

# 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

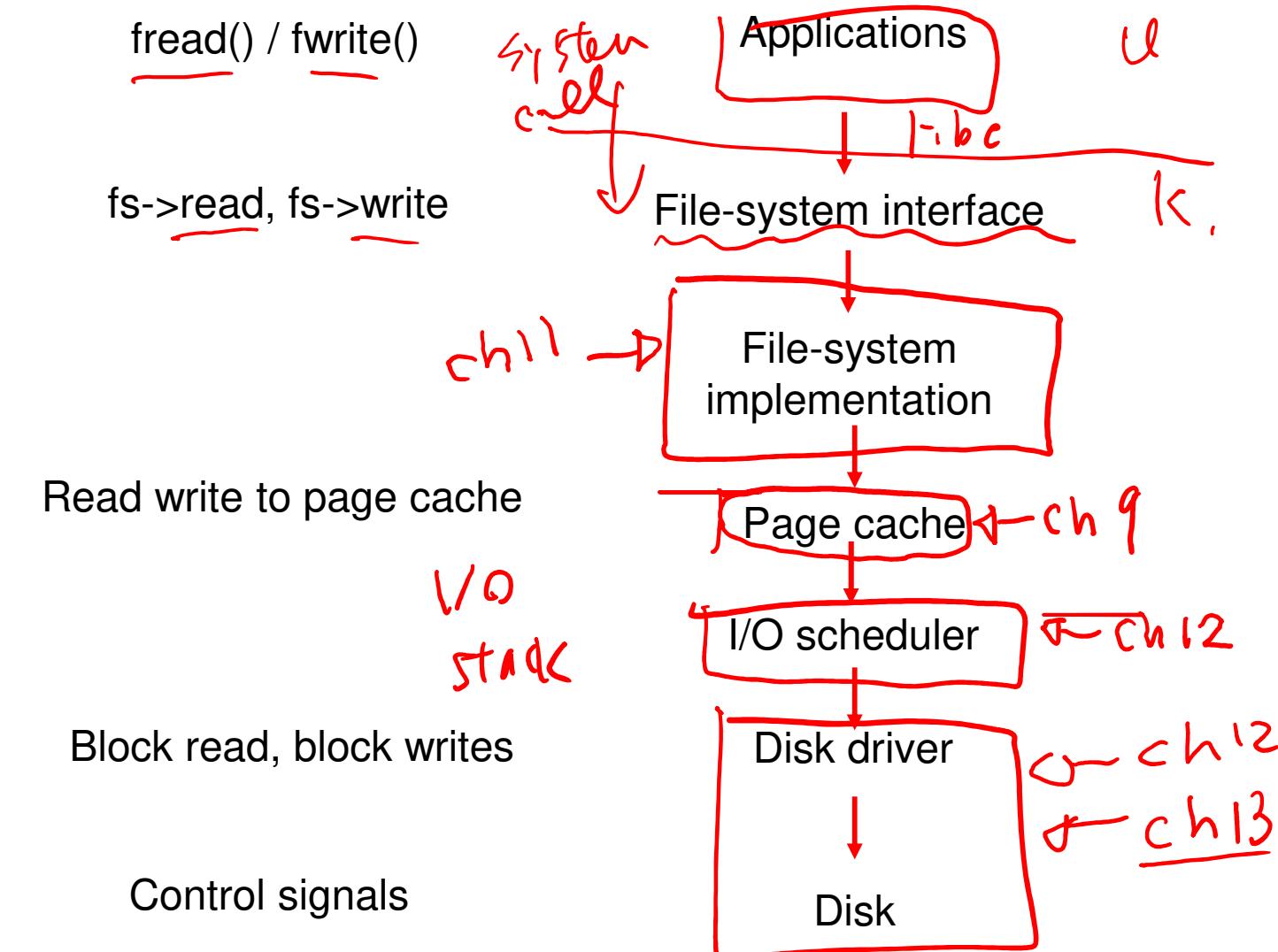
# Objectives

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

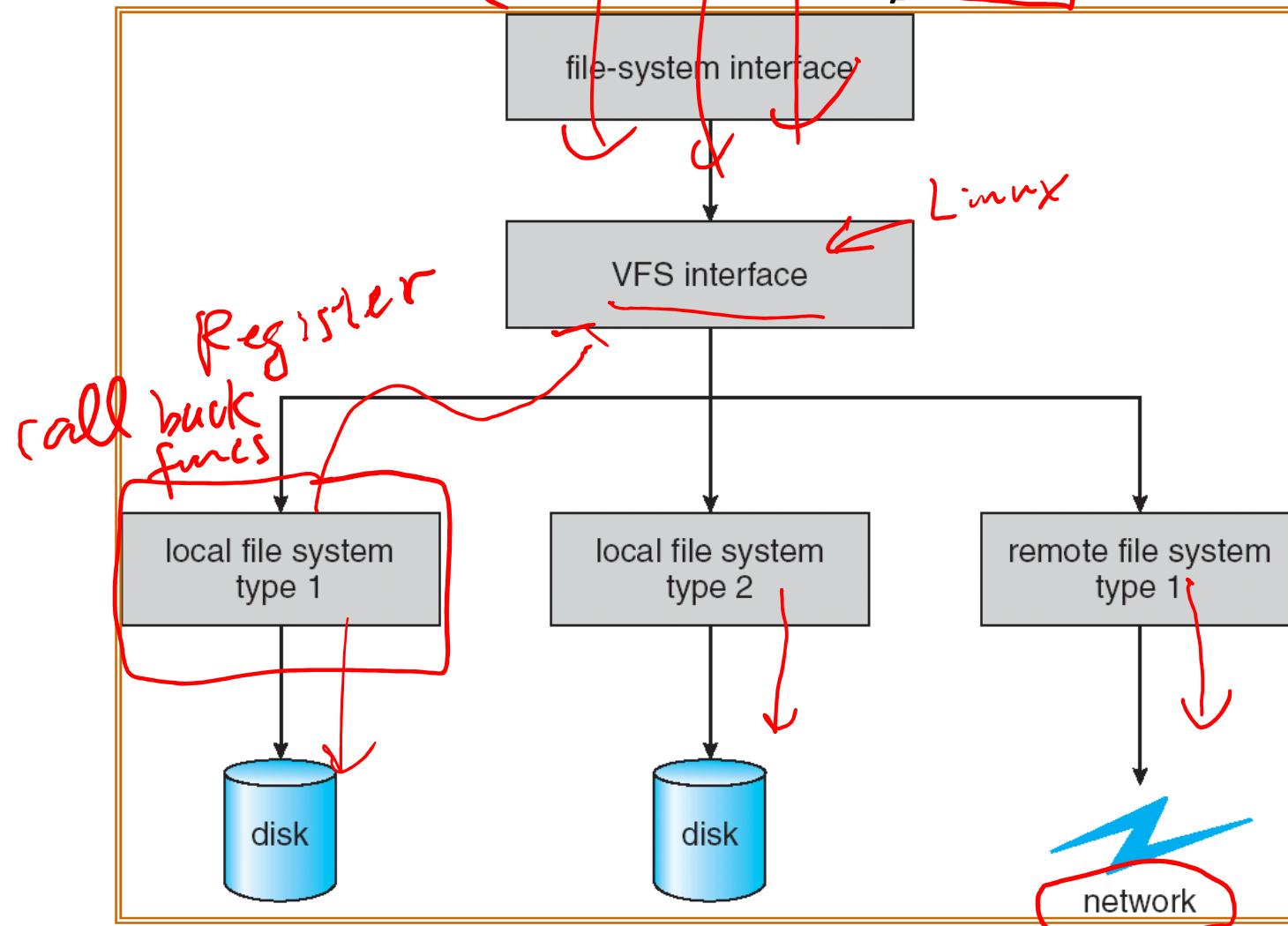
# File-System Structure

- File structure
  - Logical storage unit
  - Collection of related information
- File system resides on secondary storage (disks)
- File system organized into layers

# Layered File System

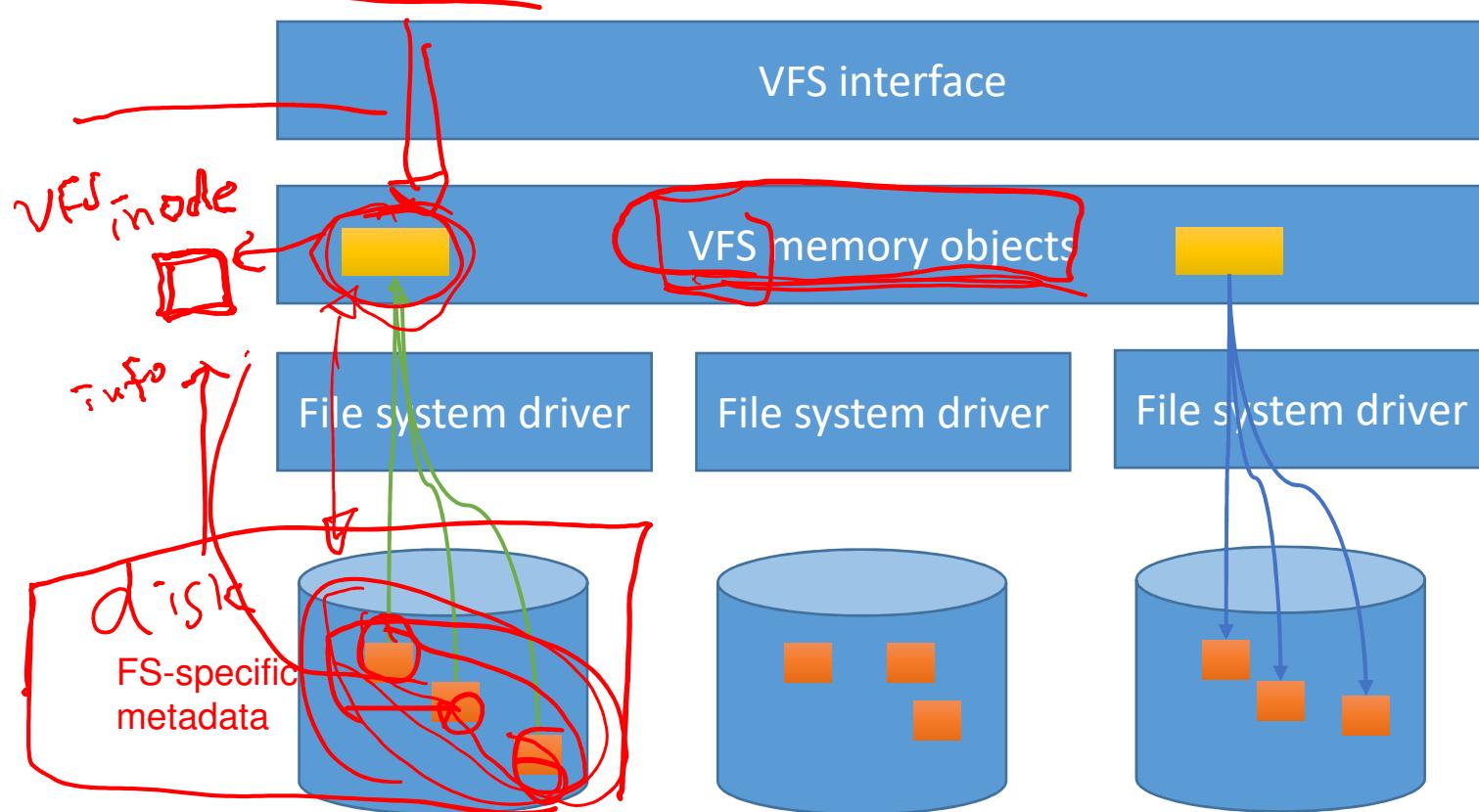


# Schematic View of Virtual File System



# 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

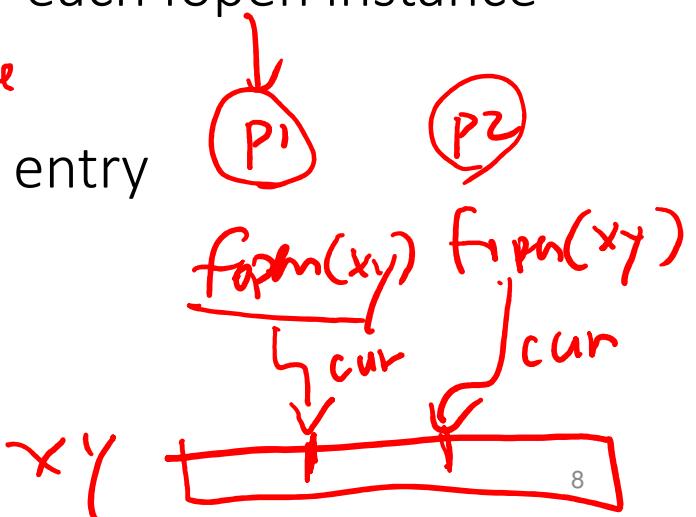
"Tree"

