

Any modern computer system will incorporate (at least) two levels of storage:

<u>primary storage:</u>	random access memory (RAM)
typical capacity	32MB to 1GB
cost per MB	\$3.00
typical access time	5ns to 60ns
burst transfer rate	??
<u>secondary storage:</u>	magnetic disk/optical devices/tape systems
typical capacity	4GB to 100GB for fixed media; $\infty$ for removable
cost per MB	\$0.01 for fixed media, more for removable
typical access time	8ms to 12ms for fixed media, larger for removable
burst transfer rate	??

**Note: all statistics here are guaranteed invalid by July 1, 2000.**

## Spatial units:

byte	8 bits
kilobyte (KB)	1024 or $2^{10}$ bytes
megabyte (MB)	1024 kilobytes or $2^{20}$ bytes*
gigabyte (GB)	1024 megabytes or $2^{30}$ bytes

**\* Most, if not all, secondary storage vendors lie about this and call 1,000,000 bytes a megabyte and 1,000,000,000 bytes a gigabyte.**

**We will ignore the issue.**

## Time units:

nanosecond (ns)	one-billionth ( $10^{-9}$ ) of a second
microsecond ( $\mu$ s)	one-millionth ( $10^{-6}$ ) of a second
millisecond (ms)	one-thousandth ( $10^{-3}$ ) of a second

While the particular values given earlier are volatile, the relative performances suggested are actually quite stable:

Primary storage:

- costs several hundred times as much per unit as secondary storage.
- has access times that are 250,000 to 1,000,000 times faster than secondary storage
- has transfer rates that are ?? times faster than secondary storage

Why do WE care (in a data structures class)?

Often data must be first read from disk into memory for processing, and then results must be written back to disk after processing.

In many cases, data sets are too large to store in memory at once.

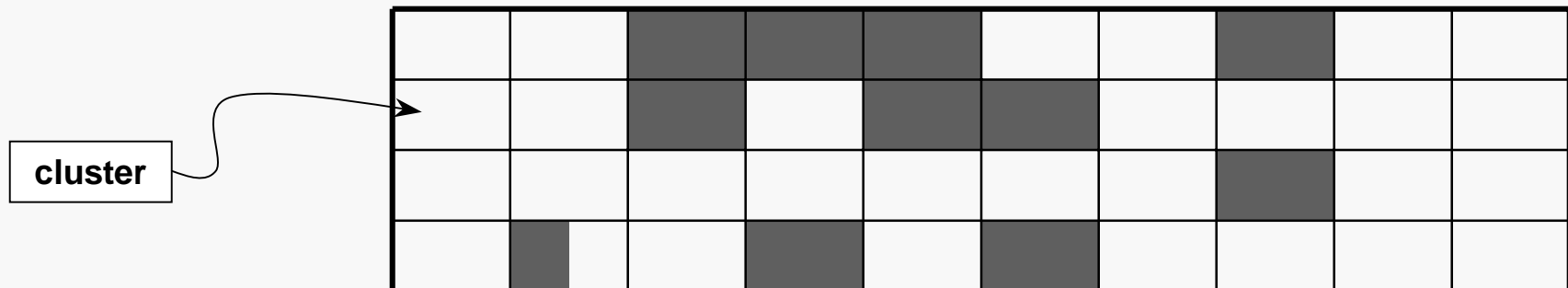
In many file systems, space is allocated in fixed-size chunks called clusters.

Typical cluster sizes range from 1KB up to 32KB, usually equaling a power of 2.

When a file is stored on disk, an integer number of clusters are allocated for the file; since files sizes are typically not a multiple of the cluster size, this means that a certain amount of space is wasted to internal fragmentation.

On average,  $\frac{1}{2}$  of a cluster is wasted per file stored. Obviously that adds up... but that's not really our concern in this course.

The clusters are also not usually stored contiguously on the disk. This external fragmentation can cause a serious degradation of performance when a file is being read from or written to disk.



To illustrate the issues, we will consider the physical organization and behavior of a typical hard disk design.

Note that the presentation here is an over-simplification and that contemporary hard drive designs incorporate control sophistication not discussed here.

Despite advances, the basic performance issues remain the same.

