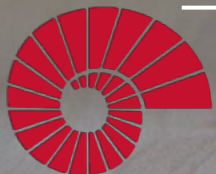


# COMP201

## Computer Systems & Programming

### Lecture #25 – Managing The Heap



KOÇ  
UNIVERSITY

Aykut Erdem // Koç University // Fall 2023

# Recap

- Static Linking
- Symbol Resolution
- Relocation
- Static Libraries
- Shared Libraries

# Plan for Today

- The heap so far
- What is a heap allocator?
- Heap allocator requirements and goals
- Bump Allocator
- Implicit Free List Allocator

**Disclaimer:** Slides for this lecture were borrowed from

—Nick Troccoli's Stanford CS107 class

—Ruth Anderson's UW CSE 351 class

# Multiple Ways to Store Program Data

- Static global data
  - *Fixed size* at compile-time
  - Entire *lifetime of the program* (loaded from executable)
  - Portion is read-only (e.g. string literals)
- Stack-allocated data
  - Local/temporary variables
    - *Can* be dynamically sized (in some versions of C)
  - *Known lifetime* (deallocated on `return`)
- **Dynamic (heap) data**
  - Size known only at runtime (*i.e.* based on user-input)
  - Lifetime known only at runtime (long-lived data structures)

```
int array[1024];

int* foo(int n) {
    int tmp;
    int local_array[n];

    int* dyn =
        (int*)malloc(n*sizeof(int));
    return dyn;
}
```

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

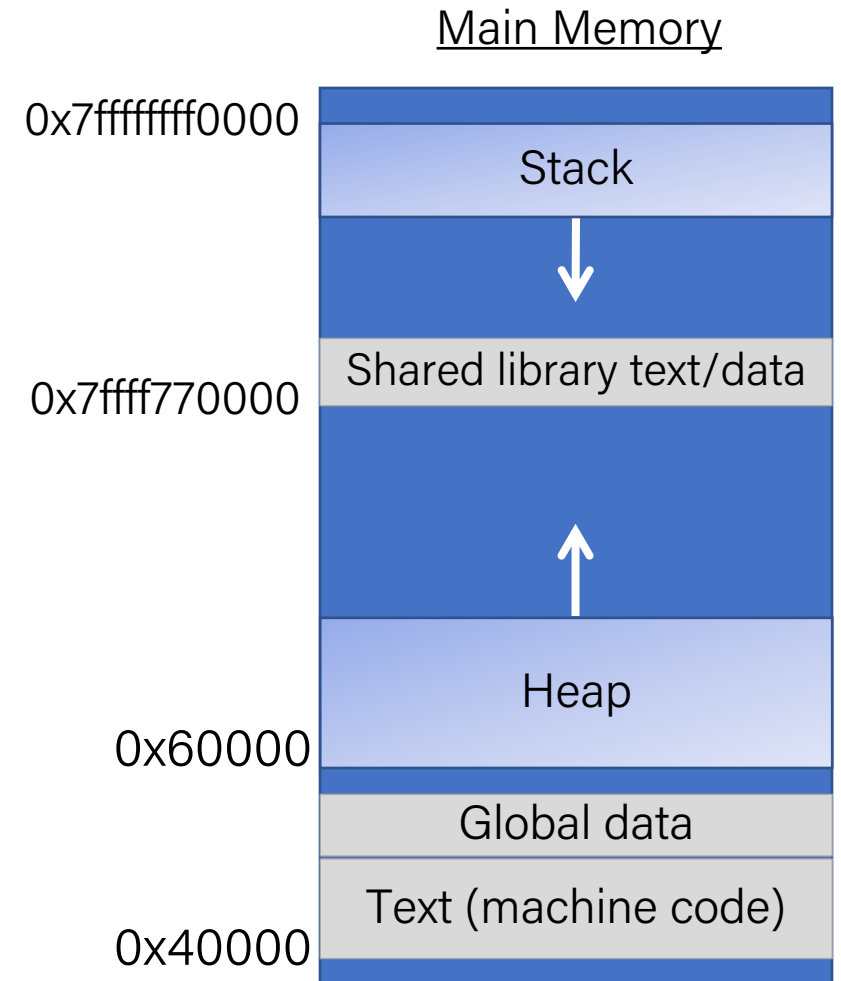
# Lecture Plan

- The heap so far
- What is a heap allocator?
- Heap allocator requirements and goals
- Bump Allocator
- Implicit Free List Allocator

