

Goals for Today



- Learning Objective:
 - Inspect Linux Disk Scheduling Algorithms
 - Survey concepts in File Systems Design
- Announcements, etc:
 - Midterm scores posted on Compass (finally)!
 - MP2 was due yesterday at 5am!
 - MP3 is now available for download on Compass!
 - ***DUE APRIL 15th (19 days from now)... get started!!***
 - MP2.5 (Extra Credit) releasing later today!
 - **Next Class:** MP3 Q&A, Midterm Debrief



Reminder: Please put away devices at the start of class



CS 423

Operating System Design: File Systems

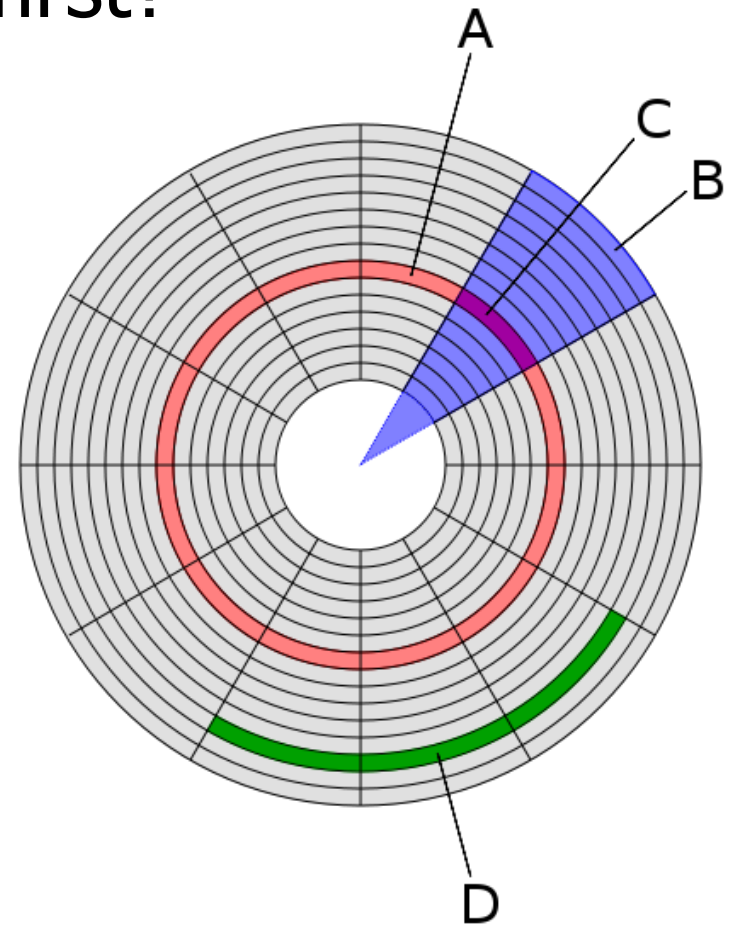
Professor Adam Bates

Disk Scheduling



- Which disk request is serviced first?
 - FCFS
 - Shortest seek time first
 - SCAN (Elevator)
 - C-SCAN (Circular SCAN)

A: Track.
B: Sector.
C: Sector of Track.
D: File



Disk Scheduling Decision — Given a series of access requests, on which track should the disk arm be placed next to maximize fairness, throughput, etc?

Linux I/O Schedulers



- What disk (I/O) schedulers are available in Linux?

```
$ cat /sys/block/sda/queue/scheduler  
noop deadline [cfq]
```

^ scheduler enabled on our VMs

- As of Linux 2.6.10, it is possible to change the IO scheduler for a given block device on the fly!
- How to enable a specific scheduler?

```
$ echo SCHEDNAME > /sys/block/DEV/queue/scheduler
```

- SCHEDNAME = Desired I/O scheduler
- DEV = device name (e.g., sda)

Linux NOOP Scheduler



- Insert all incoming I/O requests into a simple FIFO
- Merges duplicate requests (results can be cached)
- When would this be useful?

Linux NOOP Scheduler



- Insert all incoming I/O requests into a simple FIFO
- Merges duplicate requests (results can be cached)
- When would this be useful?
 - Solid State Drives! Avoids scheduling overhead
 - Scheduling is handled at a lower layer of the I/O stack (e.g., RAID Controller, Network-Attached)
 - Host doesn't actually know details of sector positions (e.g., RAID controller)



- Imposes a deadline on all I/O operations to prevent starvation of requests
- Maintains 4 queues:
 - 2 Sorted Queues (R,W), order by Sector
 - 2 Deadline Queues (R,W), order by Exp Time
- Scheduling Decision:
 - Check if 1st request in deadline queue has expired.
 - Otherwise, serve request(s) from Sorted Queue.
 - Prioritizes reads (DL=500ms) over writes (DL=5s) .Why?