We will consider three primary factors that affect the time required to read requested data from the disk into memory:

seek time      time for the appropriate I/O head to reach the desired track

latency      time for the appropriate sector to rotate to the I/O head

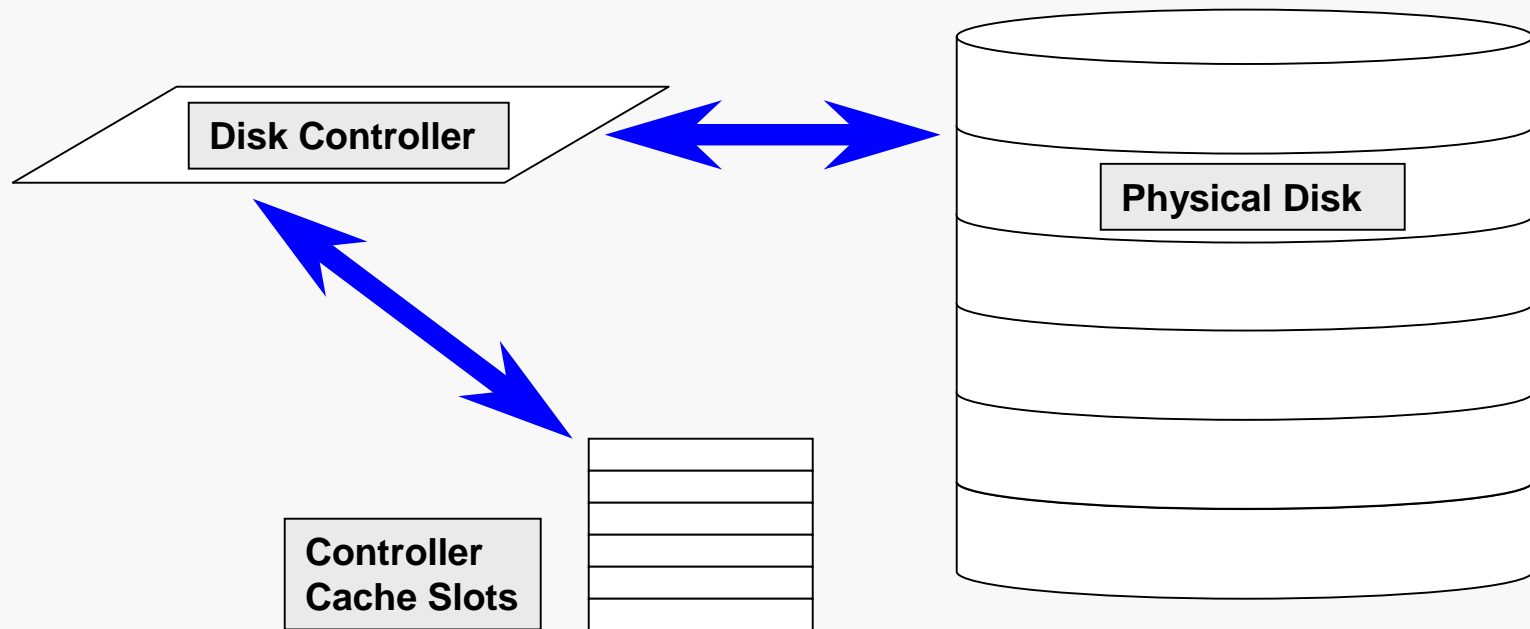
transfer time      time for the data to be read to move past the I/O head

First we must have a clear picture of the hardware...

All disk systems incorporate an on-board controller which acts as an interface between the CPU and the disk hardware.

