

# OPERATING SYSTEMS & SYSTEMS PROGRAMMING I

## MEMORY MANAGEMENT

---

Keith Bugeja / Joshua Ellul / Alessio Magro / Kevin Vella

April 30, 2018

Department of Computer Science  
Faculty of ICT, University of Malta

# OVERVIEW

---

## Introduction

- Process address space
- Pages and paging
- Memory regions

## Memory allocation

- Primitives
- Manipulation
- Common problems

## Advanced memory allocation

- Anonymous memory mappings
- Shared memory

# MEMORY MANAGEMENT

---

Linux *virtualises* physical memory

- Physical memory not directly accessible by processes
- Processes are presented with a logical view of memory
- Kernel associates a virtual address space with each process:
  - Linear** starts at zero and contiguously increases, and
  - Flat** directly accessible without the need for segmentation

In memory management, the page is the smallest addressable unit of memory that the memory management unit (MMU) can manage:

In memory management, the page is the smallest addressable unit of memory that the memory management unit (MMU) can manage:

- The virtual address space is divided into pages, the size of which is determined by the machine architecture (e.g. 4 KB for 32-bit and 8 KB for 64-bit systems)

In memory management, the page is the smallest addressable unit of memory that the memory management unit (MMU) can manage:

- The virtual address space is divided into pages, the size of which is determined by the machine architecture (e.g. 4 KB for 32-bit and 8 KB for 64-bit systems)
- A 32-bit address space contains around 1 million 4 KB pages
  - Not all pages necessarily correspond to anything and thus are either *valid* or *invalid*
  - A valid page is associated with an actual page of data, either in physical memory (RAM) or on secondary storage
  - Access to an invalid (unmapped) page results in a segmentation violation



When a page is on secondary storage, attempts to access it causes the memory management unit to generate a *page fault*

- Kernel intervenes and transparently *pages in* the data from secondary storage to physical memory
- The kernel may have to move data out of memory to make space for data paged in - a process called *paging out*
  - To minimise subsequent page ins, the kernel attempts to page out (evict) data that is least likely to be used in the near future

Multiple pages of virtual memory, even in different virtual address spaces owned by different processes, may map to a single physical page

- virtual address spaces may share data in physical memory (e.g. many processes using the standard C library map it into their virtual address space, even though only one copy exists in physical memory).
- shared data may be read-only, writable or both read and writable
- basis for copy-on-write (COW); MMU intercepts write operation and raises exception - kernel in response creates a new copy of the page for the writing process

# MEMORY REGIONS

Kernel arranges pages into blocks that share certain properties (e.g. access permissions); these blocks are called *mappings* or *memory regions*.

**text segment** contains process program code, string literals, constant variables and other read-only data. Mapped directly from object file

**stack** contains process execution stack, which grows and shrinks dynamically

**data segment/heap** contains the process dynamic memory

**bss segment** contains uninitialised global variables (all zeros, according to C standard)

## Hint

Examine `/proc/self/maps`

# MEMORY ALLOCATION

---

Foundation of memory management is the allocation, use and eventual return of *dynamic memory*

- Allocated at runtime (not compile time)
- Sizes may be unknown until the moment of allocation

# ALLOCATING DYNAMIC MEMORY

The C interface for allocating dynamic memory is `malloc`:

```
1  #include <stdlib.h>
2
3  void * malloc (size_t size);
```

# ALLOCATING DYNAMIC MEMORY

The C interface for allocating dynamic memory is `malloc`:

```
1  #include <stdlib.h>
2
3  void * malloc (size_t size);
```

Allocates `size` bytes of memory and returns a pointer to the allocated region

- contents of memory are undefined!

# ALLOCATING DYNAMIC MEMORY

The C interface for allocating dynamic memory is `malloc`:

```
1 #include <stdlib.h>
2
3 void * malloc (size_t size);
```

Allocates `size` bytes of memory and returns a pointer to the allocated region

- contents of memory are undefined!

On failure, `malloc` returns `NULL` and sets `errno` to `ENOMEM`



For instance, to allocate an 8 KB chunk of memory:

# ALLOCATING DYNAMIC MEMORY

For instance, to allocate an 8 KB chunk of memory:

```
1 char *buffer = malloc(8192);  
2  
3 if (!buffer)  
4     perror("malloc");
```

# ALLOCATING DYNAMIC MEMORY

For instance, to allocate an 8 KB chunk of memory:

```
1 char *buffer = malloc(8192);  
2  
3 if (!buffer)  
4     perror("malloc");
```

## Note

C automatically promotes pointers to `void` to any other pointer type on assignment.

## ALLOCATING DYNAMIC MEMORY

It is important to check the return value from a call to `malloc` and handle error conditions

- many programs employ a `malloc` wrapper that prints an error message and terminates the program if `malloc` returns `NULL`

# ALLOCATING DYNAMIC MEMORY

It is important to check the return value from a call to `malloc` and handle error conditions

- many programs employ a `malloc` wrapper that prints an error message and terminates the program if `malloc` returns `NULL`

```
1 void * xmalloc(size_t size)
2 {
3     void *buffer = malloc(size);
4
5     if (!buffer) {
6         perror("xmalloc");
7         exit(EXIT_FAILURE);
8     }
9
10    return buffer;
11 }
```

# ALLOCATING ARRAYS

Dynamic allocation of arrays, where the size of an element may be fixed but the number of elements is dynamic, is accomplished through the `calloc` function:

```
1 #include <stdlib.h>
2
3 void * calloc (size_t nr, size_t size);
```

# ALLOCATING ARRAYS

Dynamic allocation of arrays, where the size of an element may be fixed but the number of elements is dynamic, is accomplished through the `calloc` function:

```
1 #include <stdlib.h>
2
3 void * calloc (size_t nr, size_t size);
```

Allocates `size * nr` bytes of memory and returns a pointer to the allocated region

- contents of memory are zeroed!

