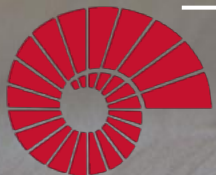


# COMP201

## Computer Systems & Programming

### Lecture #34 – Managing The Heap



**KOÇ  
UNIVERSITY**

Aykut Erdem // Koç University // Fall 2020

# Recap

- Static Linking
- Symbol Resolution
- Relocation
- Static Libraries

# Plan for Today

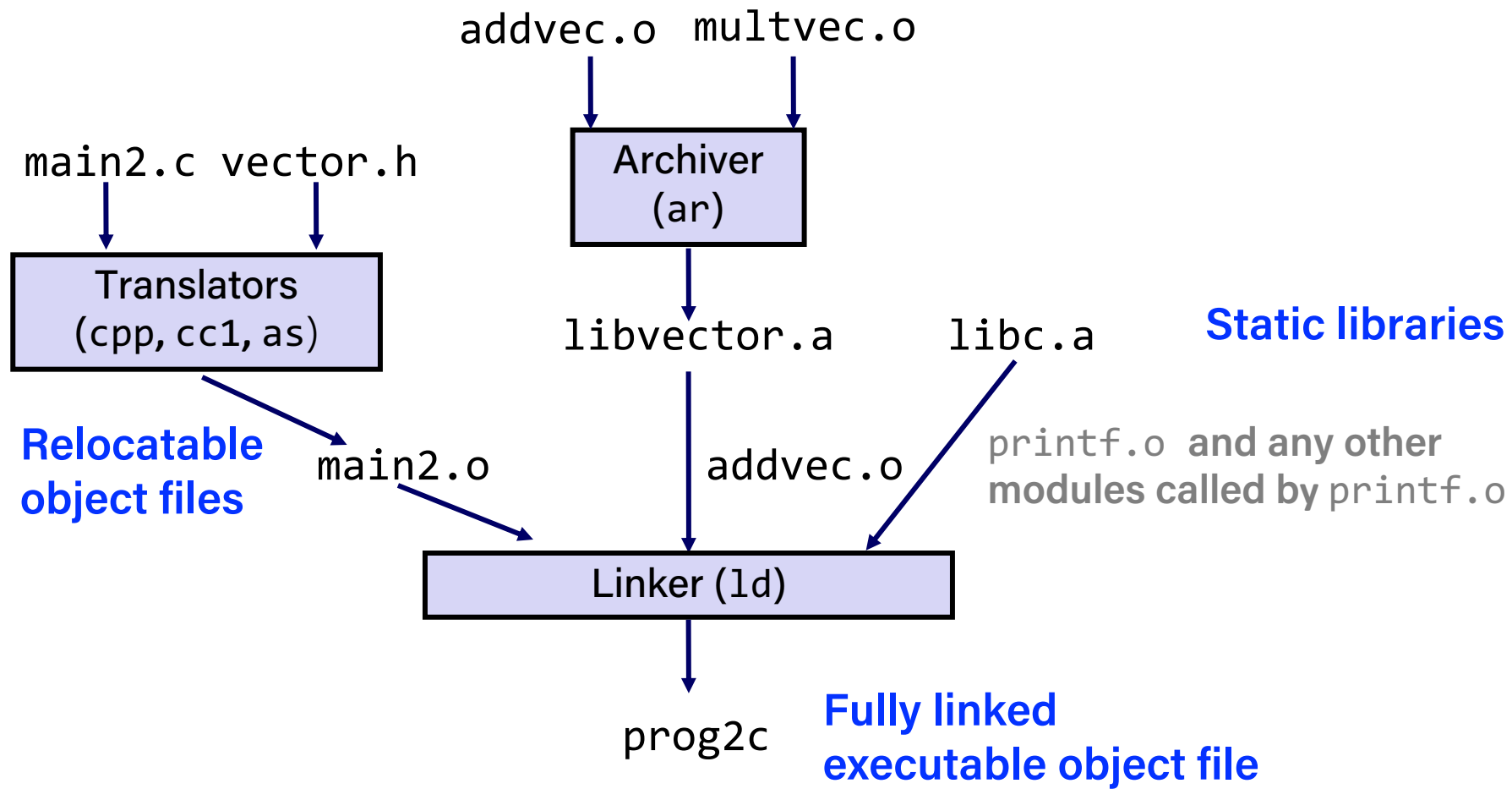
- Shared Libraries
- Case study: Library interpositioning
- The heap so far
- What is a heap allocator?
- Heap allocator requirements and goals
- Method 0: Bump Allocator
- Method 1: Implicit Free List Allocator
- Method 2: Explicit Free List Allocator

**Disclaimer:** Slides for this lecture were borrowed from  
—Nick Troccoli's Stanford CS107 class

# Lecture Plan

- Shared Libraries
- Case study: Library interpositioning
- The heap so far
- What is a heap allocator?
- Heap allocator requirements and goals
- Method 0: Bump Allocator
- Method 1: Implicit Free List Allocator
- Method 2: Explicit Free List Allocator

# Recap: Linking with Static Libraries



*"c" for "compile-time"*

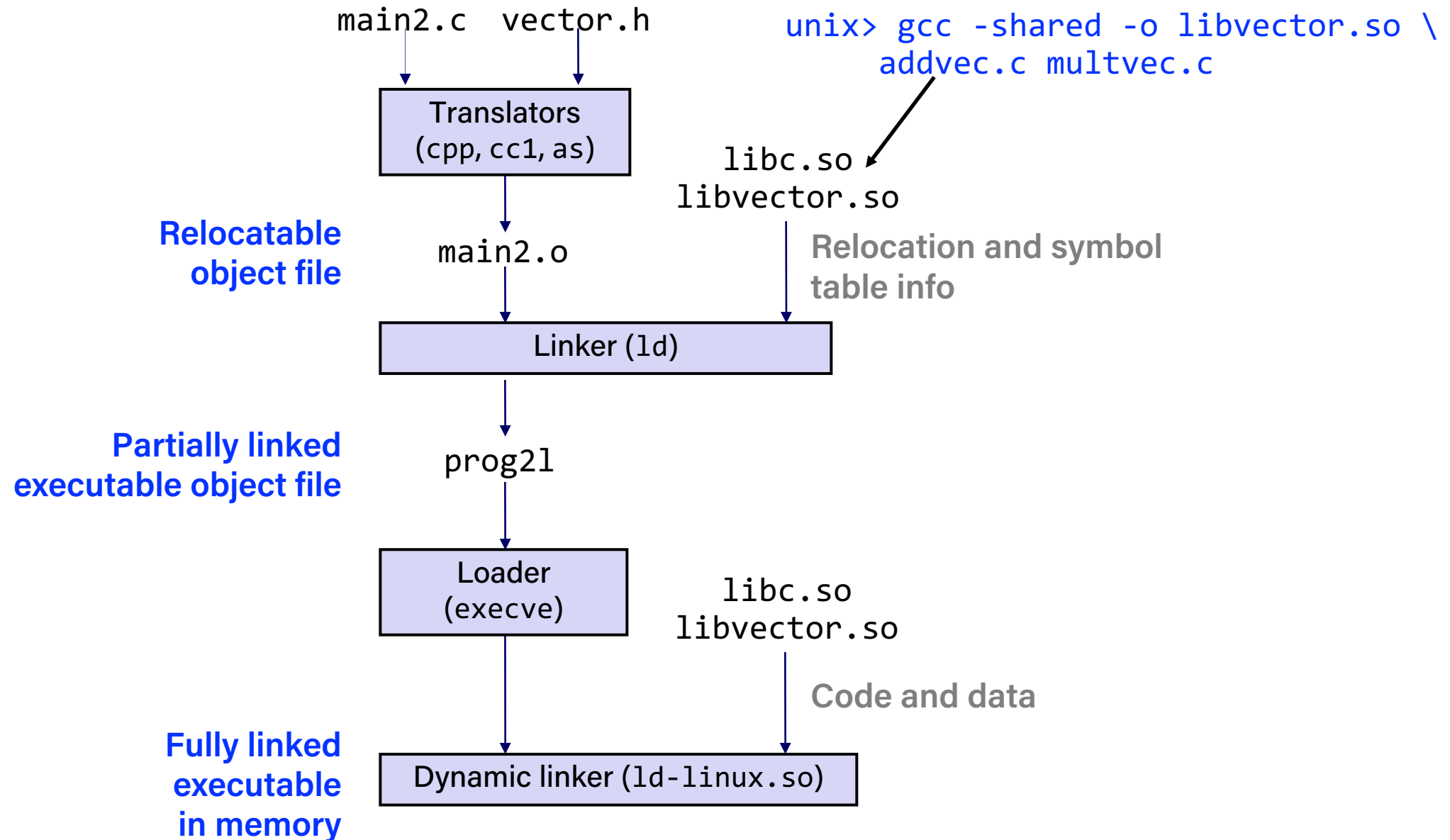
# Modern Solution: Shared Libraries

- **Static libraries have the following disadvantages:**
  - Duplication in the stored executables (every function needs libc)
  - Duplication in the running executables
  - Minor bug fixes of system libraries require each application to explicitly relink
- **Modern solution: Shared Libraries**
  - Object files that contain code and data that are loaded and linked into an application *dynamically*, at either *load-time* or *run-time*
  - Also called: dynamic link libraries, DLLs, .so files