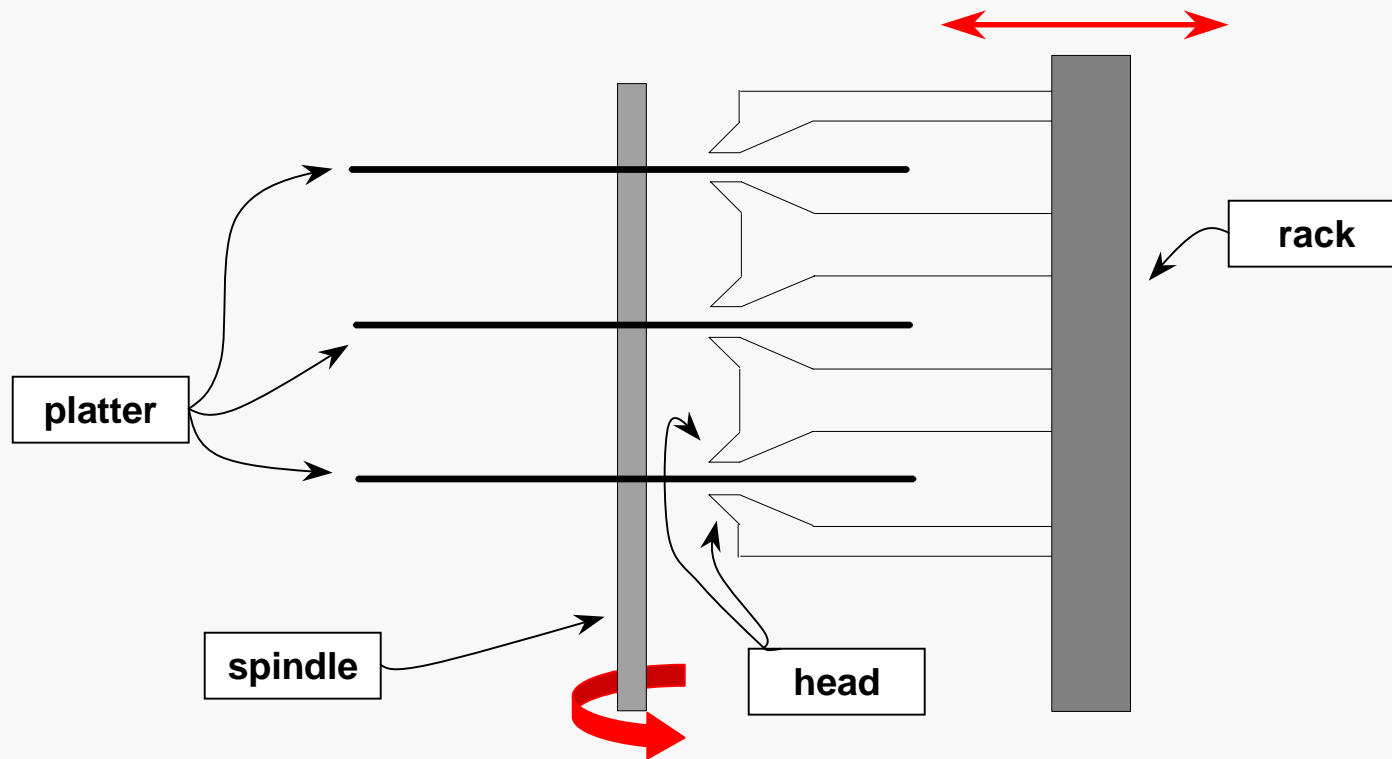
The controller masks the internal physical organization of the disk (covered in the next few slides) from the CPU and user processes.

The controller has an internal cache (typically 256KB up to several MB) that it uses to buffer data for read/write requests.



A typical hard drive contains a number of circular platters, attached to a rotating spindle. Each platter holds data on one or both of its surfaces.

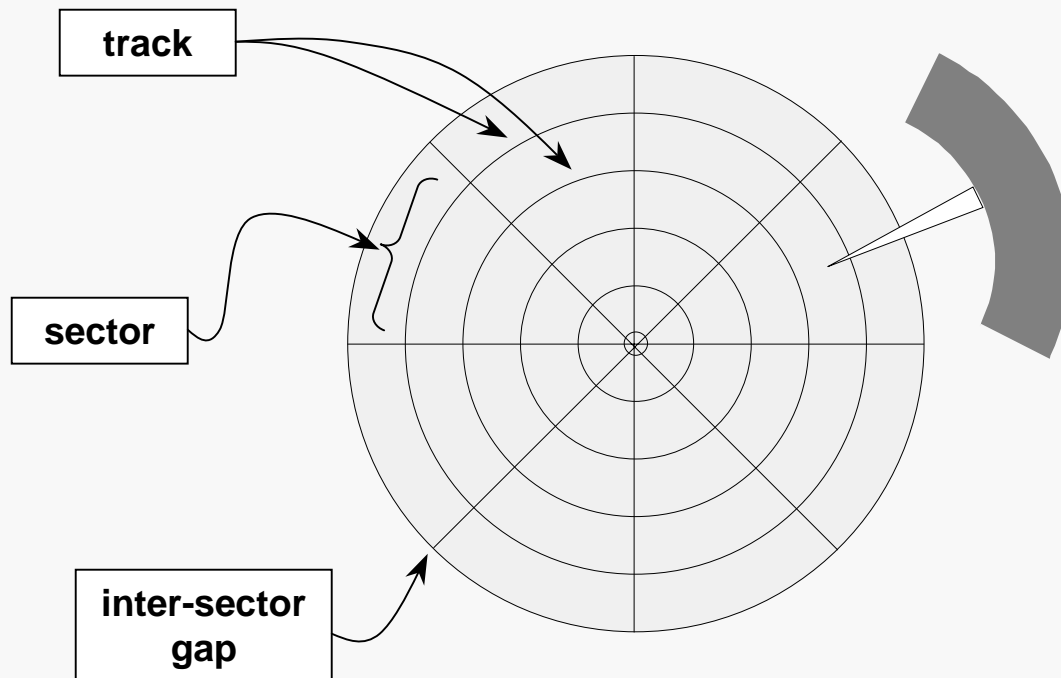
The data is read and written by I/O heads, typically one per data surface. These I/O heads are mounted in a rack and all move in and out in unison. Typically, only one I/O head can be active (reading or writing data) at any given time.