Linux CFQ Scheduler



- CFQ = Completely Fair Queueing!
- Maintain per-process queues.
- Allocate time slices for each queue to access the disk
- I/O Priority dictates time slice, # requests per queue
- Asynchronous requests handled separately — batched together in priority queues
- CFQ is often the default scheduler

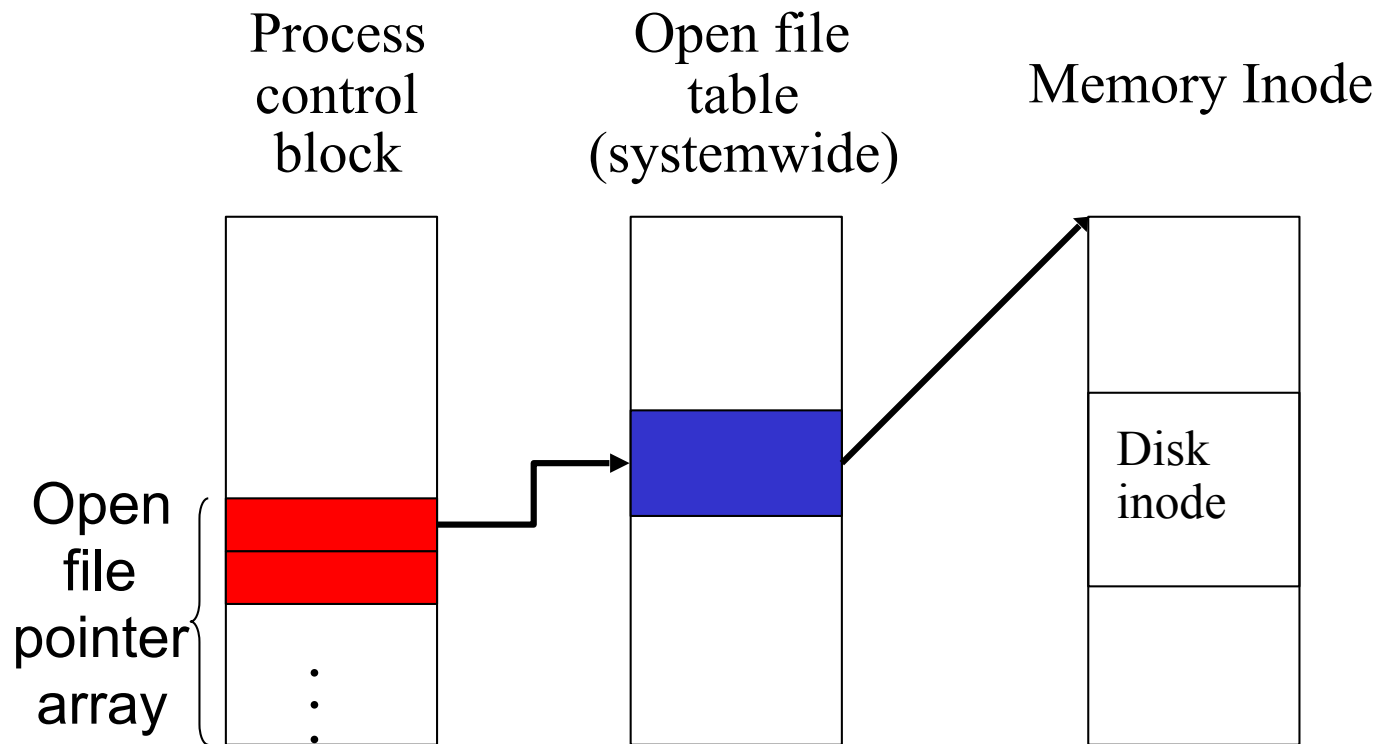


- Deceptive Idleness: A process appears to be finished reading from disk, but is actually processing data. Another (nearby) request is coming soon!
- Bad for synchronous read workloads because seek time is increased.
- Anticipatory Scheduling: Idle for a few milliseconds after a read operation in *anticipation* of another close-by read request.
- Deprecated — CFQ can approximate.

Data Structures for a FS



Data structures in a typical file system:



Disk Layout for a FS



Disk layout in a typical file system:

Boot block	Super block	File metadata (i-node in Unix)	File data blocks
---------------	----------------	-----------------------------------	------------------

- Data Structures:
 - File data blocks: File contents
 - File metadata: How to find file data blocks
 - Directories: File names pointing to file metadata
 - Free map: List of free disk blocks

Disk Layout for a FS



Disk layout in a typical file system:

Boot block	Super block	File metadata (i-node in Unix)	File data blocks
---------------	----------------	-----------------------------------	------------------

- Superblock defines a file system
 - size of the file system
 - size of the file descriptor area
 - free list pointer, or pointer to bitmap
 - location of the file descriptor of the root directory
 - other meta-data such as permission and various times
- For reliability, replicate the superblock

Design Constraints



- How can we allocate files efficiently?
 - For small files:
 - Small blocks for storage efficiency
 - Files used together should be stored together
 - For large files:
 - Contiguous allocation for sequential access
 - Efficient lookup for random access
- Challenge: May not know at file creation where our file will be small or large!!

