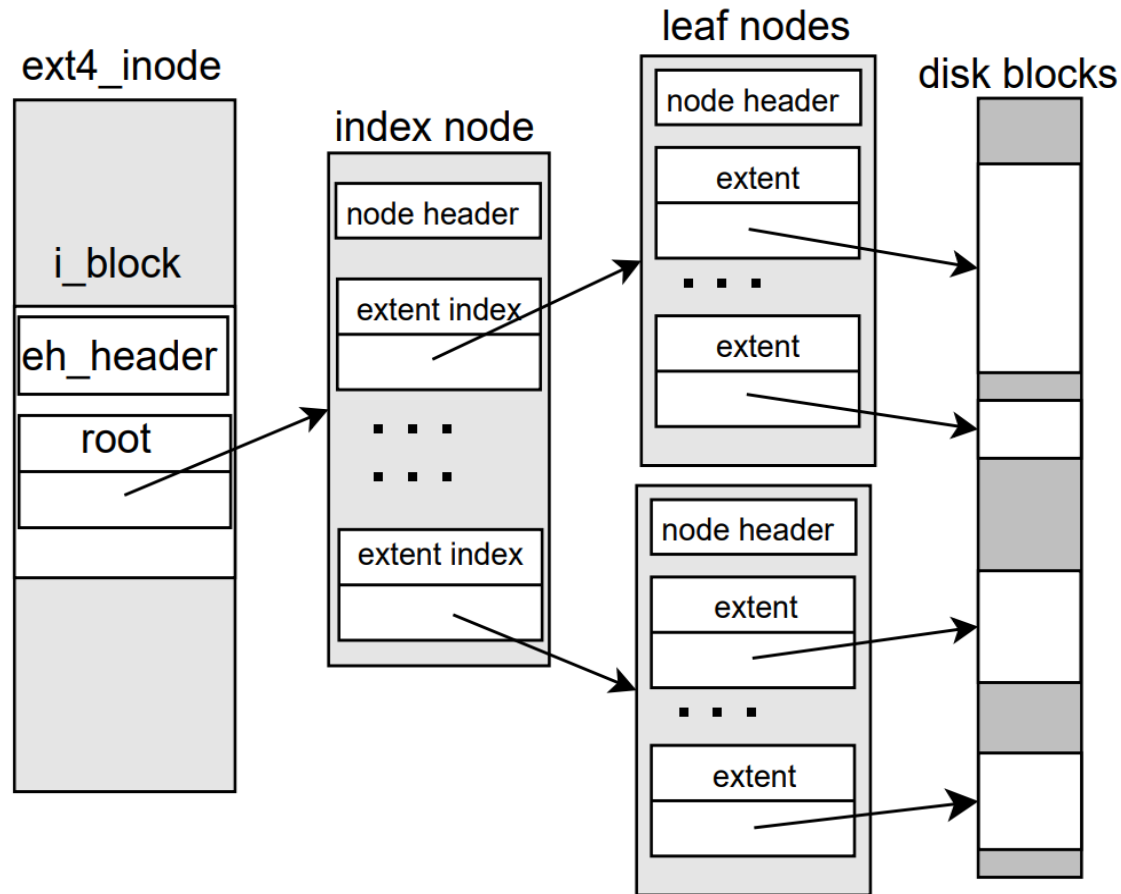


“All problems in computer science can be solved by another level of indirection” -- David Wheeler

Extent-Based Allocation

- A hybrid of contiguous allocation and linked/indexed allocation
- Extent-based file systems allocate disk blocks in **extents**
- An extent is a set of **contiguous** disk blocks
 - Extents are allocated upon file space allocation, but they are usually **larger than** the demanded size
 - Sequential access within extents
 - All extents of a file need not be of the same size
- Example: Linux ext4 file system

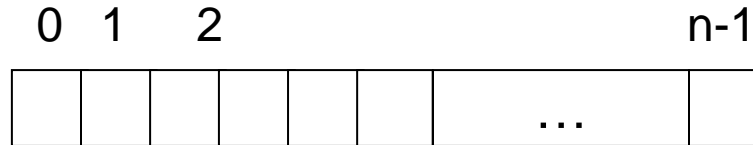
Ext4 Extent Allocation (extent tree)



1. Mathur, Avantika, et al. "The new ext4 filesystem: current status and future plans." *Proceedings of the Linux symposium*. Vol. 2. 2007.
2. <https://blogs.oracle.com/linux/post/understanding-ext4-disk-layout-part-2>

Issue 3: Free-Space Management

- Bit vector (n blocks)



