

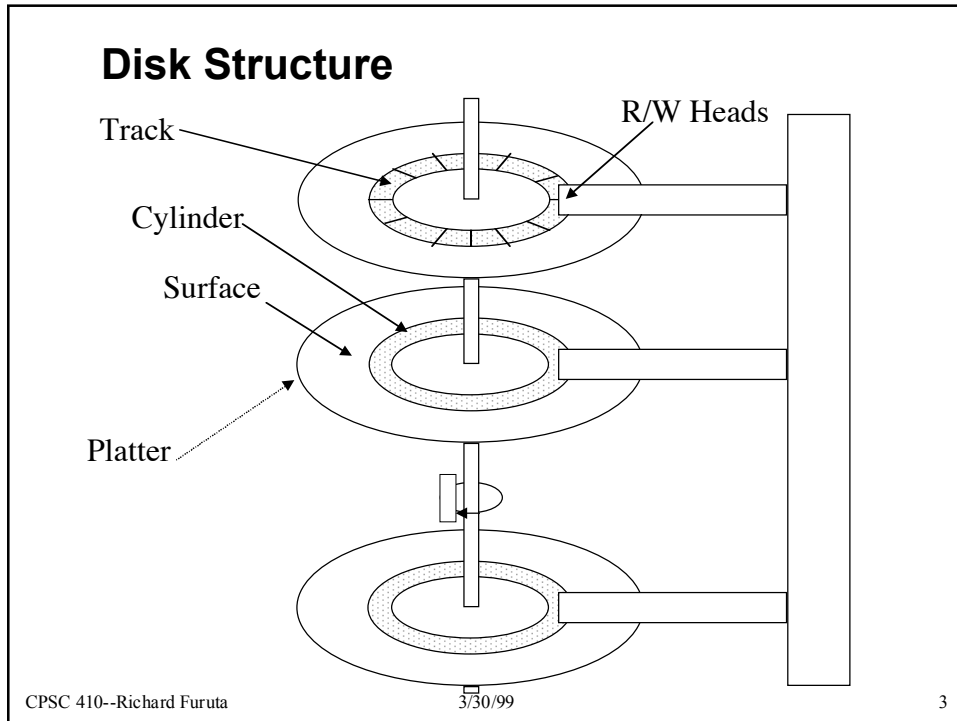
# Silberschatz and Galvin

## Chapter 13

### Secondary Storage Structure

## Secondary Storage Structure Topics

- Disk Structure
- Disk Scheduling
- Disk Management
- Swap-Space Management
- Disk Reliability
- Stable-Storage Implementation



## Disk structure

- Disk drives are addressed as large 1-dimensional arrays of logical blocks, where the logical block is the smallest unit of transfer.
- The 1-dimensional array of logical blocks is mapped onto the sectors of the disk sequentially.
  - Sector 0 is the first sector of the first track on the outermost cylinder.
  - Mapping proceeds in order through that track, then the rest of the tracks in that cylinder, and then through the rest of the cylinders from outermost to innermost.

## Disk Access Time

- seek time: time to position heads on cylinder (a fixed head disk does not require seek time but is more expensive than a moving-head disk)
- rotational latency: delay in accessing material once seek accomplished (time required to wait for data to rotate around under head)
- Transmission time: time to transfer information once it is under the head.
- access time = seek time + rotational latency  
+read/write transmission time  
seek time  $\gg$  read/write time

## Disk Scheduling

- The operating system is responsible for using hardware efficiently — for the disk drives, this means having a fast access time and disk bandwidth.
  - Disk bandwidth is the total number of bytes transferred, divided by the total time between the first request for service and the completion of the last transfer.
- Accomplish this by minimizing seek time
  - Seek time approximates seek distance