# Shared Libraries (cont.)

- **Dynamic linking can occur when executable is first loaded and run (load-time linking).**
  - Common case for Linux, handled automatically by the dynamic linker (`ld-linux.so`).
  - Standard C library (`libc.so`) usually dynamically linked.
- **Dynamic linking can also occur after program has begun (run-time linking).**
  - In Linux, this is done by calls to the `dlopen()` interface.
    - Distributing software.
    - High-performance web servers.
    - Runtime library interpositioning.
- **Shared library routines can be shared by multiple processes.**
  - More on this when you learn about virtual memory

# Dynamic Linking at Load-time





# Dynamic Linking at Run-time

```
#include <stdio.h>
#include <stdlib.h>
#include <dlfcn.h>

int x[2] = {1, 2};
int y[2] = {3, 4};
int z[2];

int main()
{
    void *handle;
    void (*addvec)(int *, int *, int *, int);
    char *error;

    /* Dynamically load the shared library that contains addvec() */
    handle = dlopen("./libvector.so", RTLD_LAZY);
    if (!handle) {
        fprintf(stderr, "%s\n", dlerror());
        exit(1);
    }
}
```

*dLL.c*

# Dynamic Linking at Run-time

...

```
/* Get a pointer to the addvec() function we just loaded */
addvec = dlsym(handle, "addvec");
if ((error = dlerror()) != NULL) {
    fprintf(stderr, "%s\n", error);
    exit(1);
}
```

```
/* Now we can call addvec() just like any other function */
addvec(x, y, z, 2);
printf("z = [%d %d]\n", z[0], z[1]);
```

```
/* Unload the shared library */
if (dlclose(handle) < 0) {
    fprintf(stderr, "%s\n", dlerror());
    exit(1);
}
return 0;
```

```
}
```

*dll.c*

# Linking Summary

- Linking is a technique that allows programs to be constructed from multiple object files.
- Linking can happen at different times in a program's lifetime:
  - Compile time (when a program is compiled)
  - Load time (when a program is loaded into memory)
  - Run time (while a program is executing)
- Understanding linking can help you avoid nasty errors and make you a better programmer.

# Lecture Plan

- Shared Libraries
- Case study: Library interpositioning
- The heap so far
- What is a heap allocator?
- Heap allocator requirements and goals
- Method 0: Bump Allocator
- Method 1: Implicit Free List Allocator
- Method 2: Explicit Free List Allocator

# Case Study: Library Interpositioning

- **Library interpositioning:** powerful linking technique that allows programmers to intercept calls to arbitrary functions
- Interpositioning can occur at:
  - Compile time: When the source code is compiled
  - Link time: When the relocatable object files are statically linked to form an executable object file
  - Load/run time: When an executable object file is loaded into memory, dynamically linked, and then executed.

# Some Interpositioning Applications

- **Security**

- Confinement (sandboxing)
- Behind the scenes encryption

- **Debugging**

- In 2014, two Facebook engineers debugged a treacherous 1-year old bug in their iPhone app using interpositioning
- Code in the SPDY networking stack was writing to the wrong location
- Solved by intercepting calls to Posix write functions (`write`, `writew`, `pwrite`)

Source: Facebook engineering blog post at

<https://code.facebook.com/posts/313033472212144/debugging-file-corruption-on-ios/>

# Some Interpositioning Applications

- **Monitoring and Profiling**

- Count number of calls to functions
- Characterize call sites and arguments to functions
- Malloc tracing
  - Detecting memory leaks
  - **Generating address traces**

# Example program

```
#include <stdio.h>
#include <malloc.h>

int main()
{
    int *p = malloc(32);
    free(p);
    return(0);
}                                     int.c
```

- **Goal:** trace the addresses and sizes of the allocated and freed blocks, without breaking the program, and without modifying the source code.
- **Three solutions:** interpose on the `lib malloc` and `free` functions at compile time, link time, and load/run time.



# Compile-time Interpositioning

```
#ifdef COMPILETIME
#include <stdio.h>
#include <malloc.h>

/* malloc wrapper function */
void *mymalloc(size_t size)
{
    void *ptr = malloc(size);
    printf("malloc(%d)=%p\n",
           (int)size, ptr);
    return ptr;
}

/* free wrapper function */
void myfree(void *ptr)
{
    free(ptr);
    printf("free(%p)\n", ptr);
}
#endif
```

*mymalloc.c*

# Compile-time Interpositioning

```
#define malloc(size) mymalloc(size)
#define free(ptr) myfree(ptr)

void *mymalloc(size_t size);
void myfree(void *ptr);
```

*malloc.h*

```
linux> make intc
gcc -Wall -DCOMPILETIME -c mymalloc.c
gcc -Wall -I. -o intc int.c mymalloc.o
linux> make runc
./intc
malloc(32)=0x1edc010
free(0x1edc010)
linux>
```

# Link-time Interpositioning

```
#ifdef LINKTIME
#include <stdio.h>
```

```
void *__real_malloc(size_t size);
void __real_free(void *ptr);
```

```
/* malloc wrapper function */
void *__wrap_malloc(size_t size)
{
    void *ptr = __real_malloc(size); /* Call libc malloc */
    printf("malloc(%d) = %p\n", (int)size, ptr);
    return ptr;
}
```

```
/* free wrapper function */
void __wrap_free(void *ptr)
{
    __real_free(ptr); /* Call libc free */
    printf("free(%p)\n", ptr);
}
#endif
```

*mymalloc.c*

# Link-time Interpositioning

```
linux> make intl
gcc -Wall -DLINKTIME -c mymalloc.c
gcc -Wall -c int.c
gcc -Wall -Wl,--wrap,malloc -Wl,--wrap,free -o intl int.o mymalloc.o
linux> make runl
./intl
malloc(32) = 0x1aa0010
free(0x1aa0010)
linux>
```

- The “-Wl” flag passes argument to linker, replacing each comma with a space.
- The “--wrap,malloc” arg instructs linker to resolve references in a special way
  - Refs to malloc should be resolved as \_\_wrap\_malloc
  - Refs to \_\_real\_malloc should be resolved as malloc

# Load/Run-time Interpositioning

```
#ifdef RUNTIME
#define _GNU_SOURCE
#include <stdio.h>
#include <stdlib.h>
#include <dlfcn.h>

/* malloc wrapper function */
void *malloc(size_t size)
{
    void *(*mallocp)(size_t size);
    char *error;

    mallocp = dlsym(RTLD_NEXT, "malloc"); /* Get addr of libc malloc */
    if ((error = dlerror()) != NULL) {
        fputs(error, stderr);
        exit(1);
    }
    char *ptr = mallocp(size); /* Call libc malloc */
    printf("malloc(%d) = %p\n", (int)size, ptr);
    return ptr;
}
```

*mymalloc.c*

# Load/Run-time Interpositioning

```
/* free wrapper function */
void free(void *ptr)
{
    void (*freep)(void *) = NULL;
    char *error;

    if (!ptr)
        return;

    freep = dlsym(RTLD_NEXT, "free"); /* Get address of libc free */
    if ((error = dlerror()) != NULL) {
        fputs(error, stderr);
        exit(1);
    }
    freep(ptr); /* Call libc free */
    printf("free(%p)\n", ptr);
}
#endif
```

*mymalloc.c*

# Load/Run-time Interpositioning

```
linux> make intr
gcc -Wall -DRUNTIME -shared -fpic -o mymalloc.so mymalloc.c -ldl
gcc -Wall -o intr int.c
linux> make runr
(LD_PRELOAD="./mymalloc.so" ./intr)
malloc(32) = 0xe60010
free(0xe60010)
linux>
```

- The `LD_PRELOAD` environment variable tells the dynamic linker to resolve unresolved refs (e.g., to `malloc`) by looking in `mymalloc.so` first.

# Interpositioning Recap

- **Compile Time**

- Apparent calls to `malloc/free` get macro-expanded into calls to `mymalloc/myfree`

- **Link Time**

- Use linker trick to have special name resolutions
  - `malloc` → `__wrap_malloc`
  - `__real_malloc` → `malloc`

- **Load/Run Time**

- Implement custom version of `malloc/free` that use dynamic linking to load library `malloc/free` under different names



COMP201 Topic 8: How do the  
core malloc/realloc/free  
memory-allocation operations  
work?