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

Scanning for 0's to find free blocks

Block number calculation

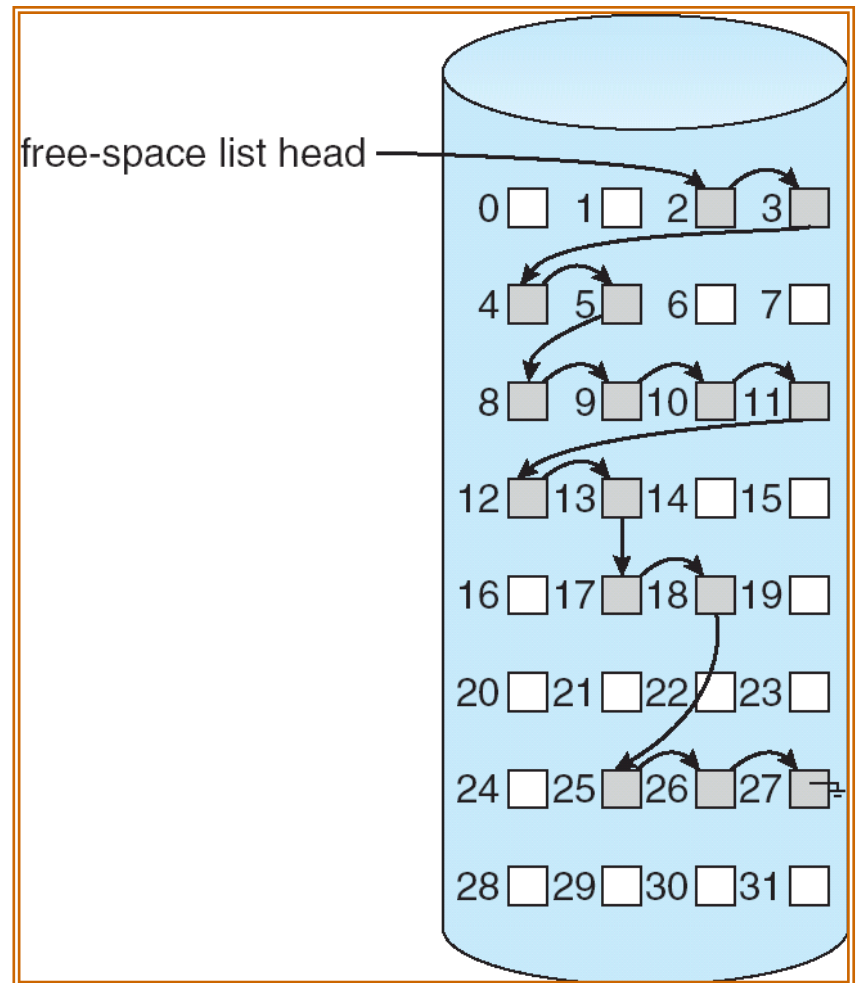
(number of bits per word) *
(number of **all**-0-value words) +
offset of first 1 bit

Free-Space Management (Cont.)

- Bit map requires extra space
- Example: 1TB hard drive with 4KB blocks
 - block size = 2^{12} bytes
 - disk size = 2^{40} bytes (1 terabyte)
 - $n = 2^{40}/2^{12} = 2^{28}$ bits (bitmap size = 32MB)
- Easier to get physically contiguous files, because neighboring bits represent adjacent blocks
- Used in UNIX FFS, Linux Ext family, ...

Linked Free Space List on Disk

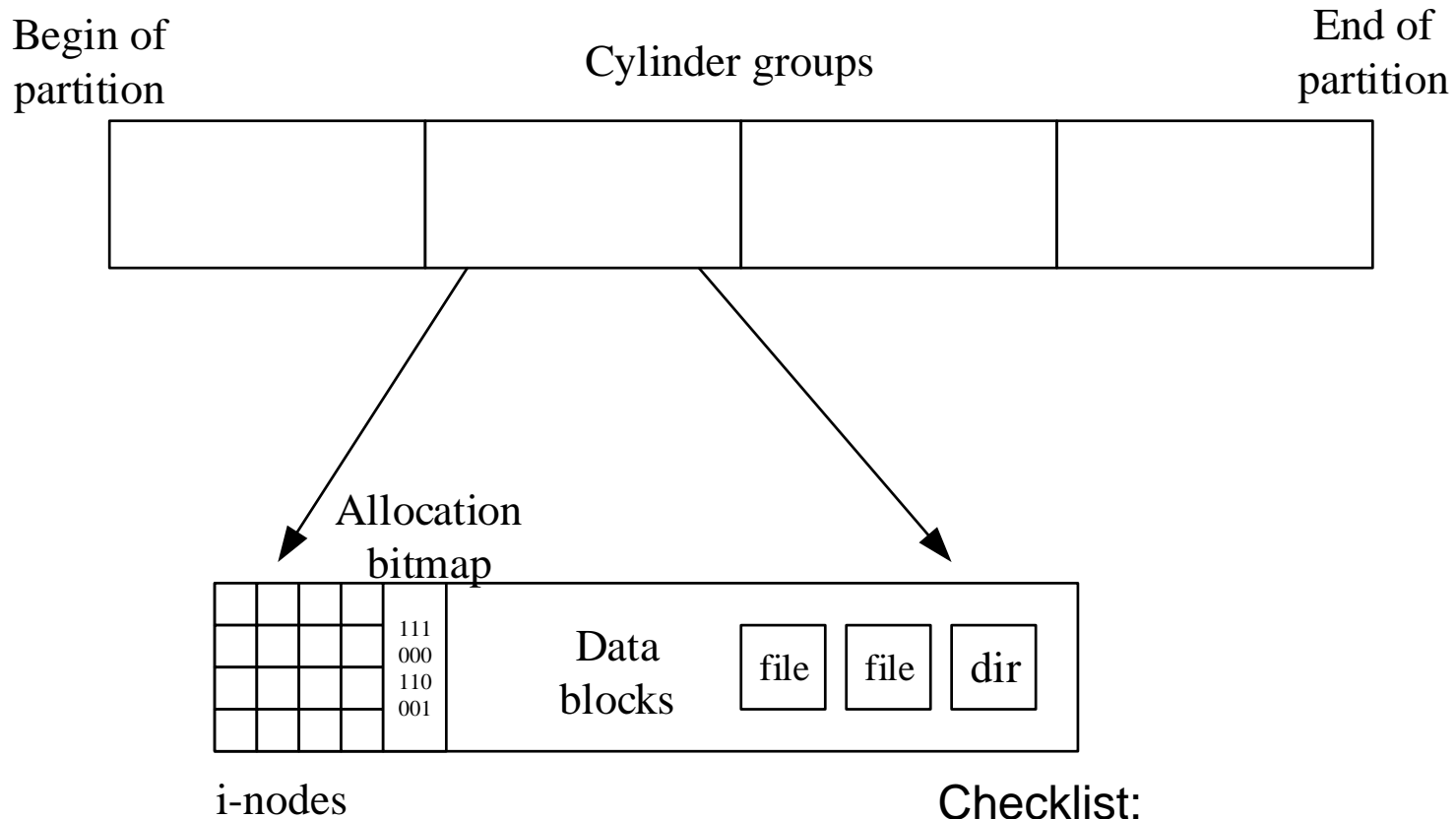
- Allocating and deallocating free blocks in a constant time
- No waste of free space
- But cannot get contiguous space easily, prone to fragmentation
- *Not* seen in modern file systems