## Disk I/O Request

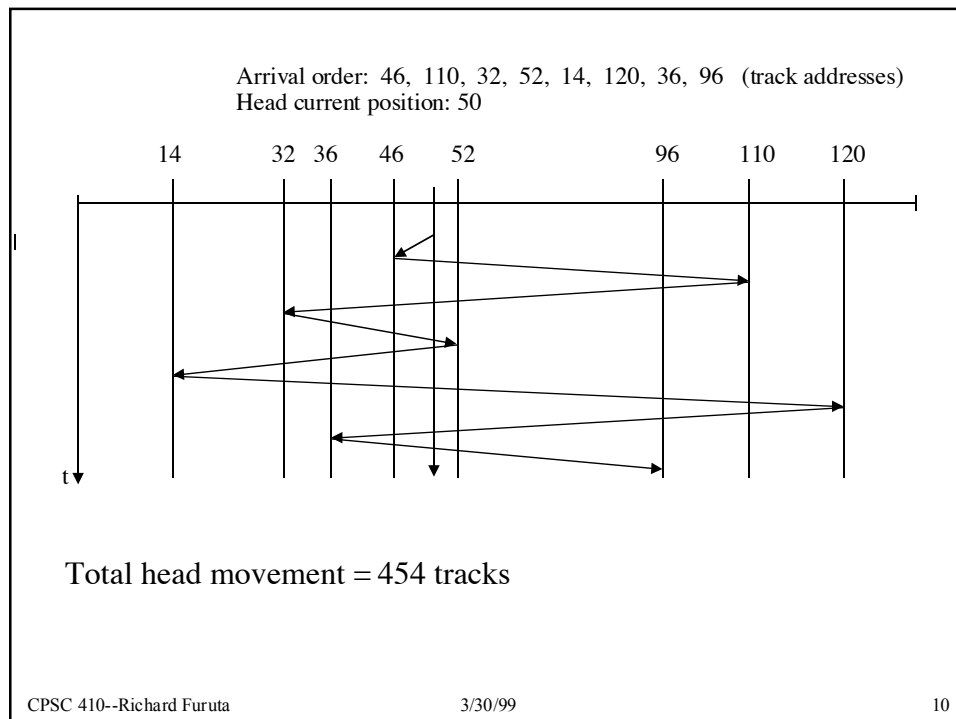
- Disk I/O request specifies
  - whether the operation is input or output
  - disk address (block number, which is translated into drive, cylinder, surface, and sector coordinates)
  - memory address to copy to or from
  - byte count giving the amount of information to be transferred

## Disk Scheduling

- Many requests may be pending at once. Which should be handled first?
- Head moving strategy developed
- Attempting to manage the overall disk seek time. Latency is not controllable and transfer time depends on the size of the transfer request
- Different strategies:
  - FCFS
  - SSTF
  - SCAN
  - LOOK

## FCFS Scheduling

- Simplest form
- First-come, first-served scheduling
- Requests served in order of arrival
- Advantage: simple queueing
- Disadvantage: does not provide the “best” seek time



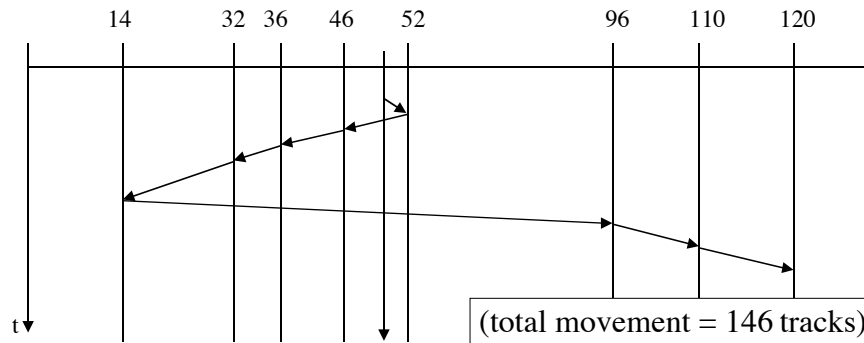
# SSTF Scheduling

- Shortest-seek-time-first

## Shortest-seek-time-first (SSTF)

- Better performance but not optimal.
- Starvation problem.