# How do malloc/realloc/free work?

Pulling together all our COMP201 topics this semester:

- Testing
- Efficiency
- Bit-level manipulation
- Memory management
- Pointers
- Generics
- Assembly
- And more...

# Learning Goals

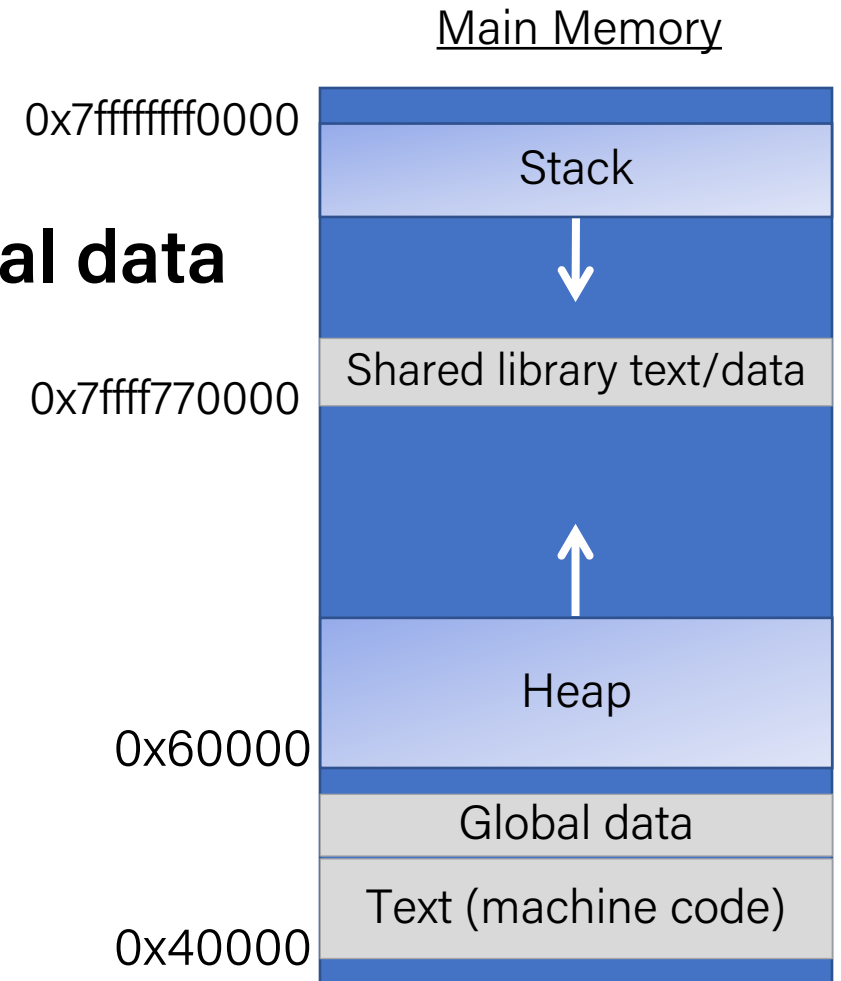
- Learn the restrictions, goals and assumptions of a heap allocator
- Understand the conflicting goals of utilization and throughput
- Learn about different ways to implement a heap allocator

# Lecture Plan

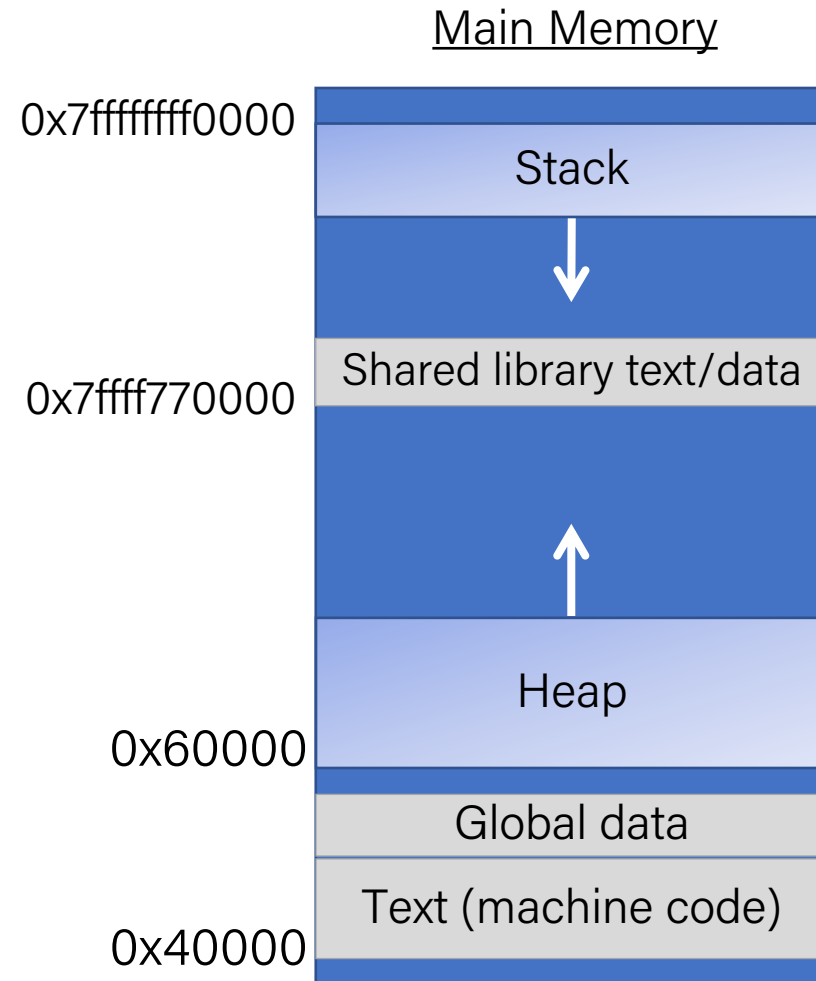
- Shared Libraries
- Case study: Library interpositioning
- The heap so far
- What is a heap allocator?
- Heap allocator requirements and goals
- Method 0: Bump Allocator
- Method 1: Implicit Free List Allocator
- Method 2: Explicit Free List Allocator

# Running a program

- **Creates new process**
- **Sets up address space/segments**
- **Read executable file, load instructions, global data**  
Mapped from file into gray segments
- **Libraries loaded on demand**
- **Set up stack**  
Reserve stack segment, init %rsp, call main
- **malloc written in C, will init self on use**  
Asks OS for large memory region,  
parcels out to service requests



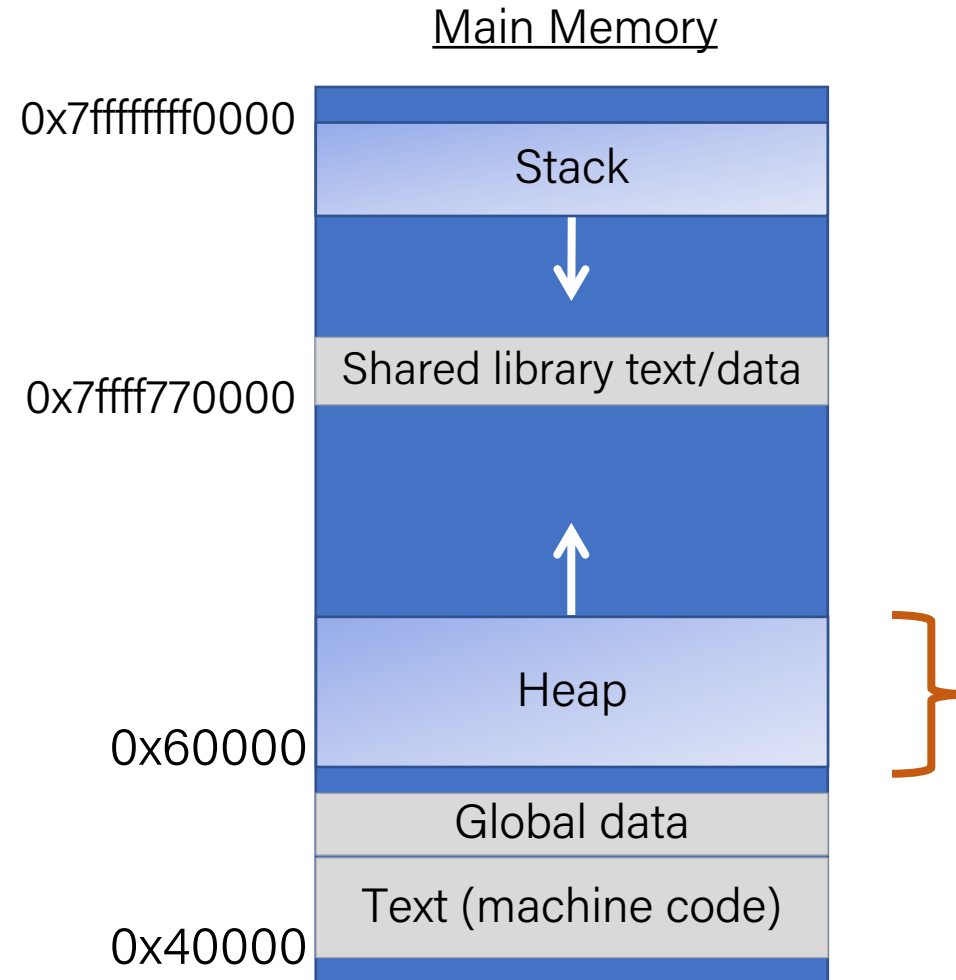
# The Stack Revisited



**Stack memory "goes away"** after function call ends.

**Automatically managed** at compile-time by gcc

# Today: The Heap



**Heap memory persists** until caller indicates it no longer needs it.

**Managed** by C standard library functions (malloc, realloc, free)

This lecture:  
How does heap management work?