FS type → ext4.  
block size → 4kB?  
volume size → 1TB

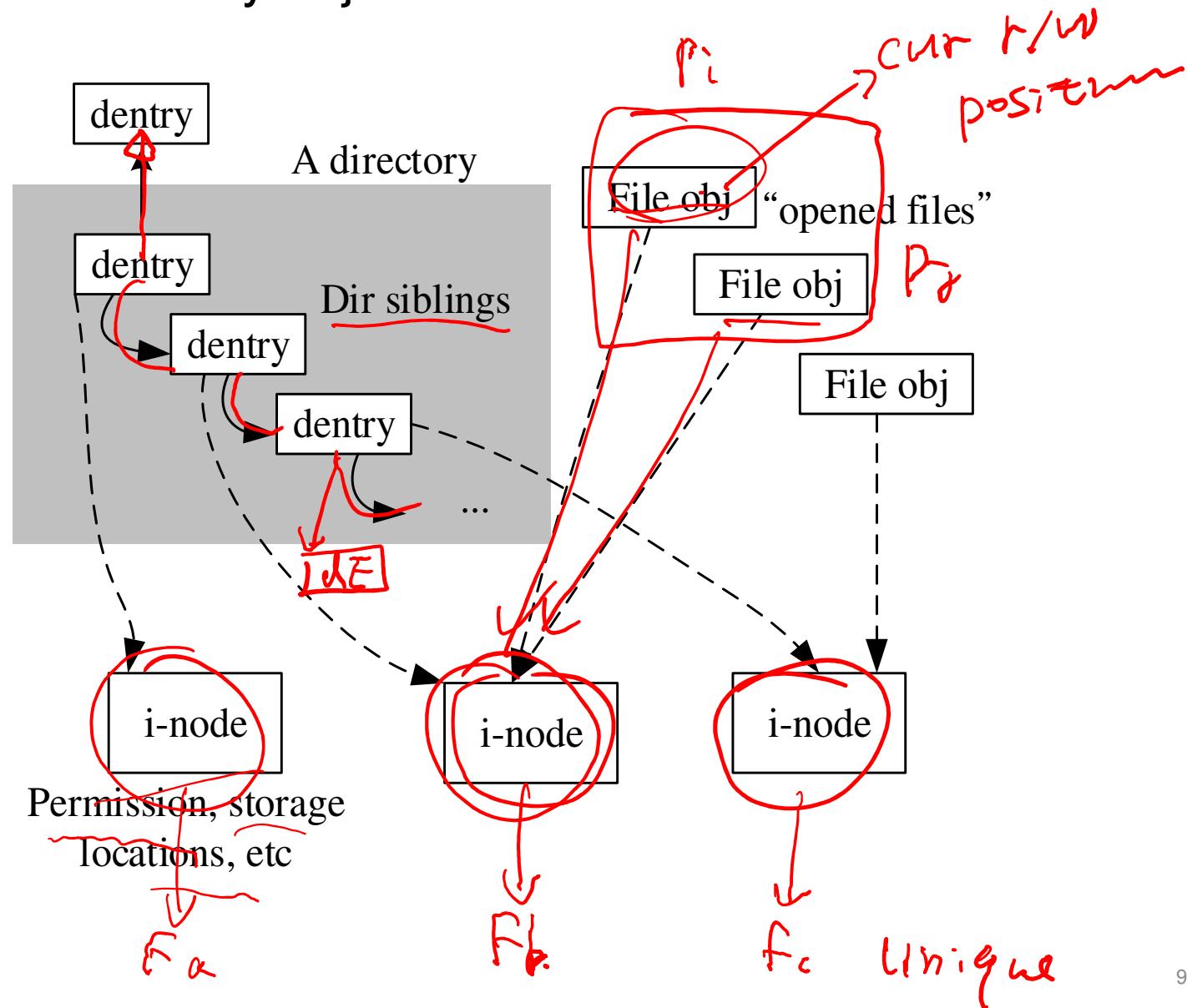
→ permission

file size  
data blocks ptr

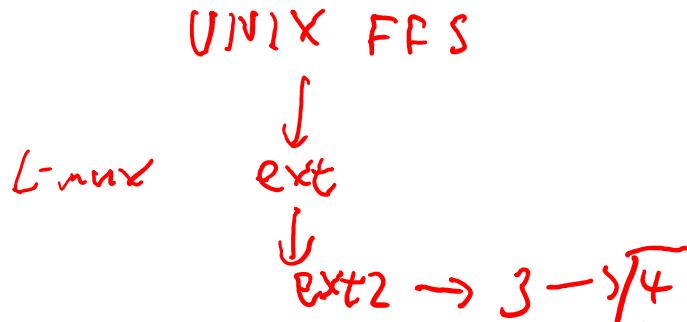
file-open-instance



# In-memory objects of Linux VFS



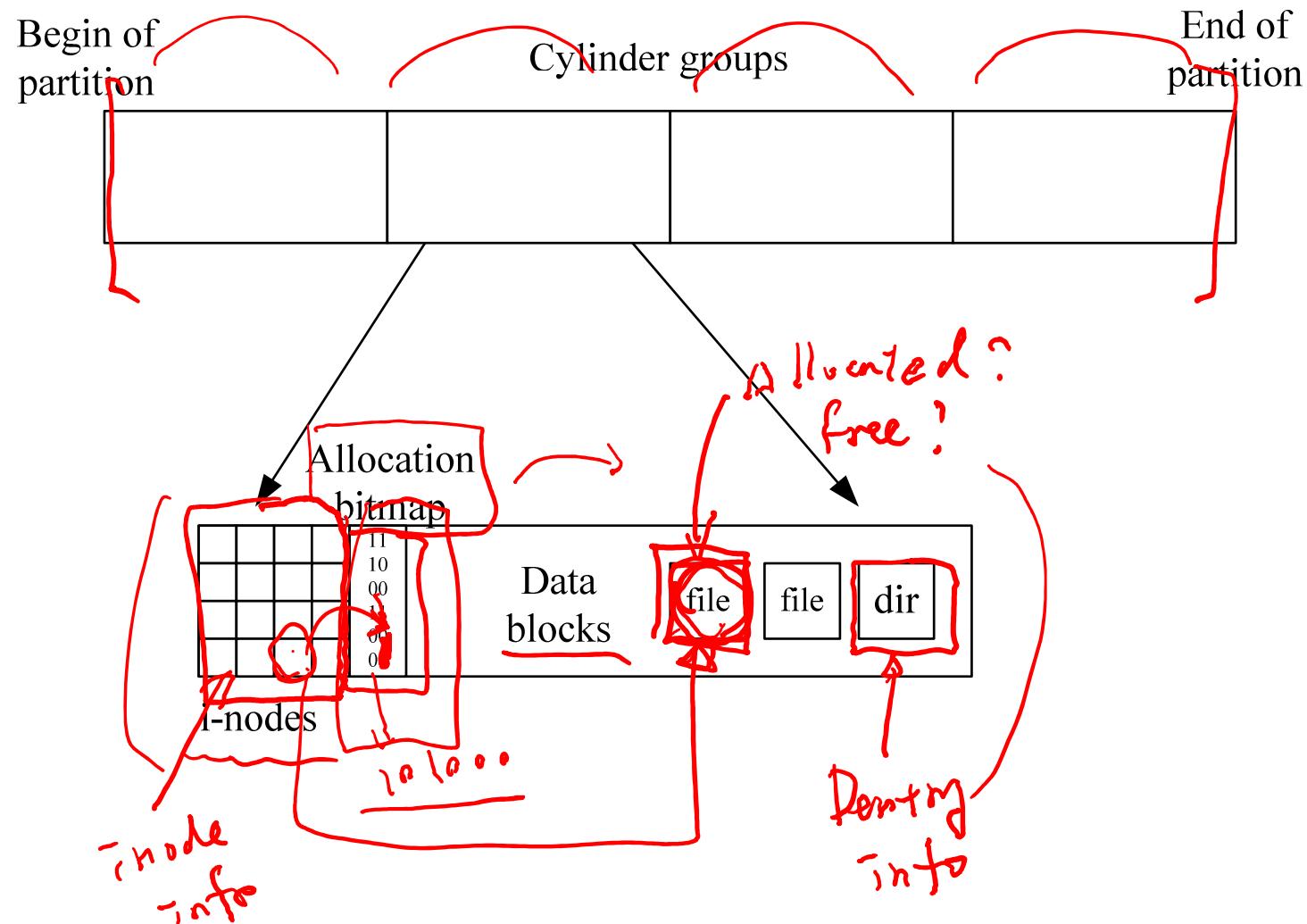
## Disk Metadata



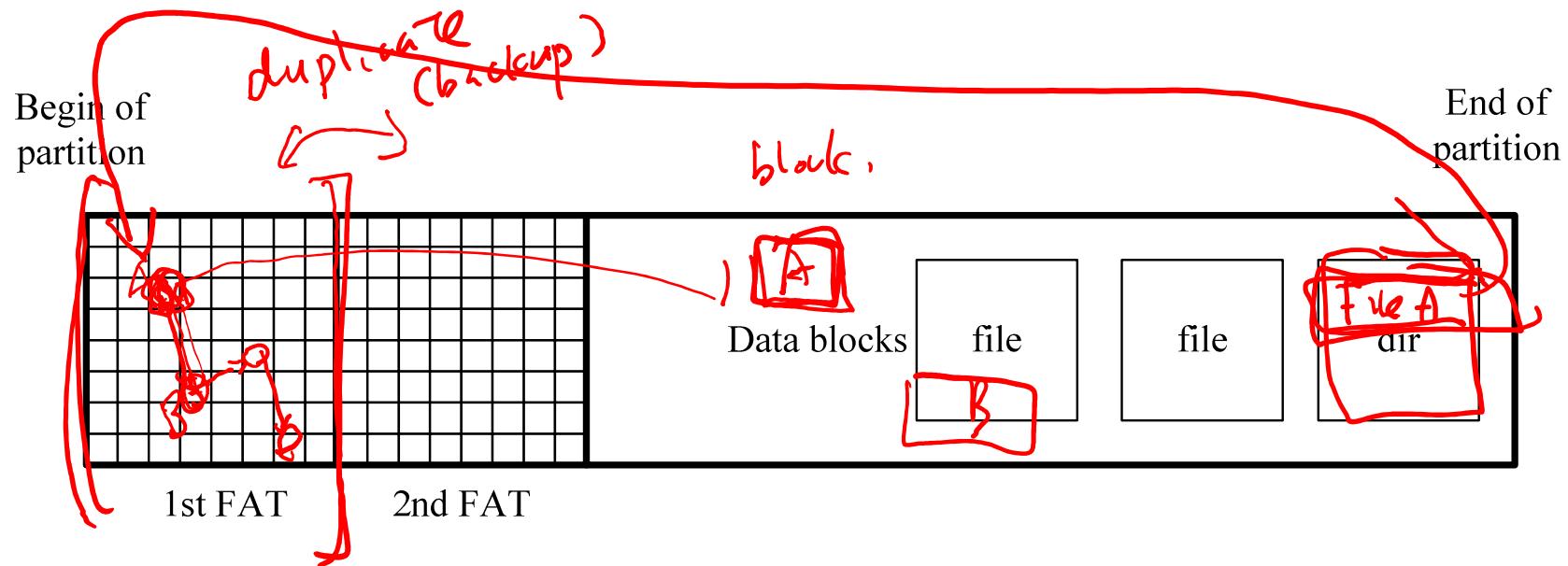
- File-system-specific; vary from file system to file system
- Linux ~~ext~~ file system
  - Super block, Inodes, Allocation bitmaps
- Microsoft ~~FAT~~ file system → ~~Simple~~ NTFS
  - File allocation tables, Directories
- File system driver must fill the in-memory objects with the information in disk 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

$\sim 128MB / 256MB$

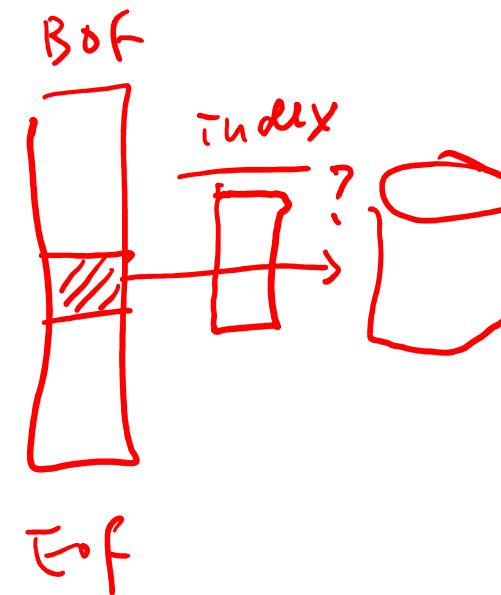


## Disk layout of FAT 12/16/32 file systems



# Key Issues of File System Implementation

- Directory implementation
- Allocation (index) methods
- 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)
  - 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 00 00 FF FF FF 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 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 00 00 FF FF 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 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 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 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 2F 00 74 00 78 00 74 00 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 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.." ..

# Example: Directory Dump in FAT

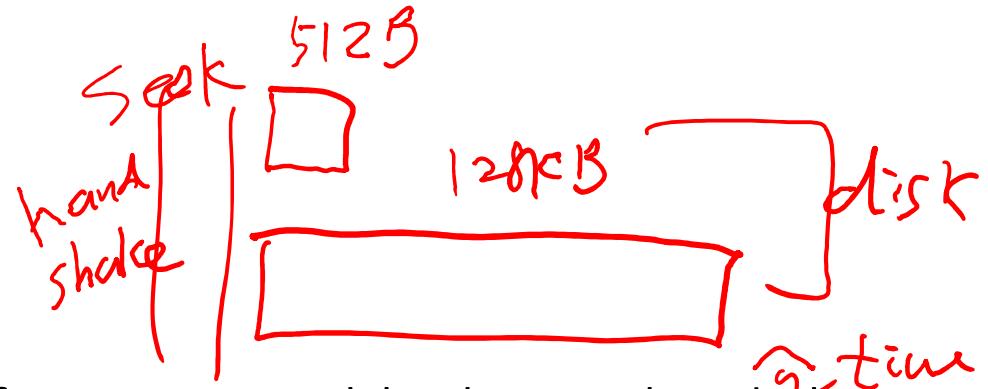
*usr data*

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	00	00	FF	FF	FF	FF	.....							
000167968	4D	59	44	49	52	32	20	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	00	00	FF	FF	FF	FF	.....							
000168032	4D	59	44	49	52	31	20	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 ✓

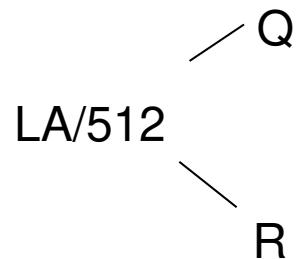
## Contiguous Allocation



- 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
  - file offset can be directly translated into sector block #
  - Less I/Os involved
  - Always sequential disk read/write
- Wasteful of space (dynamic storage-allocation problem)
  - File deletion leaves free holes (external fragmentation)
  - Needs compaction, maybe done in background or downtime