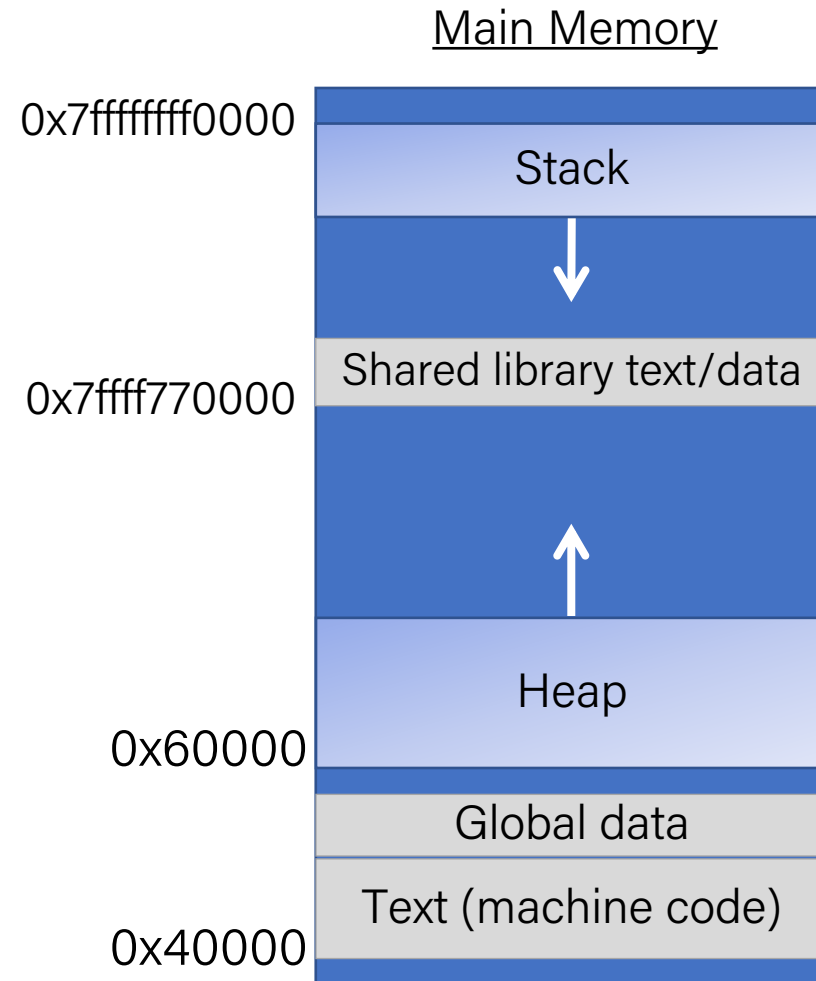


# Running a program

- **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



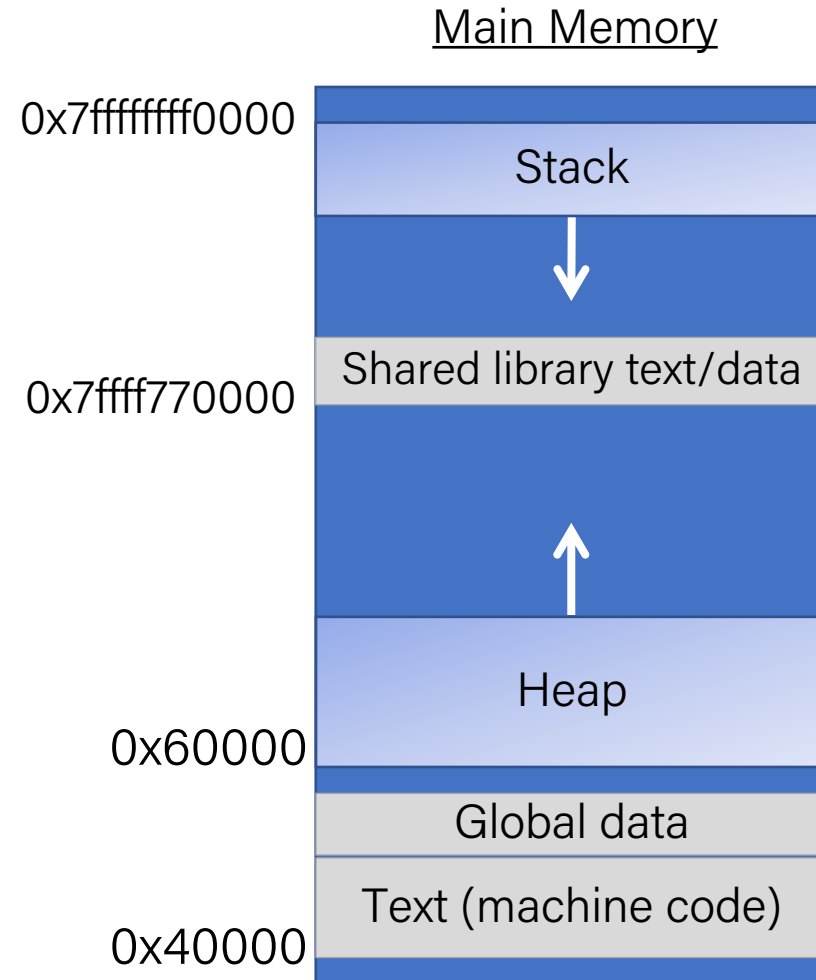
# The Stack Revisited



**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
- Bump Allocator
- Implicit Free List Allocator

# Revisited: Allocating Memory in C

- Need to `#include <stdlib.h>`
- **`void* malloc(size_t size)`**
  - Allocates a continuous block of `size` bytes of uninitialized memory
  - Returns a pointer to the beginning of the allocated block; `NULL` indicates failed request
    - Typically aligned to an 8-byte (x86) or 16-byte (x86-64) boundary
    - Returns `NULL` if allocation failed (also sets `errno`) or `size==0`
  - Different blocks not necessarily adjacent
- Good practices:
  - `ptr = (int*) malloc(n*sizeof(int));`
    - `sizeof` makes code more portable
    - `void*` is implicitly cast into any pointer type; explicit typecast will help you catch coding errors when pointer types don't match

# Revisited: Allocating Memory in C

- Need to `#include <stdlib.h>`
- **`void* malloc(size_t size)`**
  - Allocates a continuous block of `size` bytes of uninitialized memory
  - Returns a pointer to the beginning of the allocated block; NULL indicates failed request
    - Typically aligned to an 8-byte (x86) or 16-byte (x86-64) boundary
    - Returns NULL if allocation failed (also sets `errno`) or `size==0`
  - Different blocks not necessarily adjacent
- Related functions:
  - **`void* calloc(size_t nitems, size_t size)`**
    - “Zeros out” allocated block
  - **`void* realloc(void* ptr, size_t size)`**
    - Changes the size of a previously allocated block (if possible)
  - **`void* sbrk(intptr_t increment)`**
    - Used internally by allocators to grow or shrink the heap

# Revisited: Freeing Memory in C

- Need to `#include <stdlib.h>`
- **void** `free(void* p)`
  - Releases whole block pointed to by `p` to the pool of available memory
  - Pointer `p` must be the address *originally* returned by `m/c/realloc` (*i.e.* beginning of the block), otherwise system exception raised
  - Don't call `free` on a block that has already been released or on `NULL`

# Memory Allocation Example in C

```
void foo(int n, int m) {
    int i, *p;
    p = (int*) malloc(n*sizeof(int));           /* allocate block of n ints */
    if (p == NULL) {                           /* check for allocation error */
        perror("malloc");
        exit(0);
    }
    for (i=0; i<n; i++)                        /* initialize int array */
        p[i] = i;

    p = (int*) realloc(p, (n+m)*sizeof(int));  /* add space for m ints to end of p block */
    if (p == NULL) {                           /* check for allocation error */
        perror("realloc");
        exit(0);
    }
    for (i=n; i < n+m; i++)                    /* initialize new spaces */
        p[i] = i;
    for (i=0; i<n+m; i++)                      /* print new array */
        printf("%d\n", p[i]);
    free(p);                                   /* free p */
}
```



