



Amirkabir University of Technology
(Tehran Polytechnic)
Department of Computer Engineering

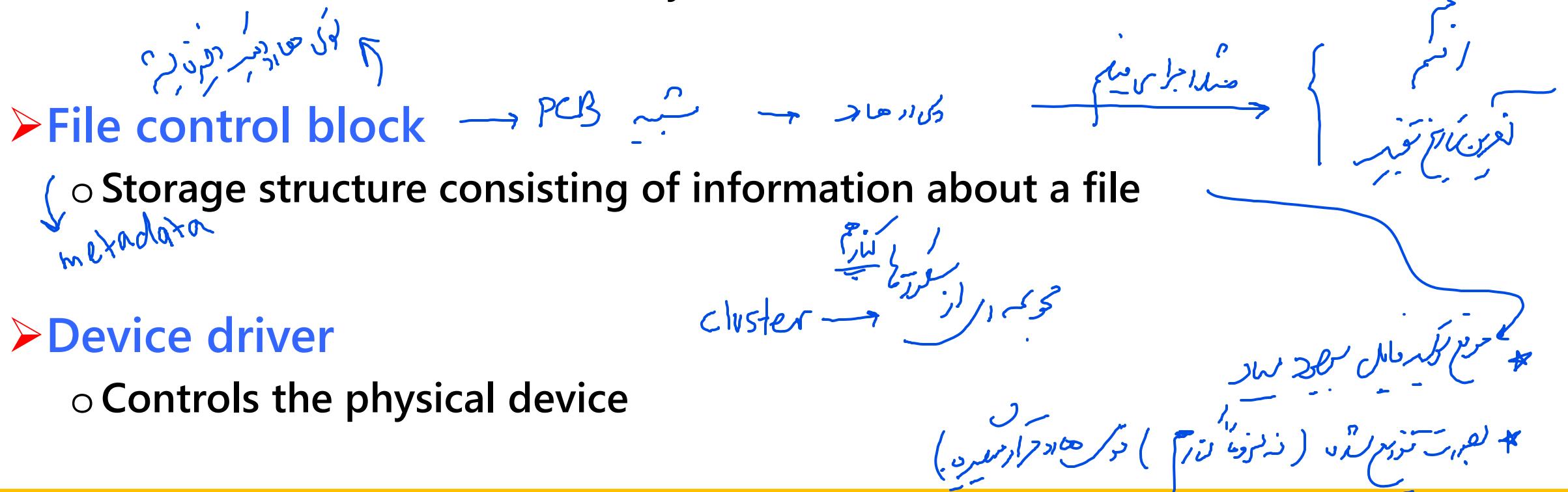
File System Implementation

Hamid R. Zarandi
h_zarandi@aut.ac.ir

File system structure

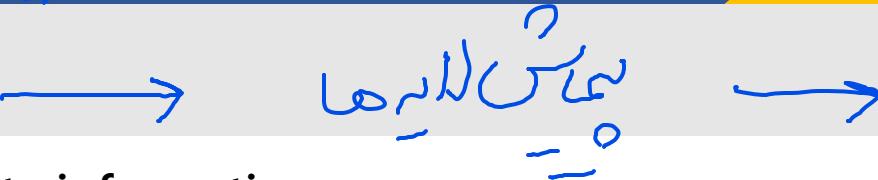
➤ Why file system?

- Provided user **interface to storage**, mapping **logical to physical**
- Provides **efficient** and **convenient** access to **disk** by allowing data to be stored, located, retrieved easily



- Controls the physical device

Layered file system



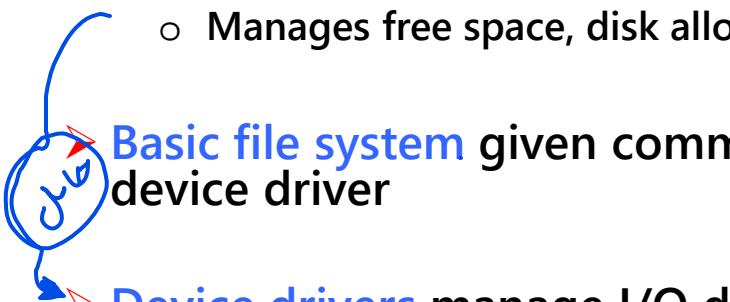
➤ Logical file system manages metadata information

- Translates file name into **file number**, **file handle**, **location** by maintaining file control blocks (**inodes** in UNIX), **FCB** and **metadata**
- Directory management
- Protection



➤ File organization module understands files, logical address, and physical blocks

- Translates logical block # to physical block #
- Manages free space, disk allocation

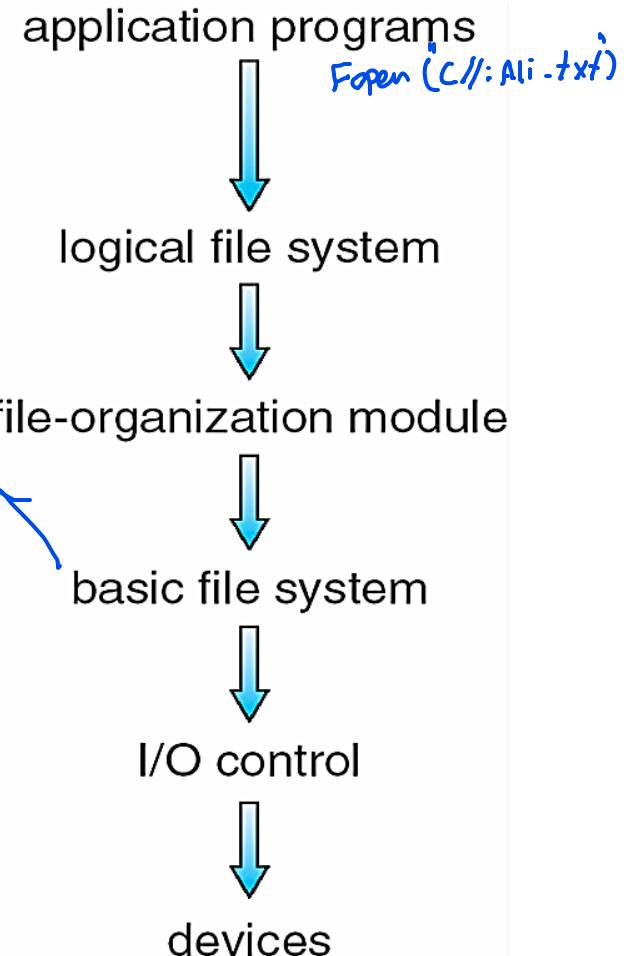


➤ Device drivers manage I/O devices at the I/O control layer

- Given commands like "read drive1, cylinder 72, track 2, sector 10, into memory location 1060" outputs low-level hardware specific commands to hardware controller
- Interrupt Service Routines

➤ Layered design

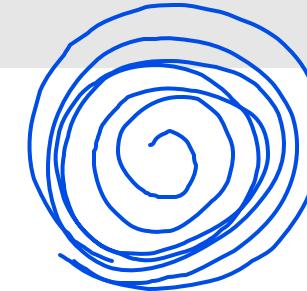
- Reduction of complexity and redundancy
- Overheads and performance penalty translation of file to disk



File system layers

➤ Each with its own format

- CD-ROM: ISO 9660
- Unix: UFS, FFS
- Windows: FAT, FAT32, NTFS as well as floppy, CD, DVD Blu-ray file systems
- Linux: more than 40 types (extended file system ext2 and ext3)
 - +distributed file systems
- New ones still arriving – ZFS, GoogleFS, Oracle ASM, FUSE



➤ File systems should have

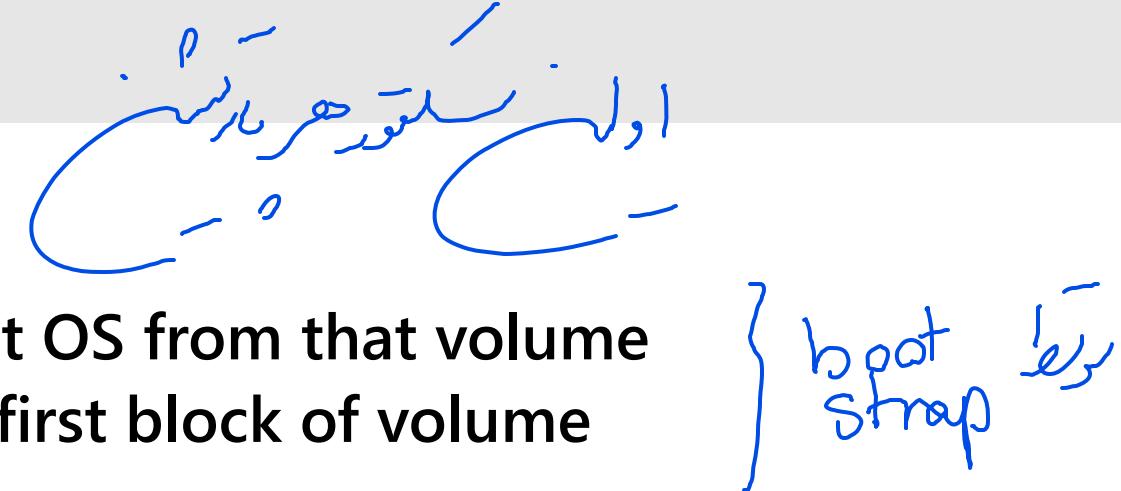
- 1) How to boot? (?)
- 2) Total # of blocks?
- 3) # and location of free blocks?
- 4) Directory structures?
- 5) Files?