# Lecture Plan

- The heap so far
- What is a heap allocator?
- Heap allocator requirements and goals
- Method 0: Bump Allocator
- Method 1: Implicit Free List Allocator
- Method 2: Explicit Free List Allocator



# Your role so far: Client

```
void *malloc(size_t size);
```

Returns a pointer to a block of heap memory of at least size bytes, or NULL if an error occurred.

```
void free(void *ptr);
```

Frees the heap-allocated block starting at the specified address.

```
void *realloc(void *ptr, size_t size);
```

Changes the size of the heap-allocated block starting at the specified address to be the new specified size. Returns the address of the new, larger allocated memory region.

# Your role now: Heap Hotel Concierge



(aka **Heap Allocator**)

# What is a heap allocator?

- A heap allocator is a set of functions that fulfills requests for heap memory.
- On initialization, a heap allocator is provided the starting address and size of a large contiguous block of memory (the heap).

0x10      0x11      0x12      0x13      0x14      0x15      0x16      0x17      0x18      0x19

AVAILABLE

# What is a heap allocator?

- A heap allocator is a set of functions that fulfills requests for heap memory.
- On initialization, a heap allocator is provided the starting address and size of a large contiguous block of memory (the heap).
- A heap allocator must manage this memory as clients request or no longer need it.

**Request 1:** Hi! May I please have 2 bytes of heap memory?

**Allocator:** Sure, I've given you address 0x10.

0x10      0x11      0x12      0x13      0x14      0x15      0x16      0x17      0x18      0x19

AVAILABLE

# What is a heap allocator?

- A heap allocator is a set of functions that fulfills requests for heap memory.
- On initialization, a heap allocator is provided the starting address and size of a large contiguous block of memory (the heap).
- A heap allocator must manage this memory as clients request or no longer need it.

**Request 1:** Hi! May I please have 2 bytes of heap memory?

**Allocator:** Sure, I've given you address 0x10.

0x10      0x11      0x12      0x13      0x14      0x15      0x16      0x17      0x18      0x19

FOR REQUEST 1

AVAILABLE

# What is a heap allocator?

- A heap allocator is a set of functions that fulfills requests for heap memory.
- On initialization, a heap allocator is provided the starting address and size of a large contiguous block of memory (the heap).
- A heap allocator must manage this memory as clients request or no longer need it.

**Request 2:** Howdy!  
May I please have 3  
bytes of heap memory?

**Allocator:** Sure, I've  
given you address 0x12.

0x10      0x11      0x12      0x13      0x14      0x15      0x16      0x17      0x18      0x19

FOR REQUEST 1

AVAILABLE

# What is a heap allocator?

- A heap allocator is a set of functions that fulfills requests for heap memory.
- On initialization, a heap allocator is provided the starting address and size of a large contiguous block of memory (the heap).
- A heap allocator must manage this memory as clients request or no longer need it.

**Request 2:** Howdy!  
May I please have 3  
bytes of heap memory?

**Allocator:** Sure, I've  
given you address 0x12.

0x10      0x11      0x12      0x13      0x14      0x15      0x16      0x17      0x18      0x19

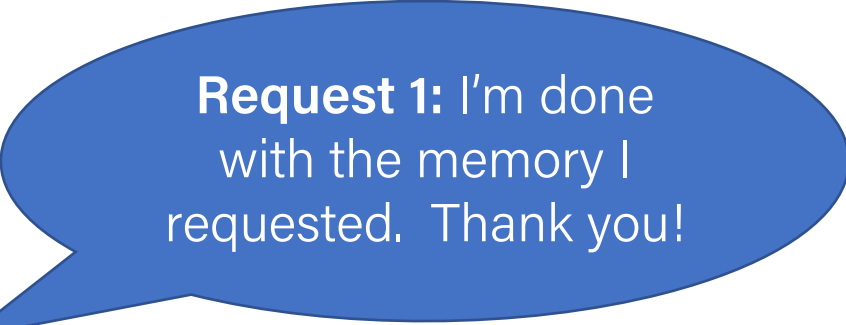
FOR REQUEST 1

FOR REQUEST 2

AVAILABLE

# What is a heap allocator?

- A heap allocator is a set of functions that fulfills requests for heap memory.
- On initialization, a heap allocator is provided the starting address and size of a large contiguous block of memory (the heap).
- A heap allocator must manage this memory as clients request or no longer need pieces of it.



**Request 1:** I'm done with the memory I requested. Thank you!



**Allocator:** Thanks. Have a good day!

0x10      0x11      0x12      0x13      0x14      0x15      0x16      0x17      0x18      0x19

FOR REQUEST 1

FOR REQUEST 2

AVAILABLE



# What is a heap allocator?

- A heap allocator is a set of functions that fulfills requests for heap memory.
- On initialization, a heap allocator is provided the starting address and size of a large contiguous block of memory (the heap).
- A heap allocator must manage this memory as clients request or no longer need pieces of it.

**Request 1:** I'm done with the memory I requested. Thank you!

**Allocator:** Thanks. Have a good day!

0x10      0x11      0x12      0x13      0x14      0x15      0x16      0x17      0x18      0x19

AVAILABLE

FOR REQUEST 2

AVAILABLE

# What is a heap allocator?

- A heap allocator is a set of functions that fulfills requests for heap memory.
- On initialization, a heap allocator is provided the starting address and size of a large contiguous block of memory (the heap).
- A heap allocator must manage this memory as clients request or no longer need pieces of it.

**Request 3:** Hello there!  
I'd like to request 2 bytes  
of heap memory, please.

**Allocator:** Sure thing.  
I've given you address  
0x10.

0x10      0x11      0x12      0x13      0x14      0x15      0x16      0x17      0x18      0x19

AVAILABLE

FOR REQUEST 2

AVAILABLE

# What is a heap allocator?

- A heap allocator is a set of functions that fulfills requests for heap memory.
- On initialization, a heap allocator is provided the starting address and size of a large contiguous block of memory (the heap).
- A heap allocator must manage this memory as clients request or no longer need pieces of it.

**Request 3:** Hello there!  
I'd like to request 2 bytes of  
heap memory, please.

**Allocator:** Sure thing.  
I've given you address  
0x10.

0x10      0x11      0x12      0x13      0x14      0x15      0x16      0x17      0x18      0x19

FOR REQUEST 3

FOR REQUEST 2

AVAILABLE

# What is a heap allocator?

- A heap allocator is a set of functions that fulfills requests for heap memory.
- On initialization, a heap allocator is provided the starting address and size of a large contiguous block of memory (the heap).
- A heap allocator must manage this memory as clients request or no longer need pieces of it.

**Request 3:** Hi again!  
I'd like to request the  
region of memory at 0x10  
be reallocated to 4 bytes.

**Allocator:** Sure thing.  
I've given you address  
0x15.

0x10      0x11      0x12      0x13      0x14      0x15      0x16      0x17      0x18      0x19

FOR REQUEST 3

FOR REQUEST 2

AVAILABLE

# What is a heap allocator?

- A heap allocator is a set of functions that fulfills requests for heap memory.
- On initialization, a heap allocator is provided the starting address and size of a large contiguous block of memory (the heap).
- A heap allocator must manage this memory as clients request or no longer need pieces of it.

**Request 3:** Hi again!  
I'd like to request the  
region of memory at 0x10  
be reallocated to 4 bytes.

**Allocator:** Sure thing.  
I've given you address  
0x15.

0x10      0x11      0x12      0x13      0x14      0x15      0x16      0x17      0x18      0x19

AVAILABLE	FOR REQUEST 2	FOR REQUEST 3	AVAILABLE
-----------	---------------	---------------	-----------

# Lecture Plan

- Shared Libraries
- Case study: Library interpositioning
- The heap so far
- What is a heap allocator?
- **Heap allocator requirements and goals**
- Method 0: Bump Allocator
- Method 1: Implicit Free List Allocator
- Method 2: Explicit Free List Allocator

# Heap Allocator Functions

```
void *malloc(size_t size);
```

```
void free(void *ptr);
```

```
void *realloc(void *ptr, size_t size);
```

# Heap Allocator Requirements

A heap allocator must...