Comparison

- Directory Implementation
 - Plain table: FAT, Ext2
 - B-tree: XFS, NTFS, Ext3&4
- Allocation methods
 - Linked list: FAT
 - Indexed allocation: Ext2&3
 - Extent: Ext 4
- Free space management
 - Bitmap: Ext family

Review: ext4 file system



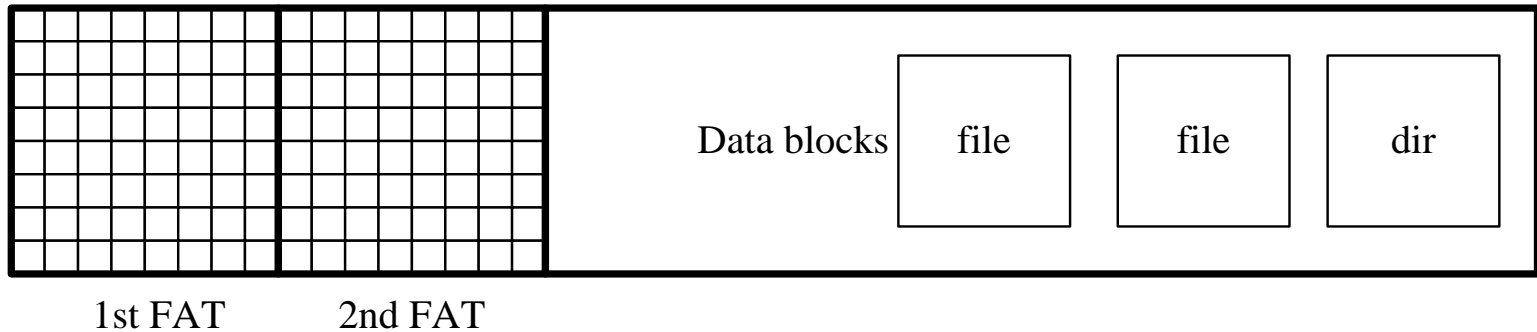
Checklist:

- Indexed(123) extent(4) allocation;
- Allocation bitmaps
- Directories (H-tree)

Review: FAT file system

Begin of
partition

End of
partition



Check list

- Linear directory table
- Linked allocation
- Scan 0 in FAT for free space (similar to bitmap)