این فایل سیستم ها باید داشته باشند.
۱) آنچه برای بوت کردن نیاز است.
۲) تعداد کلی بلک ها.
۳) مکان و تعداد بلک های آزاد.
۴) ساختار دیرکتوری.
۵) فایل ها.

File system implementation

➤ Boot control block (Boot sector)

- Contains info needed by system to boot OS from that volume
- Needed if volume contains OS, usually first block of volume



➤ Volume control block (superblock, master file table) (Master Boot Record)

- Contains volume details

- Total # of blocks, # of free blocks, block size, free block pointers or array

(0,0,0)

➤ Per-file File Control Block (FCB)

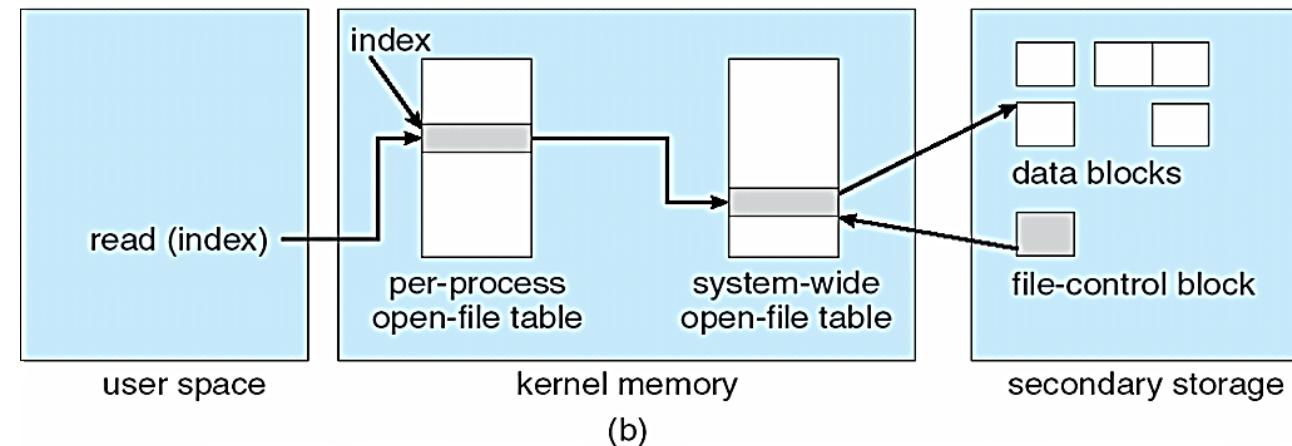
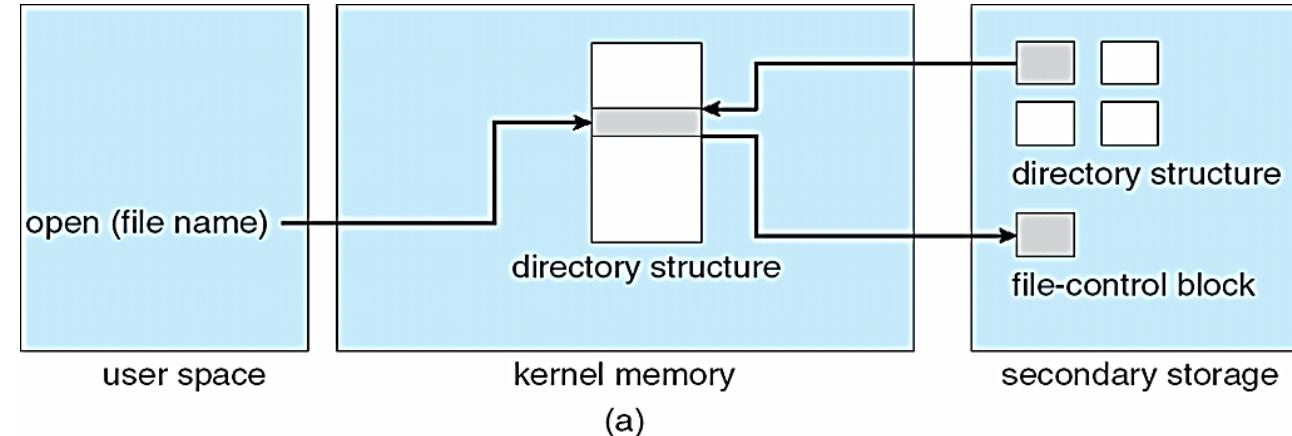
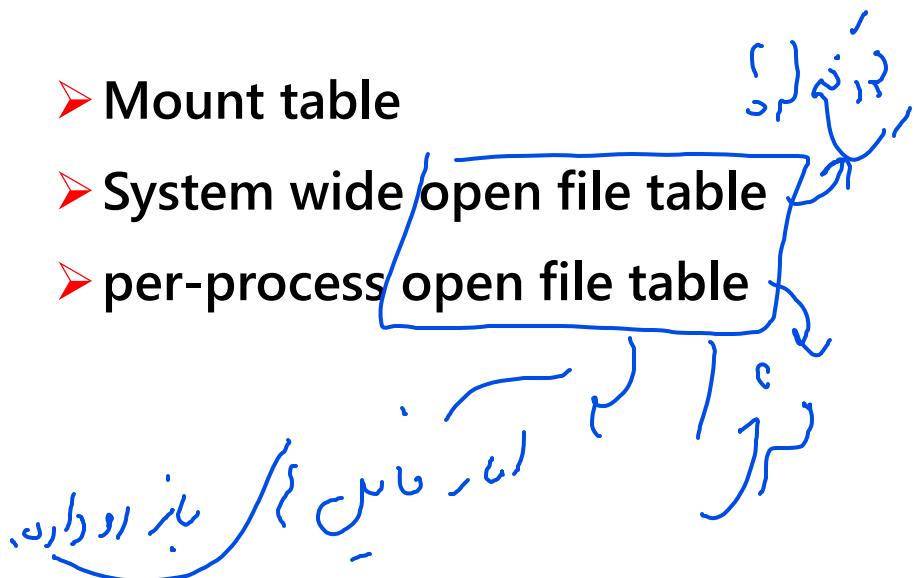
- Contains many details about the file

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



In-memory file system structures

- Mount table
- System wide open file table
- per-process open file table



Per process →

~ Kernel →

Partitions and mounting

RAM is used

➤ Boot loader

➤ Root partition

- Contains the OS
- Other partitions can hold other **Oses**, other **file systems**, or be **raw**
- Mounted at boot time
- **Other partitions can mount automatically or manually**

formatted

boot, swap, /, /home, /var, /tmp, /etc, /usr

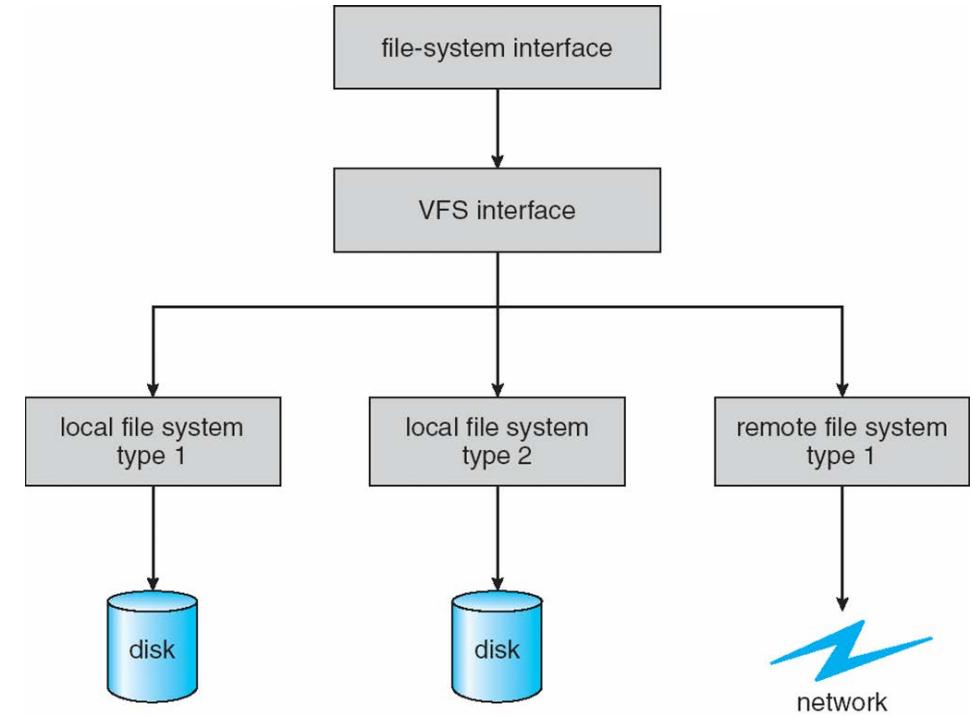
Virtual file systems (VFS)