1. Handle arbitrary request sequences of allocations and frees
2. Keep track of which memory is allocated and which is available
3. Decide which memory to provide to fulfill an allocation request
4. Immediately respond to requests without delay



# Heap Allocator Requirements

A heap allocator must...

1. **Handle arbitrary request sequences of allocations and frees**
2. Keep track of which memory is allocated and which is available
3. Decide which memory to provide to fulfill an allocation request
4. Immediately respond to requests without delay

A heap allocator cannot assume anything about the order of allocation and free requests, or even that every allocation request is accompanied by a matching free request.

# Heap Allocator Requirements

A heap allocator must...

1. Handle arbitrary request sequences of allocations and frees
2. **Keep track of which memory is allocated and which is available**
3. Decide which memory to provide to fulfill an allocation request
4. Immediately respond to requests without delay

A heap allocator marks memory regions as **allocated** or **available**. It must remember which is which to properly provide memory to clients.

# Heap Allocator Requirements

A heap allocator must...

1. Handle arbitrary request sequences of allocations and frees
2. Keep track of which memory is allocated and which is available
- 3. Decide which memory to provide to fulfill an allocation request**
4. Immediately respond to requests without delay

A heap allocator may have options for which memory to use to fulfill an allocation request. It must decide this based on a variety of factors.

# Heap Allocator Requirements

A heap allocator must...

1. Handle arbitrary request sequences of allocations and frees
2. Keep track of which memory is allocated and which is available
3. Decide which memory to provide to fulfill an allocation request
- 4. Immediately respond to requests without delay**

A heap allocator must respond immediately to allocation requests and should not e.g. prioritize or reorder certain requests to improve performance.

# Heap Allocator Requirements

A heap allocator must...

1. Handle arbitrary request sequences of allocations and frees
2. Keep track of which memory is allocated and which is available
3. Decide which memory to provide to fulfill an allocation request
4. Immediately respond to requests without delay
- 5. Return addresses that are 8-byte-aligned (must be multiples of 8).**

# Heap Allocator Goals

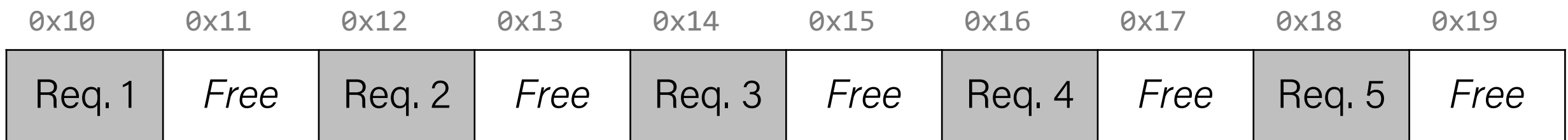
- Goal 1: Maximize **throughput**, or the number of requests completed per unit time. This means minimizing the average time to satisfy a request.
- Goal 2: Maximize memory **utilization**, or how efficiently we make use of the limited heap memory to satisfy requests.

# Utilization

- The primary cause of poor utilization is **fragmentation**. **Fragmentation** occurs when otherwise unused memory is not available to satisfy allocation requests.
- In this example, there is enough aggregate free memory to satisfy the request, but no single free block is large enough to handle the request.
- **In general:** we want the largest address used to be as low as possible.

**Request 6:** Hi! May I please have 4 bytes of heap memory?

**Allocator:** I'm sorry, I don't have a 4 byte block available...



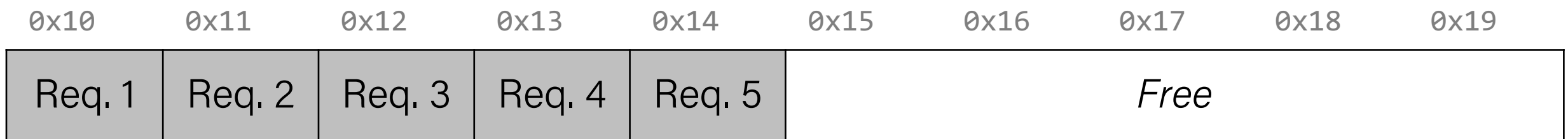
# Utilization

**Question:** what if we shifted these blocks down to make more space?  
Can we do this?

A. YES, great idea!

B. YES, it can be done, but not a good idea for some reason (e.g. not efficient use of time)

C. NO, it can't be done!

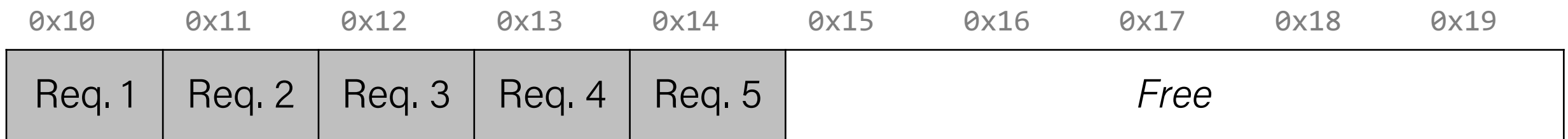




# Utilization

**Question:** what if we shifted these blocks down to make more space?  
Can we do this?

- **No** - we have already guaranteed these addresses to the client. We cannot move allocated memory around, since this will mean the client will now have incorrect pointers to their memory!



# Fragmentation

- **Internal Fragmentation:** an allocated block is larger than what is needed (e.g. due to minimum block size)
- **External Fragmentation:** no single block is large enough to satisfy an allocation request, even though enough aggregate free memory is available

# Heap Allocator Goals

- Goal 1: Maximize **throughput**, or the number of requests completed per unit time. This means minimizing the average time to satisfy a request.
- Goal 2: Maximize memory **utilization**, or how efficiently we make use of the limited heap memory to satisfy requests.

These are seemingly conflicting goals – for instance, it may take longer to better plan out heap memory use for each request.

Heap allocators must find an appropriate balance between these two goals!

# Heap Allocator Goals

- Goal 1: Maximize **throughput**, or the number of requests completed per unit time. This means minimizing the average time to satisfy a request.
- Goal 2: Maximize memory **utilization**, or how efficiently we make use of the limited heap memory to satisfy requests.

Other desirable goals:

**Locality** ("similar" blocks allocated close in space)

**Robust** (handle client errors)

**Ease of implementation/maintenance**

# Question Break

# Lecture Plan

- Shared Libraries
- Case study: Library interpositioning
- The heap so far
- What is a heap allocator?
- Heap allocator requirements and goals
- **Method 0: Bump Allocator**
- Method 1: Implicit Free List Allocator
- Method 2: Explicit Free List Allocator

# Bump Allocator

Let's say we want to entirely prioritize throughput, and do not care about utilization at all. This means we do not care about reusing memory. How could we do this?

# 1. Utilization



**Never** reuses memory

# 2. Throughput



**Ultra** fast, short routines



# Bump Allocator

- A **bump allocator** is a heap allocator design that simply allocates the next available memory address upon an allocate request and does nothing on a free request.
- Throughput: each `malloc` and `free` execute only a handful of instructions:
  - It is easy to find the next location to use
  - Free does nothing!
- Utilization: we use each memory block at most once. No freeing at all, so no memory is ever reused. ☹️

# Bump Allocator

```
void *a = malloc(8);  
void *b = malloc(4);  
void *c = malloc(24);  
free(b);  
void *d = malloc(8);
```

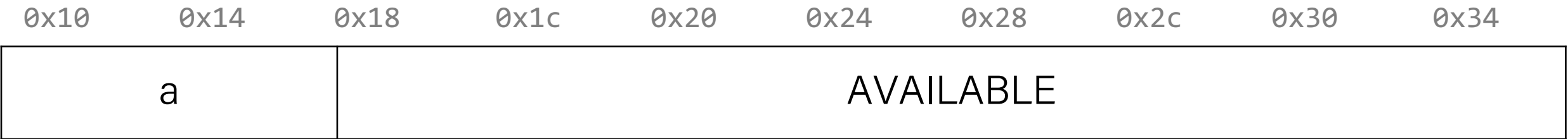
0x10      0x14      0x18      0x1c      0x20      0x24      0x28      0x2c      0x30      0x34

AVAILABLE

# Bump Allocator

```
void *a = malloc(8);  
void *b = malloc(4);  
void *c = malloc(24);  
free(b);  
void *d = malloc(8);
```