# ALLOCATING ARRAYS

Dynamic allocation of arrays, where the size of an element may be fixed but the number of elements is dynamic, is accomplished through the `calloc` function:

```
1 #include <stdlib.h>
2
3 void * calloc (size_t nr, size_t size);
```

Allocates `size * nr` bytes of memory and returns a pointer to the allocated region

- contents of memory are zeroed!

On failure, `calloc` returns `NULL` and sets `errno` to `ENOMEM`



# ALLOCATING ARRAYS

In the same vein as `xmalloc`, it is commonplace to define wrapper methods that allocate and zero memory:

```
1  #include <stdlib.h>
2
3  void * malloc0 (size_t size) {
4      return calloc(1, size);
5  }
```

# ALLOCATING ARRAYS

In the same vein as `xmalloc`, it is commonplace to define wrapper methods that allocate and zero memory:

```
1  #include <stdlib.h>
2
3  void * malloc0 (size_t size) {
4      return calloc(1, size);
5  }
```

## Exercise

Implement `xmalloc0`, a method that allocates and zeroes memory, and terminates the program with an error if the allocation fails.

## RESIZING ALLOCATIONS

The size of an existing allocation may need to change (shrink or grow), e.g. dynamic arrays or contiguous list of items, strings, etc.

## RESIZING ALLOCATIONS

The size of an existing allocation may need to change (shrink or grow), e.g. dynamic arrays or contiguous list of items, strings, etc.

C language provides **realloc** for resizing allocations:

```
1  #include <stdlib.h>
2
3  void * realloc (void *ptr, size_t size);
```

## RESIZING ALLOCATIONS

The size of an existing allocation may need to change (shrink or grow), e.g. dynamic arrays or contiguous list of items, strings, etc.

C language provides **realloc** for resizing allocations:

```
1 #include <stdlib.h>
2
3 void * realloc (void *ptr, size_t size);
```

Resizes the region pointed by **ptr** to **size** bytes

- returned pointer may or may not be the same as **ptr**
- operation may be costly due to possible copy
- if **size** is 0, effect is the same as freeing memory
- if **ptr** is **NULL**, result is the same as **malloc**

## RESIZING ALLOCATIONS

The size of an existing allocation may need to change (shrink or grow), e.g. dynamic arrays or contiguous list of items, strings, etc.

C language provides **realloc** for resizing allocations:

```
1 #include <stdlib.h>
2
3 void * realloc (void *ptr, size_t size);
```

Resizes the region pointed by **ptr** to **size** bytes

- returned pointer may or may not be the same as **ptr**
- operation may be costly due to possible copy
- if **size** is 0, effect is the same as freeing memory
- if **ptr** is **NULL**, result is the same as **malloc**

On failure, **realloc** returns **NULL** and sets **errno** to **ENOMEM**; the state of the memory pointed by **ptr** is unchanged.

## FREEING DYNAMIC MEMORY

Dynamic allocations are permanent parts of the process address space until they are manually freed.

# FREEING DYNAMIC MEMORY

Dynamic allocations are permanent parts of the process address space until they are manually freed.

Memory allocated with `malloc`, `calloc` and `realloc` must be returned to the system when no longer in use via `free`:

```
1 #include <stdlib.h>
2
3 void free (void *ptr);
```



# FREEING DYNAMIC MEMORY

Dynamic allocations are permanent parts of the process address space until they are manually freed.

Memory allocated with `malloc`, `calloc` and `realloc` must be returned to the system when no longer in use via `free`:

```
1 #include <stdlib.h>
2
3 void free (void *ptr);
```

A call to `free` releases the memory at `ptr`

- previously allocated using one of the functions above
- releases entire block of allocated memory
- if `ptr` is `NULL`, `free` does nothing

## EXAMPLE

Time for an example: a function that concatenates two strings into a third

```
1 char *concatenate(char *s1, char *s2)
2 {
3     // make sure both are valid strings
4     assert(s1 != NULL && s2 != NULL);
5
6     // assume strings are null-terminated
7     // get respective lengths
8     size_t l1 = strlen(s1),
9           l2 = strlen(s2);
10
11     // allocate a zeroed memory chunk the cumulative size of the
12     // two strings plus an additional character for null termination
13     char *s = calloc(1, l1 + l2 + 1);
14     if (!s) {
15         perror("calloc");
16         exit(EXIT_FAILURE);
17     }
18
19     // copy strings into newly allocated buffer
20     strcpy(s, s1);
21     strcpy(s + l1, s2);
22
23     return s;
24 }
```