➤ Virtual File Systems (VFS) on Unix

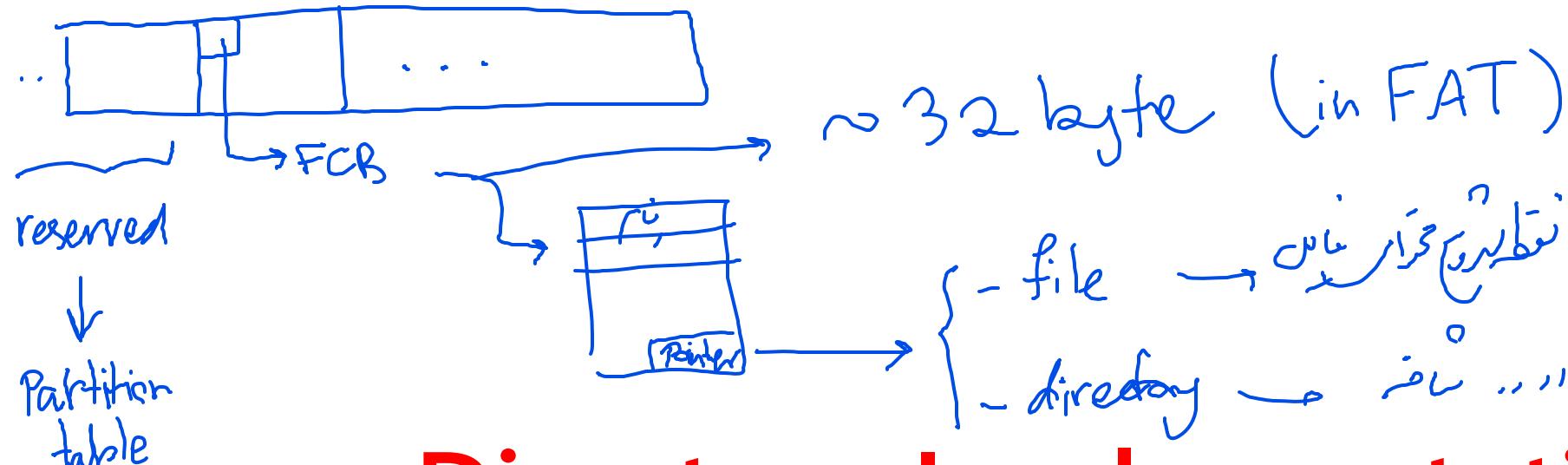
- Provide an **object-oriented** way of implementing **file systems**

➤ VFS allows the **same system call** interface (API) to be used for **different types of file systems**

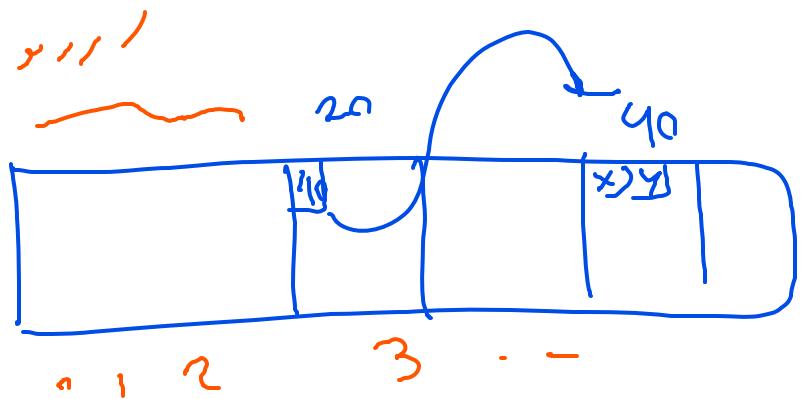
- Separates file-system generic **operations** from implementation **details**
- Implementation can be one of many file systems types, or network file system
- Then **dispatches** operation to appropriate file system implementation **routines**



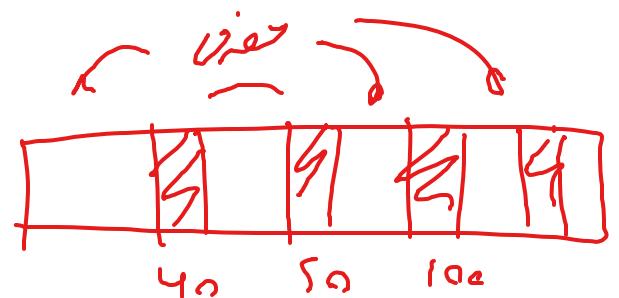
C:/ / 056 - 1
 سی دی اف اس ← پرینت



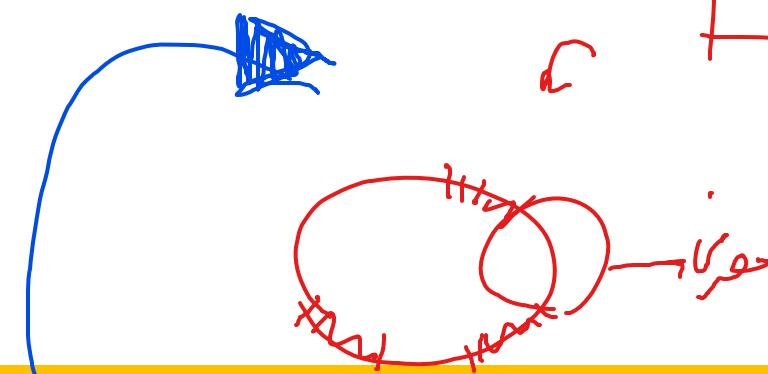
Directory Implementation



C://x/y



C:\a\z
invld wv ~

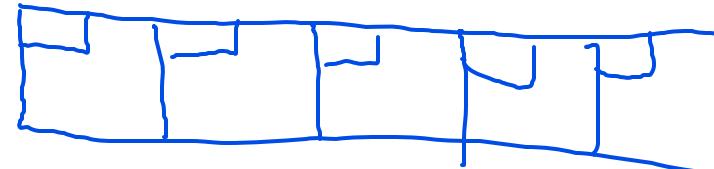


linear $\rightarrow O(n)$

Directory implementation

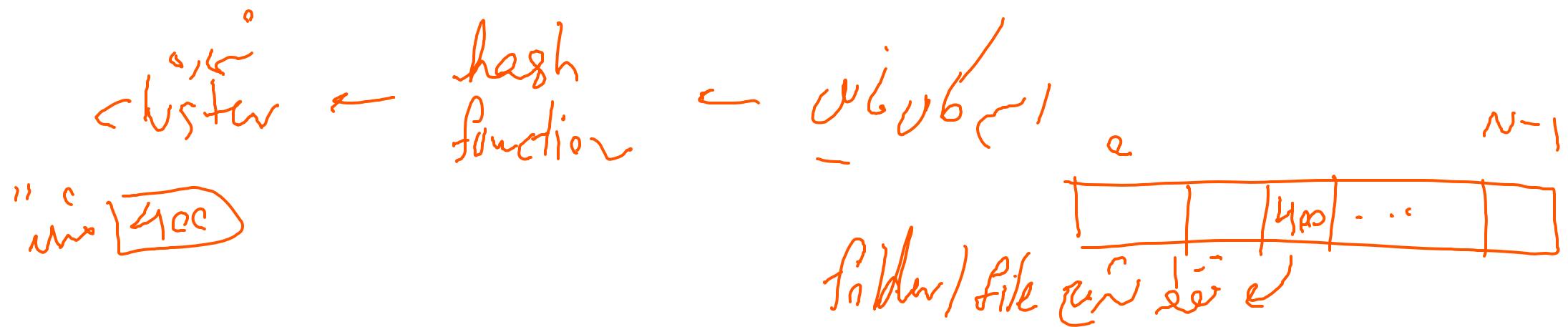
➤ 1) Linear list of file names with pointer to the data blocks

- Simple to program
- Time-consuming to execute (Why?)
 - Linear search time
 - Could keep ordered alphabetically via linked list or use B+ tree
- What would happen if a file is deleted?
 - Problem with fragmentation
- Problem of limited size