Each data surface is organized into concentric rings of data, called tracks.

Each track is divided into a number of segments, called sectors.

We will assume that each sector contains the same amount of data, regardless of which track it belongs to. (That is actually true for some older disks.)

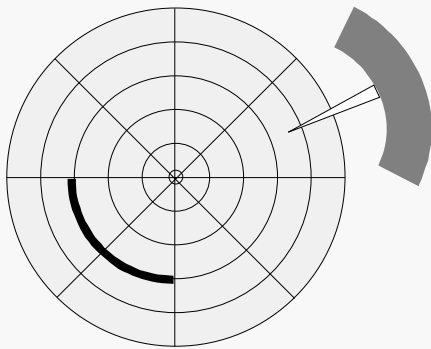


The basic unit of disk I/O is the sector.

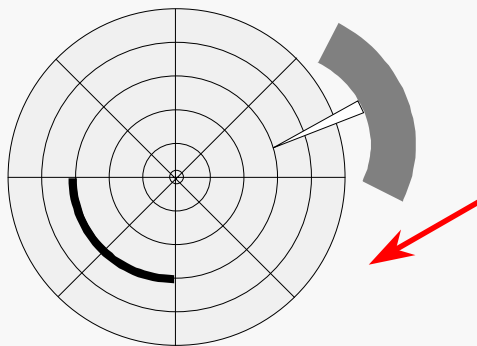
We will see shortly that reading an entire sector from disk takes only slightly longer than reading a single byte.

When a disk read is requested the following actions must be carried out:

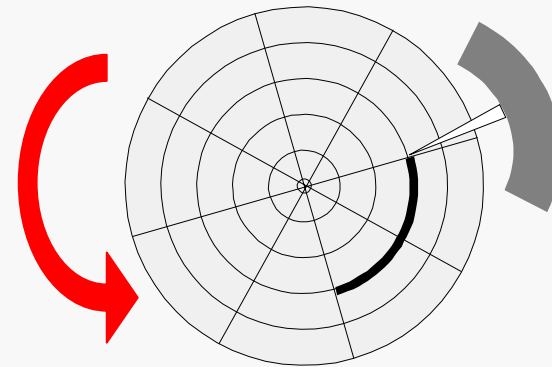
**1. Determine what sector(s) must be read — the disk controller is responsible for this.**



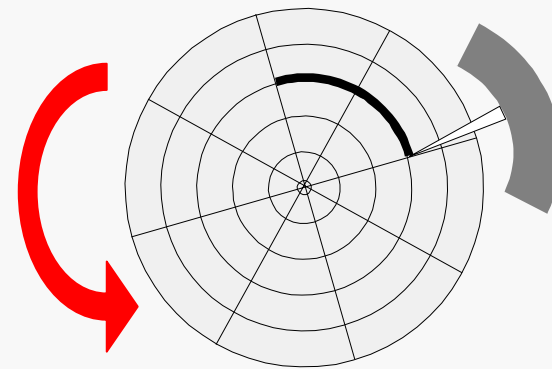
**2. Move the read head to the track containing the targeted sector(s).**



**3. Wait for the beginning of the first sector to rotate to the head.**



**4. Read the data as it rotates beneath the head.**



seek time the time required to move the head to the appropriate track

The seek time depends only on the speed with which the head rack moves, and the number of tracks that the head must move across to reach its target.