# Contiguous Allocation

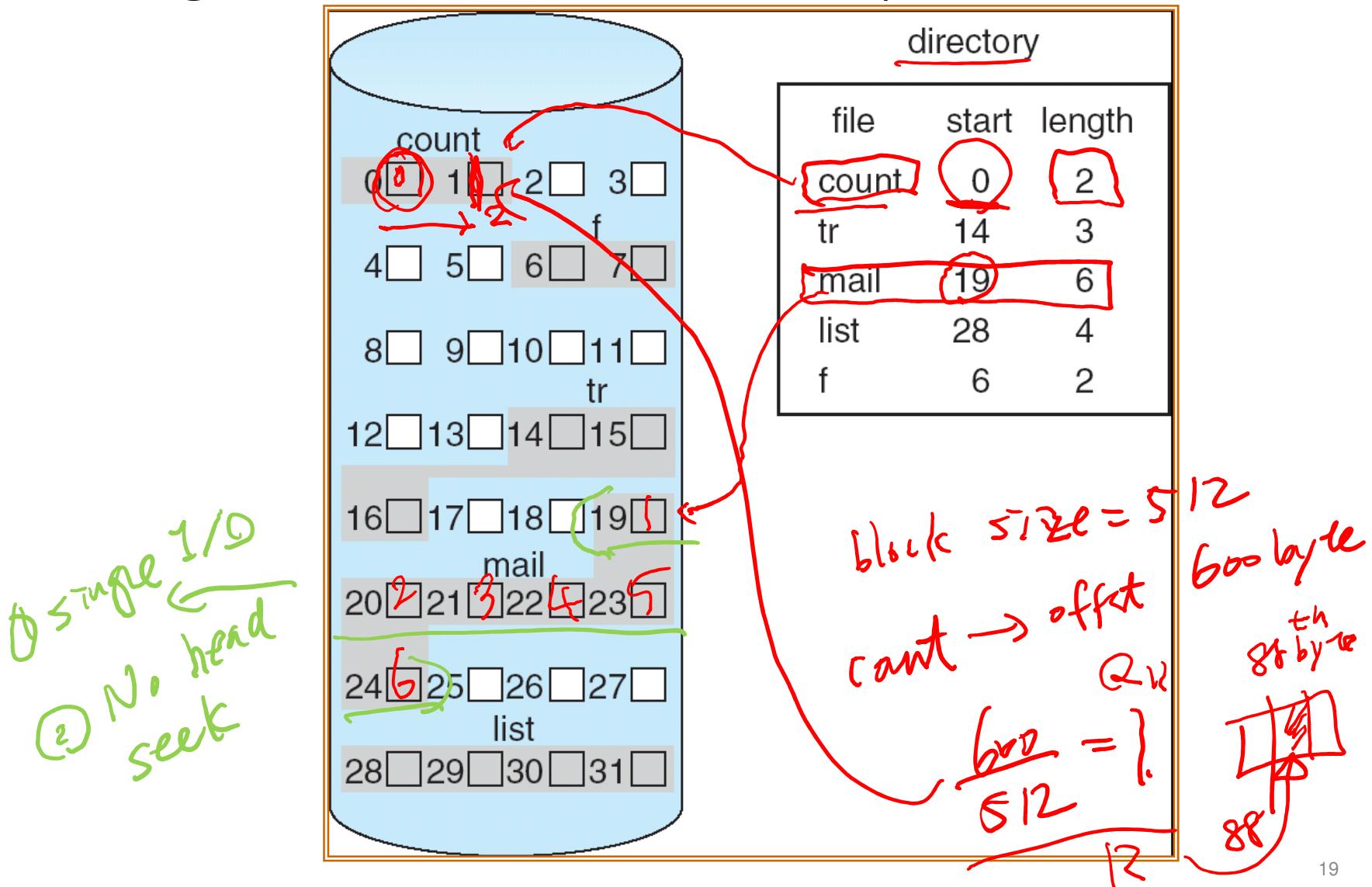
- Mapping from logical to physical



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

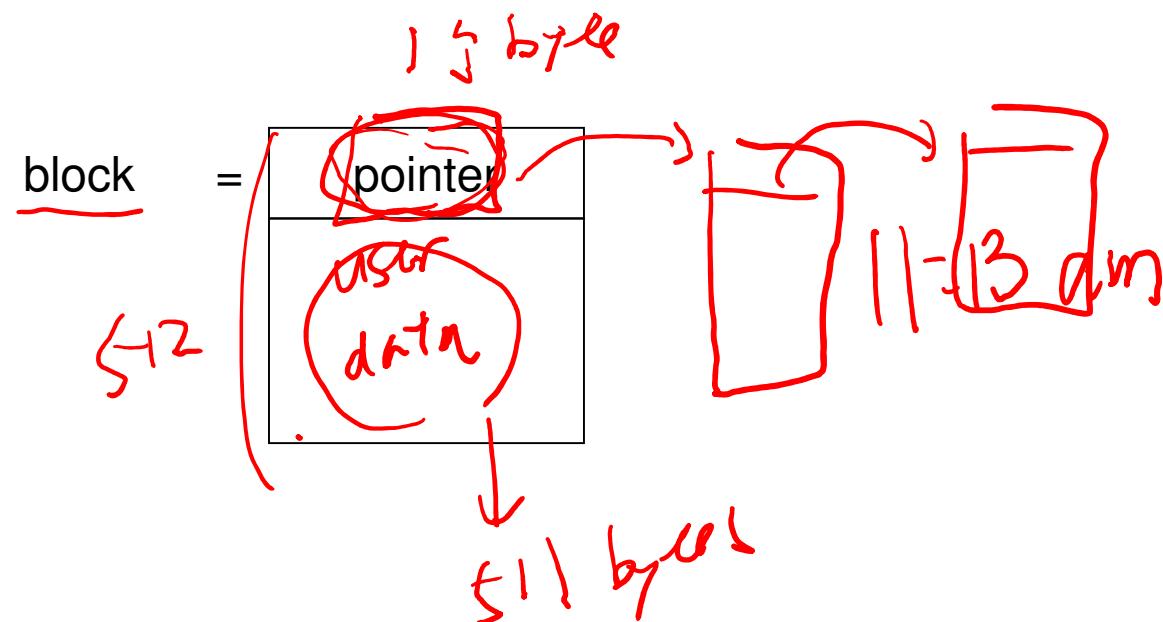
LA=byte address  
1 block=512B

# Contiguous Allocation of Disk Space



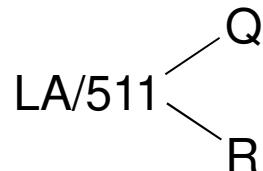
# Linked Allocation

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



## 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

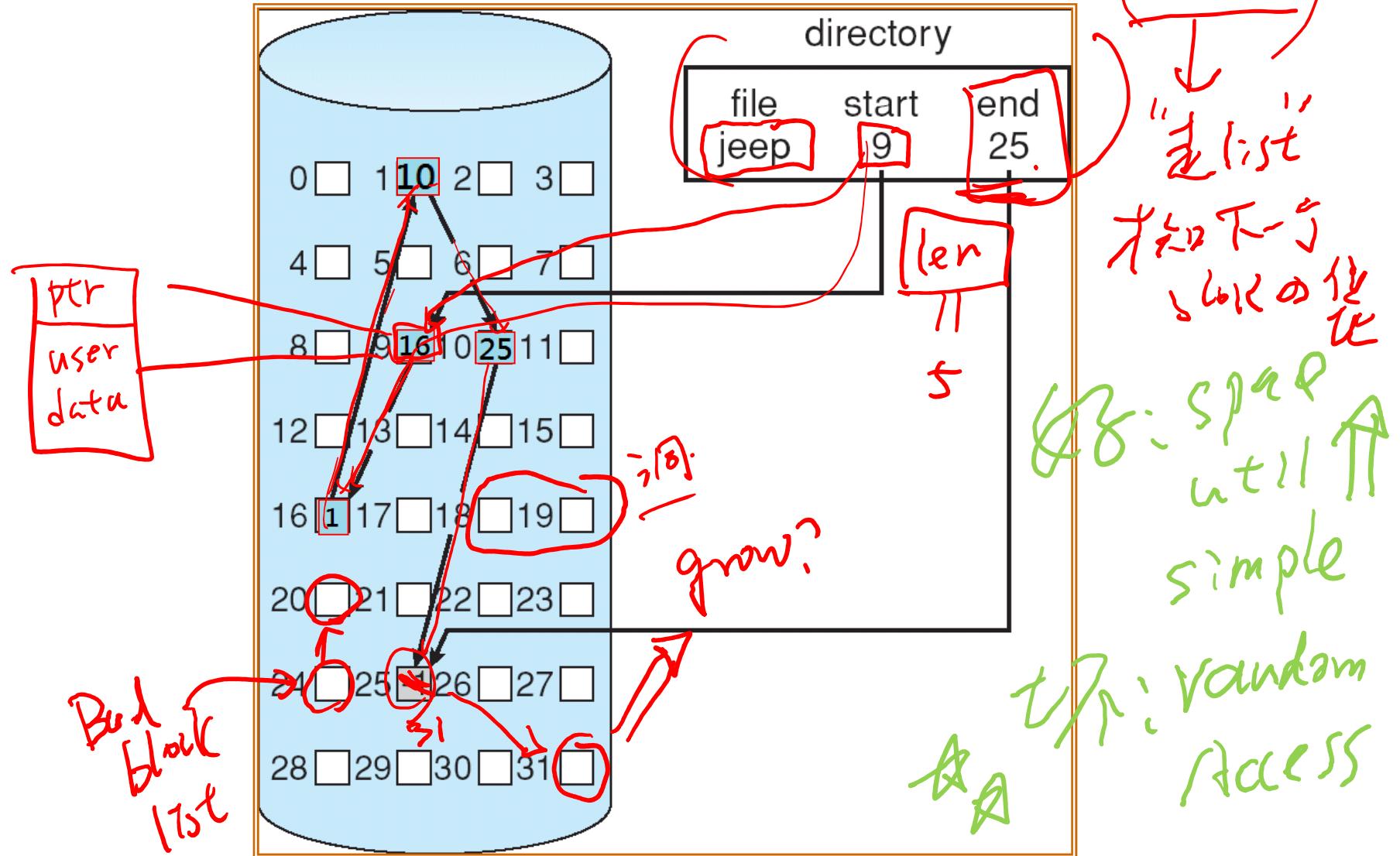


1 data block = 512B, 1 for ptr,  
So 511B for user data

Block to be accessed is the Qth block in the linked chain of blocks representing the file.

Displacement into block = R + 1 (the 0<sup>th</sup> byte is for pointer)