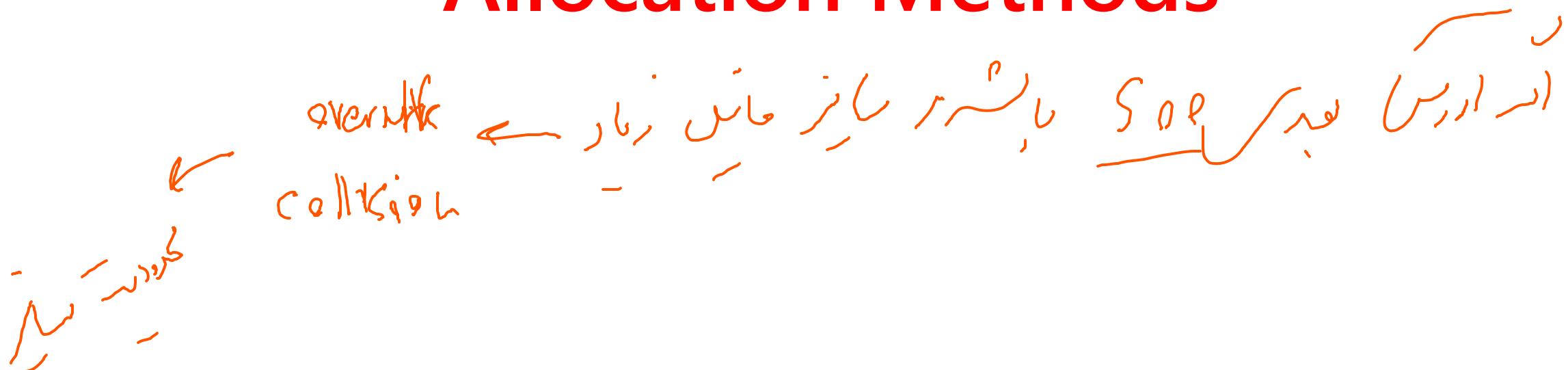


➤ 2) Hash Table

- Linear list with hash data structure
- Decreases directory search time $O(1)$
- Hash is used to map file name to a number
 - one-to-one
- Collisions – situations where two file names hash to the same location \times
- Only good if entries are fixed size, or use chained-overflow method
 - Turvies



Allocation Methods



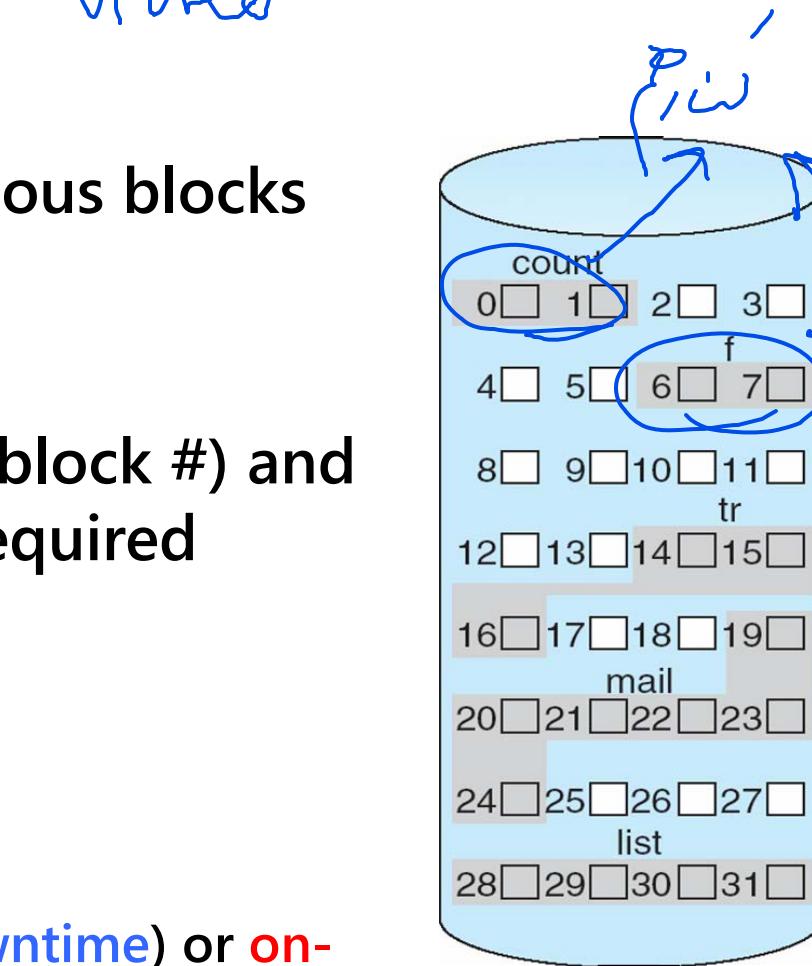
1) Contiguous allocation

ω

➤ Contiguous allocation

- Each file occupies set of contiguous blocks
- Best performance in most cases
- Simple – only starting location (block #) and length (number of blocks) are required
- Problems include
 - Finding space for file
 - Knowing file size
 - External fragmentation
 - Need for compaction off-line (downtime) or on-line

inc mB
Vi dle



➤ shift defragmentation ✓

fragmentation X

speed ✓

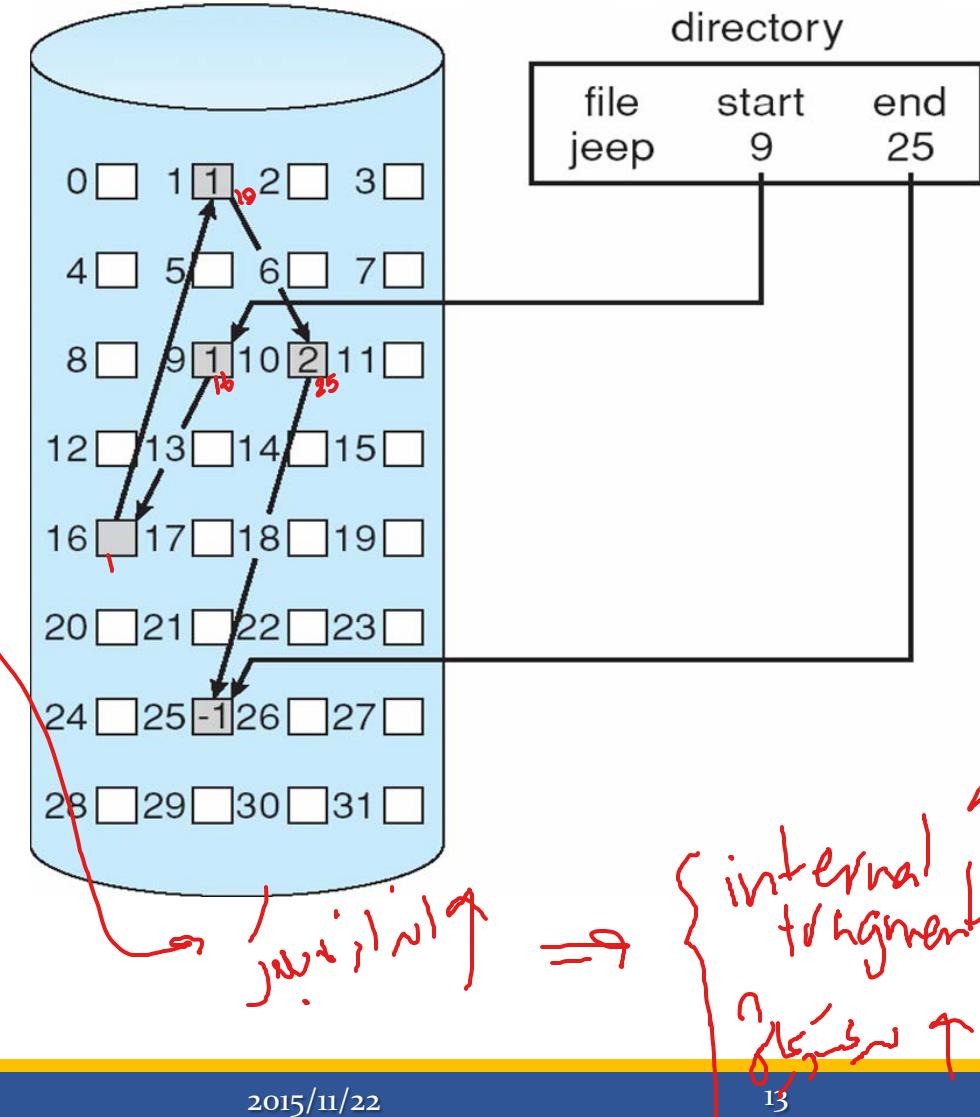
simple ✓

✓ ,
(✓)