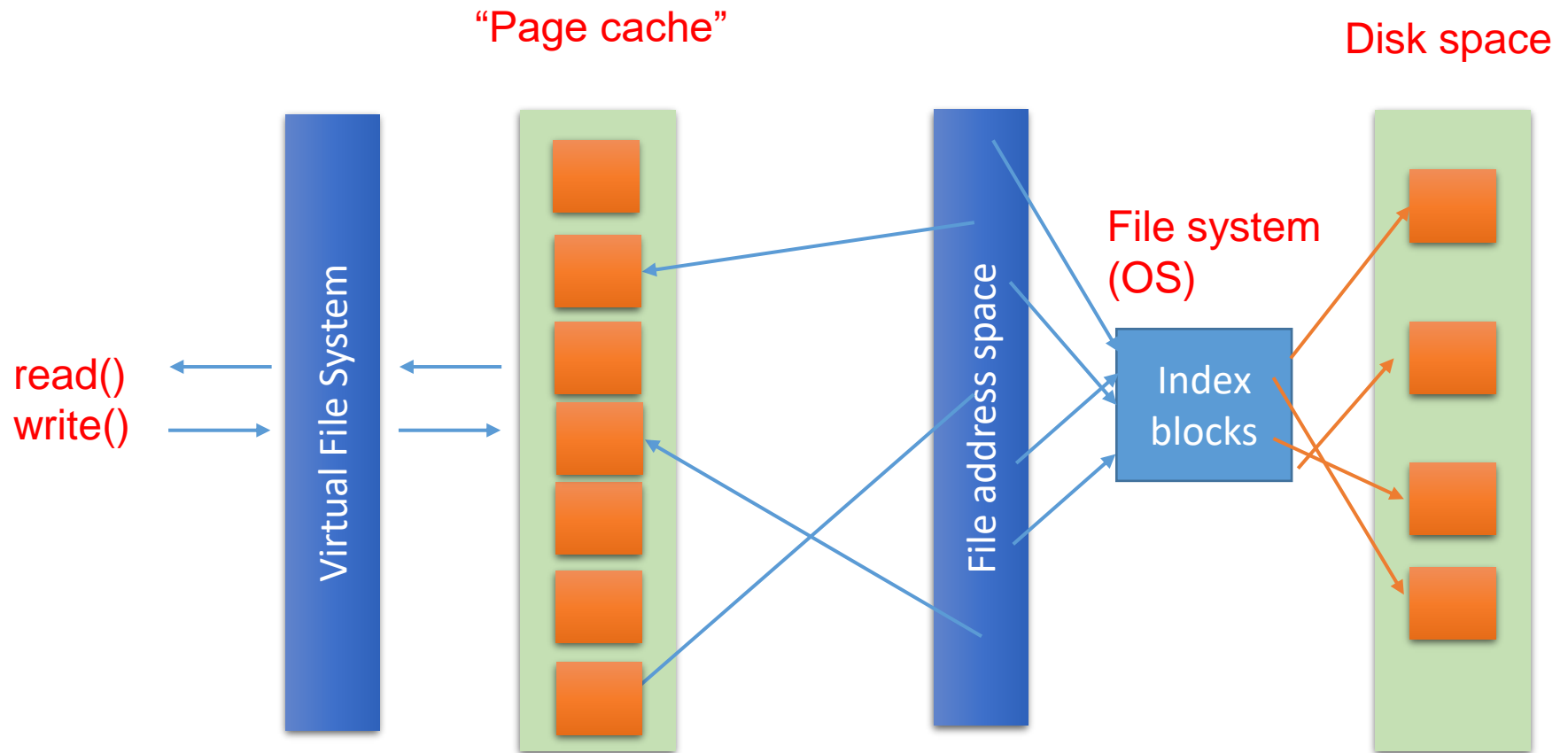
Efficiency and Performance

Generic Optimization

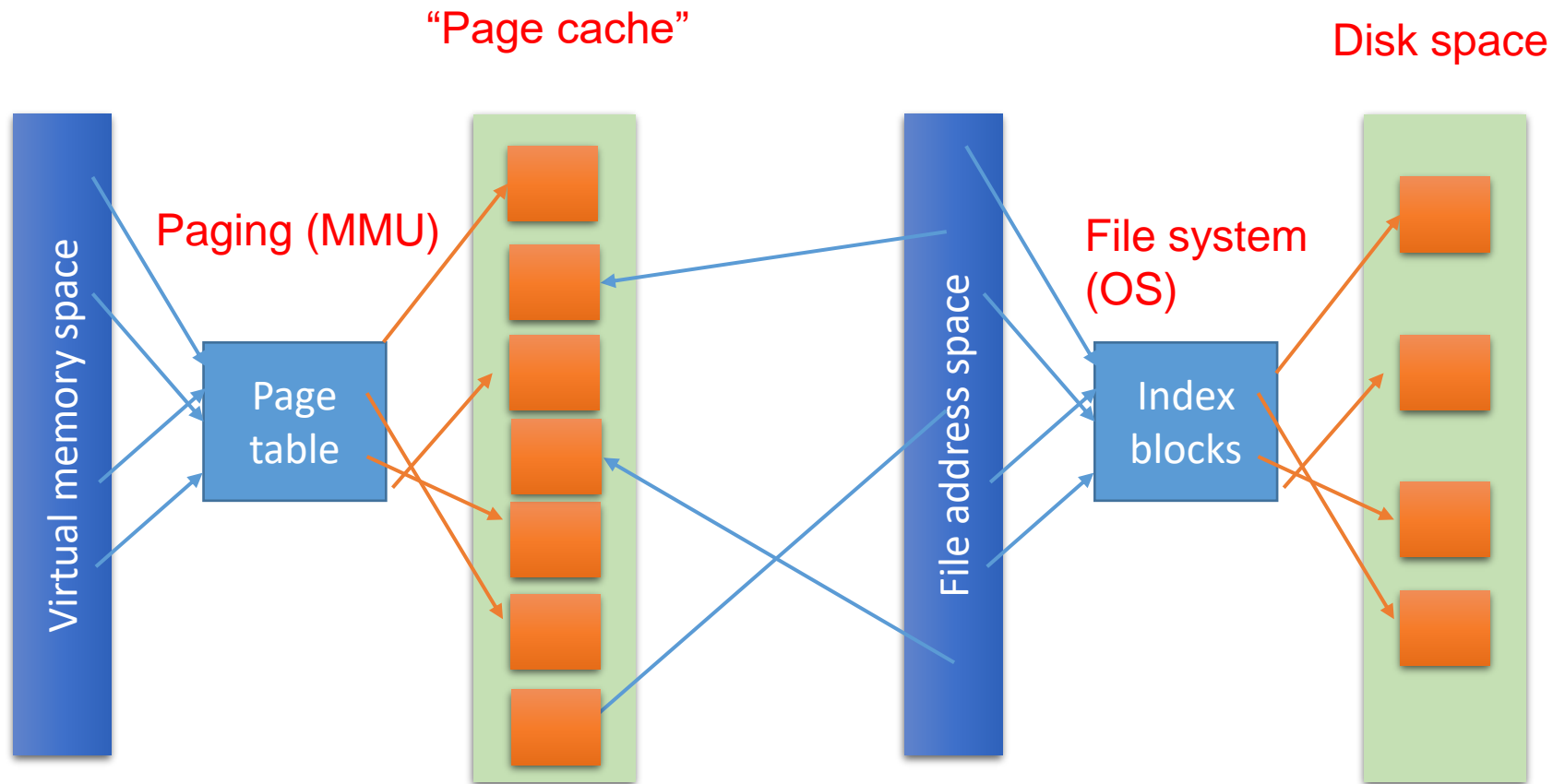
Applicable to all file systems

- Disk cache – keep recently used disk blocks in the page cache (temporal locality)
- File read-ahead aka prefetching – technique to optimize sequential access (spatial locality)
 - Similar to pre-paging. Read-ahead size doubles if prefetched data are used [[ref](#)]
 - Applications uses `fadvise()` to tell the kernel about how aggressive prefetching should be
- Defragmentation – restoring the physical contiguity of file data blocks

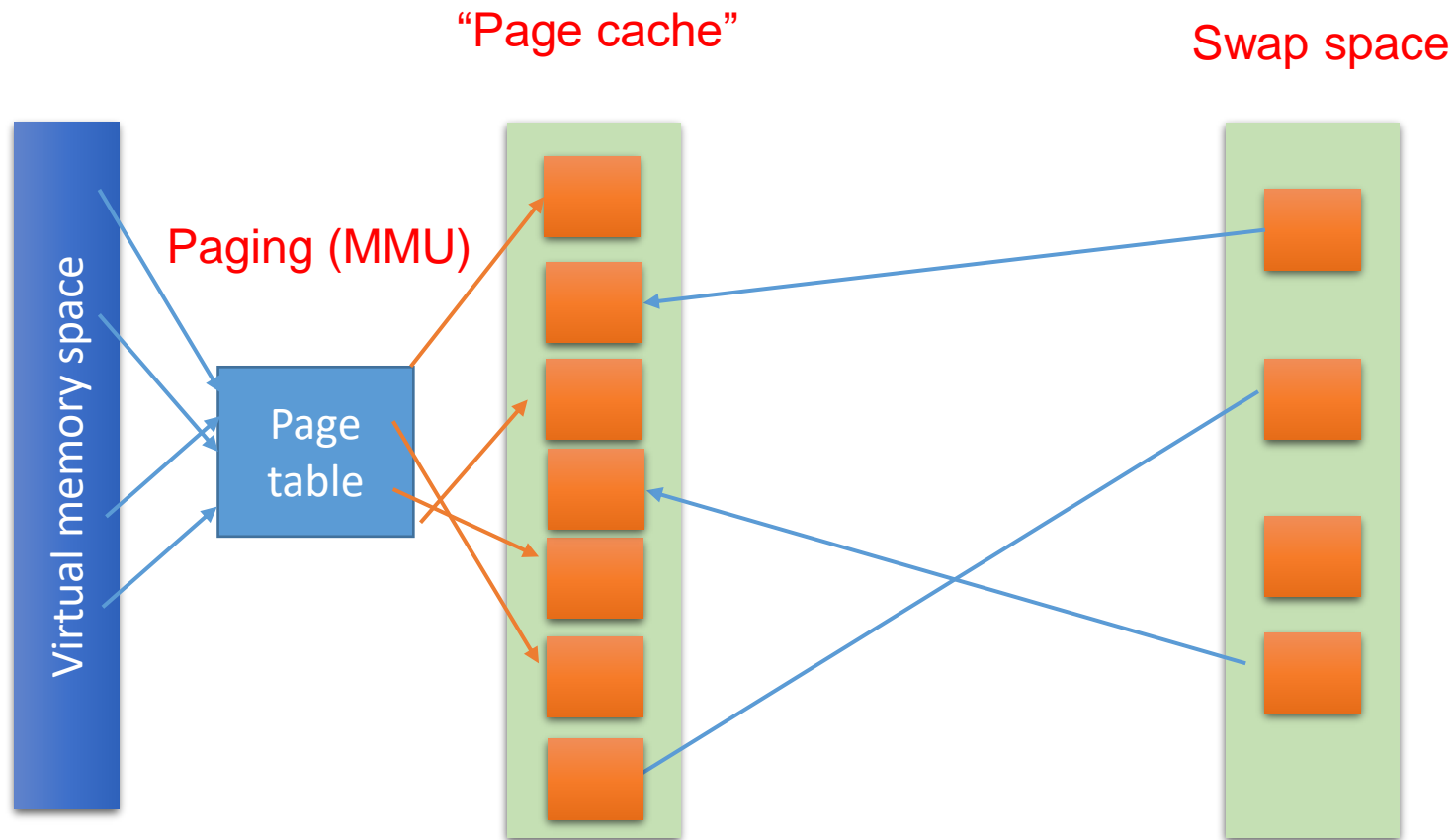
Page Caching: Regular File Pages (non-mapped)



Page Caching: mmap()'ed File Pages



Page Caching: Anonymous Pages



File Fragmentation

- File system “ages” after many creation and deletion of files
 - Free space is fragmented into small holes
 - File system cannot find contiguous free space for a new file or for an existing file to grow
- Degree of Fragmentation (DoF) of a file

$$DoF = \frac{\text{\# of extents of the file}}{\text{the ideal \# of extents for the file}}$$

- The higher the DoF of a file is, the more disk seeks are required to access the file

File Defragmentation



Making fragmented files sequential.