# 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



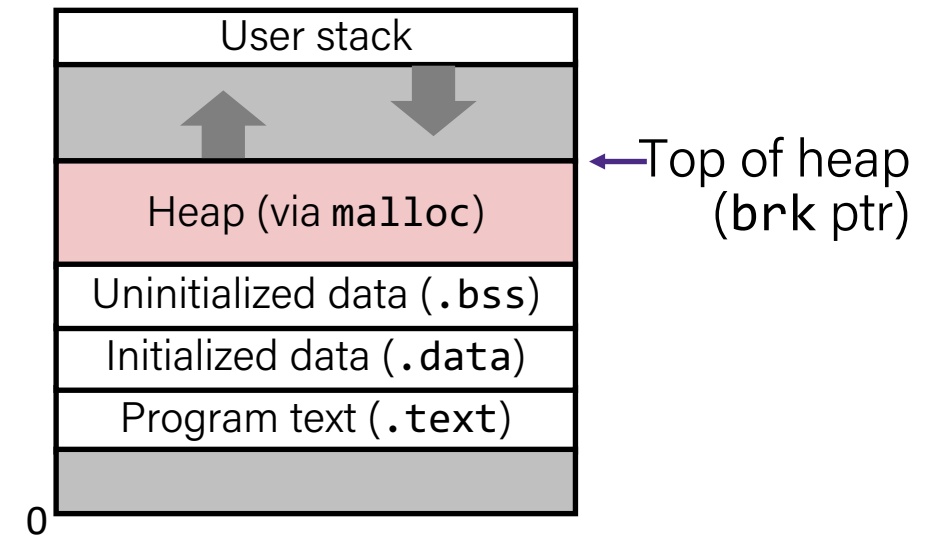
(aka **Heap Allocator**)

# Types of heap allocators

- **Explicit allocator:** programmer allocates and frees space
  - Example: `malloc` and `free` in C
- **Implicit allocator:** programmer only allocates space (no free)
  - Example: garbage collection in Java, Caml, and Lisp

# Dynamic memory allocation

- Allocator organizes heap as a collection of variable-sized **blocks**, which are either **allocated** or **free**
  - Allocator requests pages in the heap region; virtual memory hardware and OS kernel allocate these pages to the process
  - Application objects are typically smaller than pages, so the allocator manages blocks *within* pages
    - (Larger objects handled too; ignored here)



# 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 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 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 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

**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 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 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

- The heap so far
- What is a heap allocator?
- Heap allocator requirements and goals
- Bump Allocator
- Implicit 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.

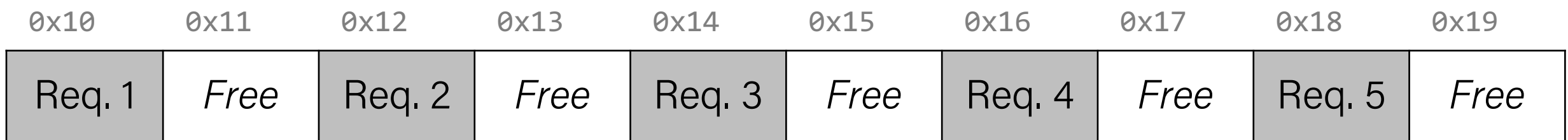


# Utilization

- 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!

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

Req. 1

Req. 2

Req. 3

Req. 4

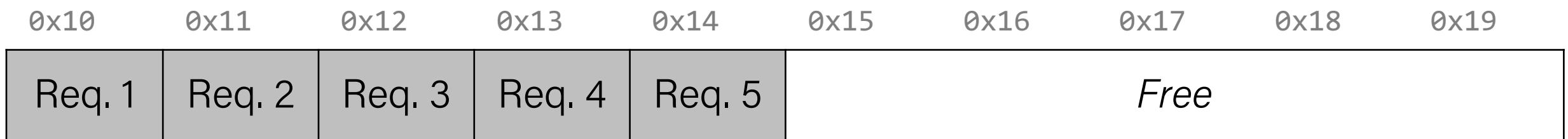
Req. 5

*Free*

# 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!

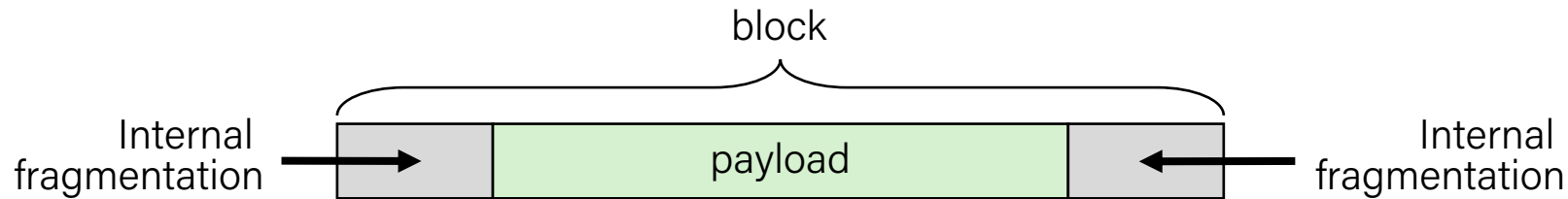


# Fragmentation

- Poor memory utilization is caused by **fragmentation**
  - Sections of memory are not used to store anything useful, but cannot satisfy allocation requests
  - Two types: **internal** and **external**
- **Recall:** Fragmentation in structs
  - Internal fragmentation was wasted space **inside** of the struct (between fields) due to alignment
  - External fragmentation was wasted space **between** struct instances (e.g. in an array) due to alignment
- Now referring to wasted space in the heap **inside** or **between** allocated blocks

# Internal Fragmentation

- For a given block, **internal fragmentation** occurs if payload is smaller than the block