2) Linked allocation

Linked allocation

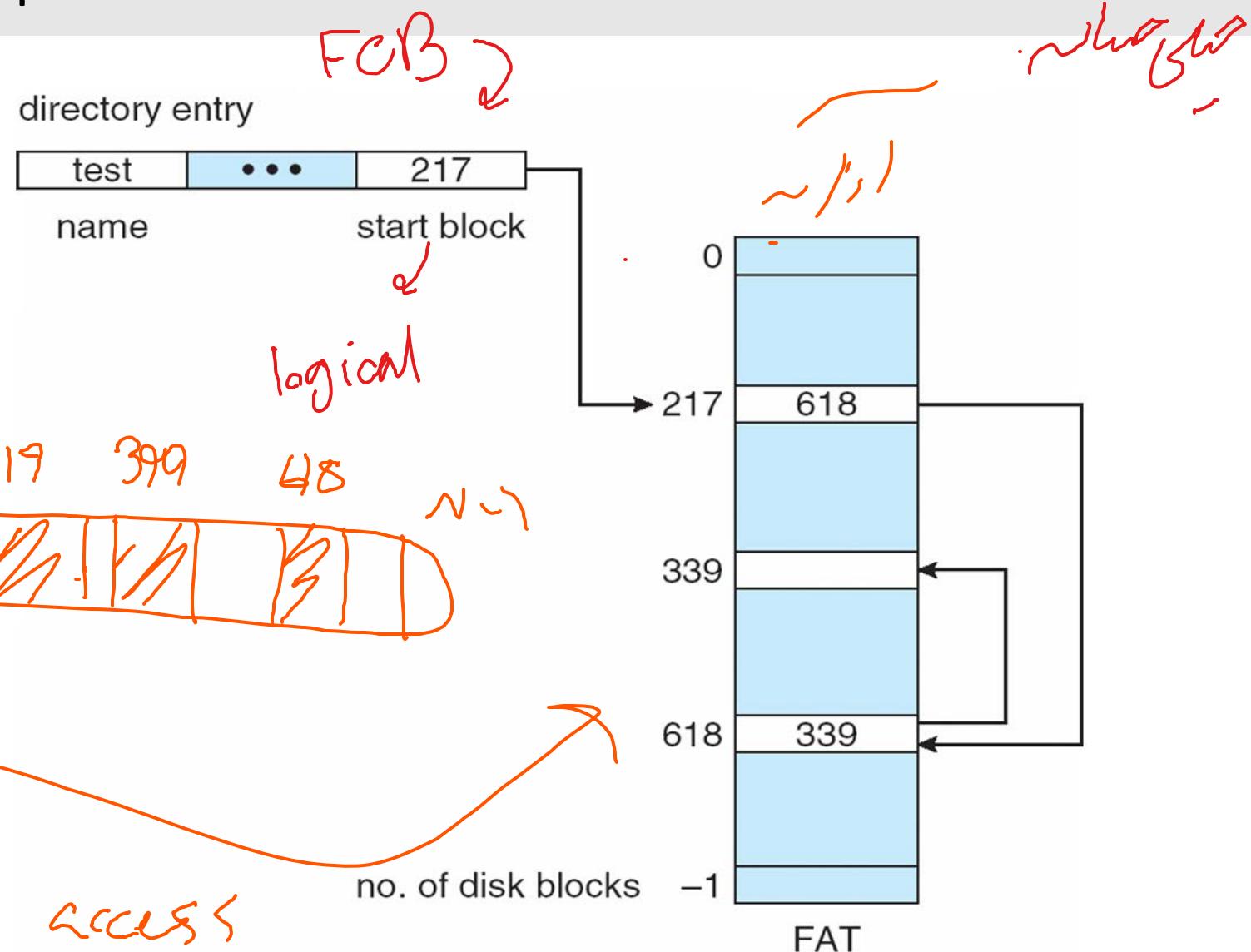
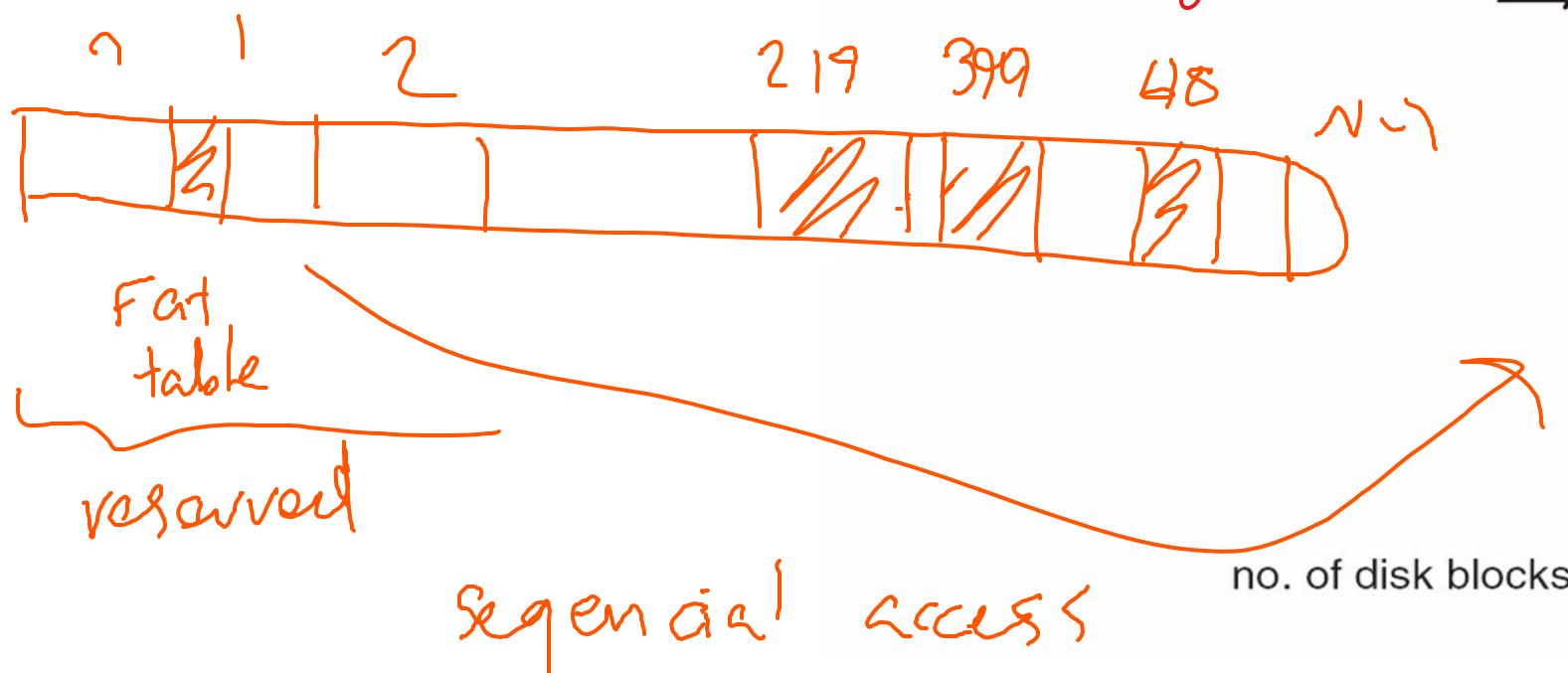
- o Each file a **linked list** of blocks
- o File ends at **nil pointer**
- o **No external fragmentation**
- o Each block contains **pointer to next block**
- o **No compaction**, external fragmentation
- o Free space management system called when new block needed
- o Improve **efficiency** by **clustering** blocks into groups but increases **internal fragmentation**
- o Reliability can be a problem (How?)
- o Locating a block can take **many I/Os** and **disk seeks**
- o 4 byte pointer in 512 byte block = **0.78% overhead**
 - Use of cluster
 - Internal fragmentation