Design Challenges



- Index structure
 - *How do we locate the blocks of a file?*
- Index granularity
 - *How much data per each index (i.e., block size)?*
- Free space
 - *How do we find unused blocks on disk?*
- Locality
 - *How do we preserve spatial locality?*
- Reliability
 - *What if machine crashes in middle of a file system op?*

File Allocation



- Contiguous
- Non-contiguous (linked)
- Tradeoffs?

Contiguous Allocation

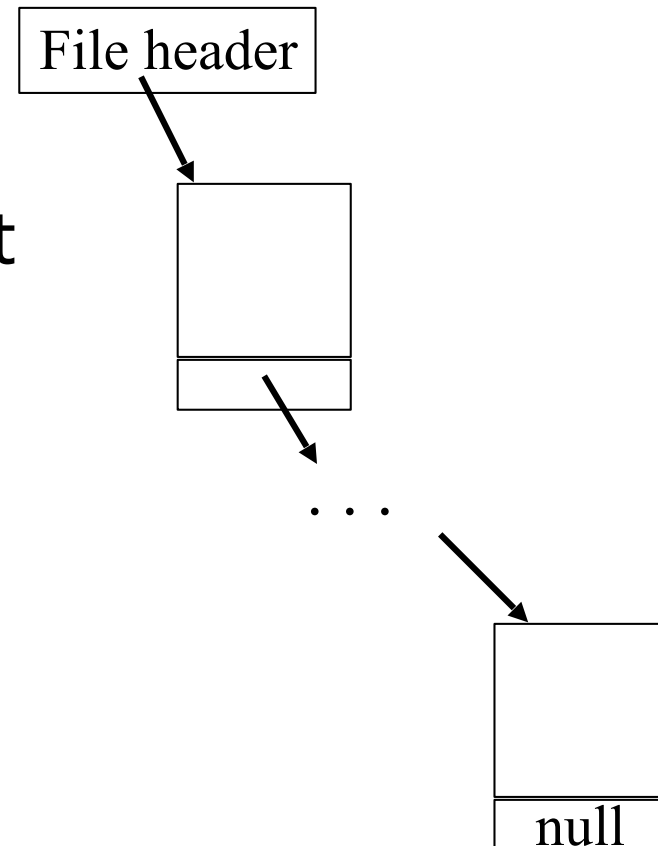


- Request in advance for the size of the file
- Search bit map or linked list to locate a space
- File header
 - first sector in file
 - number of sectors
- Pros
 - Fast sequential access
 - Easy random access
- Cons
 - External fragmentation
 - Hard to grow files

Linked Files



- File header points to 1st block on disk
- Each block points to next
- Pros
 - Can grow files dynamically
 - Free list is similar to a file
- Cons
 - random access: horrible
 - unreliable: losing a block means losing the rest