## EXAMPLE

Next we use the function to concatenate two strings:

```
1  #include <stdio.h>
2  #include <stdlib.h>
3  #include <string.h>
4  #include <assert.h>
5
6  char *concatenate(char *s1, char *s2); // concatenate goes here
7
8  int main(int argc, char *argv[]) {
9      // define strings s1 and s2
10     char *s1 = "abra", *s2 = "cadabra";
11
12     // call concatenate; returns a pointer to the newly concatenated string
13     char *s = concatenate(s1, s2);
14
15     // print output
16     printf("%s + %s = %s\n", s1, s2, s);
17
18     // free memory allocated by concatenate function
19     free(s);
20 }
```

## EXAMPLE

Running the example program displays the newly concatenated string:

```
1 $ gcc -o concat concat.c
2 $ ./concat
3 abra + cadabra = abracadabra
4 $
```

## EXAMPLE

Running the example program displays the newly concatenated string:

```
1 $ gcc -o concat concat.c
2 $ ./concat
3 abra + cadabra = abracadabra
4 $
```

### Exercise

Change the program to get the source strings **s1** and **s2** from the command line arguments.

There are a number of common pitfalls programmers fall into when manually managing memory:

- memory leak
- use-after-free
- dangling pointer

Memory leaks are among the most common and detrimental mishaps in C programming:

- typically triggered by a missing invocation to **free**
- memory is never returned to the system (during the lifetime of the process)
- process loses its only reference to the memory and becomes unable to access it again

## EXAMPLE

This program exhibits memory leaks:

```
1  ...
2
3  // concatenate prototype
4  char *concatenate(char *s1, char *s2);
5
6  char *concatenate4(char *s1, char *s2, char *s3, char *s4) {
7      char *s12 = concatenate(s1, s2); // s12 = s1s2
8      char *s34 = concatenate(s3, s4); // s34 = s3s4
9      char *s = concatenate(s12, s34); // s = s12s34
10
11     return s;
12 }
13
14 int main(int argc, char *argv[]) {
15     char *s1 = "You ", *s2 = "shall ",
16         *s3 = "not ", *s4 = "pass.";
17
18     // concatenate the four strings
19     char *s = concatenate4(s1, s2, s3, s4);
20
21     // print output
22     printf("%s + %s + %s + %s = %s\n", s1, s2, s3, s4, s);
23
24     // free memory block returned by concatenate4
25     free(s);
26 }
```



The program leaks memory in the `concatenate4` function:

- `concatenate` returns allocated memory for `s12`, `s34` and `s`
- memory pointed to by `s12` and `s34` is never freed
- being locals, `s12` and `s34` are lost once `concatenate4` returns
- memory cannot be freed once `s12` and `s34` go out of scope

Use-after-free occurs when a block of memory is freed and then subsequently accessed:

- once **free** is called on a block of memory, a program should never access its contents again

## EXAMPLE

Use-after-free is shown in the example below, where the order of `printf` and `free` is switched:

```
1  ...
2  char *concatenate(char *s1, char *s2); // concatenate goes here
3
4  int main(int argc, char *argv[]) {
5      // define strings s1 and s2
6      char *s1 = "abra", *s2 = "cadabra";
7
8      // call concatenate; returns a pointer to the newly concatenated string
9      char *s = concatenate(s1, s2);
10
11     // free memory allocated by concatenate function
12     free(s);
13
14     // print output
15     printf("%s + %s = %s\n", s1, s2, s);
16 }
```

Dangling pointers are non-NULL pointers that point to invalid (or deallocated) memory

- become invalid due to deallocation (e.g. calls to `free`, references going out of scope, etc.)
- pointer is not set to `NULL`

## EXAMPLE

Example snippet demonstrating use-after-free of a dangling pointer:

```
1  ...
2  // print s if not a NULL pointer
3  void print_string(char *s) {
4      if (s) printf("%s\n", s);
5  }
6
7  int main(int argc, char *argv[]) {
8      char *s1 = "abra", *s2 = "cadabra";
9      char *s = concatenate(s1, s2);
10     printf("%s + %s = %s\n", s1, s2, s);
11
12     // show pointer before call to free
13     printf("Pointer before free is %p\n", s);
14
15     // free memory allocated by concatenate function
16     free(s);
17
18     // show pointer after call to free
19     printf("Pointer after free is %p\n", s);
20
21     // pass the dangling pointer to function
22     print_string(s);
23 }
```

## EXAMPLE

Running the example shows that the value of the pointer before and after the call to **free** remains the same:

```
1 $ ./dangling
2 abra + cadabra = abracadabra
3 Pointer before free is 0x56446c8aa260
4 Pointer after free is 0x56446c8aa260
5
6 $
```

## EXAMPLE

Running the example shows that the value of the pointer before and after the call to **free** remains the same:

```
1 $ ./dangling
2 abra + cadabra = abracadabra
3 Pointer before free is 0x56446c8aa260
4 Pointer after free is 0x56446c8aa260
5
6 $
```

Thus, **print\_string** has no way of determining whether the pointer is valid or dangling.