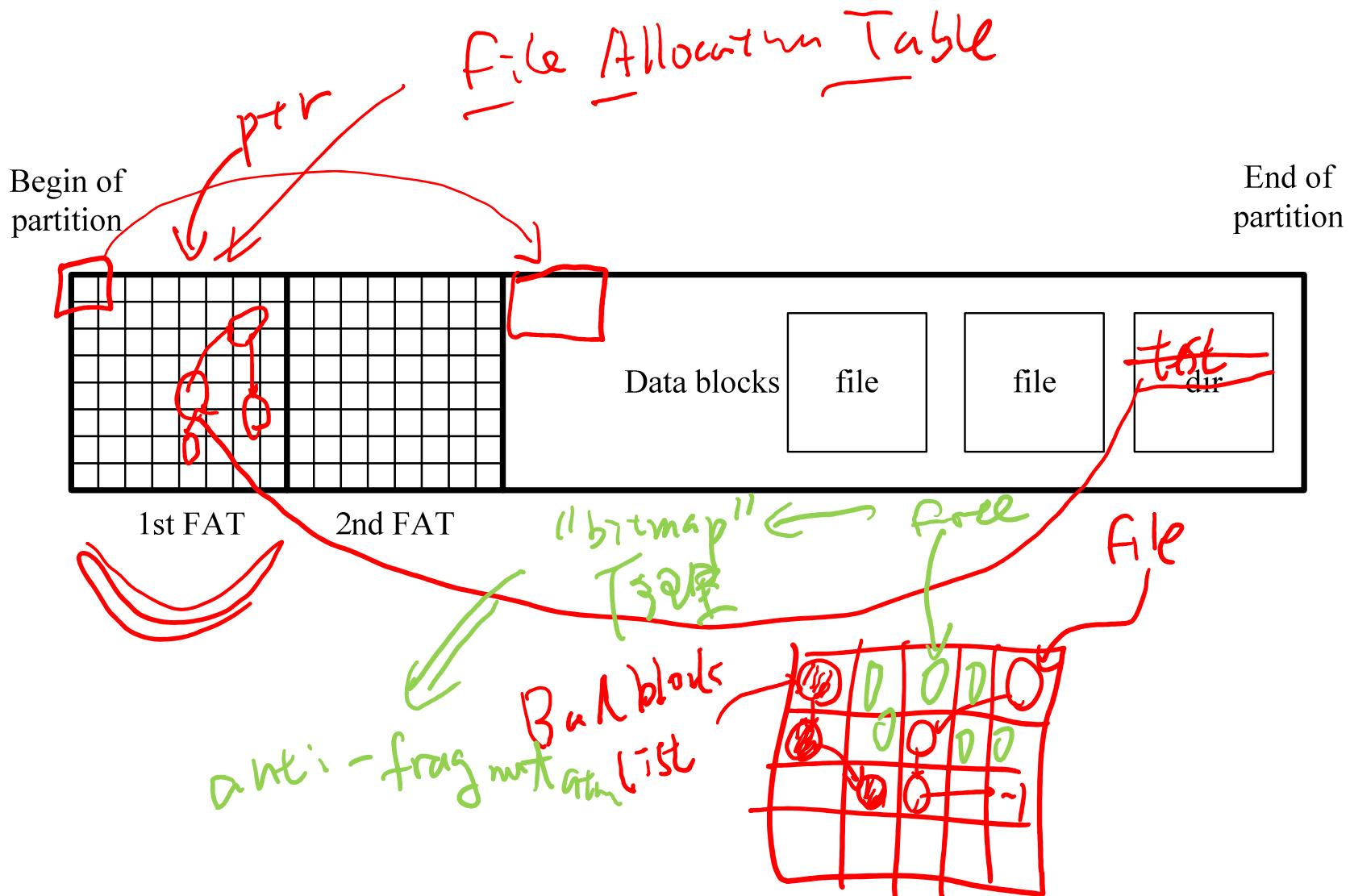
# Linked Allocation



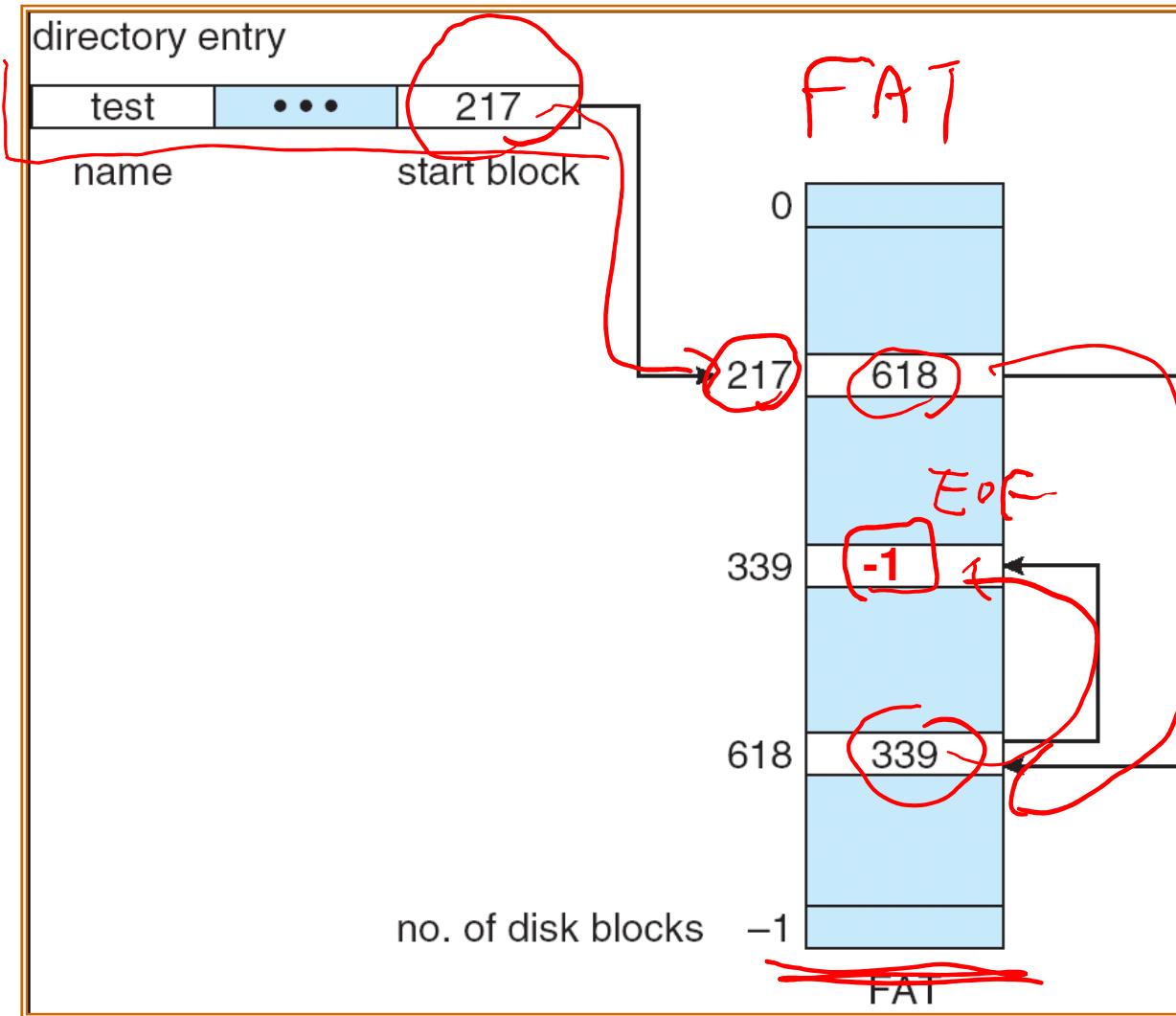
# Linked Allocation

- 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



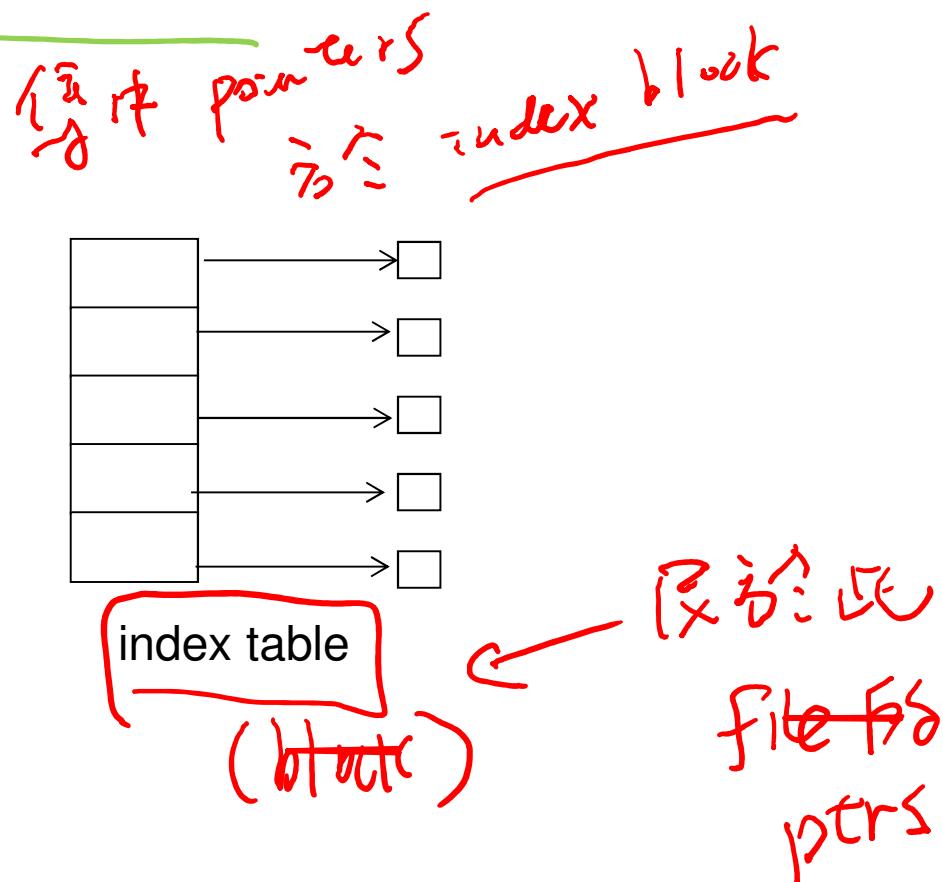
# File-Allocation Table



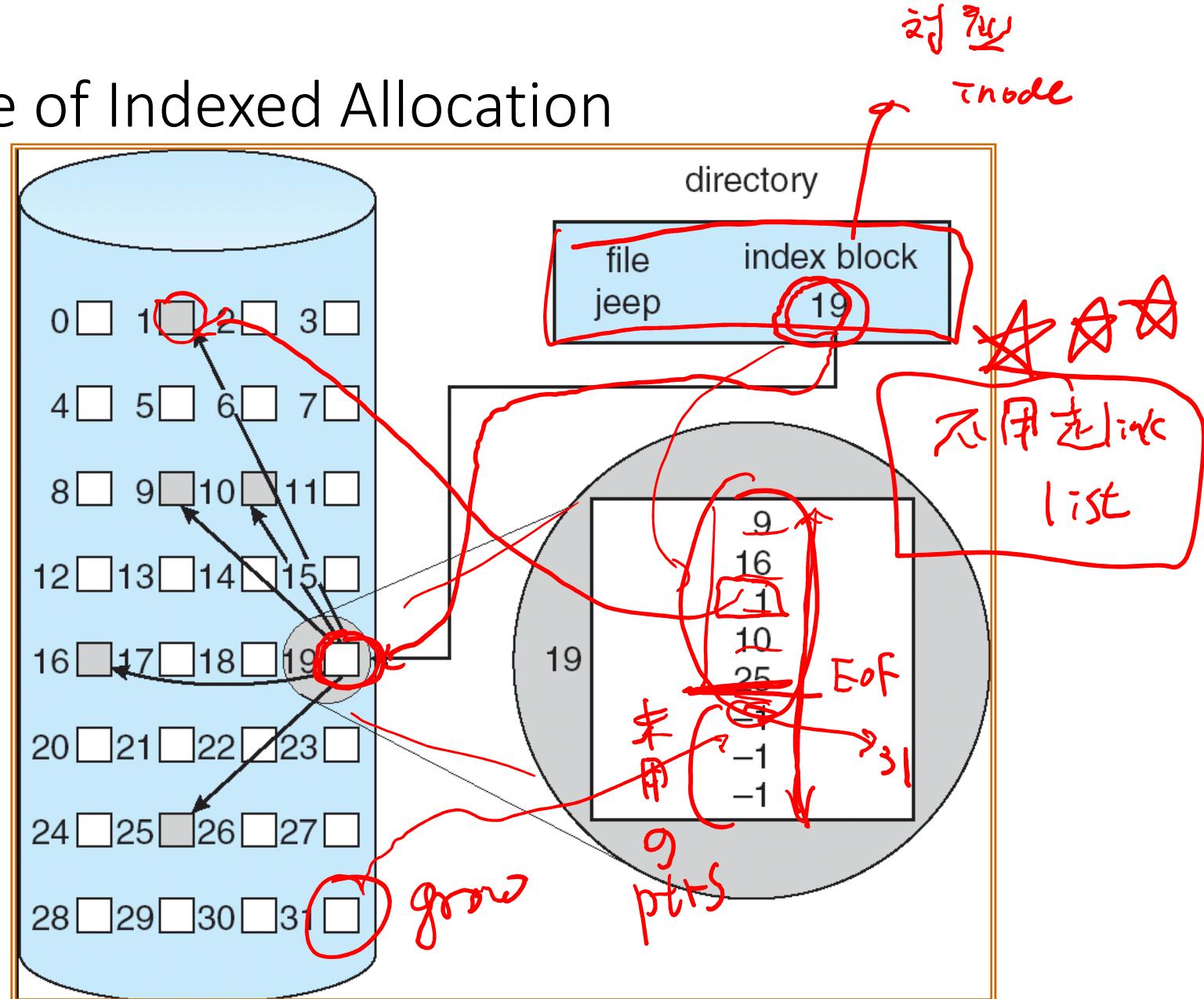
- A bad list maintains bad clusters
- Scans 0 for unallocated clusters

# Indexed Allocation

- Brings all pointers together into the index block.
  - Logical view.  
A diagram illustrating the logical view. It shows several red arrows originating from different parts of the text and pointing towards a central blue box labeled "Index block". The text "Logical view." is written in black, and the word "block" is written in red at the end of the list item.

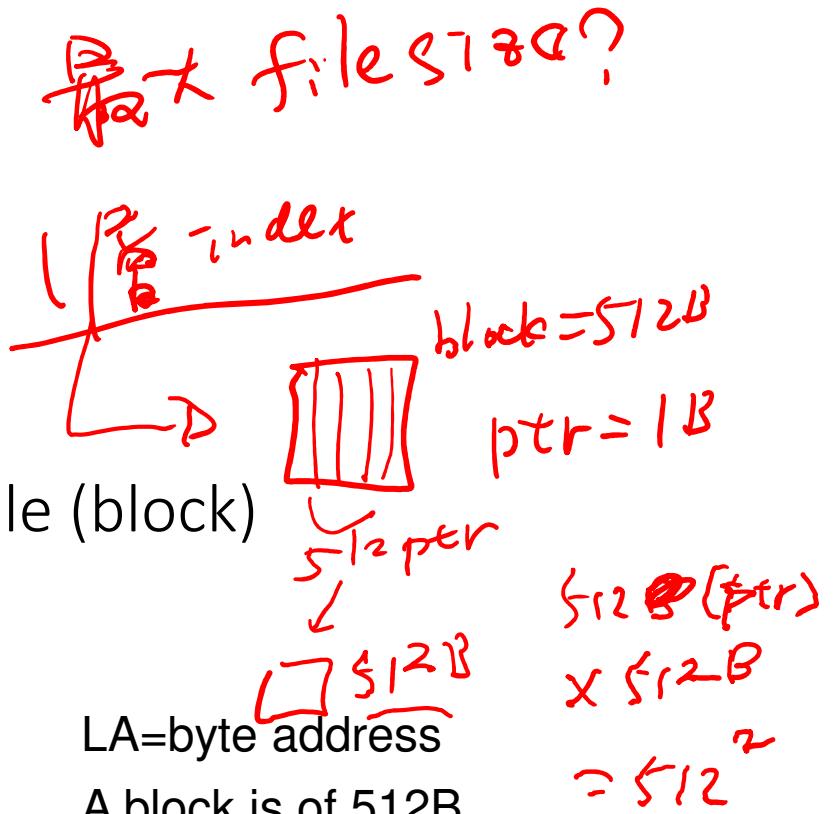


# Example of Indexed Allocation



## Indexed Allocation (Cont.)

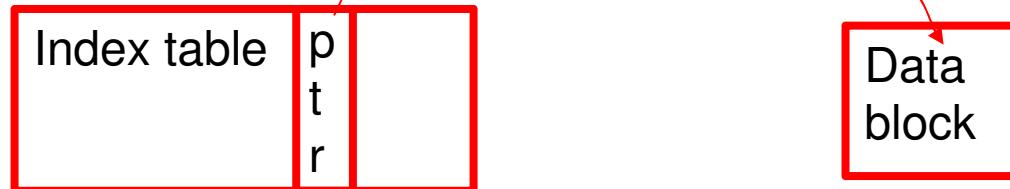
- Need a index table
- Capable of “random” access
- Per-file overhead of an index table (block)



Q = displacement into index table (entry #)

R = displacement into block

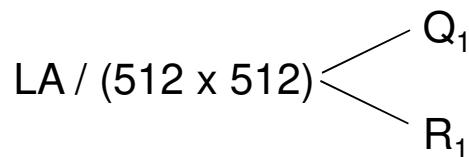
- 



# Indexed Allocation – Mapping

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

- 1<sup>st</sup> level: pointers to index tables
- 2<sup>nd</sup> level: pointers to data blocks



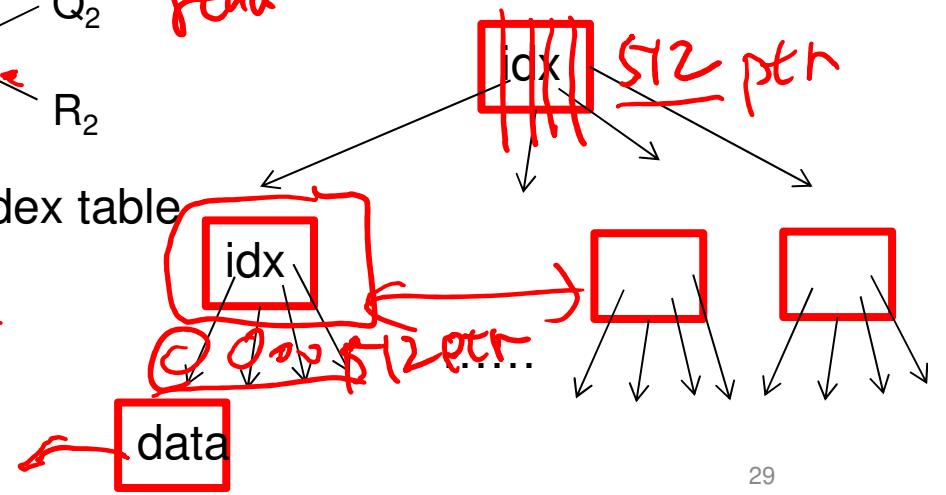
$Q_1$  = displacement into outer-index  
 $R_1$  is used as follows:



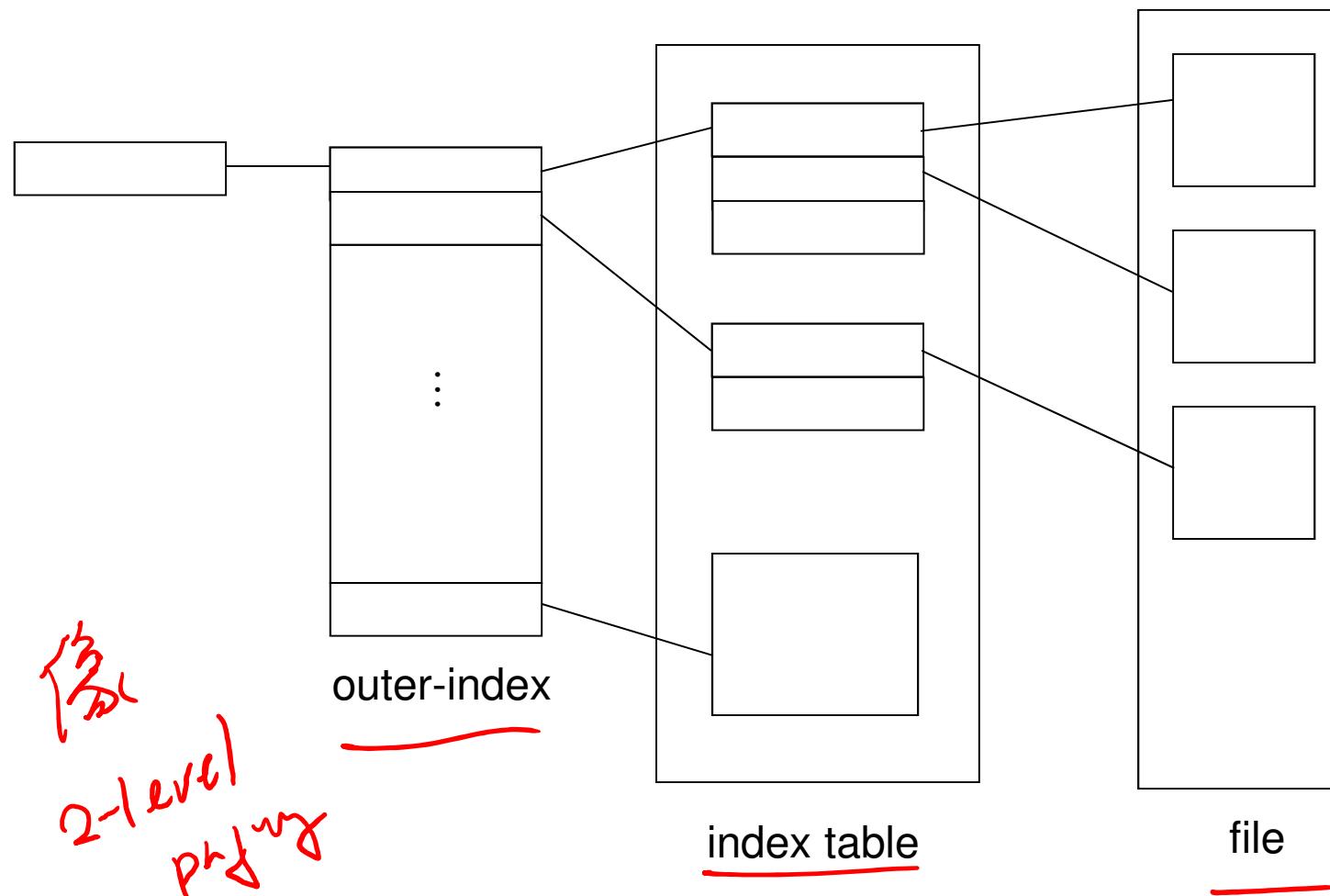
$Q_2$  = displacement into block of index table  
 $R_2$  displacement into block of file:

$\frac{1}{512 \times 512 \text{ ptr}} \times 512B$

Total addressable size =  $512 \times 512 \times 512$  bytes

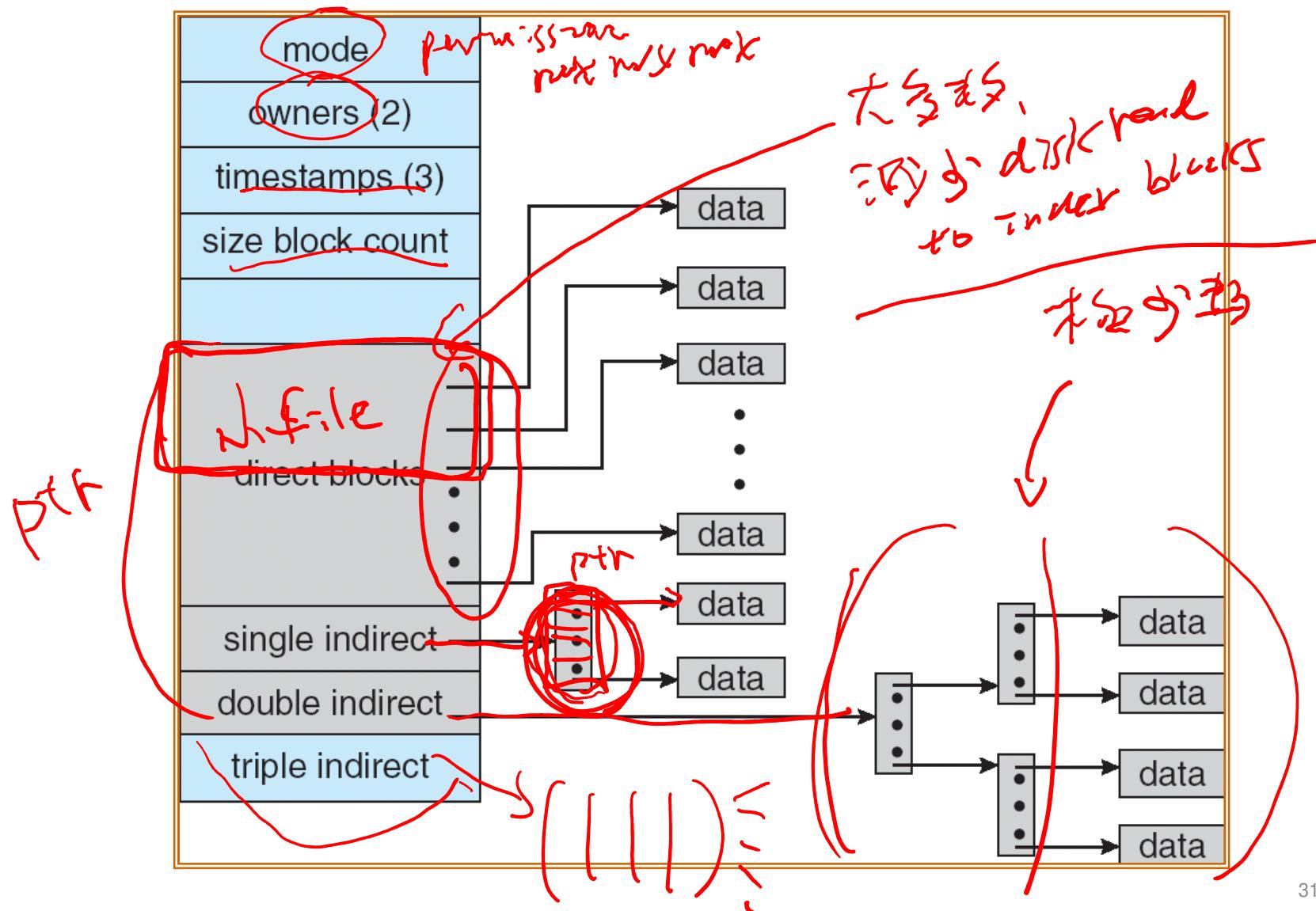


## Indexed Allocation – Mapping (Cont.)



# UNIX inode

An i-node. Small files use only direct blocks



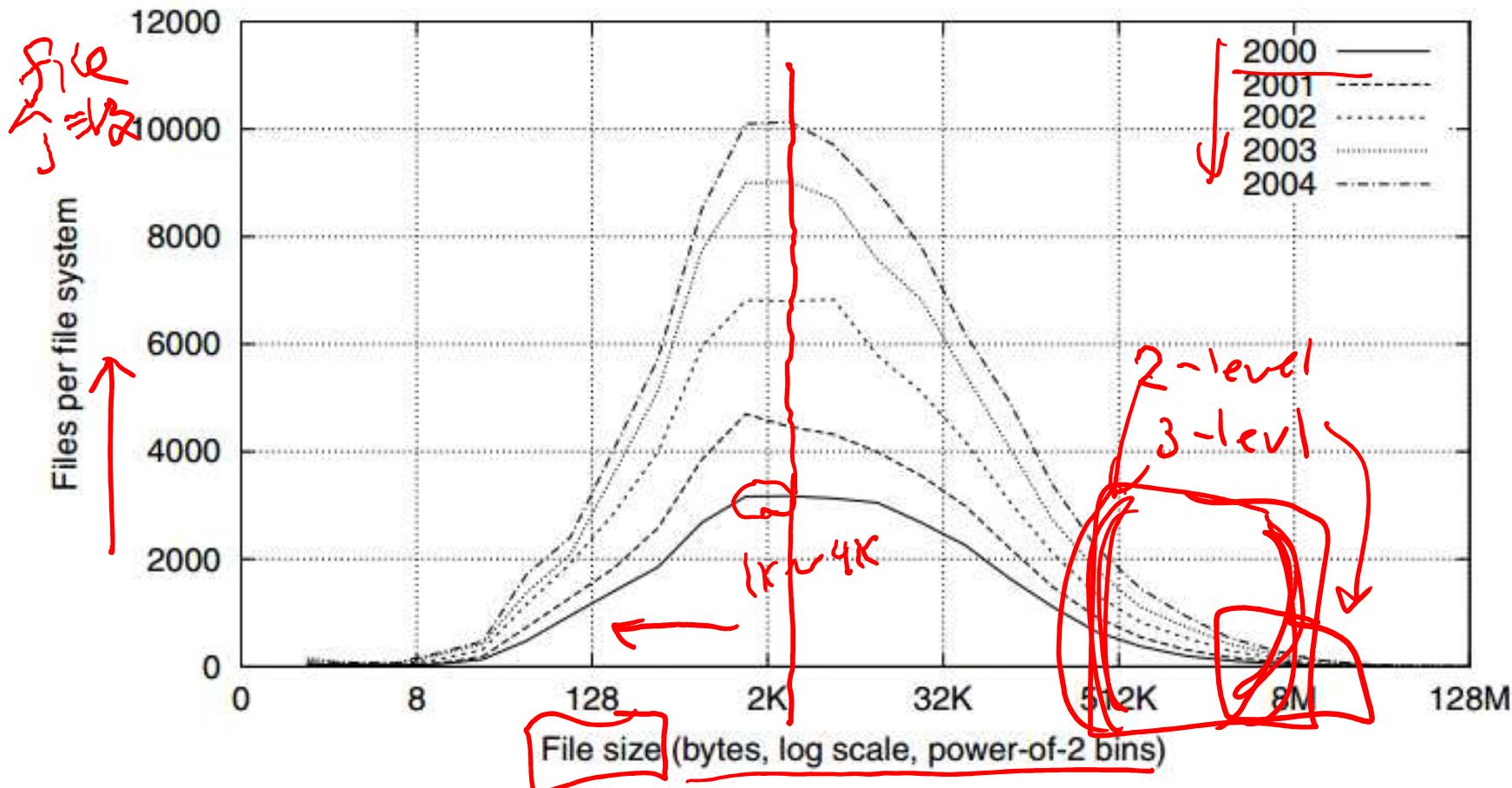
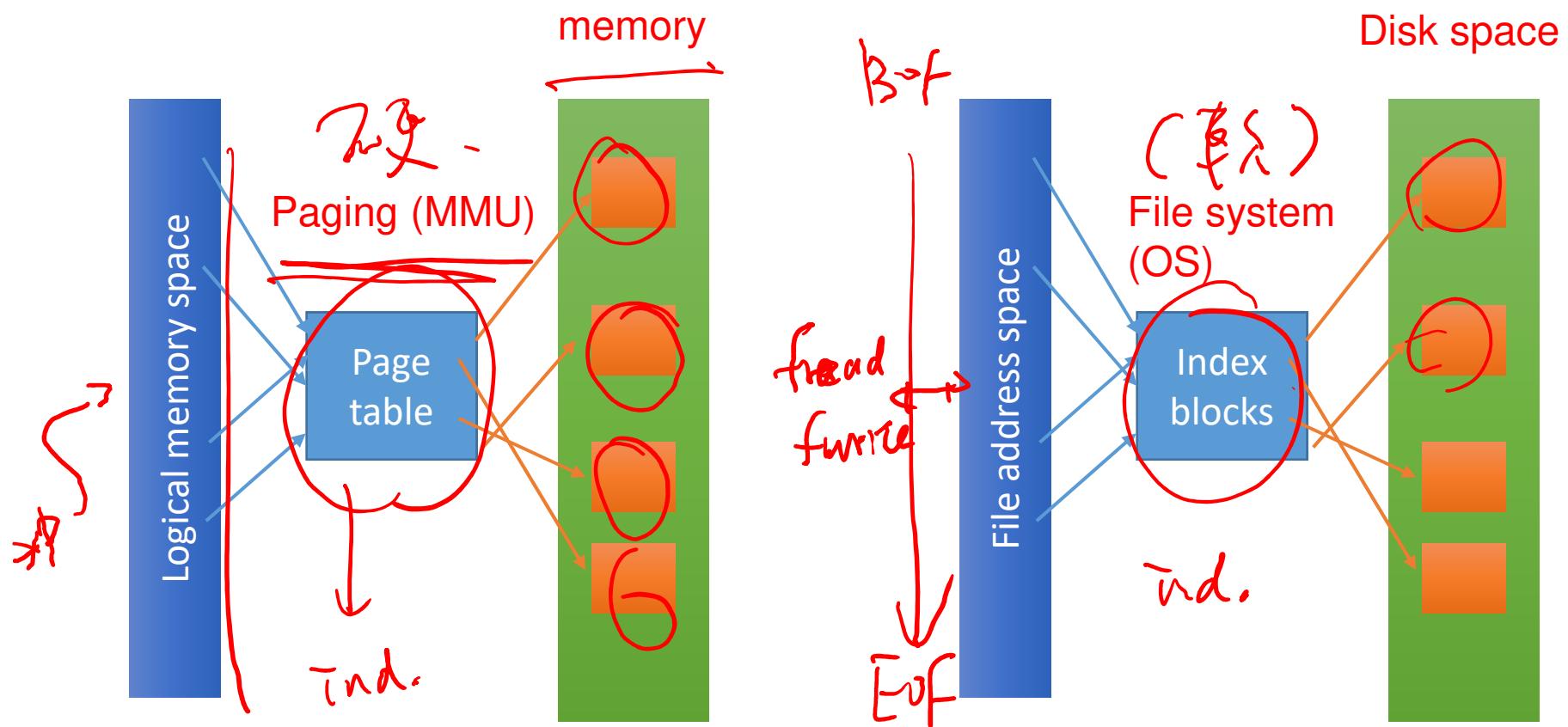


Fig. 2. Histograms of files by size.

[A. Agrawal, "A Five-Year Study of File-System Metadata"](#)

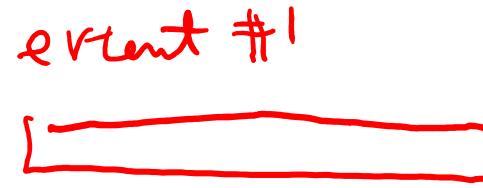
# Indirection, indirection, indirection ...



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

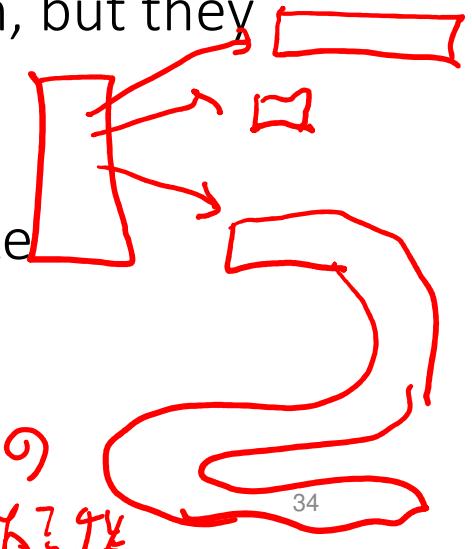
## Extent-Based Systems

- 空間の space



大小不一

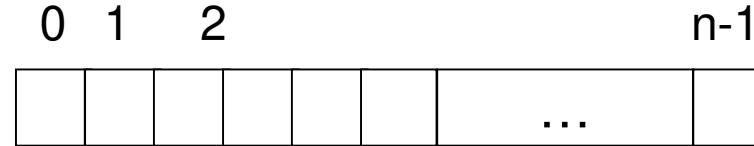
- A hybrid of contiguous allocation and linked/indexed allocation
- Extent-based file systems allocate disk blocks in  **extents**  希望.
- An extent is a set of contiguous disk blocks  **index** 
  - Extents are allocated upon file space allocation, but they are usually  **larger than**  the demanded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



領域の  
可及性

## Issue 3: Free-Space Management

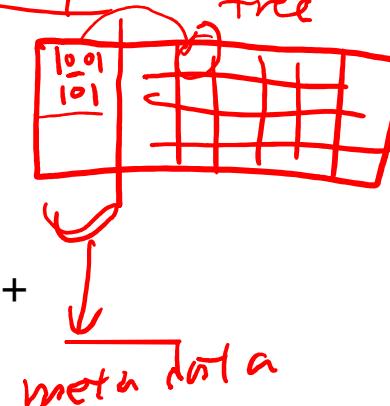
- Bit vector ( $n$  blocks)



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

Block number calculation

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



- First check whether a DWORD is not 0xffffffff
  - If not, scan for the zero bits

## Free-Space Management (Cont.)

- Bit map requires extra space
  - Example:

$$\begin{aligned} \text{block size} &= 2^{12} \text{ bytes} \\ \text{disk size} &= 2^{30} \text{ bytes (1 gigabyte)} \\ n &= 2^{30}/2^{12} = 2^{18} \text{ bits (or 32K bytes)} \end{aligned}$$

*space*  
*4KB*      *2<sup>15</sup> Bytes*      *= 32KB*      *32MB*  
*1TB?*

- Scanning for 0's to find free blocks

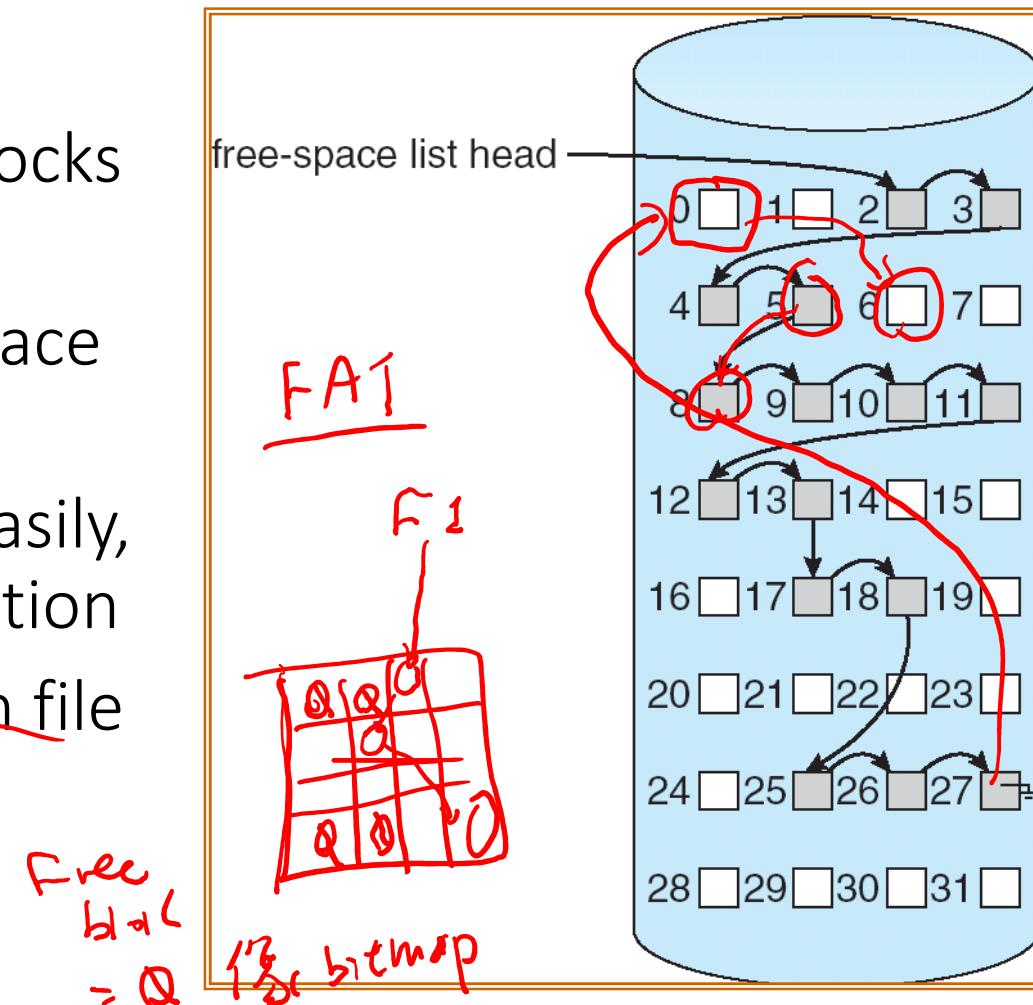
- Easy to get contiguous files

- Check whether a **DWORD** is zero (0x00000000)

- Used by UNIX FFS, Ext family, ...