FAT in linked allocation

Windows FAT

- File Allocation Table

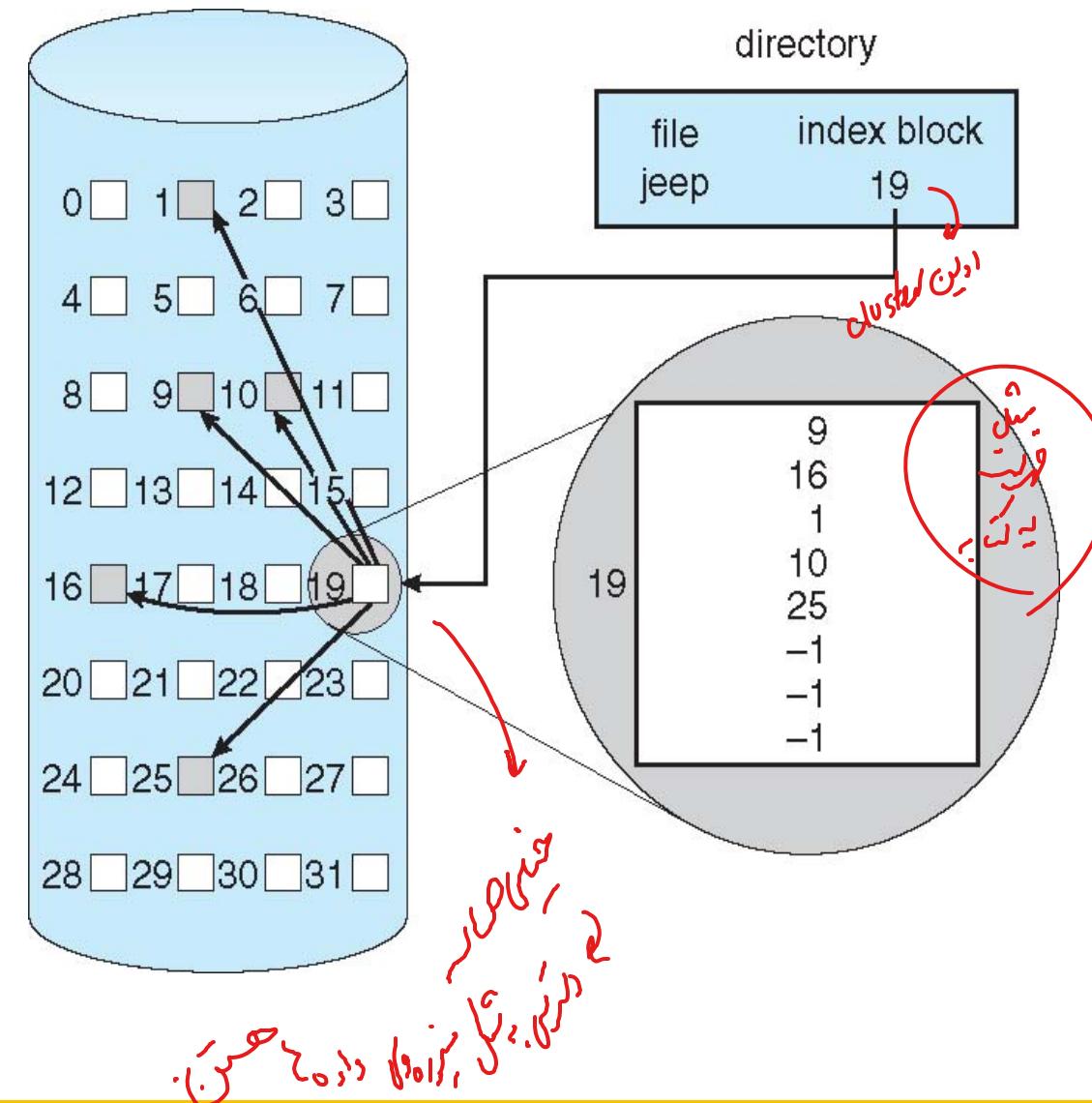


3) Indexed allocation

Direct access

Indexed allocation

- o Each file has its own **index block(s)** of pointers to its data blocks
- o Need **index table**
- o **Random access**
- o Dynamic access **without external fragmentation**, but have **overhead of index block**
- o Mapping from logical to physical in a file of **unbounded length** (block size of 512 words)
- o **Linked scheme**
 - Link blocks of index table (no limit on size)



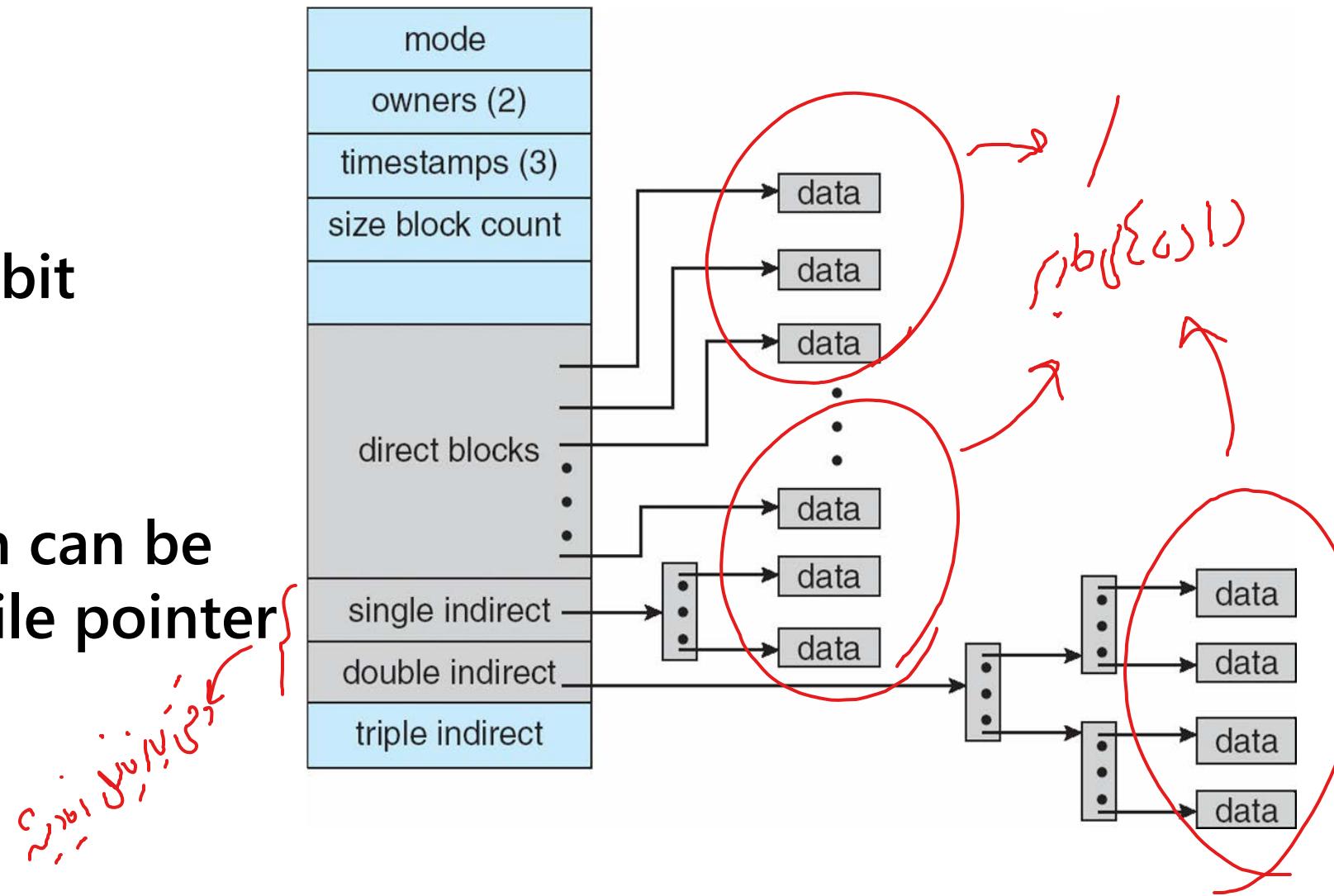
4) Combined allocation

➤ Example:

- UNIX UFS

➤ 4K bytes per block, 32-bit addresses

➤ More index blocks than can be addressed with 32-bit file pointer



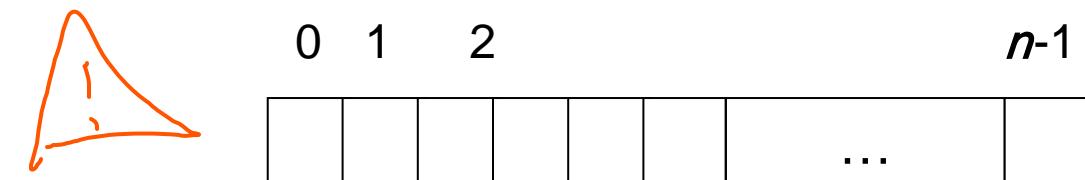
Free Space Management

clock up in invalid writing

FAT

1) Bit vector (bit map)

- File system maintains **free-space list** to track available blocks/clusters
 - (Using term “block” for simplicity)
- **Bit vector or bit map** (n blocks)



➤ Bit map requires extra space

- Example:

block size = 4KB = 2^{12} bytes

disk size = 2^{40} bytes (1 terabyte)

$n = 2^{40}/2^{12} = 2^{28}$ bits (or 32MB)