- calling **printf** on the freed memory block results in undefined behaviour

# DETECTING MEMORY ERRORS

On Linux, a number of tools, such as *Valgrind*, are available to help us detect and debug memory errors in our programs

- the *Valgrind* tool suite can automatically detect many memory management and threading bugs
- can perform detailed profiling to help spot bottlenecks in your programs
- to install, open a terminal session and type:

1

```
$ sudo apt install valgrind
```



## VALGRIND EXAMPLES

Let's run `valgrind` on the concatenation program that suffered from memory leaks:

```
1 $ valgrind ./concat4
2 ...
3 ==39702== Command: ./concat4
4 ==39702==
5 You + shall + not + pass. = You shall not pass.
6 ==39702==
7 ==39702== HEAP SUMMARY:
8 ==39702==      in use at exit: 21 bytes in 2 blocks
9 ==39702==    total heap usage: 4 allocs, 2 frees, 1,065 bytes allocated
10 ==39702==
11 ==39702== LEAK SUMMARY:
12 ==39702==    definitely lost: 21 bytes in 2 blocks
13 ==39702==    indirectly lost: 0 bytes in 0 blocks
14 ==39702==    possibly lost: 0 bytes in 0 blocks
15 ==39702==    still reachable: 0 bytes in 0 blocks
16 ==39702==    suppressed: 0 bytes in 0 blocks
```

# VALGRIND EXAMPLES

And again on the program with the dangling pointer:

```
1  $ valgrind ./dangling
2  ...
3  ==40010== Command: ./dangling
4  ==40010==
5  abra + cadabra = abracadabra
6  Pointer before free is 0x521c040
7  Pointer after free is 0x521c040
8  ...
9  ==40010==      by 0x1088A8: print_string
10 ==40010==      by 0x10894F: main
11 ==40010== Address 0x521c040 is 0 bytes inside a block of size 12
    ↪ free'd
12 ==40010==      at 0x4C30D3B: free
13 ==40010==      by 0x10892B: main
14 ==40010== Block was alloc'd at
15 ==40010==      at 0x4C31B25: calloc
16 ==40010==      by 0x1089CF: concatenate
17 ==40010==      by 0x1088E3: main
18 ...
```

# MANAGING THE DATA SEGMENT

---

# MANAGING THE DATA SEGMENT

Unix systems have traditionally provided interfaces to allow direct management of the data segment:

```
1 #include <unistd.h>
2
3 int brk (void *end);
4 void * sbrk (intptr_t increment);
```

# MANAGING THE DATA SEGMENT

Unix systems have traditionally provided interfaces to allow direct management of the data segment:

```
1 #include <unistd.h>
2
3 int brk (void *end);
4 void * sbrk (intptr_t increment);
```

Inheritance from legacy systems where the data segment and the heap lived in the same segment:

- heap grows upward from the bottom of segment
- stack grows downward from the top of segment
- line of demarcation between them is called the *break* or *break point*

# MANAGING THE DATA SEGMENT

Unix systems have traditionally provided interfaces to allow direct management of the data segment:

```
1 #include <unistd.h>
2
3 int brk (void *end);
```

A call to **brk** sets the break point address to the value specified by **end**

- if successful, the call returns 0
- on failure, it returns -1 and sets **errno** to **ENOMEM**

# MANAGING THE DATA SEGMENT

Unix systems have traditionally provided interfaces to allow direct management of the data segment:

```
1 #include <unistd.h>
2
3 void * sbrk (intptr_t increment);
```

A call to **sbrk** changes the break point address by the offset specified in **increment**:

- an increment of 0 returns the current break point

## EXAMPLE

The following example shows how calls to `malloc` affect the current break address:

```
1  #include <stdio.h>
2  #include <stdlib.h>
3  #include <unistd.h>
4
5  int main(int argc, char *argv[]) {
6      const int blocks = 50;
7      const int block_size = 32 * 1024;
8      char *block_list[blocks];
9
10     for (int j = 0; j < blocks; ++j)
11     {
12         block_list[j] = malloc(block_size);
13         printf("Current break is at %p\n", sbrk(0));
14     }
15
16     for (int j = 0; j < blocks; ++j)
17         free (block_list[j]);
18 }
```



## EXAMPLE

The program outputs the current break address after each 32 KB allocation:

```
1 $ ./sbrk
2 Current break is at 0x5588ec83e000
3 Current break is at 0x5588ec83e000
4 Current break is at 0x5588ec83e000
5 Current break is at 0x5588ec83e000
6 Current break is at 0x5588ec866000
7 ...
8 Current break is at 0x5588ec9a6000
9 Current break is at 0x5588ec9a6000
10 Current break is at 0x5588ec9ce000
11 $
```

## EXAMPLE

The program outputs the current break address after each 32 KB allocation:

```
1 $ ./sbrk
2 Current break is at 0x5588ec83e000
3 Current break is at 0x5588ec83e000
4 Current break is at 0x5588ec83e000
5 Current break is at 0x5588ec83e000
6 Current break is at 0x5588ec866000
7 ...
8 Current break is at 0x5588ec9a6000
9 Current break is at 0x5588ec9a6000
10 Current break is at 0x5588ec9ce000
11 $
```

### Note

The use of `malloc` should be preferred over `brk` and `sbrk` for its convenient and portable memory allocation interface.