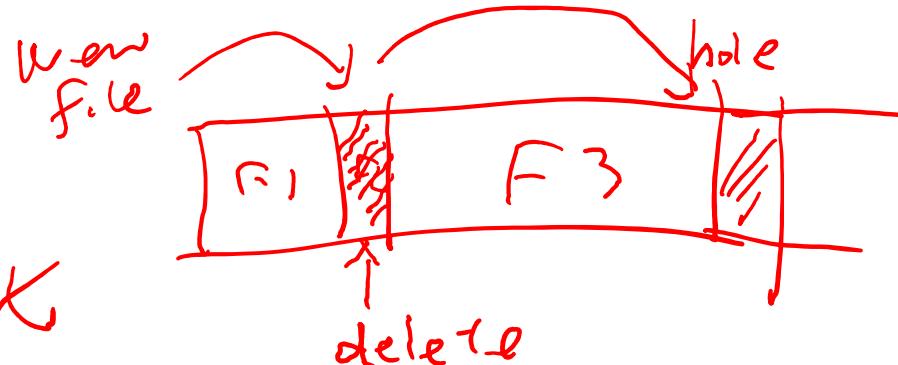
# Linked Free Space List on Disk

- 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



## File Fragmentation

FS  
文件系统  
老化和碎片化

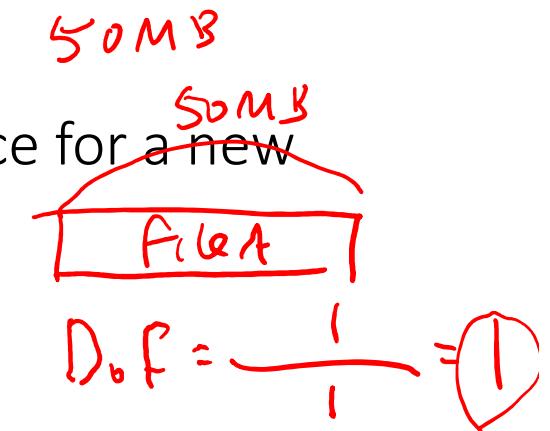


- File system "ages" after many creation and deletion of files

- 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

### 0 Degree of Fragmentation (DoF) of a file



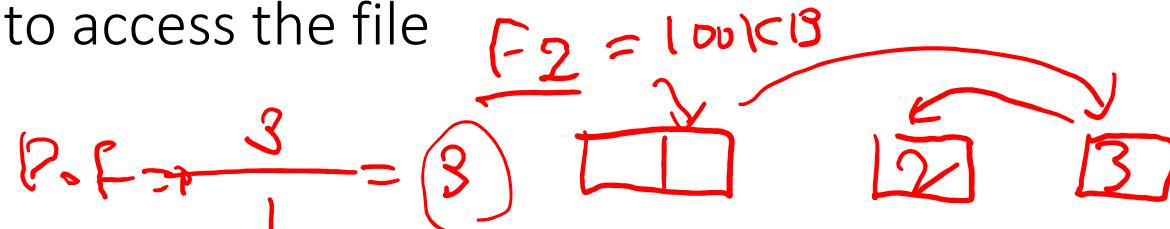
defragmentation

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

ext4  
Max. extent 128MB

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

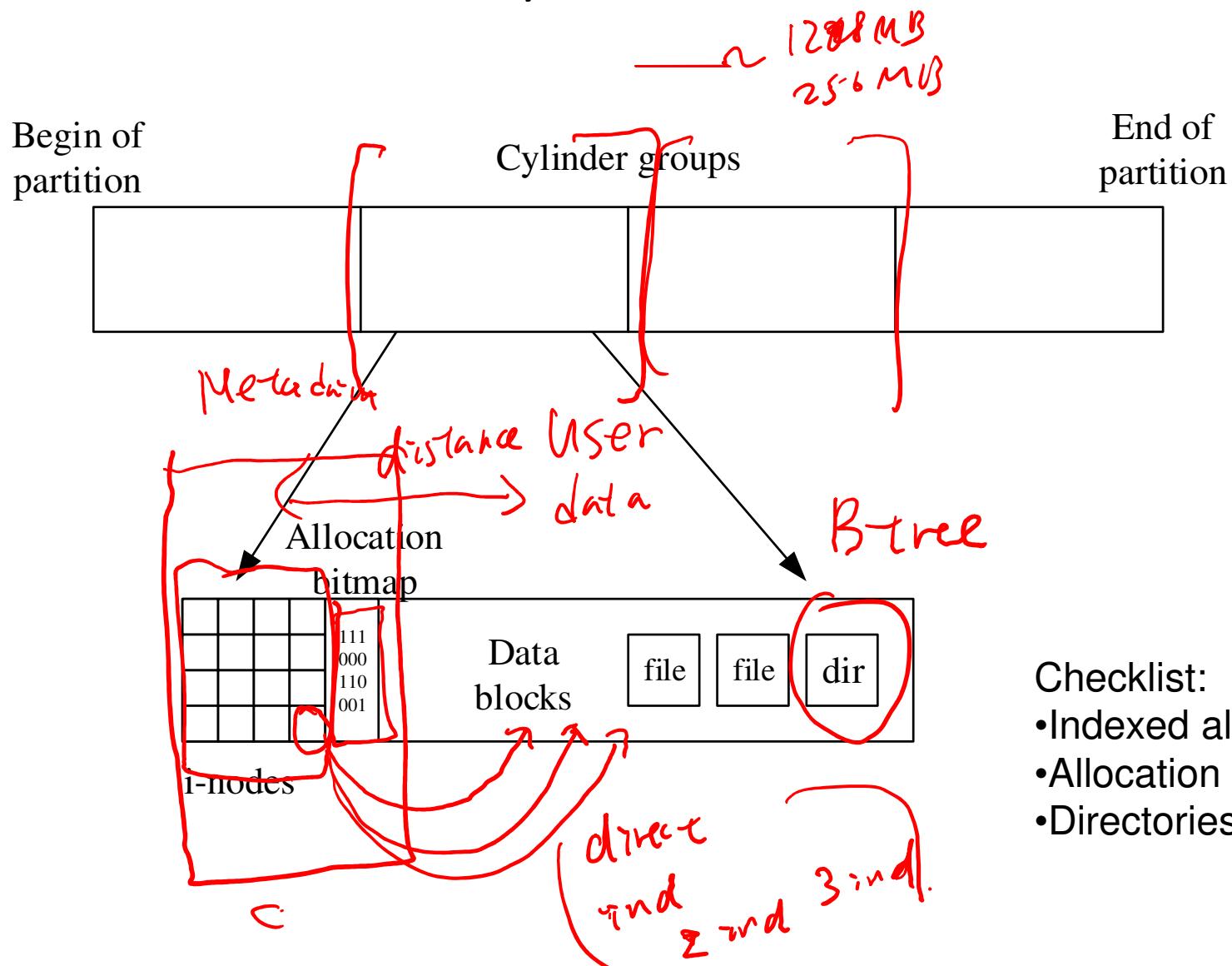
e4 defrag



# Comparison

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

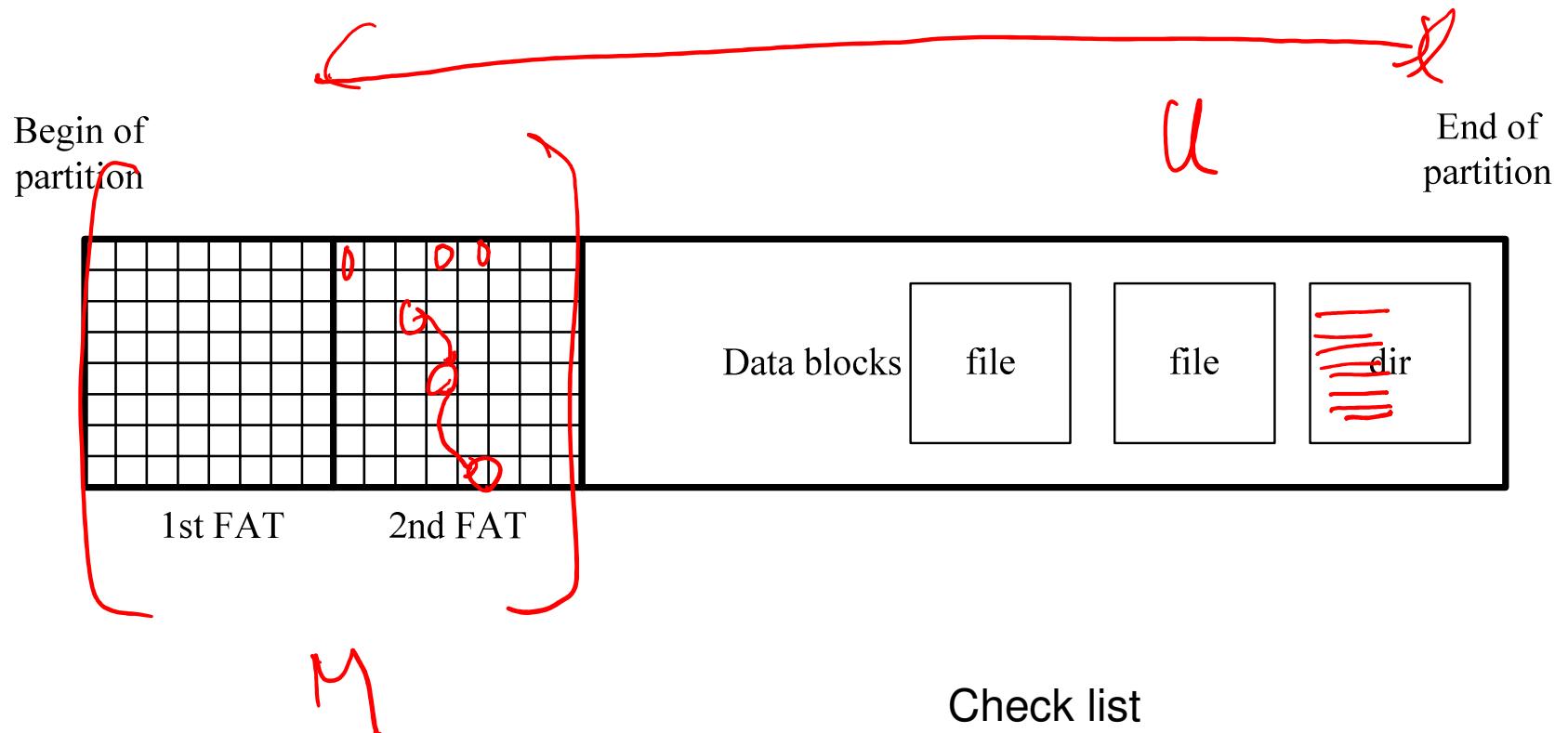
# Review: ext4 file system



## Checklist:

- Indexed allocation
- Allocation bitmaps
- Directories (H-tree)

# Review: FAT file system

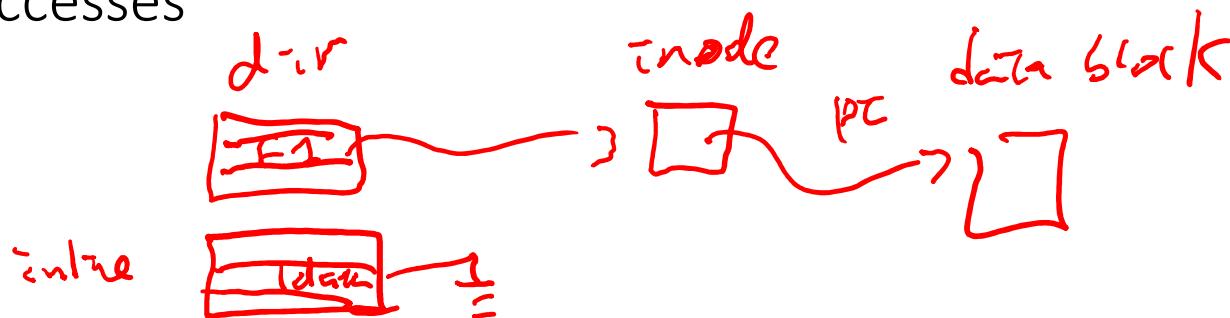


## Check list

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

# Efficiency and Performance

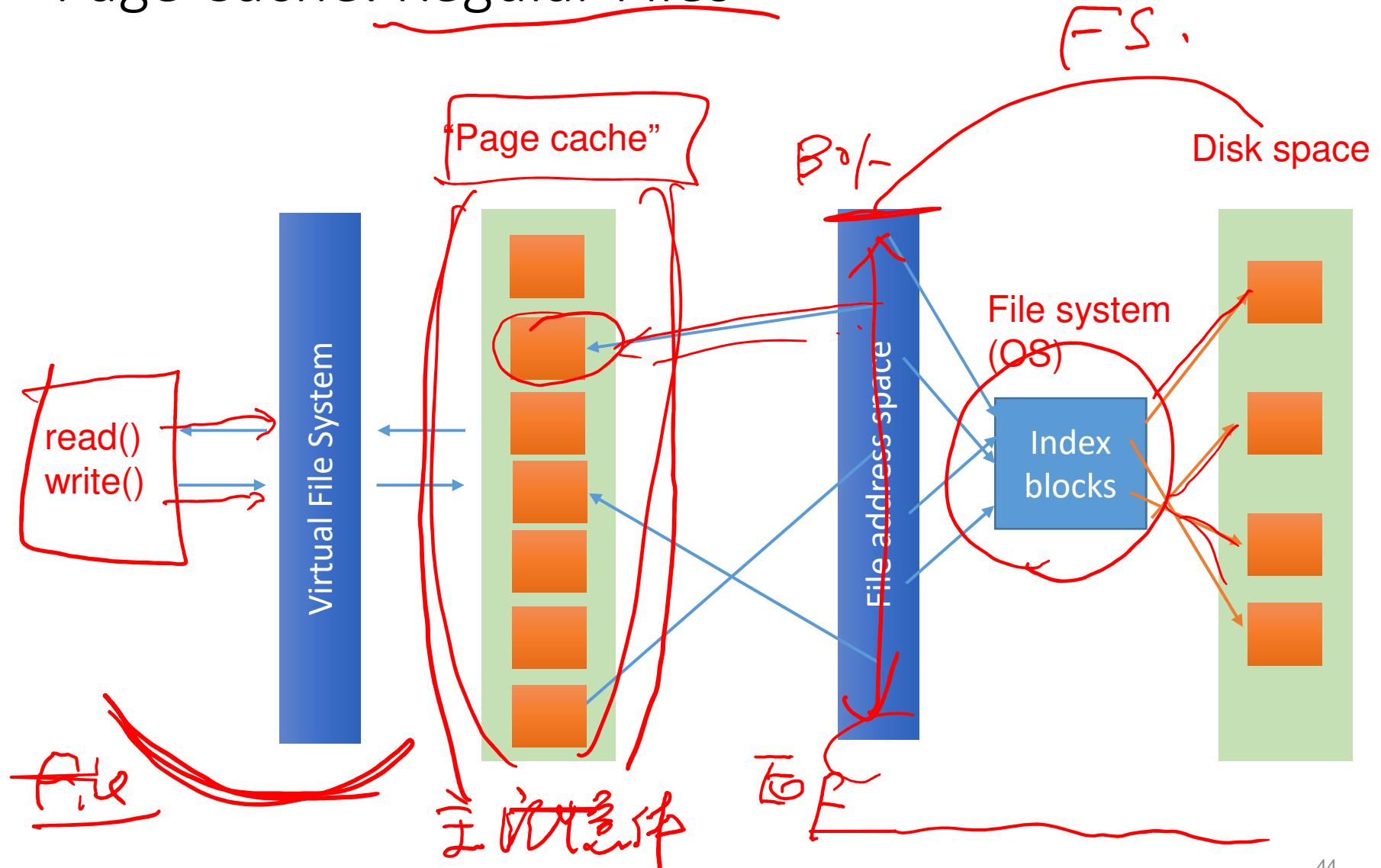
- File systems have their unique techniques for performance optimization
- For example, ext4-specific optimizations include:
  - Dividing disk space into cylinder groups to make inodes appear near to their associated data blocks inline files
  - Embedding small files into directories (<60 bytes)
  - Using extents to take advantage of sequential disk accesses