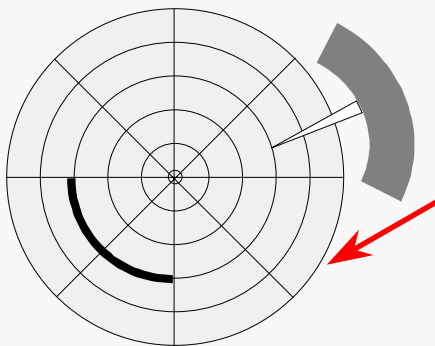
Given the following (which are constant for a particular disk):

$H_s$  = the time for the I/O head to start moving

$H_T$  = the time for the I/O head to move from one track to the next

then the time for the head to move  $n$  tracks is:

$$\text{Seek}(n) = H_s + H_T n$$



**QTP: On average, the head will move 1/3 of the way across the platter. Why?**

latency the time required for the appropriate sector to rotate to the position of the I/O head

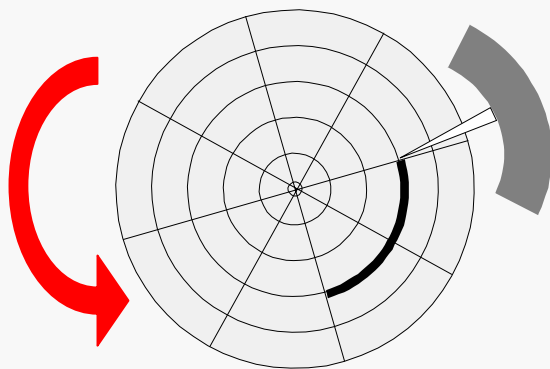
The rotational latency time depends only on the speed at which the spindle rotates, and the angle through which the track must rotate to reach the I/O head.

Given the following:

$R$  = the rotational speed of the spindle (in rotations per second)

$\theta$  = the number of radians through which the track must rotate

then the rotational latency is:



$$\text{Latency}(\theta) = \left( \frac{\theta}{2\pi} \right) \times \left( \frac{1000}{R} \right)$$

**On average, the latency time will be the time required for the platter to complete 1/2 of a full rotation.**

transfer time the time required for the appropriate sector(s) to rotate under the I/O head (for contiguous sectors on the same track)

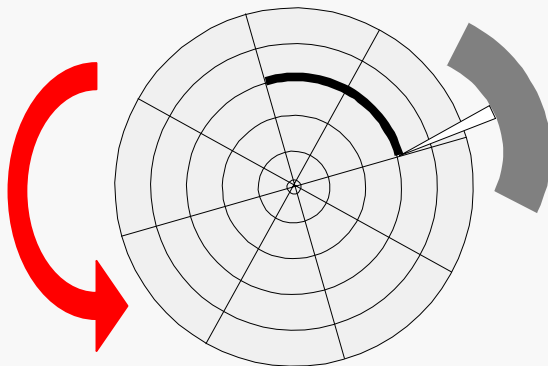
The transfer time depends only on the speed at which the spindle rotates, and the number of sectors that must be read.

Given:

$S_T$  = the total number of sectors per track

the transfer time for  $n$  contiguous sectors on the same track is:

$$\text{Transfer}(n) = \left( \frac{n}{S_T} \right) \times \left( \frac{1000}{R} \right)$$



**For sectors on different tracks, each track must be analyzed separately, allowing for seek and latency.**

The total time to read/write data is from/to disk is then the sum of the seek time, the rotational latency time, and the transfer time.

This ignores the time for controller logic, delays due to multitasking queues, etc.; that is fair because those times are normally orders of magnitude less than the times considered.

One note: we have assumed in the analysis of the data transfer time that the I/O heads are capable of reading/writing data as fast as the sector moves beneath the head. On older disk systems that was often not the case, and interleaving was necessary to optimize read/write performance. On most contemporary disk systems, the I/O heads are fast enough that such tricks are unnecessary.

## Example: Contiguous File Read

Assume a disk system with the following parameters:

<b>Total capacity:</b>	<b>675 MB</b>
<b># of platters:</b>	<b>15</b> <b>(30 data surfaces)</b>
<b># of tracks per platter:</b>	<b>576</b>
<b># of sectors per track:</b>	<b>160</b>
<b>cluster size:</b>	<b>8 KB</b>
<b>spindle speed:</b>	<b>3600 RPM</b>
<b>head start time:</b>	<b>3 ms</b>
<b>head track traversal time:</b>	<b>0.08 ms</b>

Question: how much time would be required, on average, to read a file of size 140 KB, assuming the file is stored in contiguous sectors on adjacent tracks (since one track won't hold the whole file in this case)?