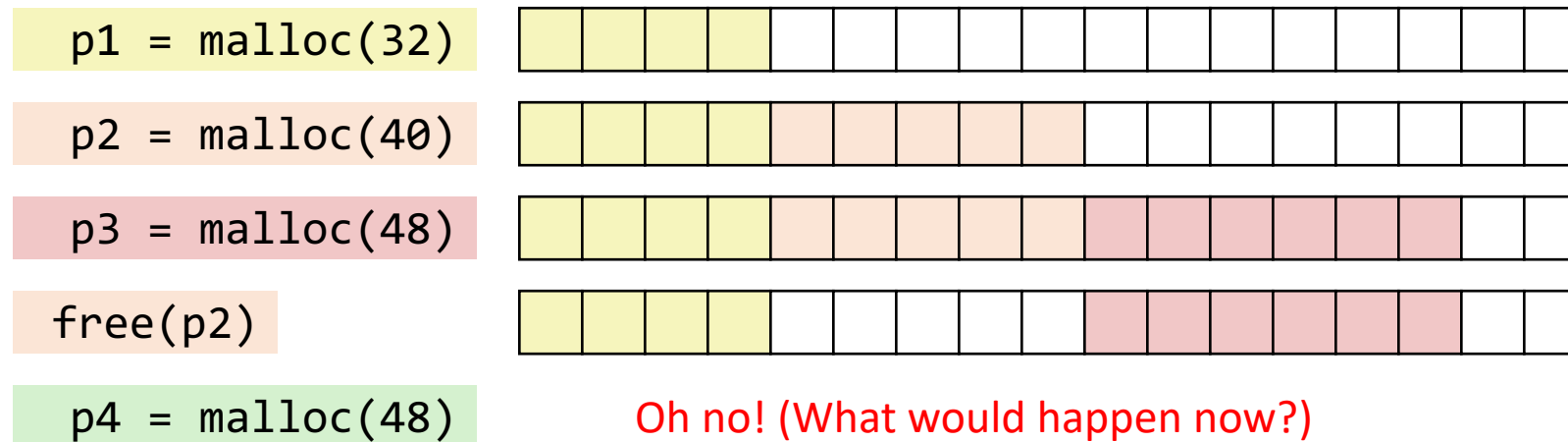


- **Causes:**
  - Padding for alignment purposes
  - Overhead of maintaining heap data structures (inside block, outside payload)
  - Explicit policy decisions (e.g. return a big block to satisfy a small request)
- Easy to measure because only depends on past requests

# External Fragmentation

□ = 8-byte word

- For the heap, **external fragmentation** occurs when allocation/free pattern leaves “holes” between blocks
  - That is, the aggregate payload is non-continuous
  - Can cause situations where there is enough aggregate heap memory to satisfy request, but no single free block is large enough



- Don't know what future requests will be
  - Difficult to impossible to know if past placements will become problematic