if clusters of 4 blocks -> 8MB of memory

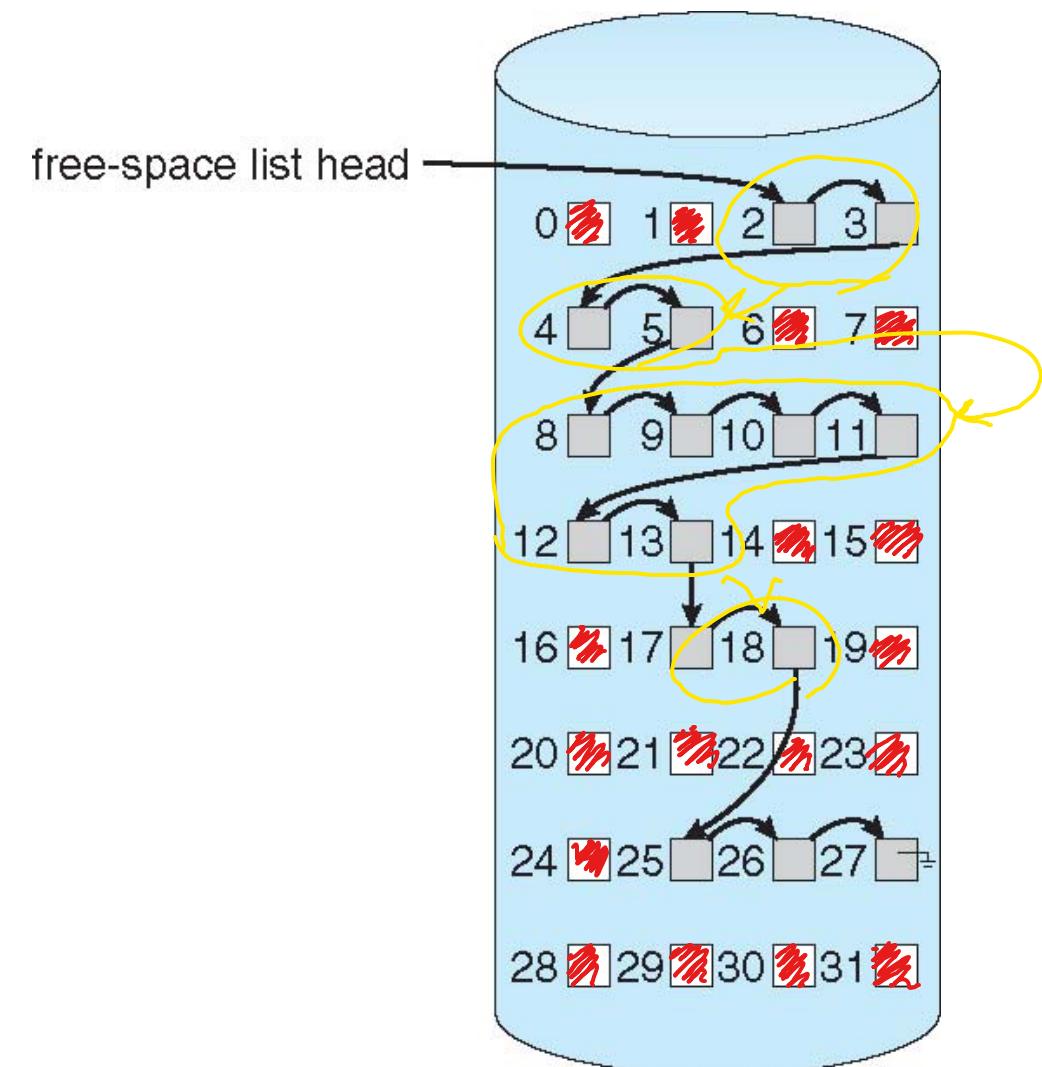
$$\text{bit}[i] = \begin{cases} 1 & \Rightarrow \text{block}[i] \text{ free} \\ 0 & \Rightarrow \text{block}[i] \text{ occupied} \end{cases}$$

➤ Easy to get contiguous files

2) Linked list

➤ Linked list (free list)

- Cannot get contiguous space easily
- No waste of space
- No need to traverse the entire list
(if # free blocks recorded)
- FAT is better to contain free blocks



3) Others: grouping, counting

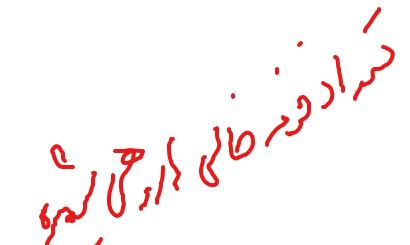
➤ Grouping

- Modify linked list to store address of **next $n-1$ free blocks** in **first free block**, plus a pointer to **next block** that contains **free-block-pointers** (like this one)
(Similar to index block)
- The **last free block** contains address of another **m free blocks**



➤ Counting

- Because space is **frequently contiguously** used and freed, with contiguous-allocation allocation, extents, or clustering
 - Keep **address** of first free block and **count** of following free blocks
 - Free space list then has entries containing addresses and counts



Performance

➤ Performance

- Keeping **data** and **metadata close together**

Buffer cache

- Separate section of main memory for **frequently** used blocks

Synchronous writes sometimes requested by apps or needed by OS

- No buffering / caching – writes must hit disk before acknowledgement
- **Asynchronous** writes more **common, buffer-able, faster**

Free-behind and **read-ahead**

- Techniques to optimize sequential access
- Reads frequently slower than writes

ویرایشی
ویرایشی

Buffer → پریفی

cache → پریفی

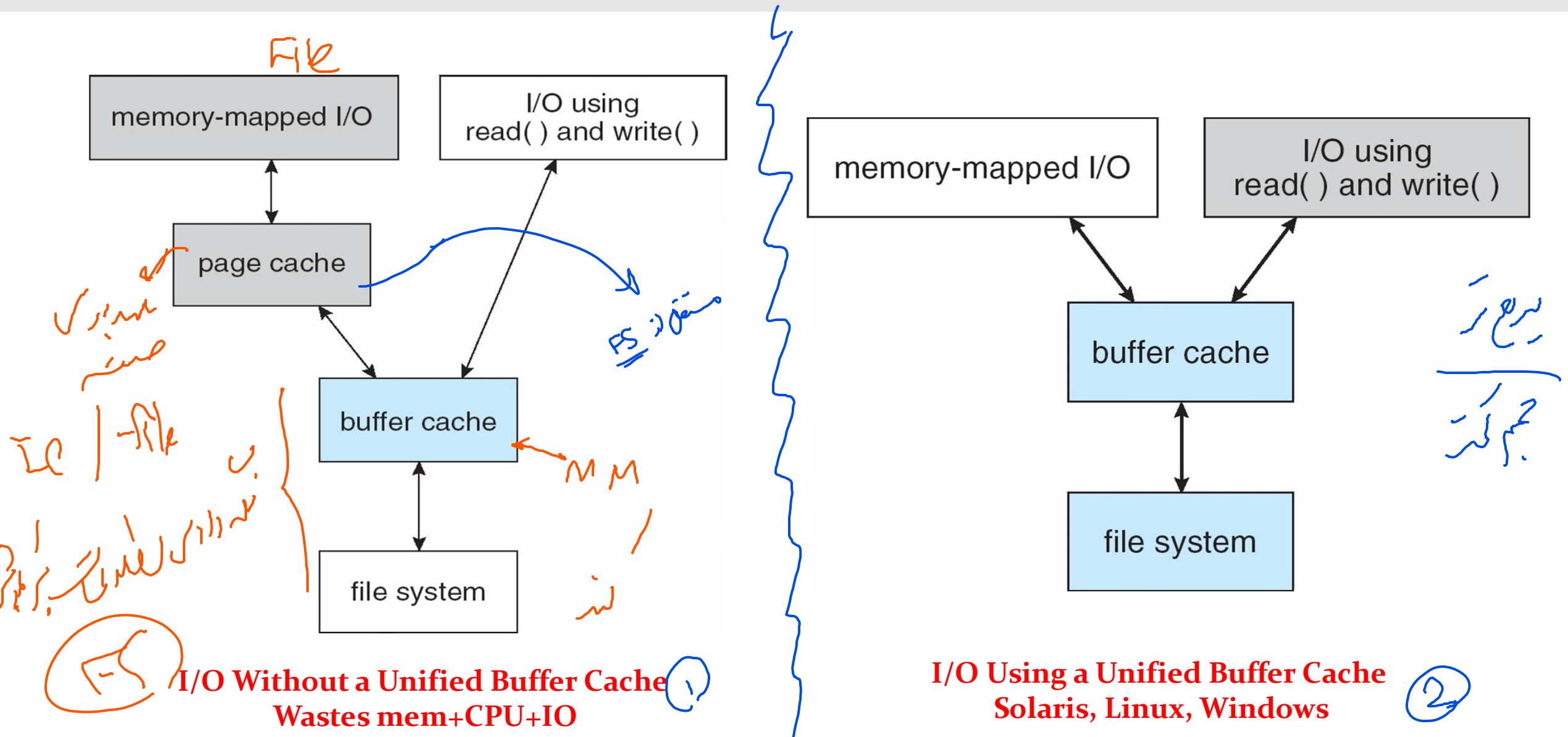
جای خالی



جای خالی (جای خالی)
جای خالی (جای خالی)

جای خالی (جای خالی)
جای خالی (جای خالی)

Unified vs. non-unified buffer cache





Failure Recovery

1) Consistency checking

→ Parity

درارس کردن سازه‌های

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

- Can be slow and sometimes fails
 - Human based!



- Use system programs to **back up** data from disk to another storage device (magnetic tape, other magnetic disk, optical)

es / / , "sectorial" / ! ?