Arrival order: 46, 110, 32, 52, 14, 120, 36, 96  
Head current position: 50



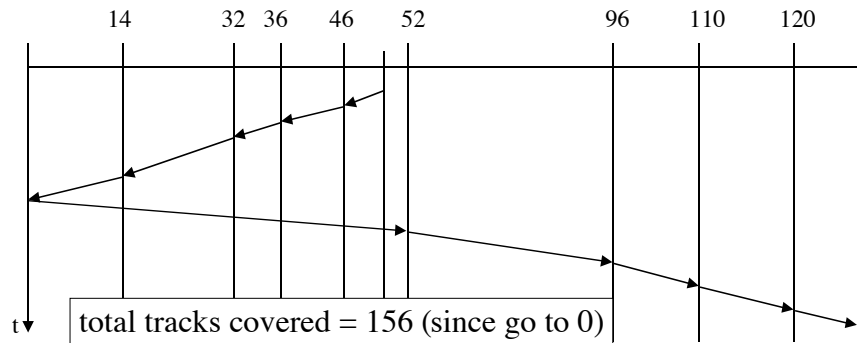
## SCAN Scheduling

- Also called the “elevator” algorithm
- Continue in the direction of first movement until reach end then reverse and head in the other direction (moves completely to the ends of the disk)
- Need to know direction of head movement

### Scan Algorithm

- Goes from one direction to another
- No starvation problem.
- Waiting time and its variance are position dependent.

Arrival order: 46, 110, 32, 52, 14, 120, 36, 96  
Head current position: 50, moving toward to 0.



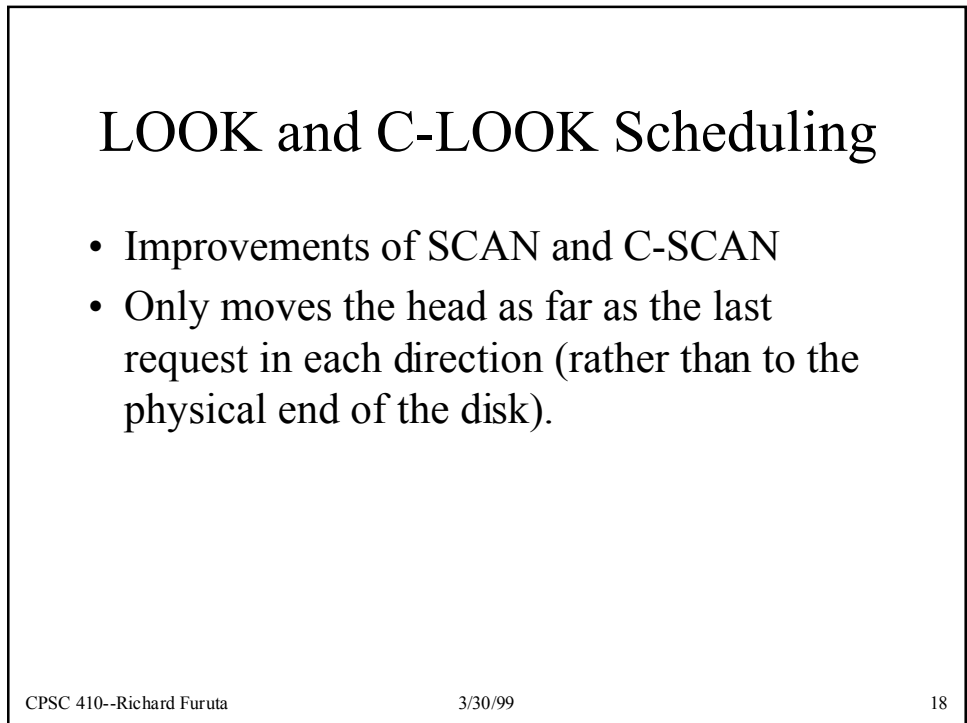
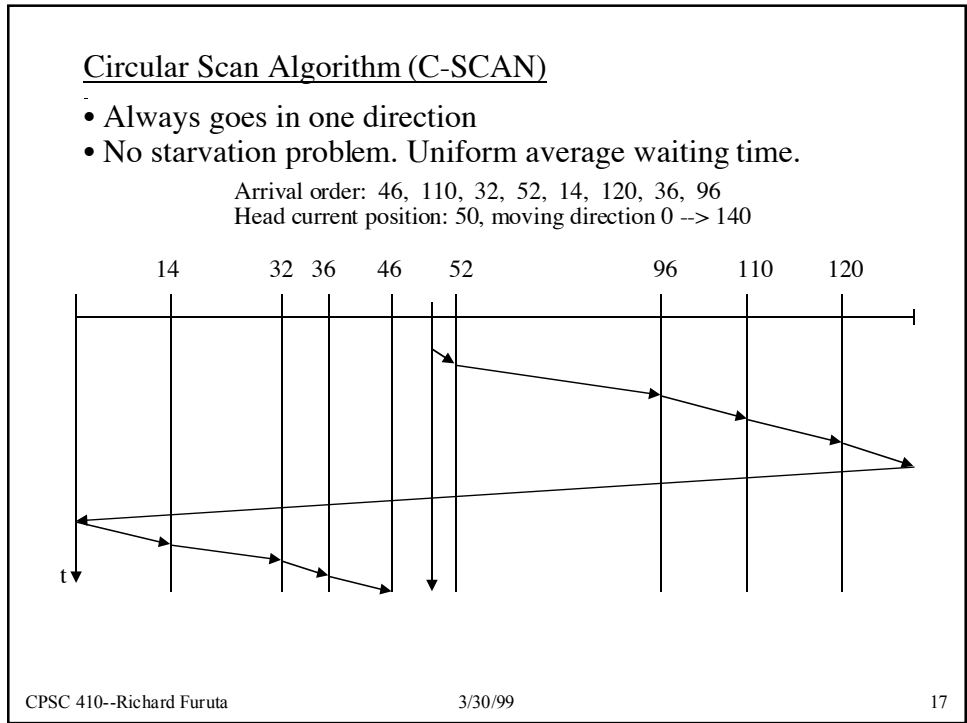
## SCAN Scheduling