→ Reducing I/O count and disk head movement on file access.

FS-Specific Optimizations

- File systems have their **unique** techniques for performance optimization

For example, Ext4 employ the following optimizations:

- Dividing disk space into cylinder groups to make inodes appear near to their associated data blocks
- Embedding small files into directories (<60 bytes); called “inline files”
- Using extents to take advantage of sequential disk accesses

Consistency and Recovery

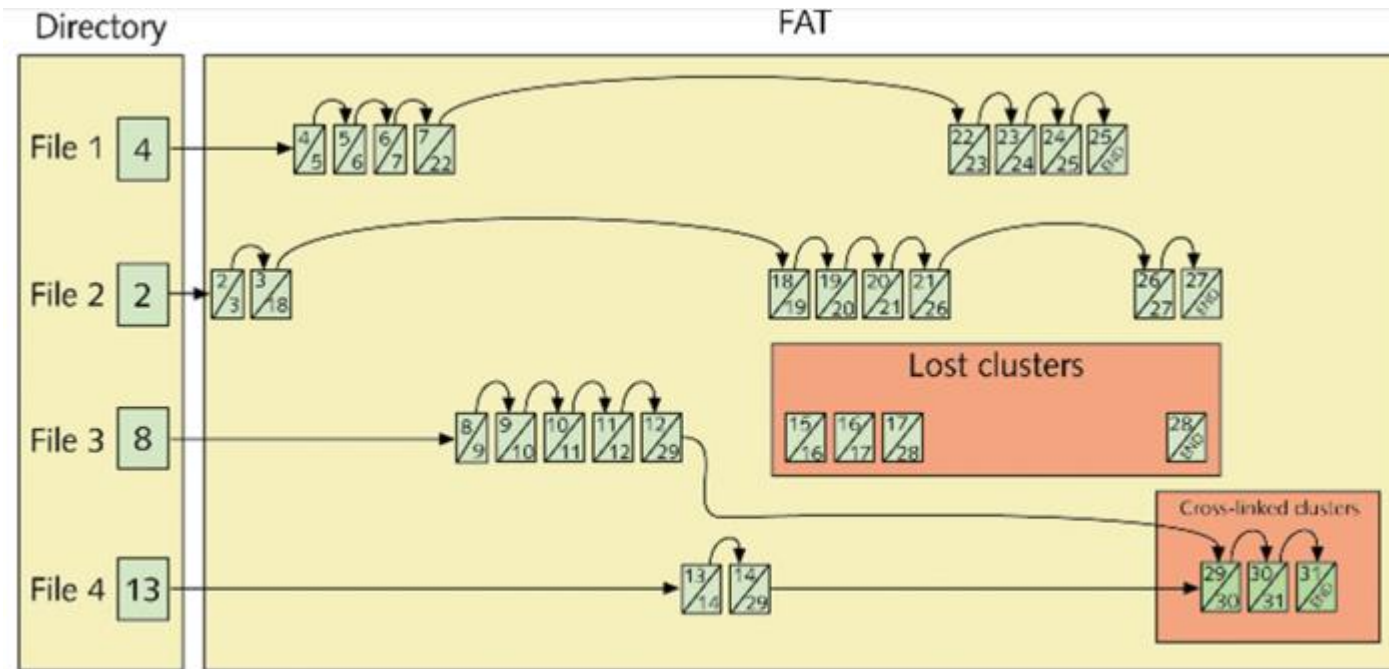
Inconsistency and Recovery

- A file operation involves **multiple block modifications**
 - To create a file in ext4 will need to modify: allocation bitmap, inode, directory, data block
- What if power fails in the middle of file creation?
 - Unwritten data/metadata will be lost
 - Loss of metadata: structural inconsistency
 - Loss of user data: undefined data (garbage) in files

Consistency checking – compares data in directory structure with data blocks on disk, and tries to fix inconsistencies

Structural Inconsistency Examples

- Ext file systems
 - A bitmap indicates that an inode has been allocated but the inode has not yet been written (and vice versa)
 - A hard link is created to a file but the file's reference count has not been incremented yet
- FAT file systems
 - A list of blocks are freed and re-allocated to another file, but the link list table has not yet been updated (cross-linked lists in FAT)



Lost and cross-linked clusters

http://faculty.salina.k-state.edu/tim/ossg/File_sys/file_system_errors.html

Recovery Utilities

- Usually a dirty bit in the super block can tell whether a volume is cleanly unmounted
- Run file system consistency check on dirty volumes
 - `fsck` (UNIX) `scandisk` (Windows)
 - A lengthy process, takes up to 1 hour on a 1 GB disk

Journaling File Systems

- The root cause of file system inconsistency
 - A file operation, which involves to modify multiple disk blocks, is interrupted
- Transactions
 - An idea borrowed from database systems
 - A set of self-contained disk block modifications
- Protecting the file system against inconsistency
 - To guarantee the atomicity of file transactions
 - **Atomicity:** “all done” or “nothing is done”
 - E.g., two `frwrite()` modifies 5 disk blocks in total. With journaling, the 5 blocks are either all modified or none is.