`malloc` is classically implemented using an algorithm called *buddy memory allocation scheme*

- data segment is divided into a series of power-of-2 partitions
- allocations are satisfied by returning the partition that is the closest fit to the requested size
- memory is freed by marking the partition as *free*
- adjacent free partitions can be coalesced into a single larger partition
- if the top of the heap is free, the system can use `brk` to lower the break point, shrinking the heap and returning memory to the kernel

The buddy memory allocation scheme has the advantage of being fast and simple; however:

- introduces internal and external fragmentation
- results in inefficient use of memory or failed memory allocations
- allows memory allocations to *pin* one another, preventing the heap from shrinking after large memory sections are freed

# ANONYMOUS MEMORY MAPPINGS

---

To avoid the problems with the buddy memory allocation scheme, `glibc` does not use the heap for large allocations:

- creates an anonymous memory mapping to satisfy a large allocation request
- large zero-filled blocks of memory, ready for use

# ANONYMOUS MEMORY MAPPING

Allocating memory via anonymous mappings has several benefits:

- no fragmentation concerns; memory is returned to the system as soon as mapping is unmapped
- mappings are resizable;
- exist in a separate memory mapping, removing any need to manage the global heap

# ANONYMOUS MEMORY MAPPING

Allocating memory via anonymous mappings has several benefits:

- no fragmentation concerns; memory is returned to the system as soon as mapping is unmapped
- mappings are resizable;
- exist in a separate memory mapping, removing any need to manage the global heap

There are also two downsides to using anonymous memory mappings rather than the heap:

- each mapping is an integer multiple of the system page size, possibly resulting in wasted space
- creating a new memory mapping incurs more overhead than satisfying an allocation from the heap



## EXAMPLE

The following example demonstrates the behaviour of `malloc` vis-a-vis small and large allocations:

```
1  #include <stdio.h>
2  #include <stdlib.h>
3  #include <unistd.h>
4
5  int main(int argc, char *argv[]) {
6      const int blocks = 4;
7      const int block_size_heap = 64 * 1024;
8      const int block_size_amm = block_size_heap * 1024;
9
10     char *block_list[blocks];
11
12     printf("Allocation size: %d\n", block_size_heap);
13     for (int j = 0; j < blocks; ++j) {
14         block_list[j] = malloc(block_size_heap);
15         printf("Address : %p, Break : %p\n", block_list[j], sbrk(0));
16     }
17
18     printf("Allocation size: %d\n", block_size_amm);
19     for (int j = 0; j < blocks; ++j) {
20         free(block_list[j]);
21         block_list[j] = malloc(block_size_amm);
22         printf("Address : %p, Break : %p\n", block_list[j], sbrk(0));
23     }
24
25     for (int j = 0; j < blocks; ++j)
26         free (block_list[j]);
27 }
```

## EXAMPLE

The program outputs the address of the newly allocated block together with the current program break address:

```
1 $ ./malloc
2 Allocation size: 65536
3 Address : 0x561a835a8670, Break : 0x561a835c9000
4 Address : 0x561a835b8680, Break : 0x561a835c9000
5 Address : 0x561a835c8690, Break : 0x561a835f9000
6 Address : 0x561a835d86a0, Break : 0x561a835f9000
7 Allocation size: 67108864
8 Address : 0x7f882981f010, Break : 0x561a835f9000
9 Address : 0x7f882581e010, Break : 0x561a835f9000
10 Address : 0x7f882181d010, Break : 0x561a835f9000
11 Address : 0x7f881d81c010, Break : 0x561a835c9000
12 $
```

## EXAMPLE

The program outputs the address of the newly allocated block together with the current program break address:

```
1 $ ./malloc
2 Allocation size: 65536
3 Address : 0x561a835a8670, Break : 0x561a835c9000
4 Address : 0x561a835b8680, Break : 0x561a835c9000
5 Address : 0x561a835c8690, Break : 0x561a835f9000
6 Address : 0x561a835d86a0, Break : 0x561a835f9000
7 Allocation size: 67108864
8 Address : 0x7f882981f010, Break : 0x561a835f9000
9 Address : 0x7f882581e010, Break : 0x561a835f9000
10 Address : 0x7f882181d010, Break : 0x561a835f9000
11 Address : 0x7f881d81c010, Break : 0x561a835c9000
12 $
```

The larger 64 MB allocations all hold addresses larger than the current program break

- allocated as anonymous memory mappings outside the global heap

## EXAMPLE

Running the program through `strace` provides us with more insight:

```
1 $ strace ./malloc
2 ...
3 write(1, "Allocation size: 65536\n", 23Allocation size: 65536) = 23
4 write(1, "Address : 0x5566bad8c680, Break "..., 49Address :
   ↪ 0x5566bad8c680, Break : 0x5566bad9d000) = 49
5 brk(0x5566badcd000)
6 ...
7 write(1, "Allocation size: 67108864\n", 26Allocation size: 67108864) =
   ↪ 26
8 mmap(NULL, 67112960, PROT_READ|PROT_WRITE, MAP_PRIVATE|MAP_ANONYMOUS,
   ↪ -1, 0) = 0x7f3d43c80000
9 write(1, "Address : 0x7f3d43c80010, Break "..., 49Address :
   ↪ 0x7f3d43c80010, Break : 0x5566badcd000) = 49
10 ...
11 munmap(0x7f3d43c80000, 67112960)          = 0
12 ...
13 exit_group(0)                             = ?
14 +++ exited with 0 +++
```

## EXAMPLE

Running the program through **strace** provides us with more insight:

```
1 $ strace ./malloc
2 ...
3 write(1, "Allocation size: 65536\n", 23Allocation size: 65536) = 23
4 write(1, "Address : 0x5566bad8c680, Break "..., 49Address :
   ↳ 0x5566bad8c680, Break : 0x5566bad9d000) = 49
5 brk(0x5566badcd000)
6 ...
7 write(1, "Allocation size: 67108864\n", 26Allocation size: 67108864) =
   ↳ 26
8 mmap(NULL, 67112960, PROT_READ|PROT_WRITE, MAP_PRIVATE|MAP_ANONYMOUS,
   ↳ -1, 0) = 0x7f3d43c80000
9 write(1, "Address : 0x7f3d43c80010, Break "..., 49Address :
   ↳ 0x7f3d43c80010, Break : 0x5566badcd000) = 49
10 ...
11 munmap(0x7f3d43c80000, 67112960)          = 0
12 ...
13 exit_group(0)                             = ?
14 +++ exited with 0 +++
```

While the smaller allocations use calls to **brk**, the larger ones use the **mmap** and **munmap** system calls.