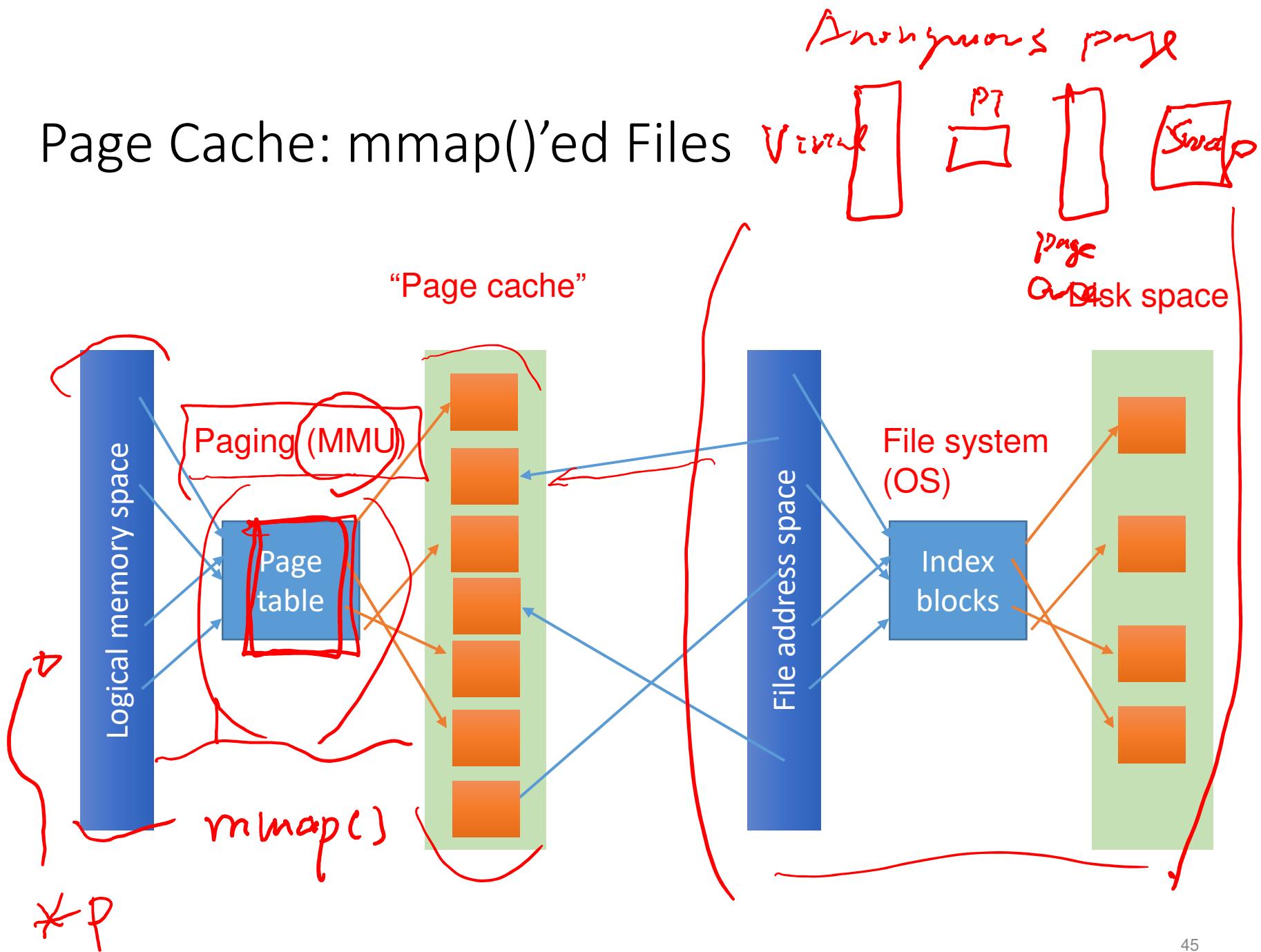
# Efficiency and Performance

- OS kernel also provides **generic** optimization methods that all file systems can use
  - **Disk cache (page cache)** – separate section of main memory for frequently used blocks (temporal locality)
  - Read-ahead (prefetch) – technique to optimize sequential access
    - Exploiting spatial locality of file access; like pre-paging

## Page Cache: Regular Files



# Page Cache: mmap()'ed Files



## Recovery

⇒ DRAM (DRAM)  
⇒ 持久性

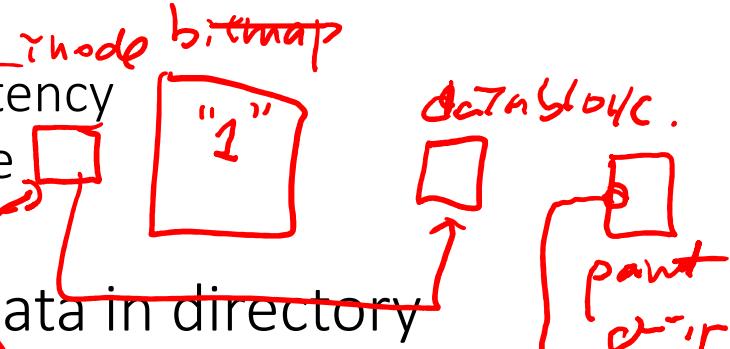
- A file operation involves **multiple block modifications**

- 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 are lost

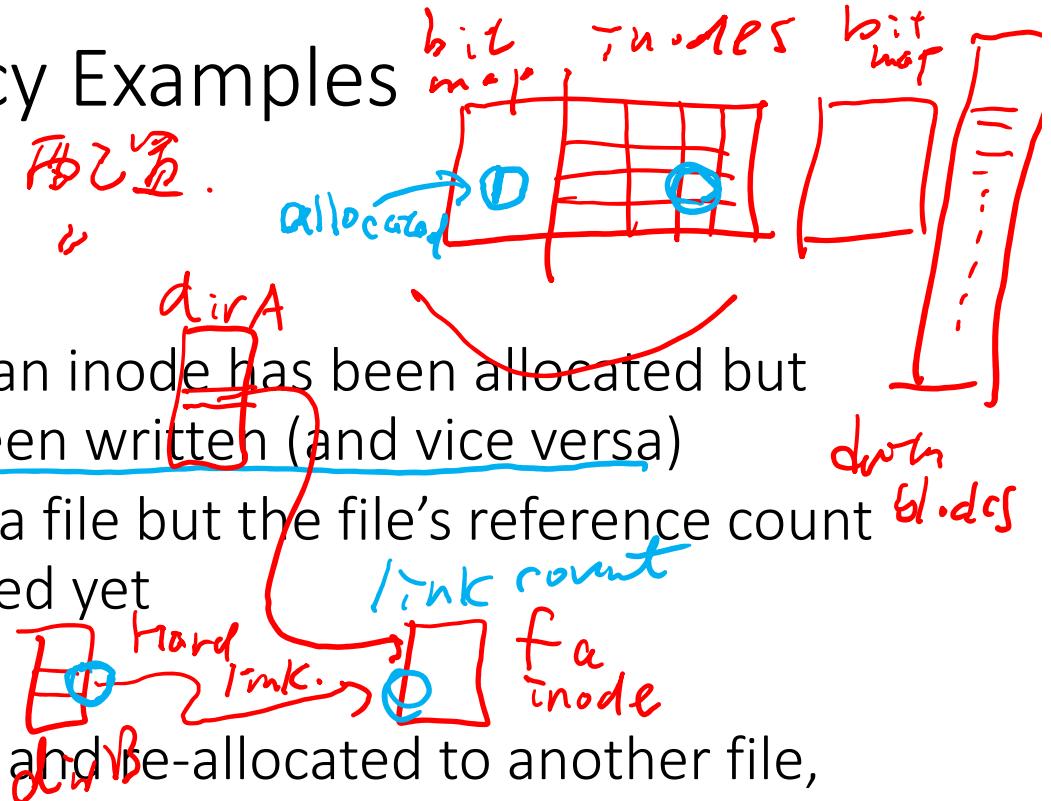
- Loss of metadata: structural inconsistency
- Loss of user data: partially written file

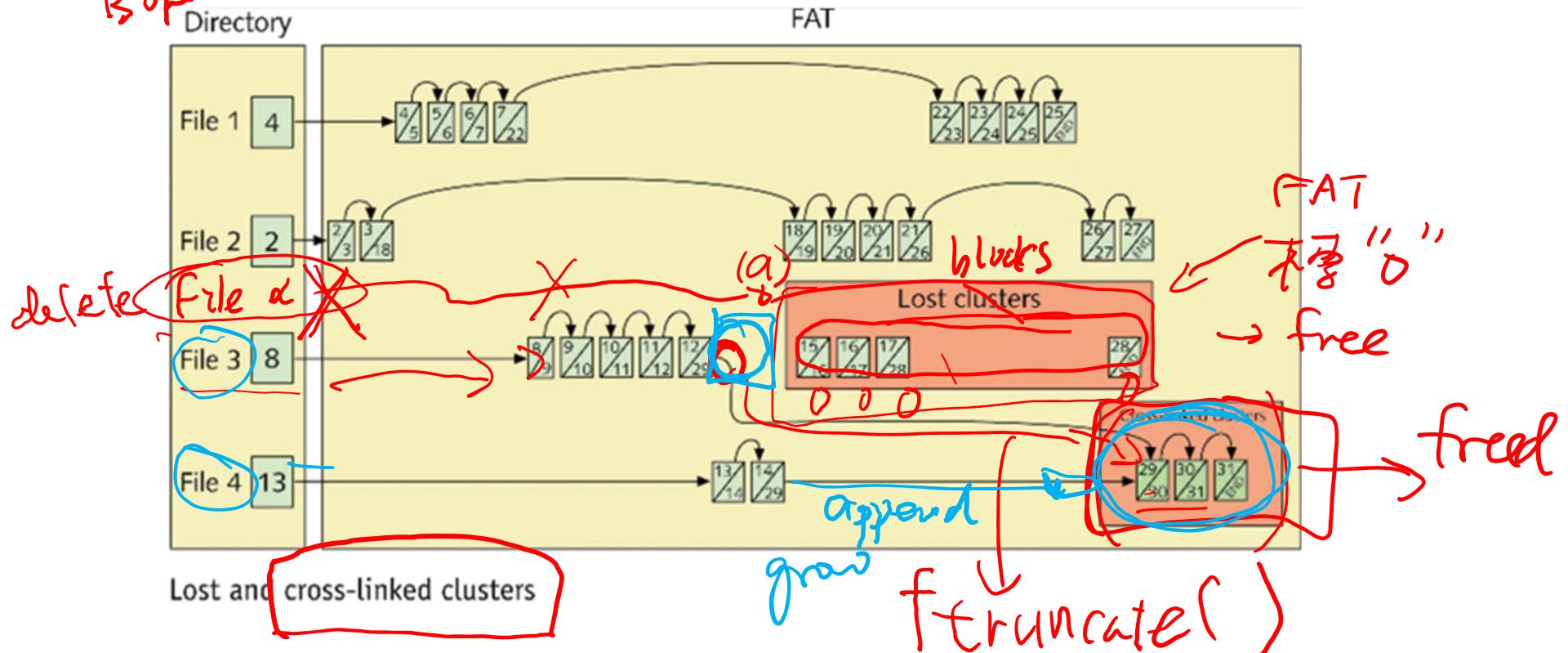
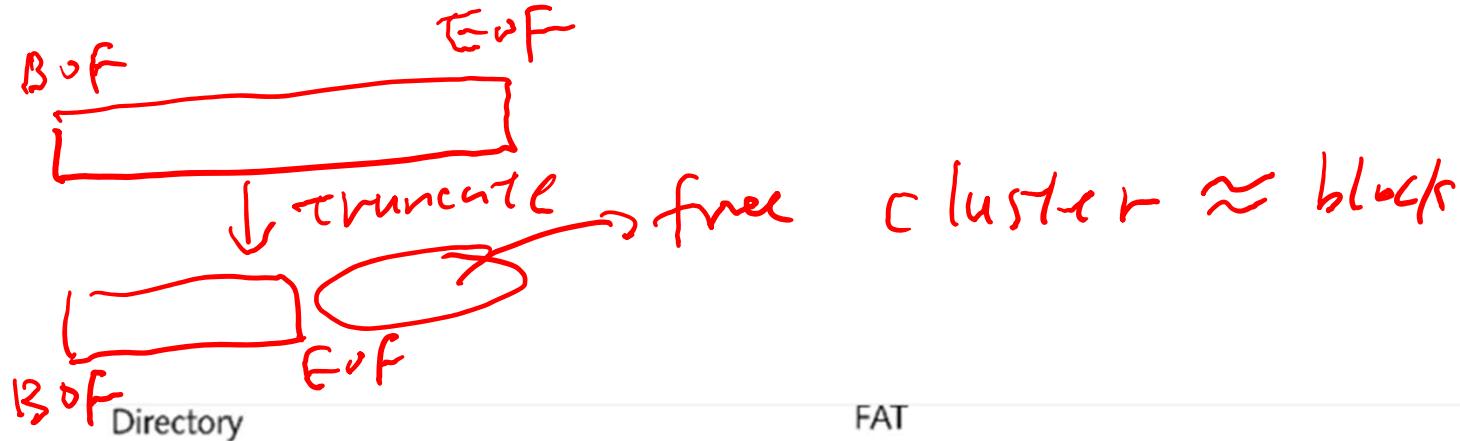
- Consistency checking – compares data in directory structure with data blocks on disk, and tries to fix inconsistencies



## 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)





[http://faculty.salina.k-state.edu/tim/ossg/File\\_sys/file\\_system\\_errors.html](http://faculty.salina.k-state.edu/tim/ossg/File_sys/file_system_errors.html)

## Recovery Utilities

↳ `mount (RW)`

`i=dirty`

Superblocks

- 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

# Journaling File Systems

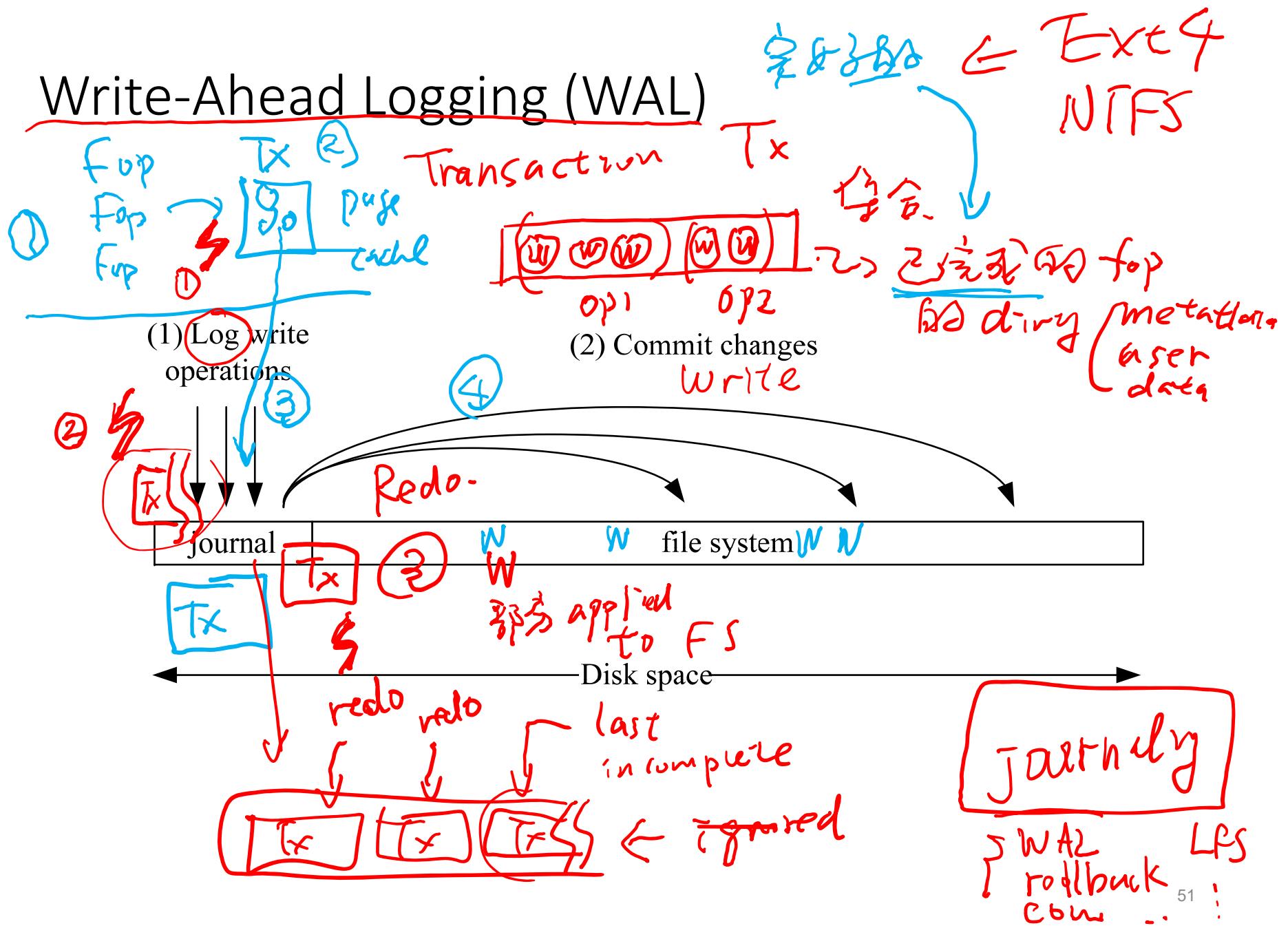
- Guarantee the atomicity of file system operations
  - Atomicity: all or none
  - Structural consistency
- Journal is a reserved disk space
- Journaling file systems collect the dirty metadata and/or dirty user data from completed file system operation(s) into a transaction
  - Write transactions to the journal first, and then modify the file system
- If system crashes, upon system reboot, the file system re-does the transactions in the journal
  - partial (incomplete) transitions will be discarded