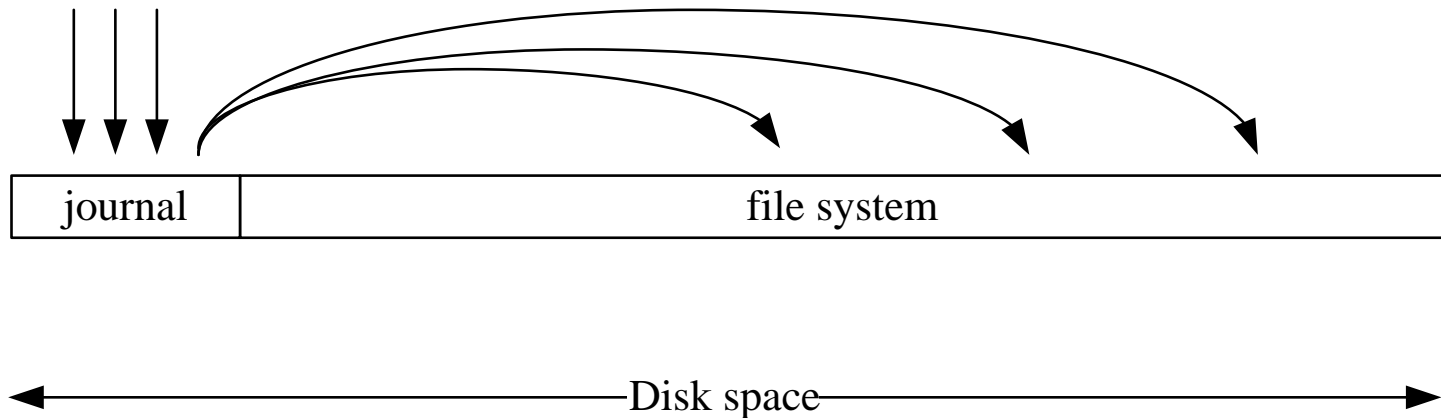
Journaling File Systems

- Journaling file systems often employ Write-Ahead Logging, WAL, to guarantee the atomicity of transactions
- WAL requires a reserved space as the journal
- Two-step approach
 - The file system commits a transaction (to the journal)
 - The file system applies the changes (to the file system)

Write-Ahead Logging (WAL)

(1) Commit a transaction

(2) Apply the changes



Crash Recovery with WAL

1. Scan the journal
 2. Found a complete transaction → redo
 3. Found a partial transaction → discard
- Transaction atomicity is thus guaranteed

Journaling File systems -- Summary

- Motivation
 - Preventing power interruptions from corrupting file systems
- Method
 - Creating a journal space for the file system
 - Collecting a set of self-contained writes as a transaction
 - Write transactions to the journal
 - Apply changes to the file system
 - On recovery, scan the disk journal. Re-do legit transactions; incomplete transactions will be discarded
- Benefit
 - Crash recovery is very fast
- Problem
 - Amplifying the write traffic

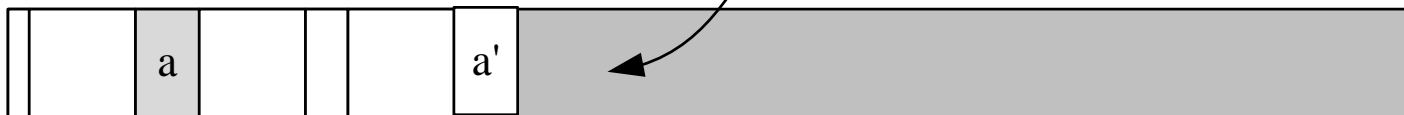
Log-Structured File System:
sequential writing always

Log-Structured File Systems

- Performance bottleneck of modern file systems
 - Read performance: not a problem with a large disk cache
 - Write performance: random writes are slow
- Key idea: out-of-place update
 - Random updates need not occur in place, they are converted to sequential writes
 - A log-structured file system can be imagined as a huge journal space without the “file system”
- Examples
 - NILFS2 (servers), F2FS (Android devices), NOVA (NVRAM)