# 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**



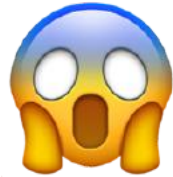
# Lecture Plan

- The heap so far
- What is a heap allocator?
- Heap allocator requirements and goals
- **Bump Allocator**
- Implicit 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);
```

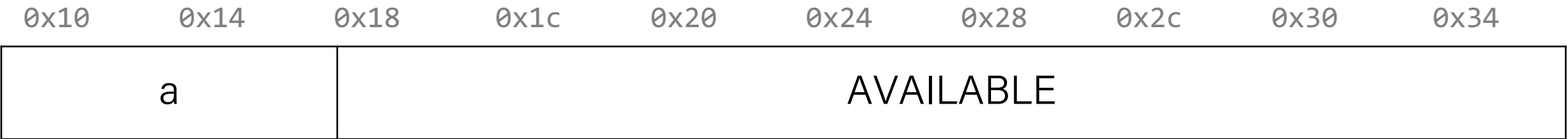
0x10      0x14      0x18      0x1c      0x20      0x24      0x28      0x2c      0x30      0x34

AVAILABLE

# Bump Allocator

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

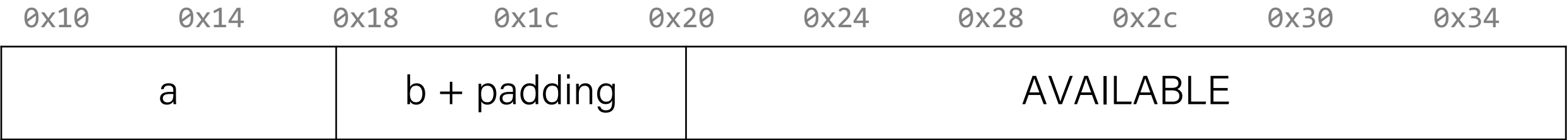
| Variable | Value |
|----------|-------|
| a        | 0x10  |



# Bump Allocator

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

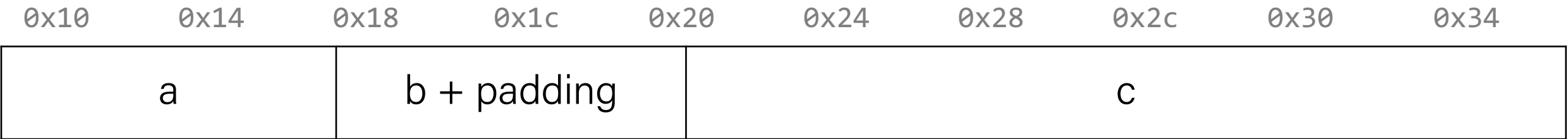
| Variable | Value |
|----------|-------|
| a        | 0x10  |
| b        | 0x18  |



# 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  |

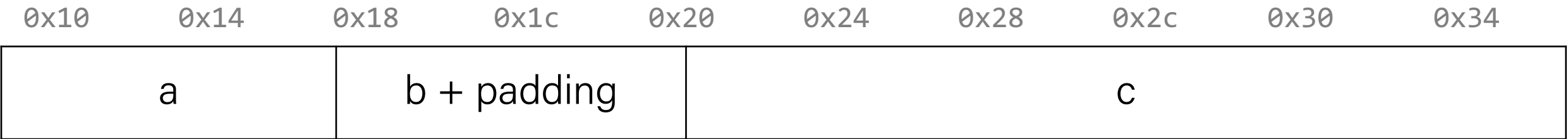




# Bump Allocator

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

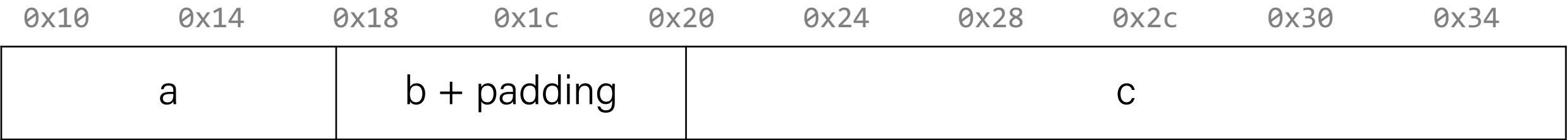
| Variable | Value |
|----------|-------|
| a        | 0x10  |
| b        | 0x18  |
| c        | 0x20  |



# 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?

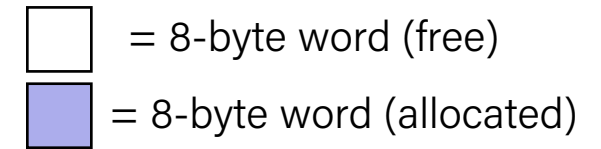
## Questions to consider:

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?

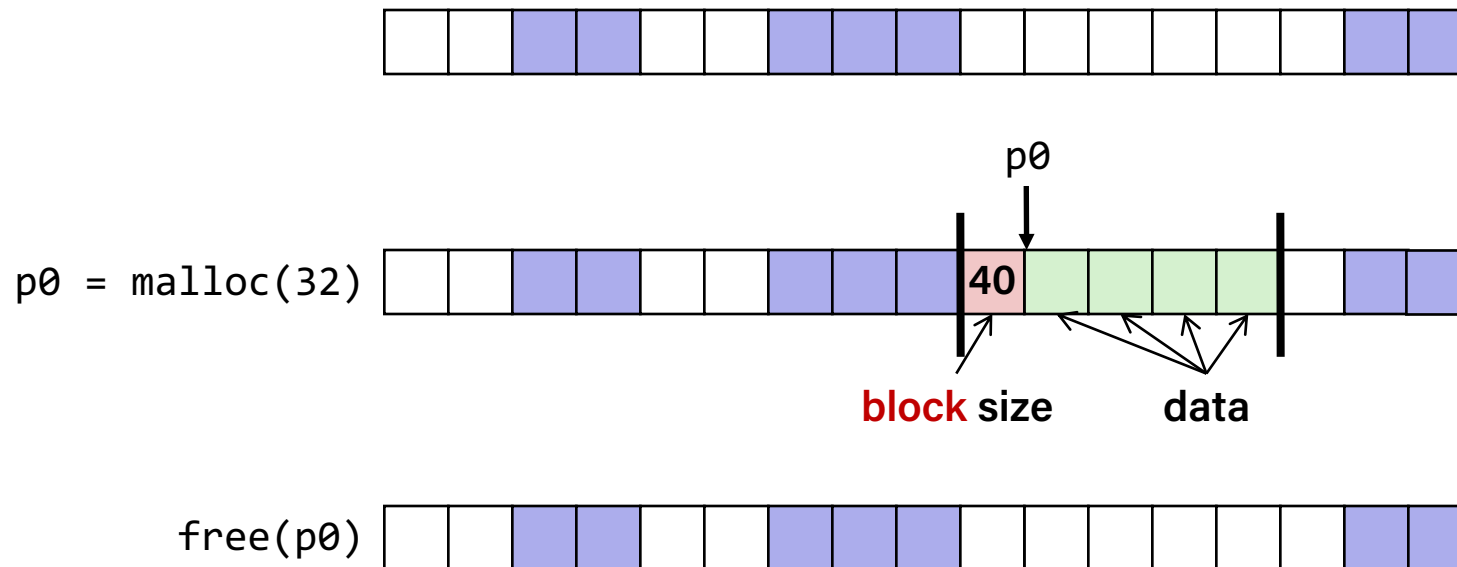
# Implementation Issues

- How do we know how much memory to free given just a pointer?
- How do we keep track of the free blocks?
- How do we pick a block to use for allocation (when many might fit)?
- What do we do with the extra space when allocating a structure that is smaller than the free block it is placed in?
- How do we reinsert a freed block into the heap?

# Knowing How Much to Free



- Standard method
  - Keep the length of a block in the word preceding the data
    - This word is often called the **header field** or **header**
  - Requires an extra word for every allocated block



# Lecture Plan

- The heap so far
- What is a heap allocator?
- Heap allocator requirements and goals
- Bump Allocator
- **Implicit 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

- **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 not.
- When we allocate a block, we look through the free list and find a block of size  $\geq$  the requested size. We update its header to reflect its allocated state and its payload size.
- When we free a block, we update its header to reflect its free state and its payload size.
- The header should be 8 bytes (or larger).
- By storing the block size of each block, we *implicitly* have a *list* of free blocks.

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!.



# 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);
```

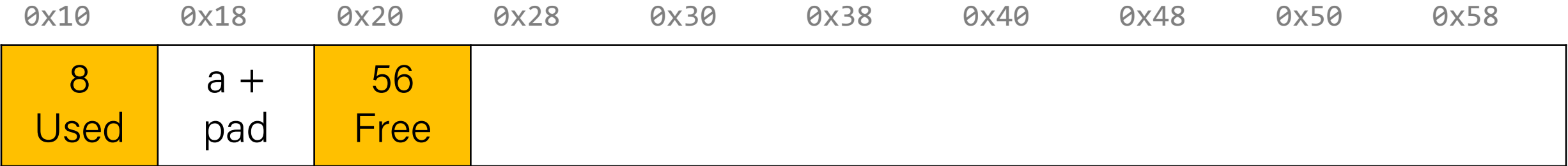
0x10      0x18      0x20      0x28      0x30      0x38      0x40      0x48      0x50      0x58