## Example: Contiguous File Read

We must compute some values from the given disk parameters:

<b>rotations per second:</b>	<b>60</b>
<b>capacity of one platter:</b>	<b>45 MB</b>
<b>capacity of one track:</b>	<b>80 KB</b>
<b>capacity of one sector:</b>	<b>512 bytes</b>
<b>sectors per cluster:</b>	<b>16</b>
<b>clusters per track:</b>	<b>10</b>

So a 140 KB file would occupy 18 clusters, or one full track plus 8 clusters on an adjacent track. The last cluster is only partially filled, but we will ignore that in our calculations.

For reading the first (full) track, assuming average values:

- seek time =  $3 + (576/3)*0.08 = 18.36$  ms
  - latency time =  $(1/2)*(1000/60) = 8.33$  ms
  - transfer time =  $(160/160)*(1000/60) = 16.67$  ms
- } total: 43.36 ms

## Example: Contiguous File Read

For reading the relevant sectors from the second track, assuming average values:

- seek time =  $3 + (1) * 0.08 = 3.36$  ms
  - latency time =  $(1/2) * (1000/60) = 8.33$  ms
  - transfer time =  $(128/160) * (1000/60) = 13.33$  ms
- } total: 24.74 ms

So the total time to read the file into memory\* would be about 68.1 ms.

**\* Actually this is the time to read the file into the disk buffer memory.**

## Example: Fragmented File Read

What if the file is scattered all over the disk?

Then it will take a lot longer... for each cluster:

- seek time =  $3 + (576/3)*0.08 = 18.36$  ms
  - latency time =  $(1/2)*(1000/60) = 8.33$  ms
  - transfer time =  $(16/160)*(1000/60) = 1.67$  ms
- } total: 28.36 ms

Since the file occupies (part or all of) 18 clusters, the total read time for a completely fragmented copy of the file would be about 510.48 ms.

# Sector Read vs Single Byte Read

The programs you have implemented simply perform read/write actions, usually on single variables which may store from 1 to a few hundred bytes of data.

Why do typical disk controllers perform read/write operations in sector-sized chunks? Performance... given the disk system described earlier:

time to read one full track: 43.36 ms (slide 14)

time to read one full sector: 26.79 ms

$$18.36 \text{ ms} + 8.33 \text{ ms} + (1/160) * (1000/60)$$

time to read one byte: 26.69 ms

$$18.36 \text{ ms} + 8.33 \text{ ms} + (1/512) * (1/160) * (1000/60)$$

**The time to read one sector is only slightly more than the time to read a single byte, largely because the seek and latency times dominate the transfer time.**

In view of the previous slide, it makes sense to design programs so that data is read from and written to disk in relatively large chunks... but there is more.

## Spatial Locality of Reference

In many cases, if a program accesses one part of a file, there is a high probability that the program will access nearby parts of the file in the near future.