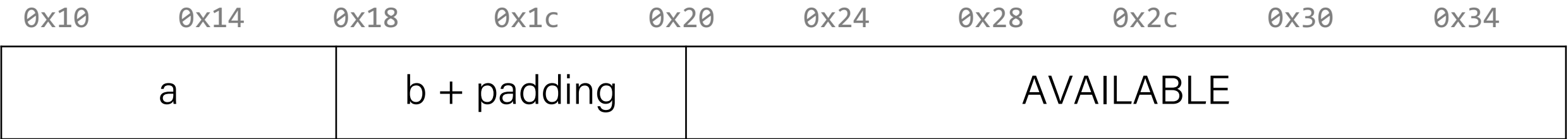
Variable	Value
a	0x10



# Bump Allocator

```
void *a = malloc(8);  
void *b = malloc(4);  
void *c = malloc(24);  
free(b);  
void *d = malloc(8);
```

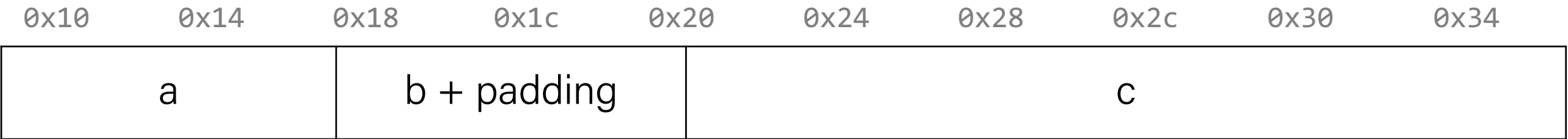
Variable	Value
a	0x10
b	0x18



# Bump Allocator

```
void *a = malloc(8);  
void *b = malloc(4);  
void *c = malloc(24);  
free(b);  
void *d = malloc(8);
```

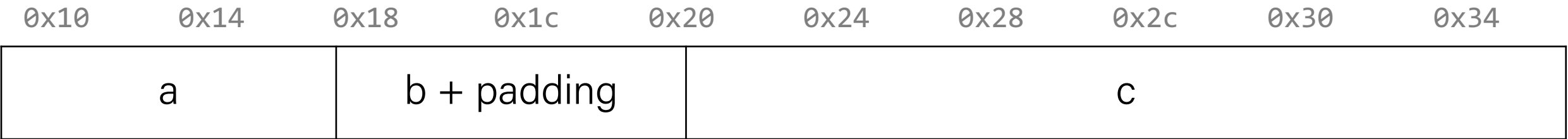
Variable	Value
a	0x10
b	0x18
c	0x20



# Bump Allocator

```
void *a = malloc(8);  
void *b = malloc(4);  
void *c = malloc(24);  
free(b);  
void *d = malloc(8);
```

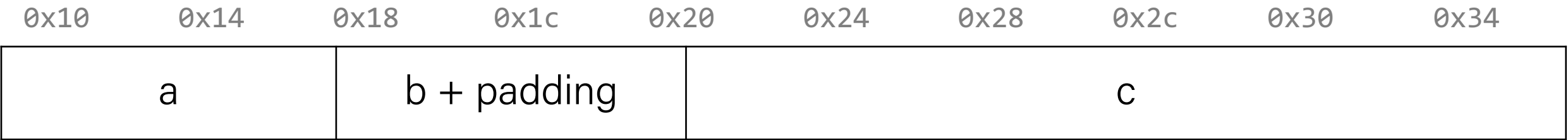
Variable	Value
a	0x10
b	0x18
c	0x20



# Bump Allocator

```
void *a = malloc(8);  
void *b = malloc(4);  
void *c = malloc(24);  
free(b);  
void *d = malloc(8);
```

Variable	Value
a	0x10
b	0x18
c	0x20
d	NULL



# Summary: Bump Allocator

- A bump allocator is an extreme heap allocator – it optimizes only for **throughput**, not **utilization**.
- Better allocators strike a more reasonable balance. How can we do this?

## Questions to consider:

1. How do we keep track of free blocks?
2. How do we choose an appropriate free block in which to place a newly allocated block?
3. After we place a newly allocated block in some free block, what do we do with the remainder of the free block?
4. What do we do with a block that has just been freed?



# Lecture Plan

- Shared Libraries
- Case study: Library interpositioning
- The heap so far
- What is a heap allocator?
- Heap allocator requirements and goals
- Method 0: Bump Allocator
- **Method 1: Implicit Free List Allocator**
- Method 2: Explicit Free List Allocator

# Implicit Free List Allocator

- **Key idea:** in order to reuse blocks, we need a way to track which blocks are allocated and which are free.
- We could store this information in a separate global data structure, but this is inefficient.
- Instead: let's allocate extra space before each block for a **header** storing its payload size and whether it is allocated or free.
- When we allocate a block, we look through the blocks to find a free one, and we update its header to reflect its allocated size and that it is now allocated.
- When we free a block, we update its header to reflect it is now free.
- The header should be 8 bytes (or larger).
- By storing the block size of each block, we *implicitly* have a *list* of free blocks.

# Implicit Free List Allocator

```
void *a = malloc(4);  
void *b = malloc(8);  
void *c = malloc(4);  
free(b);  
void *d = malloc(8);  
free(a);  
void *e = malloc(24);
```

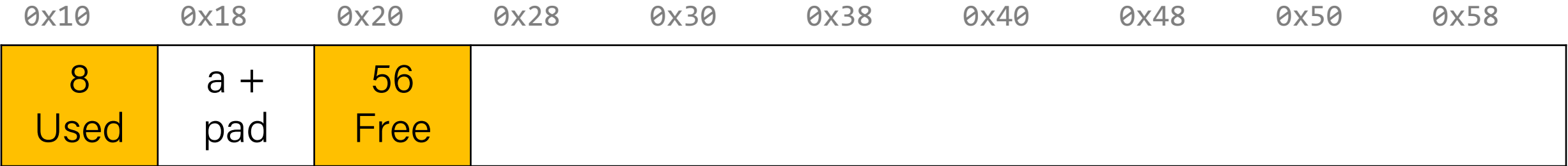
0x10      0x18      0x20      0x28      0x30      0x38      0x40      0x48      0x50      0x58

72  
Free

# Implicit Free List Allocator

```
void *a = malloc(4);
void *b = malloc(8);
void *c = malloc(4);
free(b);
void *d = malloc(8);
free(a);
void *e = malloc(24);
```

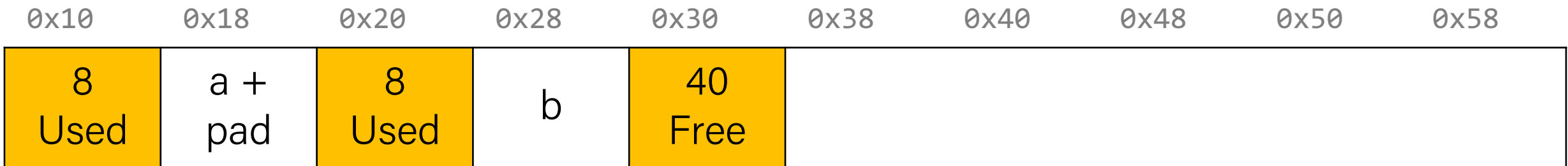
Variable	Value
a	0x18



# Implicit Free List Allocator

```
void *a = malloc(4);  
void *b = malloc(8);  
void *c = malloc(4);  
free(b);  
void *d = malloc(8);  
free(a);  
void *e = malloc(24);
```

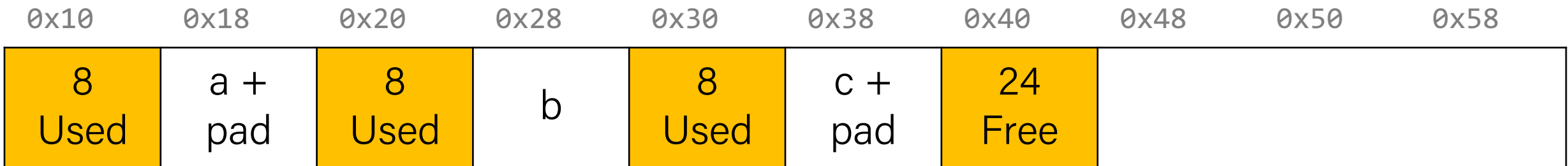
Variable	Value
a	0x18
b	0x28



# Implicit Free List Allocator

```
void *a = malloc(4);  
void *b = malloc(8);  
void *c = malloc(4);  
free(b);  
void *d = malloc(8);  
free(a);  
void *e = malloc(24);
```

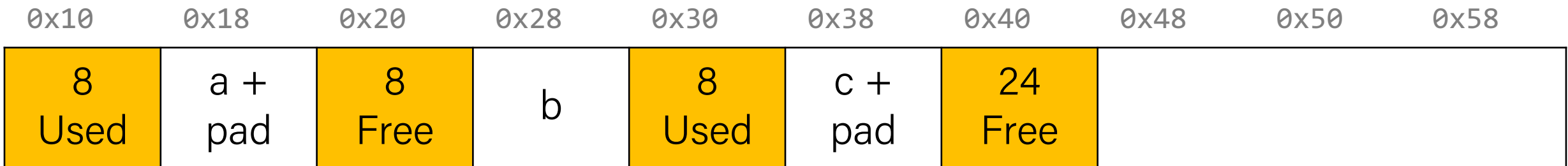
Variable	Value
a	0x18
b	0x28
c	0x38



# Implicit Free List Allocator

```
void *a = malloc(4);  
void *b = malloc(8);  
void *c = malloc(4);  
free(b);  
void *d = malloc(8);  
free(a);  
void *e = malloc(24);
```