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);
```

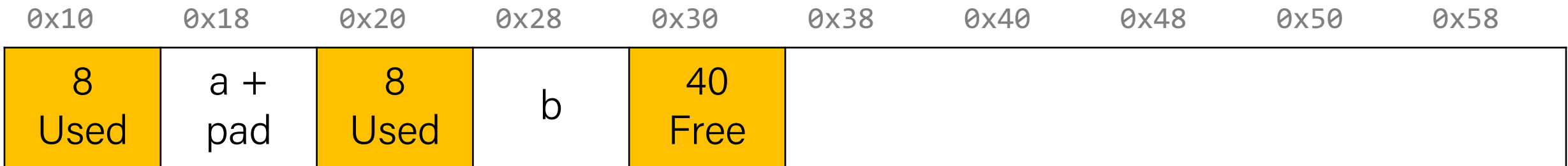
| Variable | Value |
|----------|-------|
| a        | 0x18  |



# 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);
```

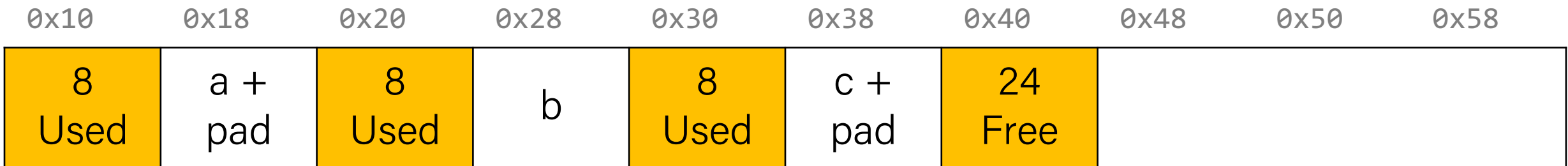
| Variable | Value |
|----------|-------|
| a        | 0x18  |
| b        | 0x28  |



# 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);
```

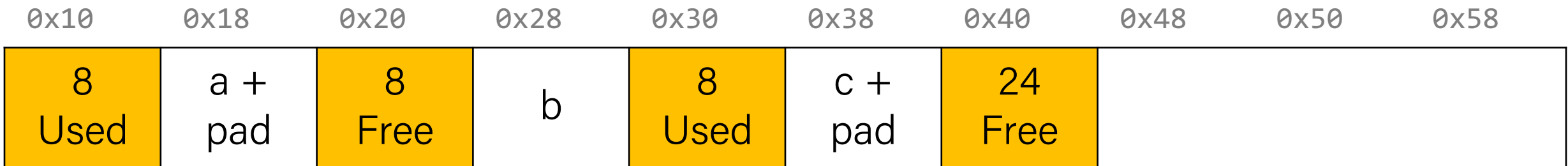
| Variable | Value |
|----------|-------|
| a        | 0x18  |
| b        | 0x28  |
| c        | 0x38  |



# 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  |

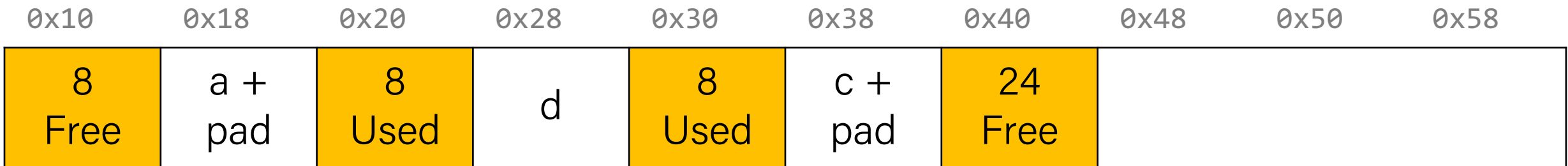
0x10      0x18      0x20      0x28      0x30      0x38      0x40      0x48      0x50      0x58

|           |            |           |   |           |            |            |  |
|-----------|------------|-----------|---|-----------|------------|------------|--|
| 8<br>Used | a +<br>pad | 8<br>Used | d | 8<br>Used | c +<br>pad | 24<br>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);
```

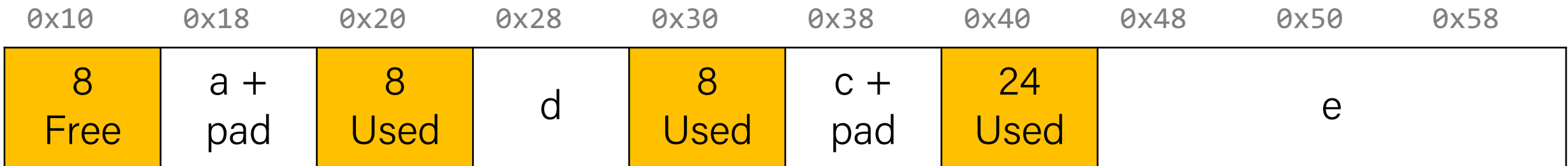
| Variable | Value |
|----------|-------|
| a        | 0x18  |
| b        | 0x28  |
| c        | 0x38  |
| d        | 0x28  |