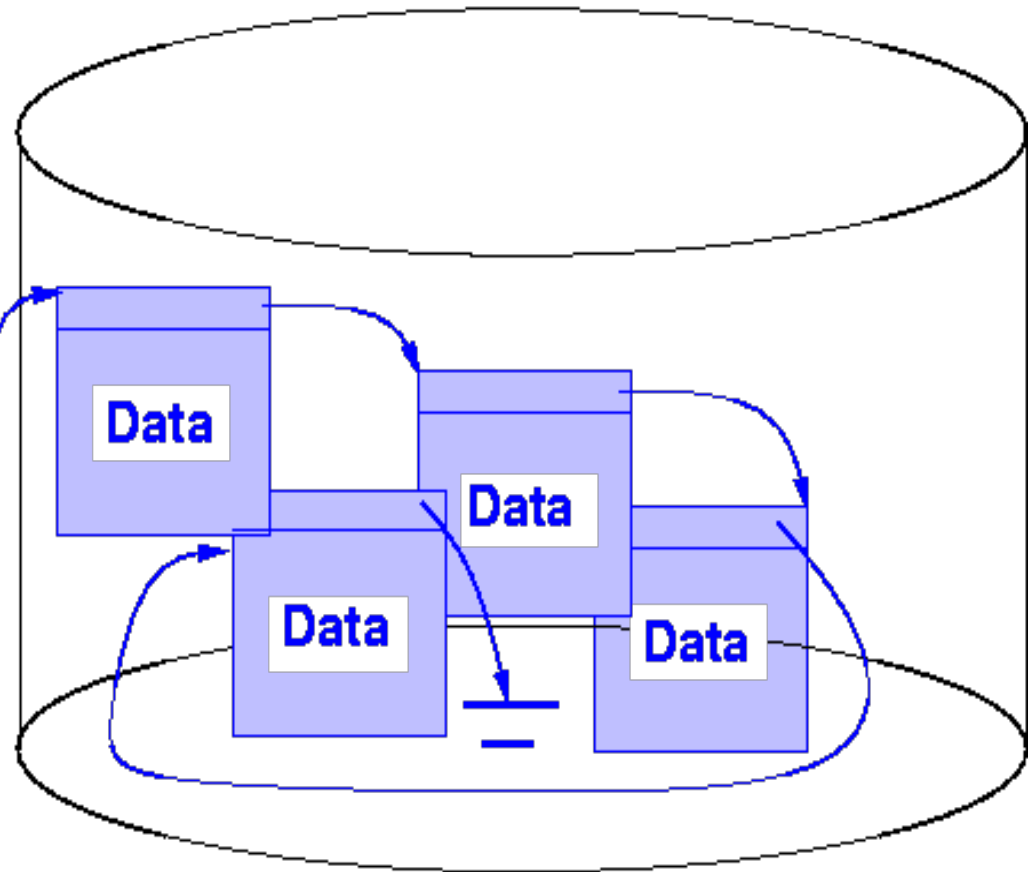


Linked Allocation

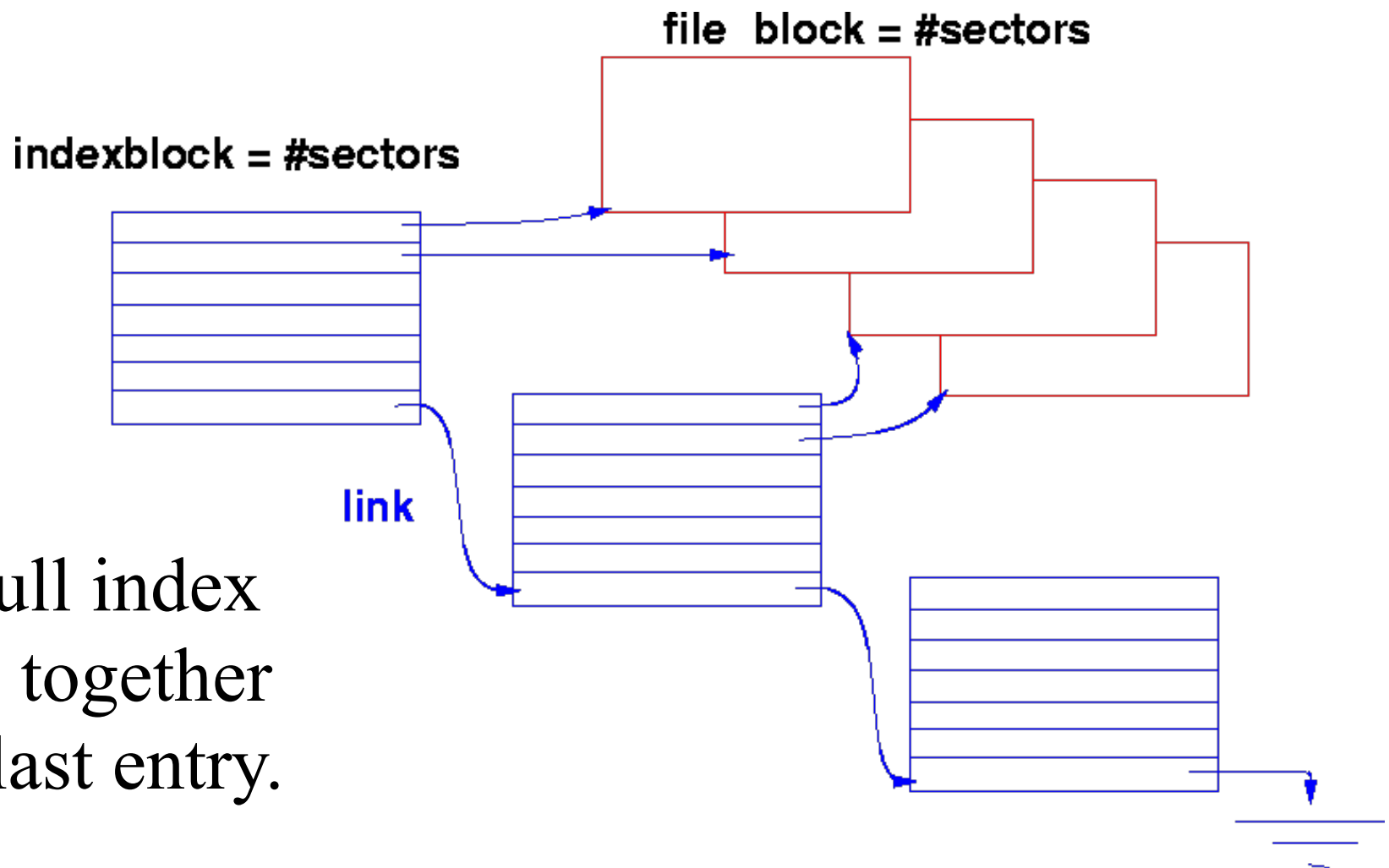


Directory

File	Address

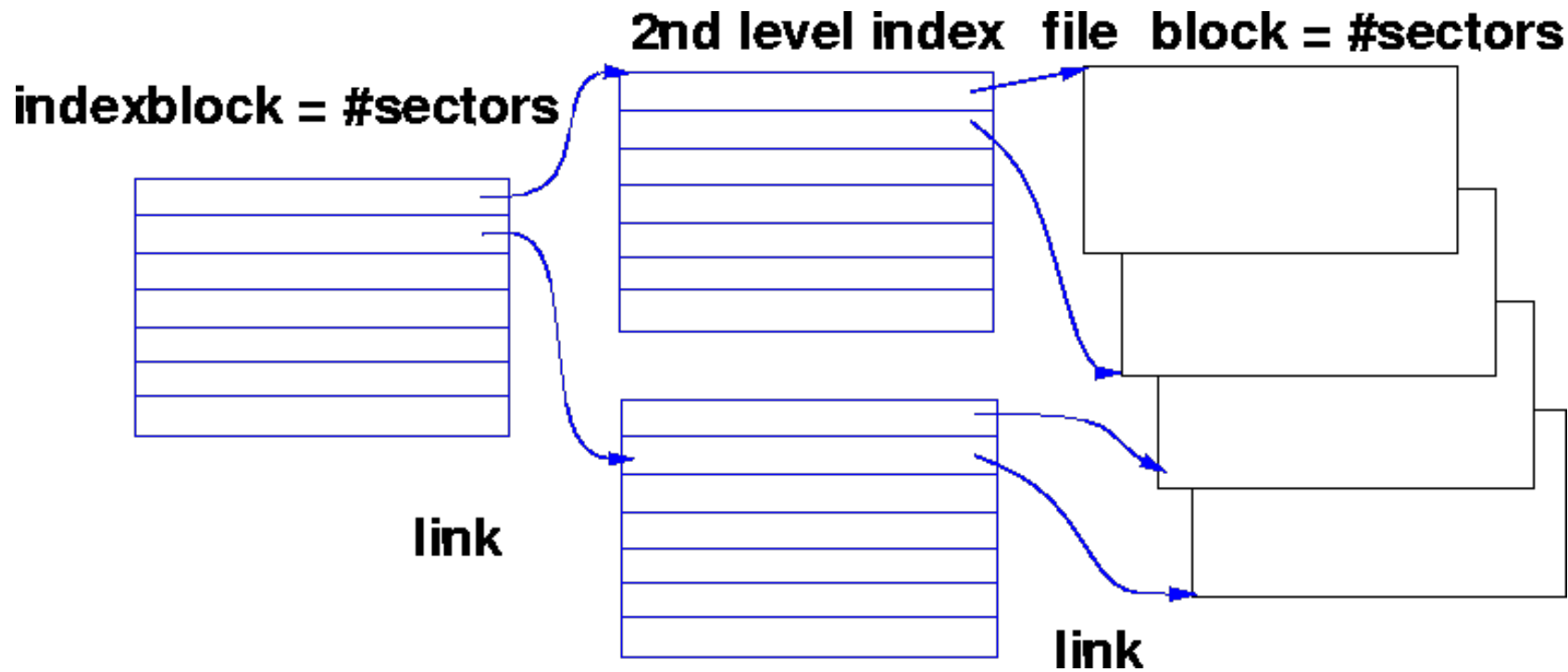


Indexed File Allocation



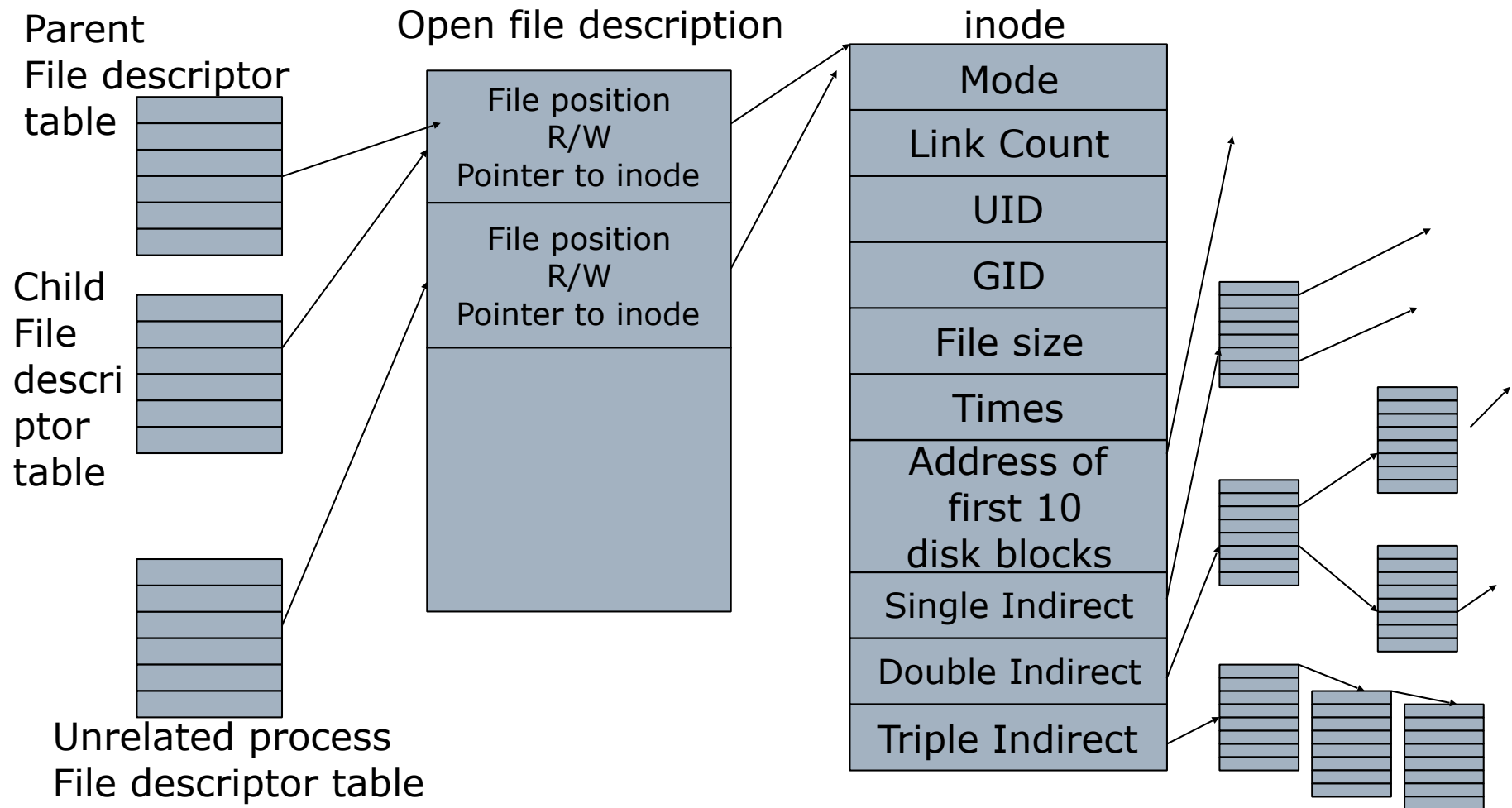