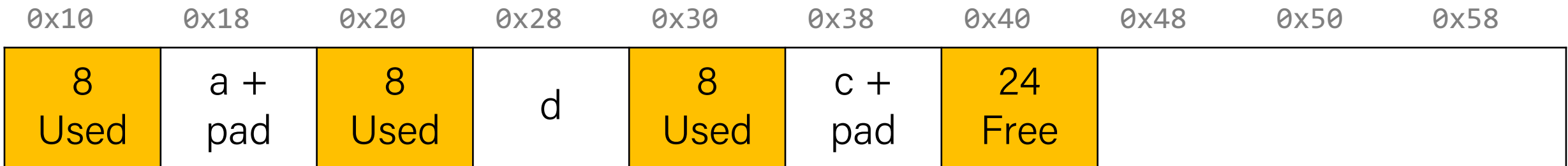
Variable	Value
a	0x18
b	0x28
c	0x38



# Implicit Free List Allocator

```
void *a = malloc(4);  
void *b = malloc(8);  
void *c = malloc(4);  
free(b);  
void *d = malloc(8);  
free(a);  
void *e = malloc(24);
```

Variable	Value
a	0x18
b	0x28
c	0x38
d	0x28

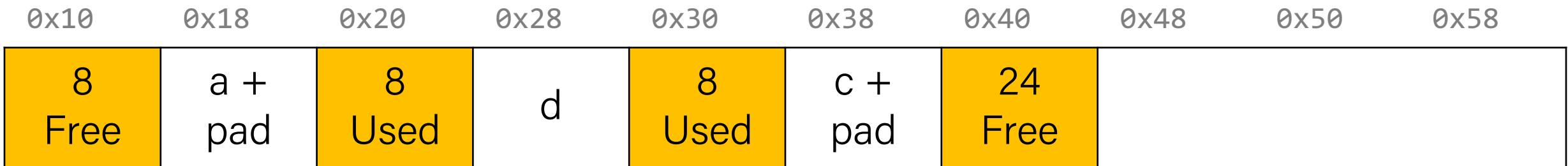




# Implicit Free List Allocator

```
void *a = malloc(4);  
void *b = malloc(8);  
void *c = malloc(4);  
free(b);  
void *d = malloc(8);  
free(a);  
void *e = malloc(24);
```

Variable	Value
a	0x18
b	0x28
c	0x38
d	0x28



# Implicit Free List Allocator

```
void *a = malloc(4);  
void *b = malloc(8);  
void *c = malloc(4);  
free(b);  
void *d = malloc(8);  
free(a);  
void *e = malloc(24);
```

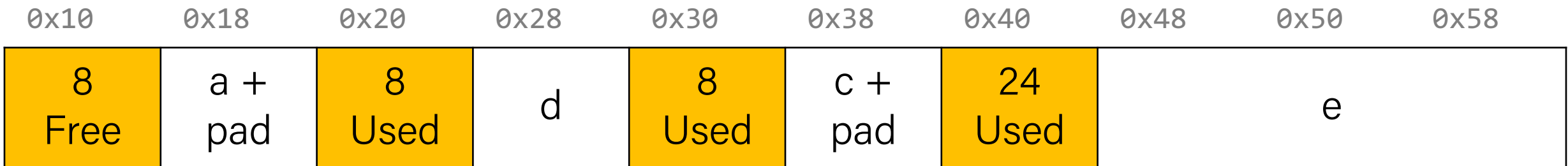
Variable	Value
a	0x18
b	0x28
c	0x38
d	0x28
e	0x48

0x10	0x18	0x20	0x28	0x30	0x38	0x40	0x48	0x50	0x58
8 Free	a + pad	8 Used	d	8 Used	c + pad	24 Used	e		

# Implicit Free List Allocator

```
void *a = malloc(4);  
void *b = malloc(8);  
void *c = malloc(4);  
free(b);  
void *d = malloc(8);  
free(a);  
void *e = malloc(24);
```

Variable	Value
a	0x18
b	0x28
c	0x38
d	0x28
e	0x48



# Representing Headers

How can we store both a size and a status (Free/Allocated) in 8 bytes?

Int for size, int for status? **no! malloc/realloc use size\_t for sizes!**

**Key idea:** block sizes will *always be multiples of 8*. (Why?)

- Least-significant 3 bits will be unused!
- *Solution:* use one of the 3 least-significant bits to store free/allocated status

# Implicit Free List Allocator

How can we choose a free block to use for an allocation request?

- **First fit:** search the list from beginning each time and choose first free block that fits.
- **Next fit:** instead of starting at the beginning, continue where previous search left off.
- **Best fit:** examine every free block and choose the one with the smallest size that fits.
- First fit/next fit easier to implement
- What are the pros/cons of each approach?

# Practice 1

- For the following heap layout, what would the heap look like after the following request is made, assuming we are using an **implicit** free list allocator with a **first-fit** approach?

[24 byte payload, free] [16 byte payload, free] [8 byte payload, allocated for A]

```
void *b = malloc(8);
```

[8 byte payload, allocated for B] [8 byte payload, free] [16 byte payload, free]  
[8 byte payload, allocated for A]

# Practice 2

- For the following heap layout, what would the heap look like after the following request is made, assuming we are using an **implicit** free list allocator with a **best-fit** approach?

[24 byte payload, free] [8 byte payload, free] [8 byte payload, allocated for A]

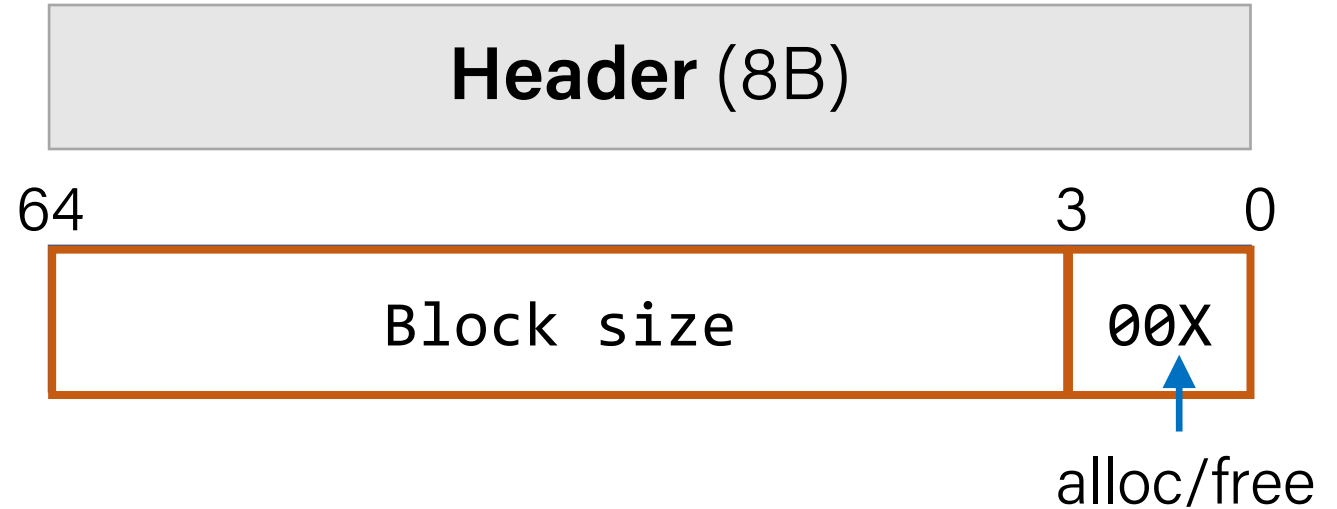
```
void *b = malloc(8);
```

[24 byte payload, free] [8 byte payload, allocated for B] [8 byte payload, allocated for A]

# Implicit Free List Summary

For **all blocks**,

- Have a header that stores size and status.
- Our list links *all* blocks, allocated (A) and free (F).