# 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  |

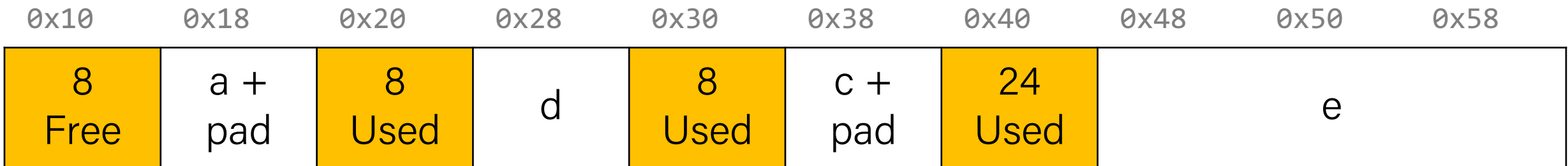




# Implicit Free List Allocator

```
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

- For each block we need: **size, is-allocated?**
  - Could store using two words, but wasteful
- Standard trick
  - If blocks are aligned, some low-order bits of **size** are always 0
  - Use lowest bit as an **allocated/free flag** (fine as long as aligning to  $K > 1$ )
  - When reading **size**, must remember to mask out this bit!

e.g. with 8-byte alignment,  
possible values for size:

00001000 = 8 bytes

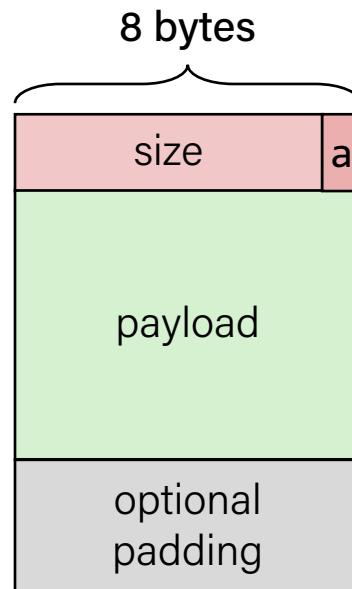
00010000 = 16 bytes

00011000 = 24 bytes

...



Format of  
allocated and  
free blocks:



**a = 1:** allocated block

**a = 0:** free block

**size:** block size (in bytes)

**payload:** application data  
(allocated blocks only)

If **x** is first word (header):

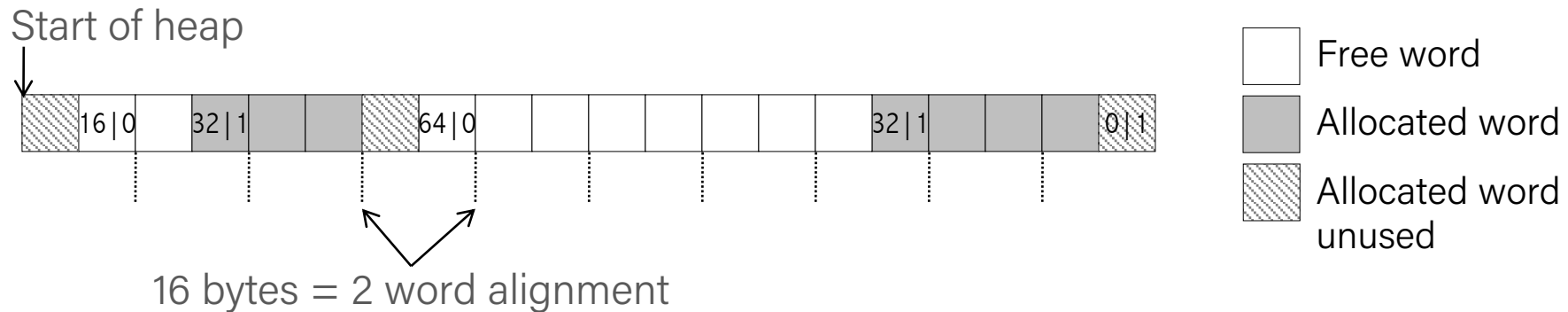
```
x = size | a;
```

```
a = x & 1;
```

```
size = x & ~1;
```

# Implicit Free List Example

- Each block begins with header (size in bytes and allocated bit)
- Sequence of blocks in heap (size|allocated): 16|0, 32|1, 64|0, 32|1



- 16-byte alignment for *payload*
  - May require initial padding (internal fragmentation)
  - Note **size**: padding is considered part of *previous* block
- Special one-word marker (0|1) marks end of list
  - Zero **size** is distinguishable from all other blocks

# 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?

# Implicit List: Finding a Free Block

- **First fit**

- Can take time linear in total number of blocks
- In practice can cause “splinters” at beginning of list

- **Next fit**

- Like first-fit, but **search list starting where previous search finished**
- Should often be faster than first-fit: avoids re-scanning unhelpful blocks
- Some research suggests that fragmentation is worse

- **Best fit**

- Search the list, choose the *best* free block: large enough AND with fewest bytes left over
- Keeps fragments small—usually helps fragmentation
- Usually worse throughput

# 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]

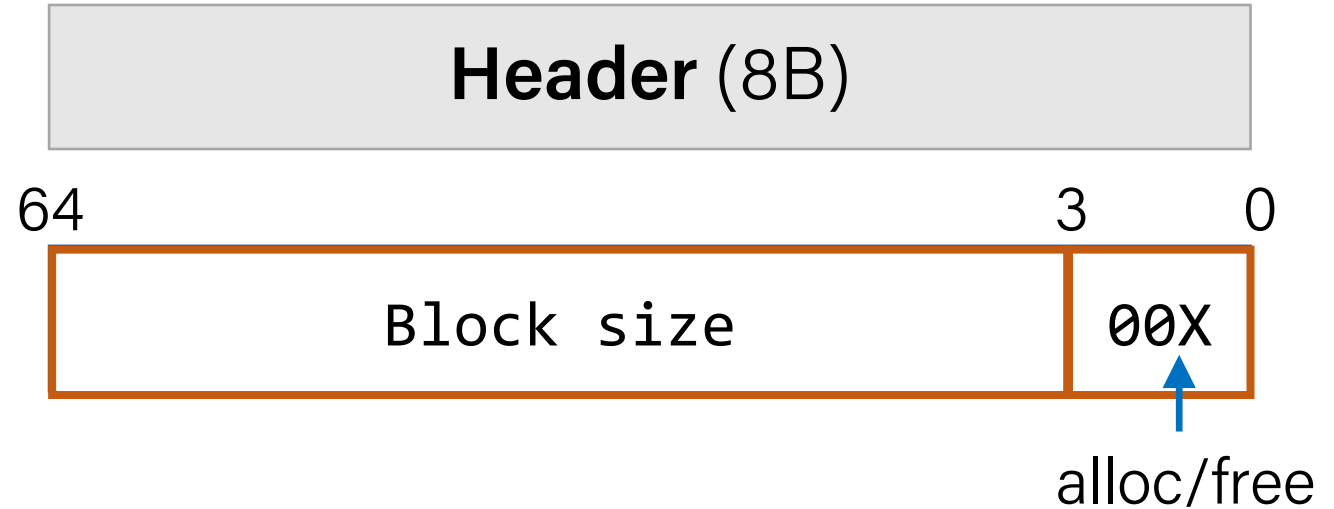
```
void *b = malloc(8);
```

[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 ☺

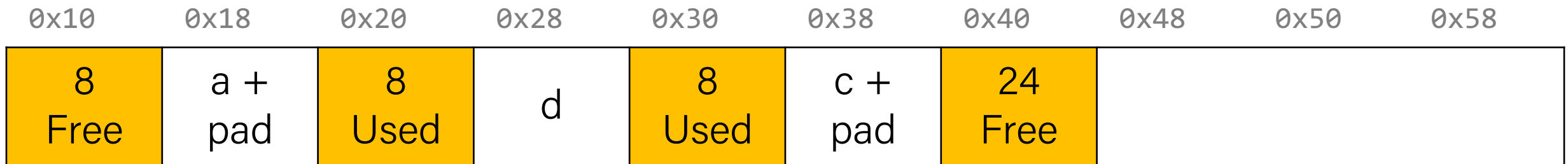


# Splitting Policy

...