The Concept of Log-Structured File Systems

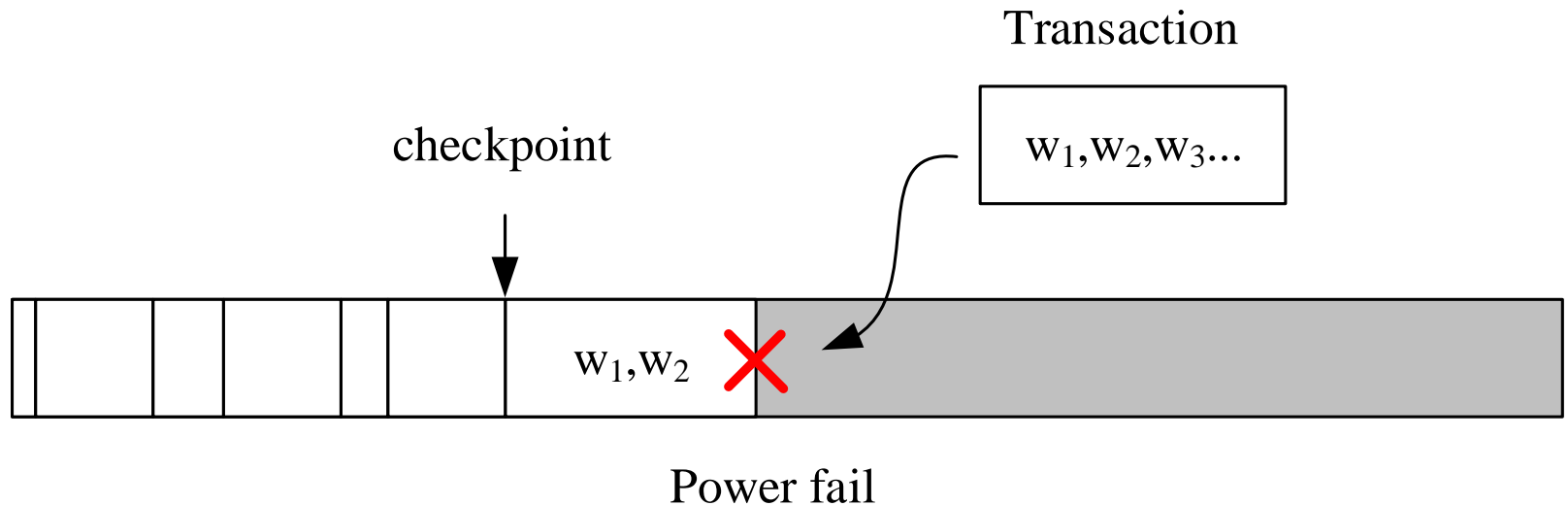
(1) Buffer multiple write operations as logs

$$W_1, W_2, W_3 \dots$$


(2) Write logs to sequential disk space

- Writes are always sequential and thus are highly efficient
- Out-of-place updates leaves garbage in the storage

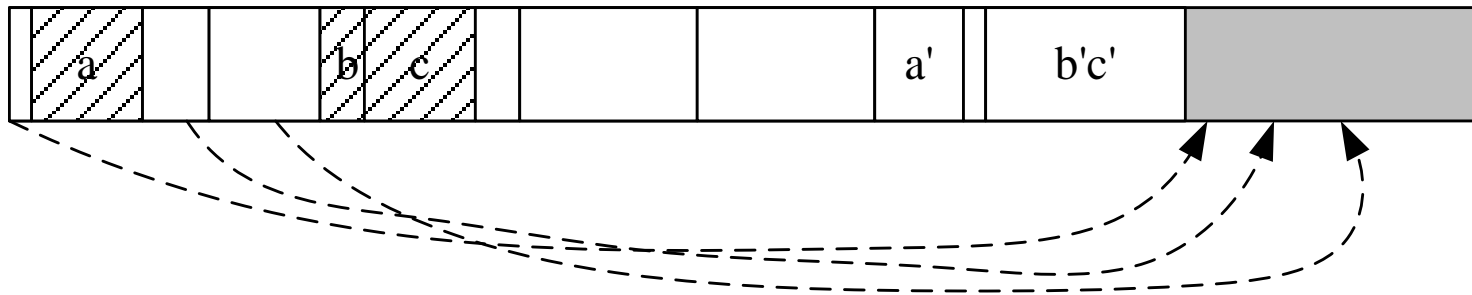
LFS Checkpoint and Recovery



- Recovery is surprisingly simple because writes are chronologically ordered

LFS Cleaning

(3) Out-of-place updates produce invalid data



(4) Reclaim contiguous disk space with compaction (garbage collection)



(5) compaction produces contiguous free space

Also known as Compaction or Garbage Collection

Log-Structured File Systems -- Summary

- Motivation:
 - RAM will be cheap and random reads are not a problem with a large page cache
 - Random writes must eventually hit the disk and they are slow
- Methods:
 - Collecting random writes (updates) into a long write burst
 - **Out-of-place** updates via sequential writing
- Benefits:
 - Great random write performance
 - Easy recovery
- Problems:
 - Need cleaning (i.e., compaction or garbage collection) to regenerate sequential space for new writes

End of Chapter 11

Review Questions

1. What is an inode? For a file whose hard-link count is two, how many inode(s) does the file have?
2. Bitmap-based allocation is more commonly seen in file-system designs than list-based allocation. Why?
3. Why is extent allocation increasingly popular in modern file systems?
4. Discuss the benefit of having a few direct pointers in inodes of a UNIX file system
5. Optional study: Ext4 extent tree (allocation) and hash tree (directory management) [[ref](#)]
6. How file defragmentation improves performance, especially for hard drives? Why is file defragmentation often considered harmful to SSDs?
7. Optional study: the Linux disk prefetching algorithm [[ref](#)]

Review Questions

6. Why it is very time consuming to checking and fixing inconsistencies for file systems?
7. What is the root cause of inconsistencies in file systems? How supporting transaction semantic avoids this problem?
8. For file system journaling, what is a transaction?
9. Power fails can occur at anytime. Discuss how write-ahead logging (WAL) guarantees the atomicity of transactions.
10. What is the design motivation of log-structured file system? (in the aspect of performance)
11. Log-structured file systems deliver extremely high write performance, but when the space utilization is higher than 70%, it begins to suffer from performance degradation. Survey why this happens.