**Moral: grab a larger chunk than you immediately need.**

## Temporal Locality of Reference

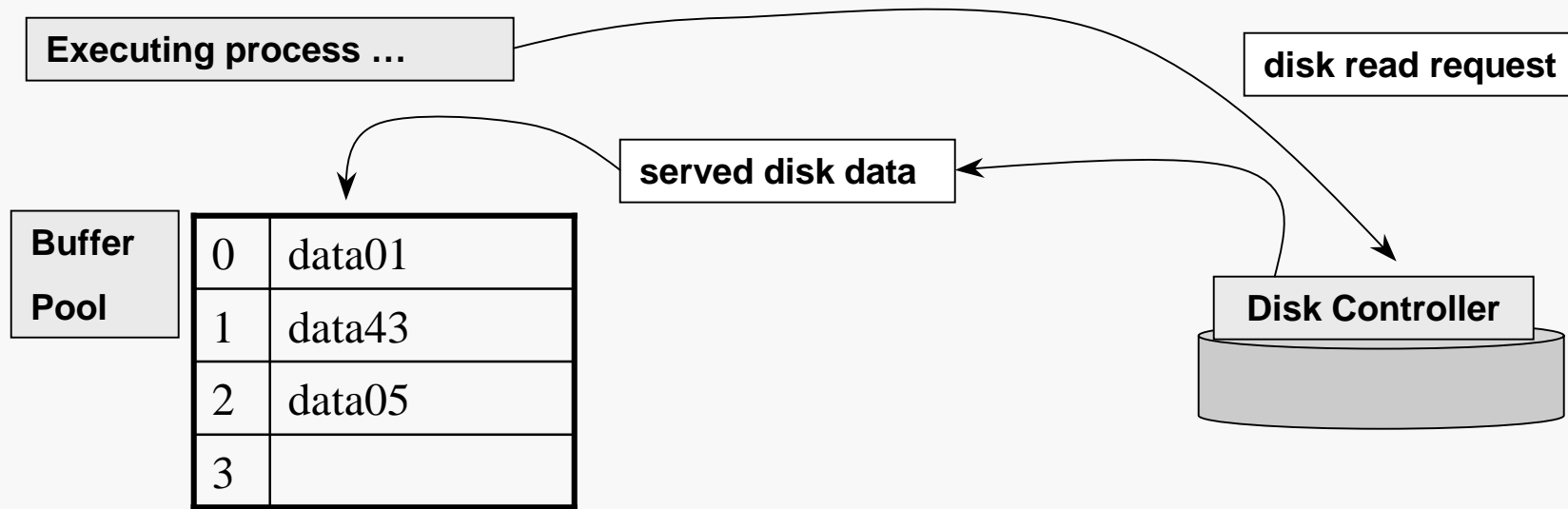
In many cases, if a program accesses one part of a file, there is a high probability that the program will access the same part of the file again in the near future.

**Moral: once you've grabbed a chunk, keep it around.**

buffer pool a series of buffers (memory locations) used by a program to cache disk data

A program that does much disk I/O can often improve its performance by employing a buffer pool to take advantage of locality of reference.

Basically, the buffer pool is just a collection of data chunks. The program reads and writes data in buffer-sized chunks, storing newly-read data chunks into the pool, replacing currently stored chunks as necessary.



The buffer pool must be organized physically and logically.

The physical organization is generally an ordered list of some sort.

The logical organization depends upon how the buffer pool deals with the issue of replacement — if a new data chunk must be added to the pool and all the buffers are currently full, one of the current elements must be replaced.

If the replaced element has been modified, it (usually) must be written back to disk or the changes will be lost.

Common buffer replacement strategies:

FIFO (first-in is first-out) organize buffers as a queue

LFU (least frequently used) replace the least-accessed buffer

LRU (least recently used) replace the longest-idle buffer

# FIFO Replacement

Logically the buffer pool is treated as a queue:

Requested "chunks":

0	0			
4	0	4		
17	0	4	17	
4	0	4	17	
13	0	4	17	13
7	4	17	13	7
8	17	13	7	8
17	17	13	7	8
0	13	7	8	0
...				

Takes no notice of the access pattern exhibited by the program. Consider what would happen with the sequence:

0 1 0 2 0 3 0 4 0 5 0 6 ...

# LFU Replacement

For LFU we must maintain an access count for each element of the buffer pool. It is also useful to keep the elements sorted by that count.

Requested "chunks":

0	0,1			
4	0,1	4,1		
17	0,1	4,1	17,1	
4	0,1	4,2	17,1	
13	0,1	4,2	17,1	13,1
7	4,2	17,1	13,1	7,1
8	4,2	13,1	7,1	8,1
17	4,2	7,1	8,1	17,1
0	4,2	8,1	17,1	0,1
...				



Aside from cost of storing and maintaining counter values, and searching for least value, consider the sequence:

- 0 (500 times)
- 1 (500 times)
- 2 3 4 2 3 4 ...

With LRU, we may use a simple list structure. On an access, we move the targeted element to the front of the list. That puts the least recently used element at the tail of the list.

Requested "chunks":

0	0			
4	4	0		
17	17	4	0	
4	4	17	0	
13	13	4	17	0
7	7	13	4	17
8	8	7	13	4
17	17	8	7	13
0	0	17	8	7
...	...	...	...	...



Consider what would happen with the sequence:

0 1 2 3 4 0 1 2 3 4 ...

In general, none of these replacement strategies is best in all cases.

All are used with some frequency.

Intuitively, LRU and LFU make more sense than FIFO.

The performance you get is determined by the access pattern exhibited by the running program, and that is often impossible to predict.

Optimal buffer replacement strategy:

replace the buffer element whose next access lies furthest in the future.