# ANONYMOUS MEMORY MAPPING

The `mmap` function creates a new mapping in the virtual address space of the calling process, while `munmap` undoes this mapping:

```
1  #include <sys/mman.h>
2
3  void * mmap (void *start,
4              size_t length,
5              int prot,
6              int flags,
7              int fd,
8              off_t offset);
9
10 int munmap (void *start, size_t length);
```

# ANONYMOUS MEMORY MAPPING

The `mmap` function creates a new mapping in the virtual address space of the calling process, while `munmap` undoes this mapping:

```
1  #include <sys/mman.h>
2
3  void * mmap (void *start,
4              size_t length,
5              int prot,
6              int flags,
7              int fd,
8              off_t offset);
9
10 int munmap (void *start, size_t length);
```

A call to `mmap` to establish an anonymous memory mapping requires `fd` be set to -1 and `offset` to 0.

- These arguments are used for memory-mapped files (more on this later)

# ANONYMOUS MEMORY MAPPING

The `mmap` function creates a new mapping in the virtual address space of the calling process, while `munmap` undoes this mapping:

```
1  #include <sys/mman.h>
2
3  void * mmap (void *start,
4              size_t length,
5              int prot,
6              int flags,
7              int fd,
8              off_t offset);
9
10 int munmap (void *start, size_t length);
```

`start` determines the start address of the mapping

- setting this value to `NULL` lets the kernel choose the address at which to create the mapping



# ANONYMOUS MEMORY MAPPING

The `mmap` function creates a new mapping in the virtual address space of the calling process, while `munmap` undoes this mapping:

```
1  #include <sys/mman.h>
2
3  void * mmap (void *start,
4              size_t length,
5              int prot,
6              int flags,
7              int fd,
8              off_t offset);
9
10 int munmap (void *start, size_t length);
```

`length` determines the size of the mapping

- value must be greater than 0

# ANONYMOUS MEMORY MAPPING

The `mmap` function creates a new mapping in the virtual address space of the calling process, while `munmap` undoes this mapping:

```
1  #include <sys/mman.h>
2
3  void * mmap (void *start,
4              size_t length,
5              int prot,
6              int flags,
7              int fd,
8              off_t offset);
9
10 int munmap (void *start, size_t length);
```

`prot` describes the desired memory protection of the mapping:

- `PROT_NONE` stipulates that pages may not accessed
- `PROT_READ`, `PROT_WRITE`, and `PROT_EXEC`, bitwise OR'd depending on requirements

# ANONYMOUS MEMORY MAPPING

The `mmap` function creates a new mapping in the virtual address space of the calling process, while `munmap` undoes this mapping:

```
1  #include <sys/mman.h>
2
3  void * mmap (void *start,
4              size_t length,
5              int prot,
6              int flags,
7              int fd,
8              off_t offset);
9
10 int munmap (void *start, size_t length);
```

`flags` describes the behaviour of the map

- for use as an anonymous map (not backed by a file), `flags` is set to `MAP_ANONYMOUS | MAP_PRIVATE`

# SHARED MEMORY

---

# SHARED MEMORY

Systems with virtual memory allow physical memory to be mapped into multiple process address spaces

- e.g. one copy of a program or library binary in memory mapped into multiple virtual address spaces

This ability can be exploited for *inter-process communication* (IPC), by establishing a region of memory that can be accessed by two or more processes

- processes read from and write to shared memory as any other allocated memory
- writes are reflected in the address space of each process attaching to the shared memory
- access by several processes may give rise to contention and concurrency problems such as race hazards

# CREATING SHARED MEMORY SEGMENTS

A shared memory segment may be allocated through the System V interface:

```
1 #include <sys/ipc.h>
2 #include <sys/shm.h>
3
4 int shmget(key_t key, size_t size, int shmflg);
```

The function **shmget** returns the identifier of the shared memory segment associated with the value of the argument **key**.

- if a new segment is being created, **size** is the minimum length of the segment rounded up to a multiple of page size; segment is initialised with zeroes
- if a client is attaching to an existing segment, **size** can be 0

## CREATING SHARED MEMORY SEGMENTS

A shared memory segment may be allocated through the System V interface:

```
1 #include <sys/ipc.h>
2 #include <sys/shm.h>
3
4 int shmget(key_t key, size_t size, int shmflg);
```

A new segment is always created if the key has value `IPC_PRIVATE`, or no shared memory segment corresponding to `key` exists and `IPC_CREAT` is specified in `shmflg`.