FS → disk blocks

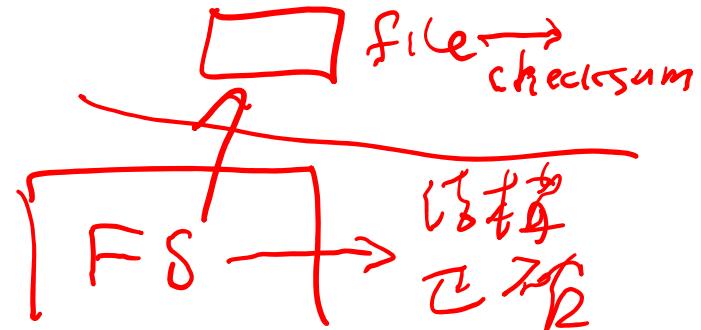
1  
2

步骤

## Write-Ahead Logging (WAL)



# Journaling File Systems



## ✓• Transactions

- An idea borrowed from database systems
- ACID properties (Atomicity, Consistency, Isolation, Durable)
- All or none (no partial)

## ✓• Journaling

- Based on write-ahead logging (WAL)
- Guarantees that file systems are structurally consistent
- Does not guarantee no loss of data
- When powering up after crash
  - Scan the journal
  - Found a complete transaction → redo
  - Found a partial transaction → discard

# Journaling File systems -- Summary

## ✓ Motivation

- Preventing power interruptions from corrupting file systems

## ✓ Method

- Adding a journal space to the file system
- Collecting a set of self-contained writes as a transaction
- Write transactions to the journal
- Apply changes to the file system (in background)
- On recovery, scan the disk journal. Re-do legit transactions; incomplete transactions will be discarded

## ○ Benefit

- Crash recovery is very fast

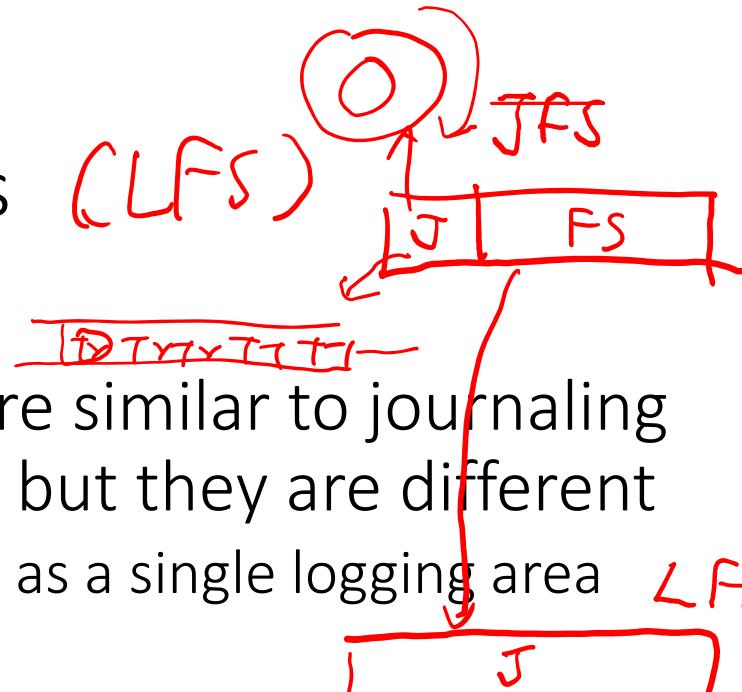
## ■ Problem

- Degraded performance due to amplification of write traffic

↳ FS write in background  
↳ Journal metadata

# Log-Structured File Systems

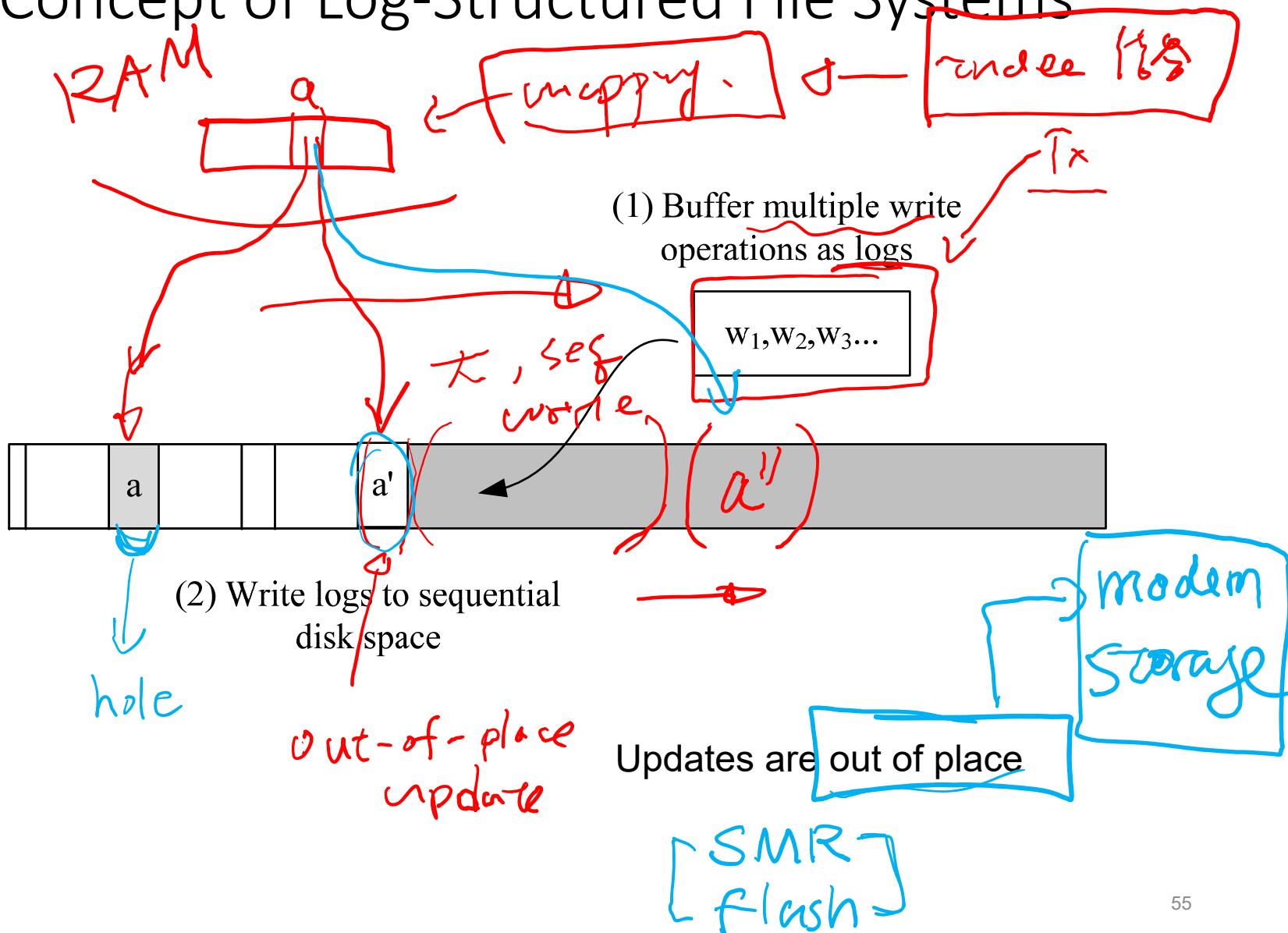
11: 12 am



- Log-structured file systems are similar to journaling file systems in many aspects; but they are different
  - LFSs treat the entire disk space as a single logging area
  - No need to “copy back”
- The primary objective is to optimize random write
  - Convert random writes into sequential writes through out-of-place updates
  - Read performance optimization is left to page caching
- Examples
  - NILFS2
  - F2FS for Android devices

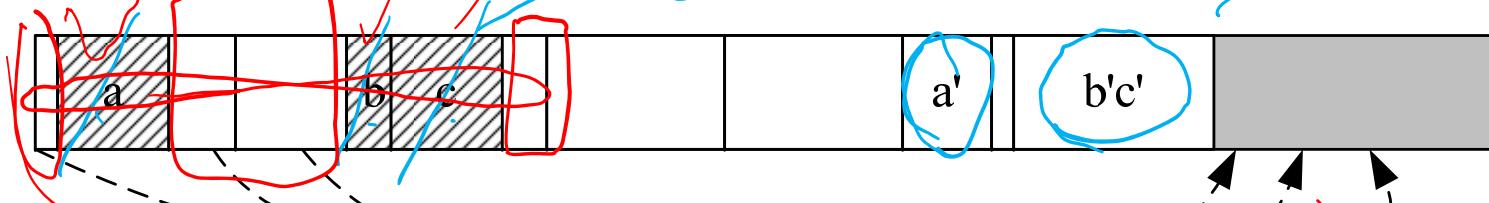


# The Concept of Log-Structured File Systems

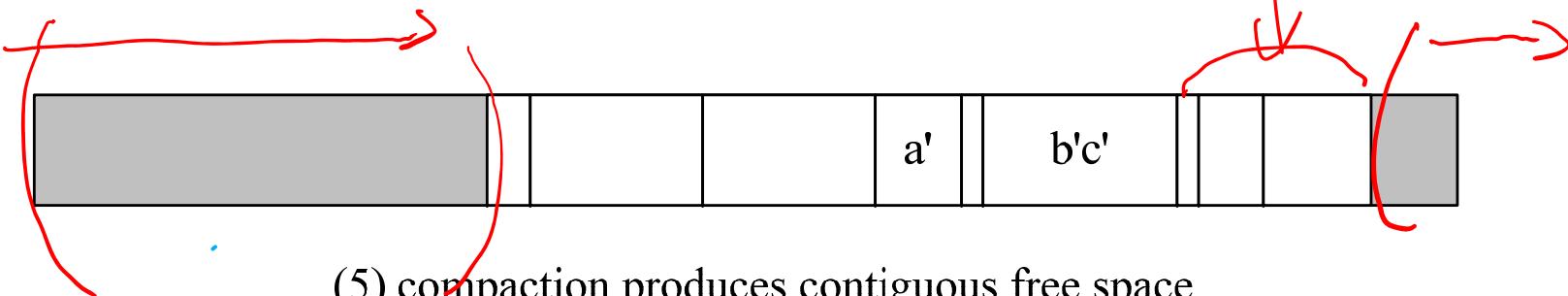


# Compaction (Garbage Collection) in LFS

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

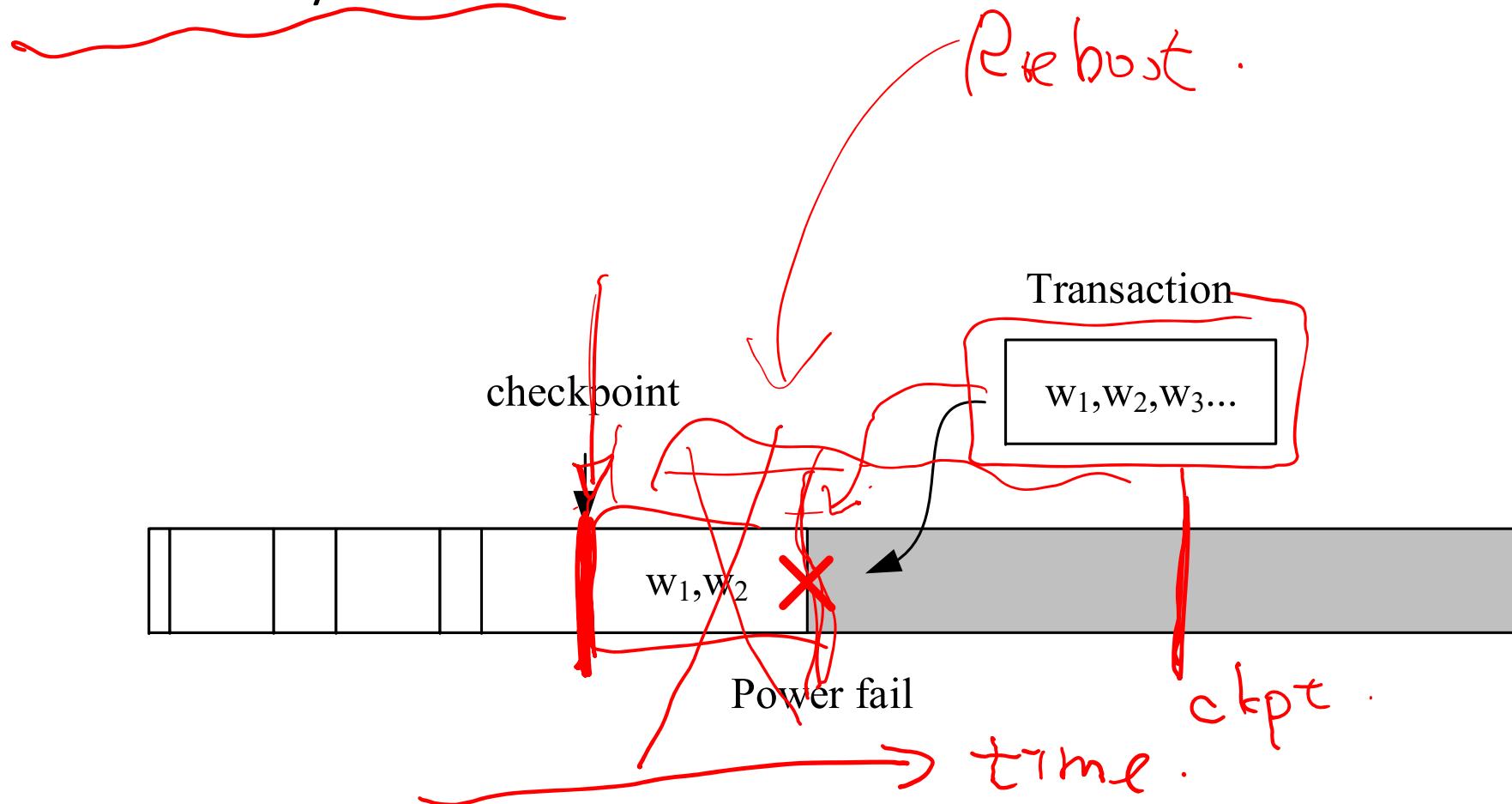


(4) Reclaim contiguous disk space with compaction (garbage collection)



(5) compaction produces contiguous free space

## Recovery in LFS



Surprisingly easy because writes are ordered chronologically

# Log-Structured File Systems -- Summary

- Motivation:
  - RAM will be cheap and random reads are efficiently handled by a large page cache
  - Random writes arrive at disk eventually and they are slow
- Methods:
  - Collecting random writes into long write bursts (logs)
  - Out-of-place updates
- Benefits:
  - Optimized random write performance
  - Easy recovery
- Problems:
  - Need compaction (garbage collection) to produce sequential space for new writes

Ind. ← ~~RAM~~



# End of Chapter 11