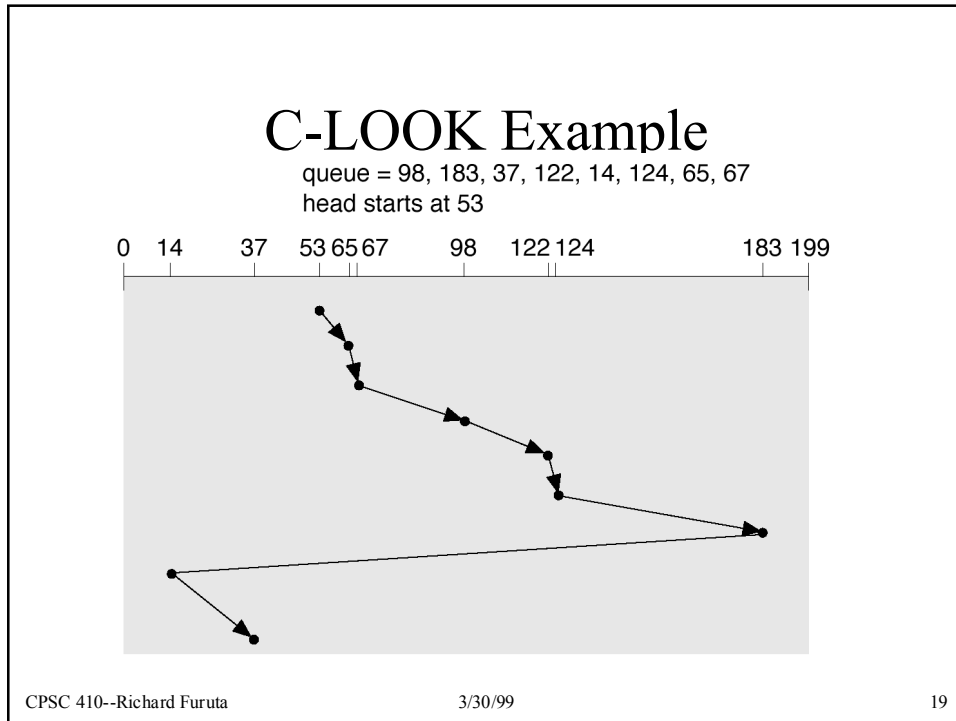
- Note that if there is a uniform distribution of (arriving) requests, the density of the requests near the head is the lowest (just been serviced). Hence when reverse direction the heaviest density of requests is at the other end of the disk and these have also waited the longest. (Uneven waiting time based on position.)