## CONTROLLING SHARED MEMORY SEGMENTS

The `shmctl` method provides an interface for controlling aspects of the shared memory segment, including its deletion:

```
1 #include <sys/ipc.h>
2 #include <sys/shm.h>
3
4 int shmctl(int shmid, int cmd, struct shm_id *buf);
```

The `cmd` entry may take one of these values (among others):

`IPC_STAT` Copy information from kernel datastructure related to segment `shmid` into `buf`



## CONTROLLING SHARED MEMORY SEGMENTS

The `shmctl` method provides an interface for controlling aspects of the shared memory segment, including its deletion:

```
1 #include <sys/ipc.h>
2 #include <sys/shm.h>
3
4 int shmctl(int shmid, int cmd, struct shm_id *buf);
```

The `cmd` entry may take one of these values (among others):

`IPC_STAT` Copy information from kernel datastructure related to segment `shmid` into `buf`

`IPC_SET` Write the values of some of the `shm_id` data structure to the kernel copy associated with segment `shmid`

## CONTROLLING SHARED MEMORY SEGMENTS

The `shmctl` method provides an interface for controlling aspects of the shared memory segment, including its deletion:

```
1 #include <sys/ipc.h>
2 #include <sys/shm.h>
3
4 int shmctl(int shmid, int cmd, struct shm_id *buf);
```

The `cmd` entry may take one of these values (among others):

**IPC\_STAT** Copy information from kernel datastructure related to segment `shmid` into `buf`

**IPC\_SET** Write the values of some of the `shm_id` data structure to the kernel copy associated with segment `shmid`

**IPC\_RMID** Remove shared memory segment; the segment is deleted only if its reference count is zero

## CONTROLLING SHARED MEMORY SEGMENTS

The `shmctl` method provides an interface for controlling aspects of the shared memory segment, including its deletion:

```
1 #include <sys/ipc.h>
2 #include <sys/shm.h>
3
4 int shmctl(int shmid, int cmd, struct shmid_ds *buf);
```

The `cmd` entry may take one of these values (among others):

**SHM\_LOCK** Lock pages of shared memory segment in memory, preventing them from being paged out

**SHM\_UNLOCK** Unlock pages of shared memory segment, allowing them to be paged out when needed

## CONTROLLING SHARED MEMORY SEGMENTS

The `shmctl` method provides an interface for controlling aspects of the shared memory segment, including its deletion:

```
1 #include <sys/ipc.h>
2 #include <sys/shm.h>
3
4 int shmctl(int shmid, int cmd, struct shmid_ds *buf);
```

The `cmd` entry may take one of these values (among others):

**SHM\_LOCK** Lock pages of shared memory segment in memory, preventing them from being paged out

**SHM\_UNLOCK** Unlock pages of shared memory segment, allowing them to be paged out when needed

### Note

These commands may only be executed by superuser.

## CONTROLLING SHARED MEMORY SEGMENTS

The `shmctl` method provides an interface for controlling aspects of the shared memory segment, including its deletion:

```
1 #include <sys/ipc.h>
2 #include <sys/shm.h>
3
4 int shmctl(int shmid, int cmd, struct shm_id *buf);
```

A successful call to `shmctl` with `SHM_STAT` returns the identifier of the shared memory segment whose index was given in `shmid`. Other operations return 0 on success.

On error, -1 is returned, and `errno` is set appropriately.

## ATTACHING SHARED MEMORY SEGMENTS

Once a shared memory segment has been created, a process may attach it to its address space via the `shmat` interface:

```
1 #include <sys/types.h>
2 #include <sys/shm.h>
3
4 void *shmat(int shmid, const void *shmaddr, int shmflg);
```

The attaching address is specified by `shmaddr`:

## ATTACHING SHARED MEMORY SEGMENTS

Once a shared memory segment has been created, a process may attach it to its address space via the `shmat` interface:

```
1 #include <sys/types.h>
2 #include <sys/shm.h>
3
4 void *shmat(int shmid, const void *shmaddr, int shmflg);
```

The attaching address is specified by `shmaddr`:

- if `shmaddr` is `NULL`, the system chooses a suitable (unused) page-aligned address to attach the segment (recommended)

## ATTACHING SHARED MEMORY SEGMENTS

Once a shared memory segment has been created, a process may attach it to its address space via the `shmat` interface:

```
1 #include <sys/types.h>
2 #include <sys/shm.h>
3
4 void *shmat(int shmid, const void *shmaddr, int shmflg);
```

The attaching address is specified by `shmaddr`:

- if `shmaddr` is `NULL`, the system chooses a suitable (unused) page-aligned address to attach the segment (recommended)
- if `SHM_RND` is specified in `shmflg`, the address is rounded down to the nearest multiple of `SHMLBA`, the low boundary address multiple (always a power of 2, currently equal to page-size)



## ATTACHING SHARED MEMORY SEGMENTS

Once a shared memory segment has been created, a process may attach it to its address space via the `shmat` interface:

```
1 #include <sys/types.h>
2 #include <sys/shm.h>
3
4 void *shmat(int shmid, const void *shmaddr, int shmflg);
```

The attaching address is specified by `shmaddr`:

- if `shmaddr` is `NULL`, the system chooses a suitable (unused) page-aligned address to attach the segment (recommended)
- if `SHM_RND` is specified in `shmflg`, the address is rounded down to the nearest multiple of `SHMLBA`, the low boundary address multiple (always a power of 2, currently equal to page-size)
- otherwise `shmaddr` must be a page-aligned address at which the attach occurs