```
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?**

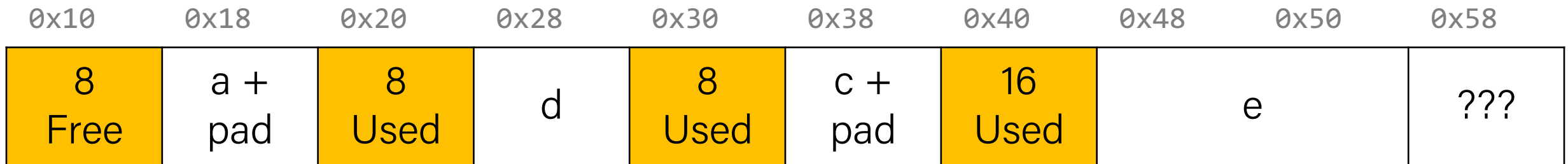


# Splitting Policy

...

```
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?**



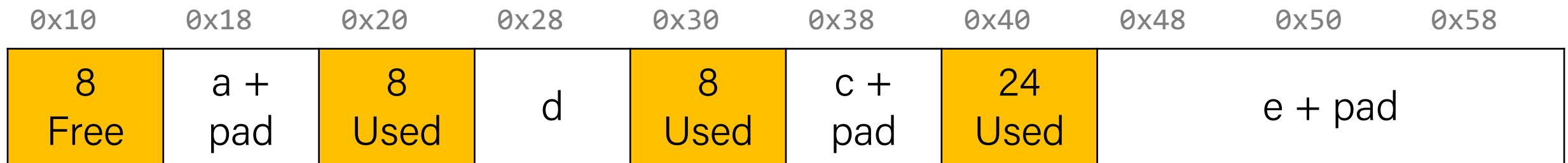
# Splitting Policy

...

```
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

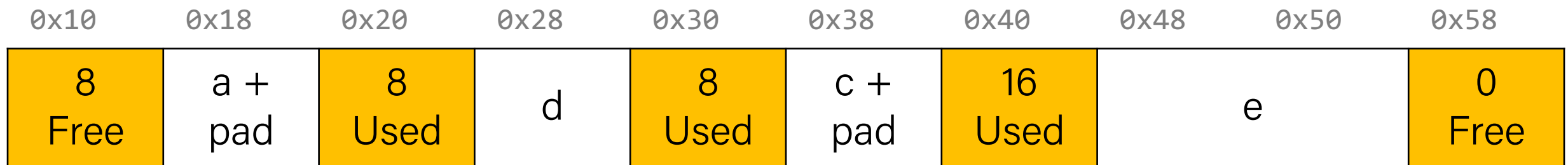
...

```
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

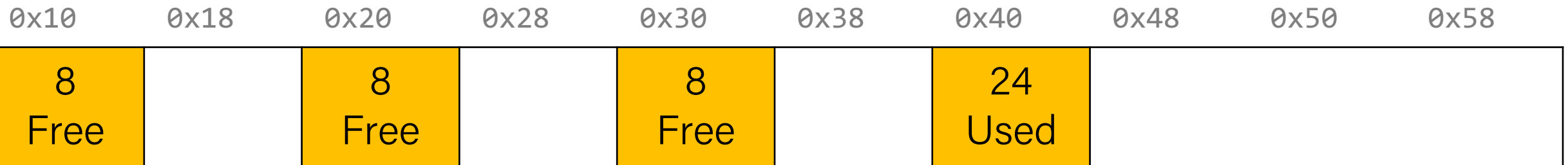
Questions we considered:

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!**

# Coalescing

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

We do not need to worry about this problem for the implicit allocator, but investigate this for the *explicit* allocator! (More about this later).



# In-Place Realloc

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

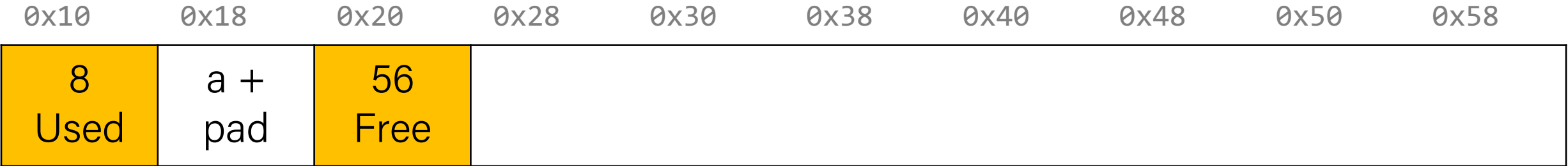
0x10      0x18      0x20      0x28      0x30      0x38      0x40      0x48      0x50      0x58

72  
Free

# In-Place Realloc

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

| Variable | Value |
|----------|-------|
| a        | 0x10  |



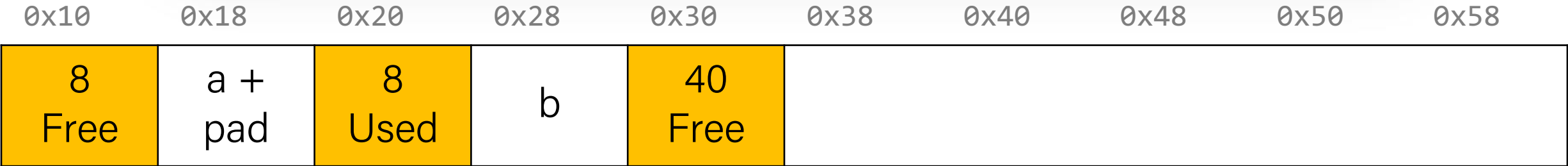


# In-Place Realloc

```
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

- The heap so far
- What is a heap allocator?
- Heap allocator requirements and goals
- Bump Allocator
- Implicit Free List Allocator

- **Next time:** *More on heap allocators*