Link full index
blocks together
using last entry.

Multilevel Indexed Files



Multiple levels of index blocks

UNIX FS Implementation





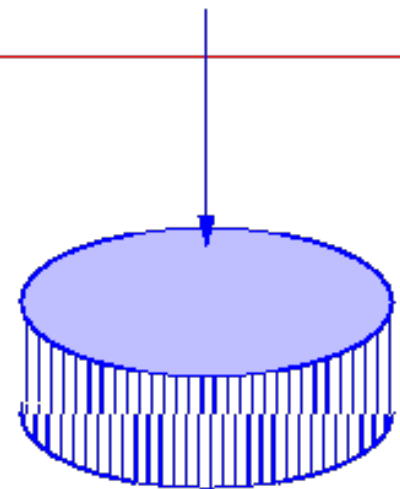
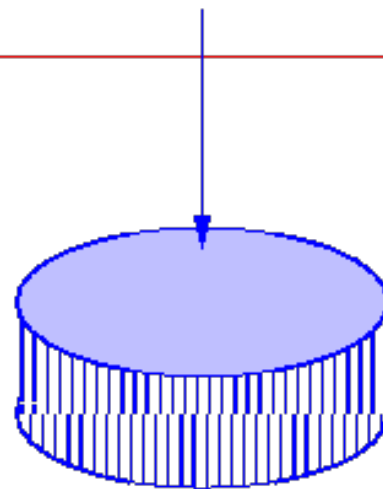