## ATTACHING SHARED MEMORY SEGMENTS

Once a shared memory segment has been created, a process may attach it to its address space via the `shmat` interface:

```
1 #include <sys/types.h>
2 #include <sys/shm.h>
3
4 void *shmat(int shmid, const void *shmaddr, int shmflg);
```

On success, `shmat` returns the address of the attached shared memory segment

- on error, (`void*`) `-1` is returned, and `errno` is set to indicate the cause of the error

## DETACHING SHARED MEMORY SEGMENTS

When done with a shared memory segment, it can be detached from the address space using the `shmdt` interface:

```
1 #include <sys/types.h>
2 #include <sys/shm.h>
3
4 int shmdt(const void *shmaddr);
```

## DETACHING SHARED MEMORY SEGMENTS

When done with a shared memory segment, it can be detached from the address space using the `shmdt` interface:

```
1 #include <sys/types.h>
2 #include <sys/shm.h>
3
4 int shmdt(const void *shmaddr);
```

The `shmaddr` argument is the value that was previously returned by a call to `shmat`

- a successful call to `shmdt` updates members of the associated `shmid_ds` structure
- decrements `shm_nattach` by 1; if the segment is marked for deletion and this value is 0, then the segment is deleted

## DETACHING SHARED MEMORY SEGMENTS

When done with a shared memory segment, it can be detached from the address space using the `shmdt` interface:

```
1 #include <sys/types.h>
2 #include <sys/shm.h>
3
4 int shmdt(const void *shmaddr);
```

A successful call to `shmdt` returns 0

- on error -1 is returned, and `errno` is set to indicate the cause of the error

To exemplify the use of shared memory segments, we're going to write a small up-time server:

- creates a shared memory segment with id `0x666`, 256 bytes in size
- every second increments a timer monotonically by 1
- stores a textual representation of the timer in the shared memory buffer

# SHARED MEMORY UP-TIME SERVER

Code listing for the shared memory up-time server:

```
1  #include <stdio.h>
2  #include <stdlib.h>
3  #include <unistd.h>
4  #include <sys/ipc.h>
5  #include <sys/shm.h>
6
7  int main(int argc, char *argv[]) {
8      int timer = 0; // monotonically increasing timer
9      int shmid; // shared memory segment id
10     char *shm; // shared memory attached address
11
12     if ((shmid = shmget(0x666, 256, IPC_CREAT | 0666)) == -1) { // create shared memory segment
13         perror("shmget() failed:");
14         exit(EXIT_FAILURE);
15     }
16
17     if ((shm = shmat(shmid, NULL, 0)) == (char *)-1) { // attach to segment
18         perror("shmat() failed:");
19         exit(EXIT_FAILURE);
20     }
21
22     for(;; ++timer) { // loop forever, increment timer by 1 every iteration
23         // convert timer to string; buffer is shared memory
24         sprintf(shm, "%02d:%02d:%02d", (timer / 3600) % 24, (timer / 60) % 60, timer % 60);
25         sleep(1); // sleep for 1 second
26     }
27 }
```

We also provide an up-time client that connects with the shared memory segment to read the current up-time:

- attaches to shared memory segment with id `0x666`
- prints the up-time value set by the up-time server



# SHARED MEMORY UP-TIME CLIENT

Code listing for client:

```
1  #include <stdio.h>
2  #include <stdlib.h>
3  #include <unistd.h>
4  #include <sys/ipc.h>
5  #include <sys/shm.h>
6
7  int main(int argc, char *argv[]) {
8      int shmid; // shared memory segment id
9      char *shm; // shared memory attached address
10
11     if ((shmid = shmget(0x666, 0, 0666)) == -1) { // acquires id of segment with key 0x666
12         perror("shmget() failed:");
13         exit(EXIT_FAILURE);
14     }
15
16     if ((shm = shmat(shmid, NULL, 0)) == (char *)-1) { // attaches to segment
17         perror("shmat() failed:");
18         exit(EXIT_FAILURE);
19     }
20
21     for(;;) { // loop forever
22         printf("Timer is : [%s]\n", shm); // print value in shared memory buffer
23         sleep(1); // sleep for 1 second
24     }
25 }
```

## EXAMPLE

Run the up-time server in the background `./shm_srv &` and start the client `shm_cli` immediately after:

```
1 $ ./shm_srv &
2 $ ./shm_cli
3 Timer is : [00:00:03]
4 Timer is : [00:00:04]
5 Timer is : [00:00:05]
6 Timer is : [00:00:06]
7 ^C
8 $ ./shm_cli
9 Timer is : [00:00:15]
10 Timer is : [00:00:16]
11 Timer is : [00:00:17]
12 ^C
13 $
```

## EXAMPLE

Run the up-time server in the background `./shm_srv &` and start the client `shm_cli` immediately after:

```
1  $ ./shm_srv &
2  $ ./shm_cli
3  Timer is : [00:00:03]
4  Timer is : [00:00:04]
5  Timer is : [00:00:05]
6  Timer is : [00:00:06]
7  ^C
8  $ ./shm_cli
9  Timer is : [00:00:15]
10 Timer is : [00:00:16]
11 Timer is : [00:00:17]
12 ^C
13 $
```

Stopping and restarting the client shows that the displayed time is independent of it.

QUESTIONS?