- # ► Recover lost file or disk by restoring data from backup ◀

cal) ۰۵ هریسون

2) Log structured file systems

نحوه ساختاری فایل سیستم

- Similar to database log-based recovery algorithm
- Log structured (or journaling) file systems record each metadata update to the file system as a transaction
- All transactions are written to a log
- Transactions can be
 - Committed → ✓
 - Rolled back → *undo* ↴
- If the file system crashes, all remaining transactions in the log must still be performed
- Where logs should be stored?



transaction (جیوال)

جیوال / نوکری / درجی (log / journal)

3) Backup and restore

- **Back up** data from disk to another storage device, such as a magnetic tape or other hard disk
- **Recovery** from the loss of an individual file, or of an entire disk, may then be a matter of restoring the data from backup
- A typical backup schedule: *Periodic*
 - Day 1. **full backup**: copy all files from the disk to a backup medium.
 - Day 2. **incremental backup**: copy all files changed since **day 1** to another medium.
 - Day 3. **incremental backup**: copy all files changed since **day 2** to another medium.
 - ...
 - Day N. **incremental backup**: copy all files changed since **day N-1** to another medium. Then go back to day 1

Network File System

The Sun Network File System (NFS)

↳ Mount (Disk FS (s))

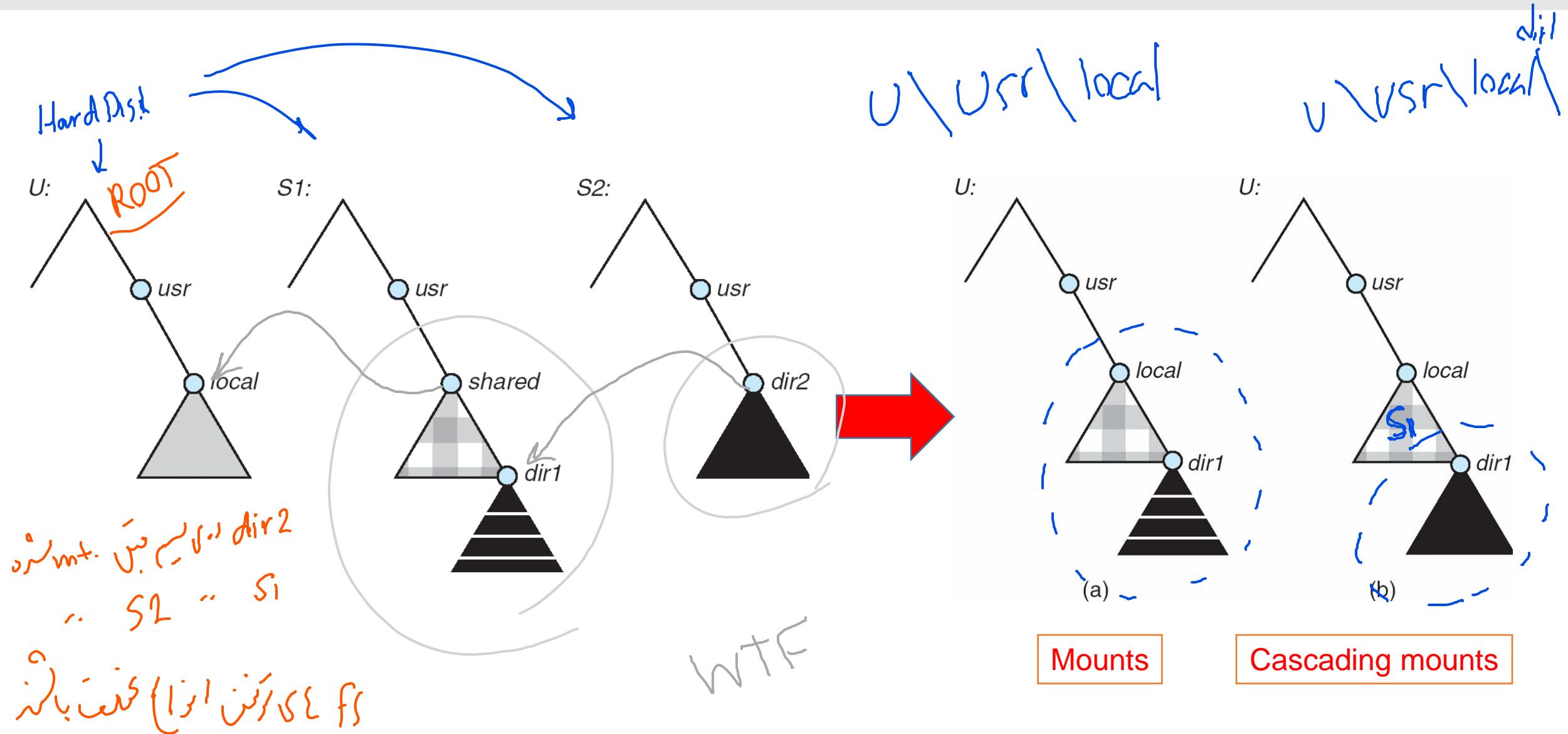
- An implementation & a specification of a software system for accessing **remote files** across **LANs** (or **WANs**)
- The implementation is part of the **Solaris** and **SunOS** operating systems running on **Sun workstations** using an **unreliable datagram protocol** (UDP/IP protocol and Ethernet)

NFS

نیافرینی از سایر سیستم‌های عامل

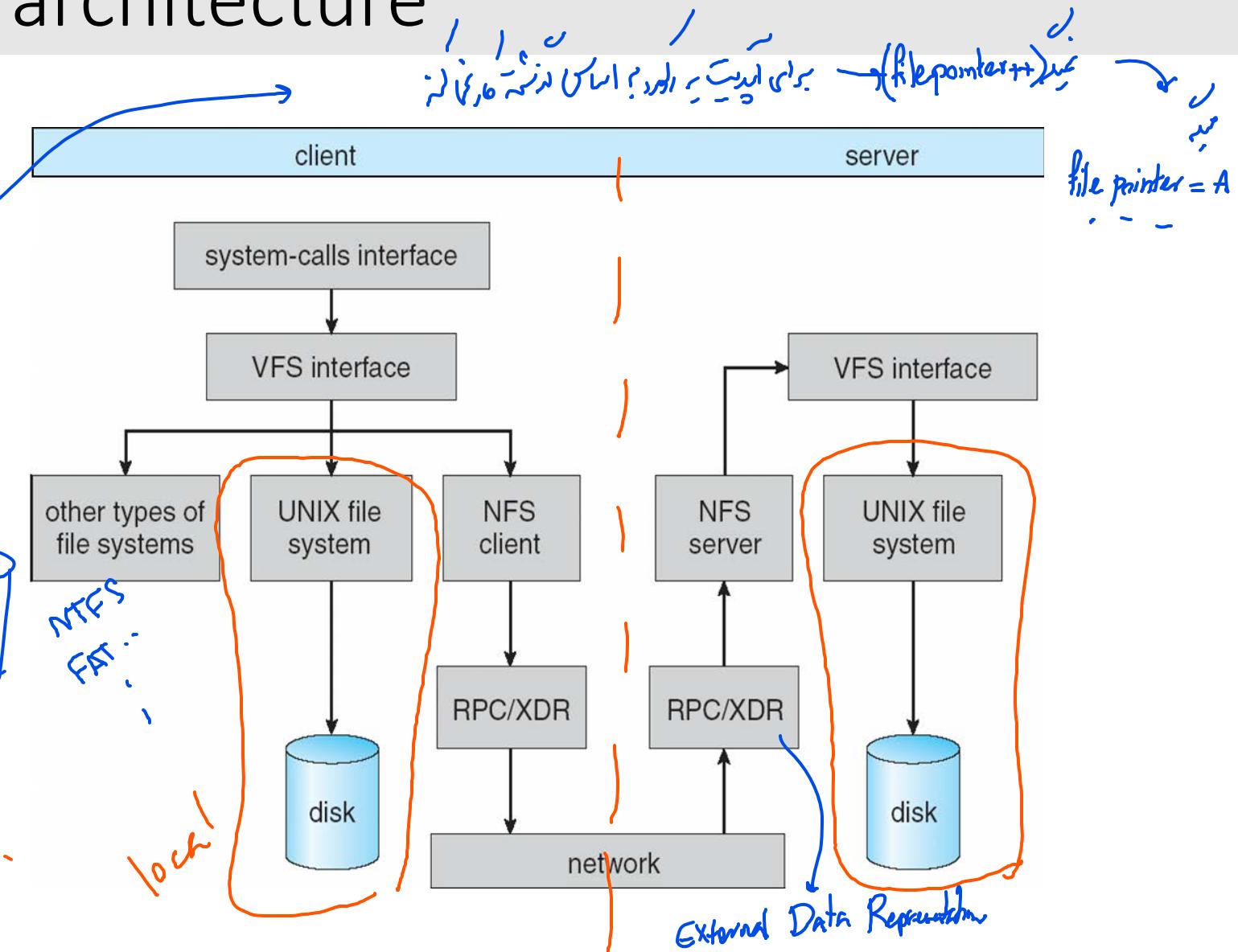
- A **remote directory** is **mounted** over a **local file system** directory
 - The mounted directory looks like an integral subtree of the local file system, replacing the subtree descending from the local directory
- Specification of the remote directory for the mount operation is **nontransparent**; 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
 - Remote Procedure call
- 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

Three independent file systems & mounting



Schematic view of NFS architecture

- NFS write procedure call is atomic
- NFS servers are state-less (why?)
- No open() and close() routines
- All routines are RPC for robustness
- Writes are synchronous (NVM used)



Questions?