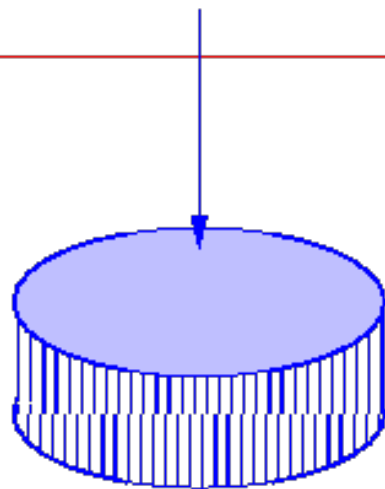
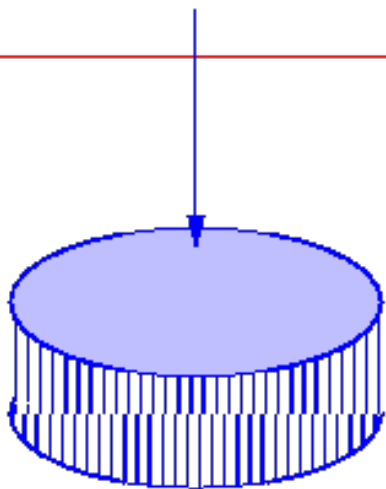
- maps symbolic names into logical file names
 - search
 - create file
 - list directory
 - backup, archival, file migration