Keeping track of free blocks:

- **Improves memory utilization** (vs bump allocator)
- **Decreases throughput** (worst case allocation request has  $O(A + F)$  time)
- Increases design complexity ☺



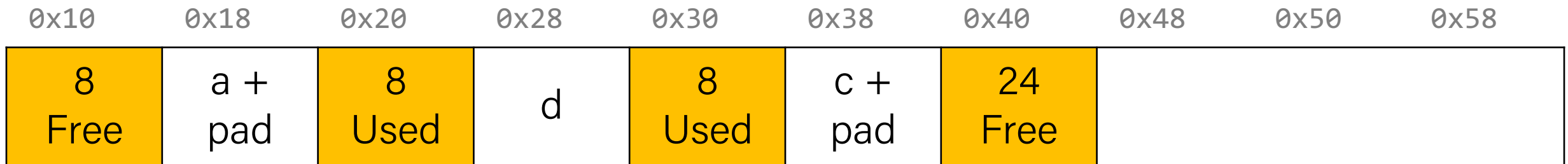
# Splitting Policy

Up to you!

...

```
void *e = malloc(16);
```

So far, we have seen that a reasonable allocation request splits a free block into an allocated block and a free block with remaining space. **What about edge cases?**



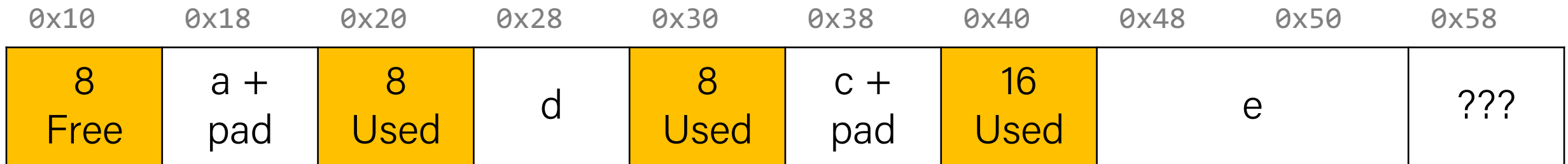
# Splitting Policy

Up to you!

...

```
void *e = malloc(16);
```

So far, we have seen that a reasonable allocation request splits a free block into an allocated block and a free block with remaining space. **What about edge cases?**



# Splitting Policy

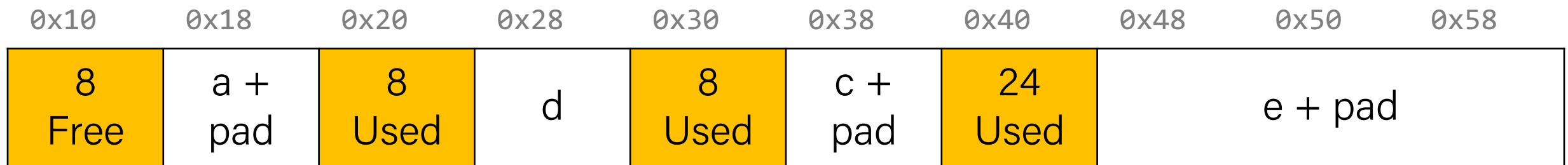
Up to you!

...

```
void *e = malloc(16);
```

So far, we have seen that a reasonable allocation request splits a free block into an allocated block and a free block with remaining space. **What about edge cases?**

**A. Throw into allocation for e as extra padding?** *Internal fragmentation – unused bytes because of padding*



# Splitting Policy

Up to you!

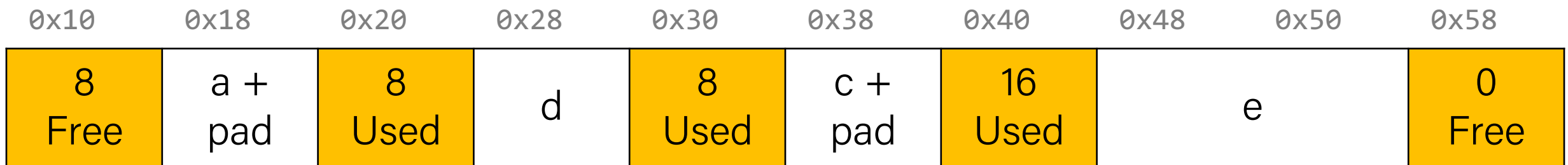
...

```
void *e = malloc(16);
```

So far, we have seen that a reasonable allocation request splits a free block into an allocated block and a free block with remaining space. **What about edge cases?**

A. Throw into allocation for e as extra padding?

**B. Make a "zero-byte free block"?** *External fragmentation – unused free blocks*



# Revisiting Our Goals

Questions we considered:

1. How do we keep track of free blocks? **Using headers!**
2. How do we choose an appropriate free block in which to place a newly allocated block? **Iterate through all blocks.**
3. After we place a newly allocated block in some free block, what do we do with the remainder of the free block? **Try to make the most of it!**
4. What do we do with a block that has just been freed? **Update its header!**

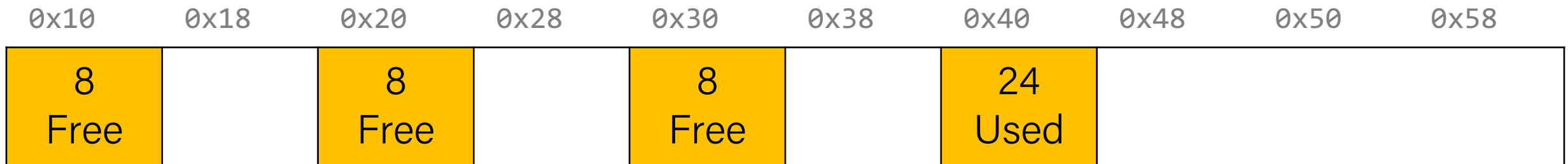
# Implicit Allocator

- **Must have** headers that track block information (size, status in-use or free) –you must use the 8 byte header size, storing the status using the free bits (this is larger than the 4 byte headers specified in the book, as this makes it easier to satisfy the alignment constraint and store information).
- **Must have** free blocks that are recycled and reused for subsequent malloc requests if possible
- **Must have** a malloc implementation that searches the heap for free blocks via an implicit list (i.e. traverses block-by-block).
- **Does not need to** have coalescing of free blocks
- **Does not need to** support in-place realloc

# Coalescing

```
void *e = malloc(24);    // returns NULL!
```

You do not need to worry about this problem for the implicit allocator, but this is a requirement for the *explicit* allocator! (More about this later).



# In-Place Realloc

```
void *a = malloc(4);  
void *b = realloc(a, 8);
```

0x10      0x18      0x20      0x28      0x30      0x38      0x40      0x48      0x50      0x58

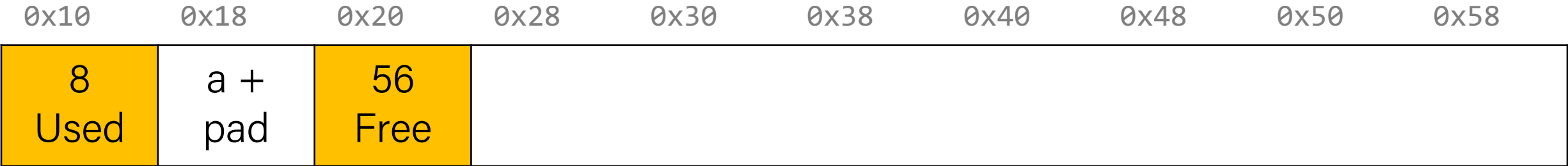
72  
Free



# In-Place Realloc

```
void *a = malloc(4);  
void *b = realloc(a, 8);
```

Variable	Value
a	0x10



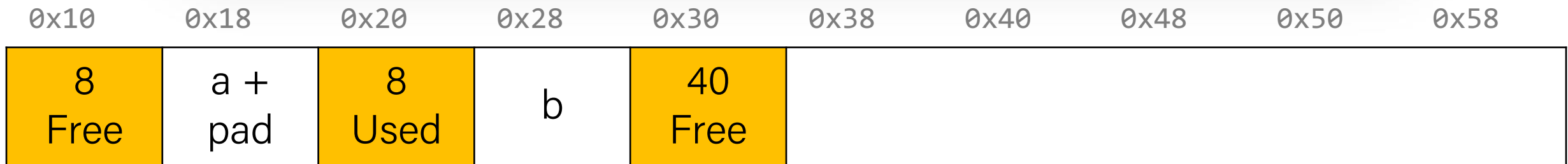
# In-Place Realloc

```
void *a = malloc(4);
```

```
void *b = realloc(a, 8);
```

Variable	Value
a	0x10
b	0x28

The implicit allocator can always move memory to a new location for a `realloc` request. The *explicit* allocator must support in-place `realloc` (more on this later).



# Summary: Implicit Allocator

- An implicit allocator is a more efficient implementation that has reasonable **throughput** and **utilization** due to its recycling of blocks.

Can we do better?

1. Can we avoid searching all blocks for free blocks to reuse?
2. Can we merge adjacent free blocks to keep large spaces available?
3. Can we avoid always copying/moving data during `realloc`?

# Recap

- Shared Libraries
- Case study: Library interpositioning
- The heap so far
- What is a heap allocator?
- Heap allocator requirements and goals
- Method 0: Bump Allocator
- Method 1: Implicit Free List Allocator