## C-SCAN Scheduling

- Circular Scan scheduling
- Treats disk as if it were circular with the last track adjacent to the first one
- As it reaches the end of the disk, restarts again on the other side
- Does move head all the way to the end of the disk







## Disk Scheduling Algorithms

- If disk queue seldom has more than one request, all scheduling algorithms are effectively equivalent. FCFS attractive because of low overhead in implementing.
- SCAN and C-SCAN are appropriate for heavy load situations
- Relationship with allocation method: are file's blocks clustered or dispersed across disk.
- Location of directories in middle of disk (rather than on edges) can reduce amount of head movement.
- Different algorithms best for different situations--modular design helps designer adjust algorithm used.

## Disk management

- Low-level formatting, or physical formatting
  - Dividing a disk into sectors that the disk controller can read and write.
- To use a disk to hold files, the operating system still needs to record its own data structures on the disk.
  - Partition the disk into one or more groups of cylinders.
  - Logical formatting or “making a file system”.

## Disk management

- Boot block initializes system.
  - The bootstrap is stored in ROM.
  - Bootstrap loader program brings in full bootstrap program from disk.
    - Bootstrap program stored at fixed location on disk (boot blocks)
    - Allows updating bootstrap program.

## Disk management: Bad blocks

- Methods such as sector sparing (also known as forwarding) used to handle bad blocks.
  - Spare sectors set aside on low-level formatting
  - Controller told to replace a bad sector logically with one of the spare sectors
  - To retain effectiveness of disk-scheduling optimization, provide spare sectors in each cylinder and also provide some spare cylinders. Use spare sector from same cylinder if possible.

## Disk management: Bad blocks

- Sector slipping
  - moves blocks following bad block downward (occupying spare sector) to free up block following bad block
  - skips bad block, using freed up block to hold that sector's information.

## Swap space management

- Swap-space — Virtual memory uses disk space as an extension of main memory.
- Swap space can be carved out of the normal file system, or, more commonly, it can be in a separate disk partition.

## Swap space management

- 4.3BSD allocates swap space when process starts; holds text segment (the program) and data segment.
- Kernel uses swap maps to track swap-space use.
- Solaris 2 allocates swap space only when a page is forced out of physical memory, not when the virtual memory page is first created.

## Disk reliability

- Several improvements in disk-use techniques involve the use of multiple disks working cooperatively.
- *Disk striping* uses a group of disks as one storage unit.
- *RAID: Redundant Array of Independent Disks*

## Disk reliability

- *RAID* schemes improve performance and improve the reliability of the storage system by storing redundant data.
  - Mirroring or shadowing keeps duplicate of each disk.
  - Block interleaved parity uses much less redundancy.
    - lost blocks can be recomputed from remaining blocks plus parity block

## Stable storage implementation

- Information stored in stable storage must *never* be lost
  - Failure during an update does not leave all copies in a damaged state
  - Recovery from failure brings all copies to a consistent and correct state, even if there is another failure during the recovery.

## Stable storage implementation

- To implement stable storage:
  - Replicate information on more than one nonvolatile storage media with independent failure modes.
  - Update information in a controlled manner to ensure that we can recover the stable data after any failure during data transfer or recovery.

## Stable storage implementation example

- Disk write results in one of these outcomes
  - Successful completion
  - Partial failure (middle of transfer)
  - Total failure (previous values remain intact)
- Maintain two copies of block and follow this protocol for writes
  - Write information on first physical block
  - When write completes successfully, write same information onto second physical block
  - Declare the operation complete only after the second write completes successfully

## Stable storage implementation example

- Recovery from failure
  - inspect each pair of physical blocks
    - both the same and no detectable error exists: no further action required
    - one block contains detectable error: replace with other
    - both blocks have no detectable error but differ in content: replace content of first with second
  - write to stable storage either succeeds completely or results in no change