Single-level Directory



Directory

Name of File			
--------------------	--	--	--



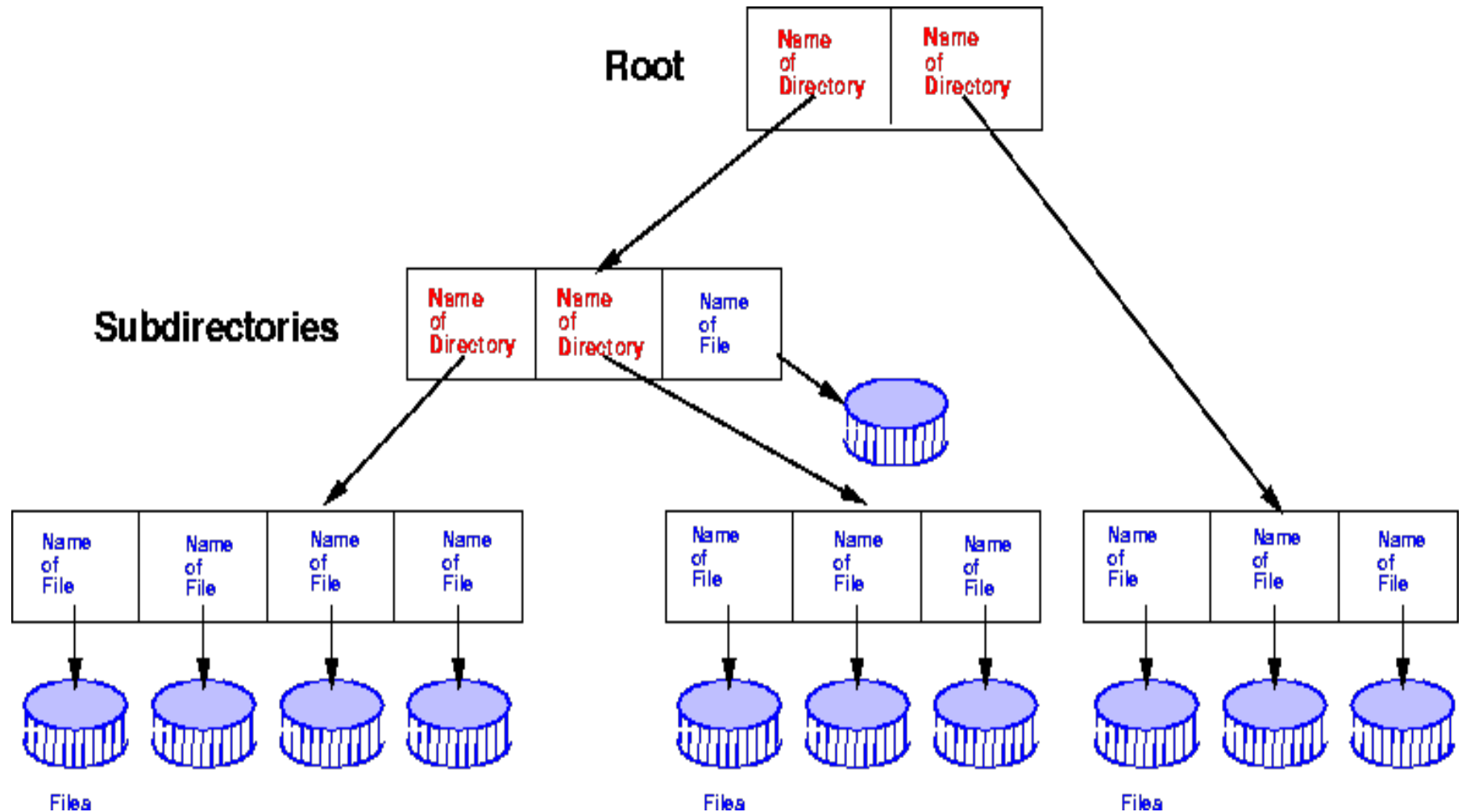
Files

Tree-Structured Directories

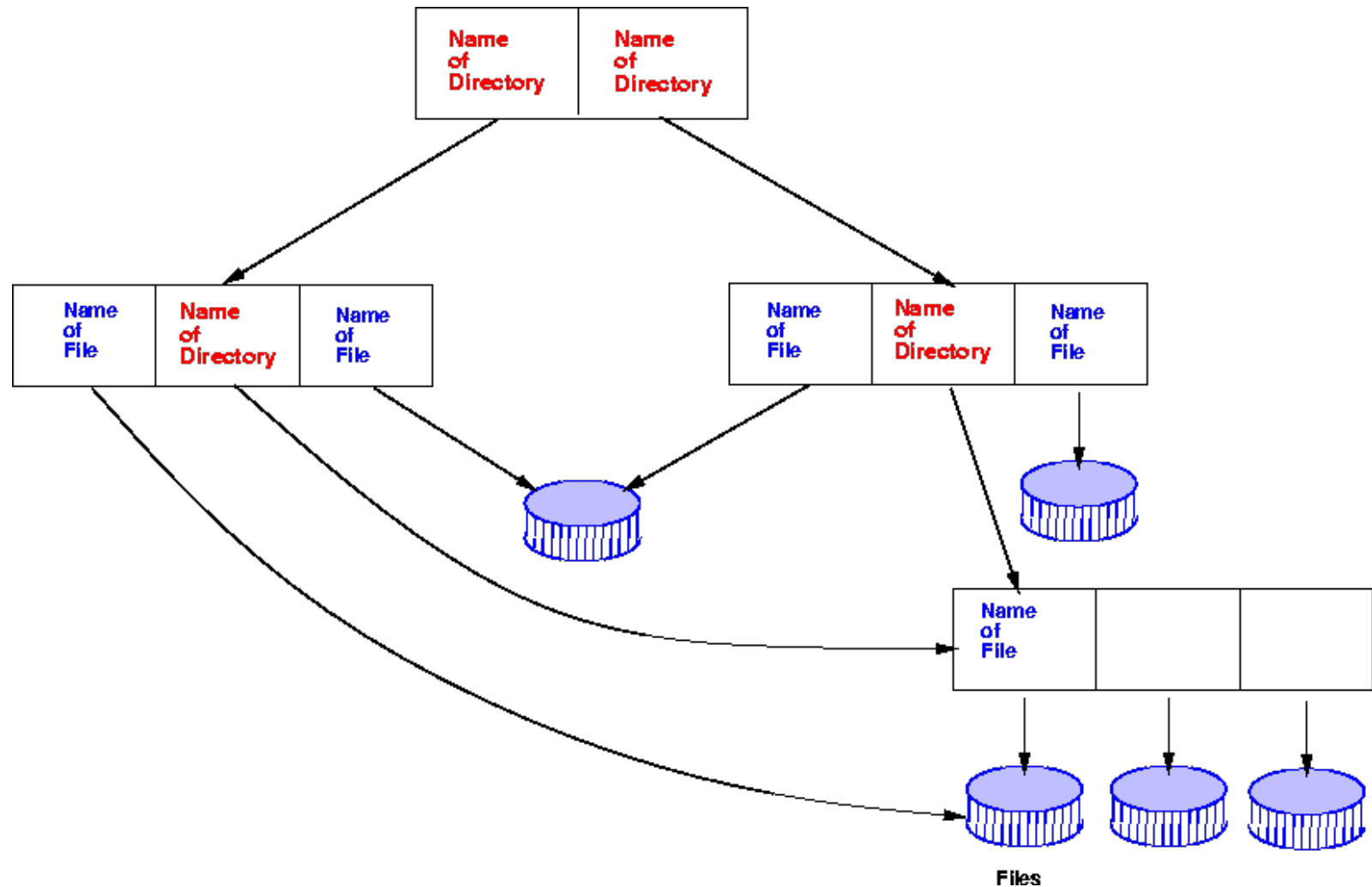


- arbitrary depth of directories
- leaf nodes are files
- interior nodes are directories
- path name lists nodes to traverse to find node
- use absolute paths from root
- use relative paths from current working directory pointer

Tree-Structured Directories



Acyclic Graph Structured Dir's



Symbolic Links



- **Symbolic** links are different than regular links (often called **hard links**). Created with **ln -s**
- Can be thought of as a directory entry that points to the name of another file.
- Does not change link count for file
 - When original deleted, symbolic link remains
- They exist because:
 - Hard links don't work across file systems
 - Hard links only work for regular files, not directories

