

CS107, Lecture 23

Managing The Heap, Continued

Reading: B&O 9.9, 9.11

This document is copyright (C) Stanford Computer Science and Nick Troccoli, licensed under Creative Commons Attribution 2.5 License. All rights reserved.

Based on slides created by Cynthia Lee, Chris Gregg, Jerry Cain, Lisa Yan and others.

NOTICE RE UPLOADING TO WEBSITES: This content is protected and may not be shared, uploaded, or distributed. (without expressed written permission)

CS107 Topic 6

How do the core malloc/realloc/free memory-allocation operations work?

Why is answering this question important?

- Combines techniques from across the quarter (bits/bytes, pointers, memory, generics, assembly, efficiency, testing, and more) to understand a real-world system that you have relied on all quarter!
- Learning about the design and tradeoffs in a real-world large system gives us a great example of how to evaluate different designs when there's no one "right" answer.

assign6: implement two different possible designs for a heap allocator, implementing malloc/realloc/free.

Learning Goals

- Learn about different ways to implement a heap allocator
- Understand the tradeoffs between bump, implicit and explicit free list allocators

Lecture Plan

- **Recap:** heap allocators so far
- Method 0: Bump Allocator
- Method 1: Implicit Free List Allocator
- Method 2: Explicit Free List Allocator

Lecture Plan

- **Recap: heap allocators so far**
- Method 0: Bump Allocator
- Method 1: Implicit Free List Allocator
- Method 2: Explicit Free List 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).
- A heap allocator must manage this memory as clients request or no longer need pieces of 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

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).

Utilization

Question: Can we / should we shift these blocks down to make more space?

- **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!

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

Req. 1

Req. 2

Req. 3

Req. 4

Req. 5

Free

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

- **Recap:** heap allocators so far
- **Method 0: Bump Allocator**
- Method 1: Implicit Free List Allocator
- Method 2: Explicit Free List Allocator

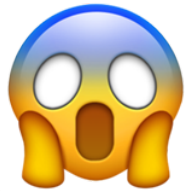
Bump Allocator

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?

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.

Bump Allocator Performance

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. 😞
- We provide a bump allocator implementation as part of the final assignment as a code reading exercise.

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

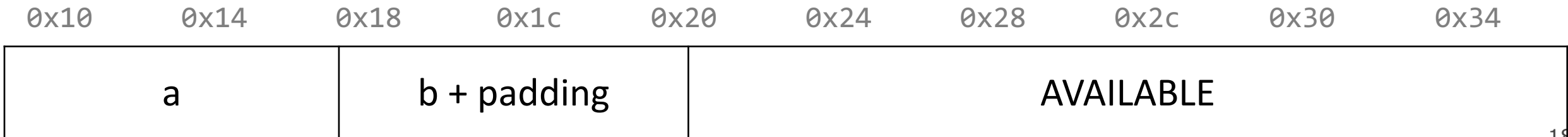
0x10 0x14 0x18 0x1c 0x20 0x24 0x28 0x2c 0x30 0x34



Bump Allocator

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

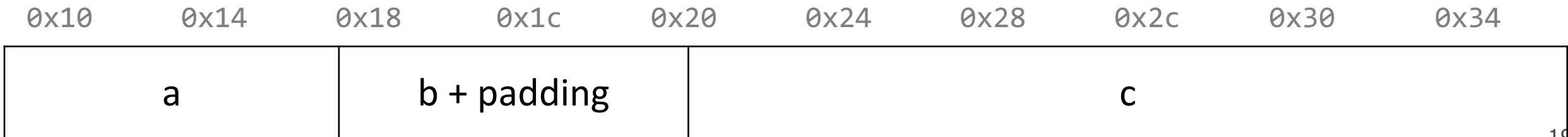
Variable	Value
a	0x10
b	0x18



Bump Allocator

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

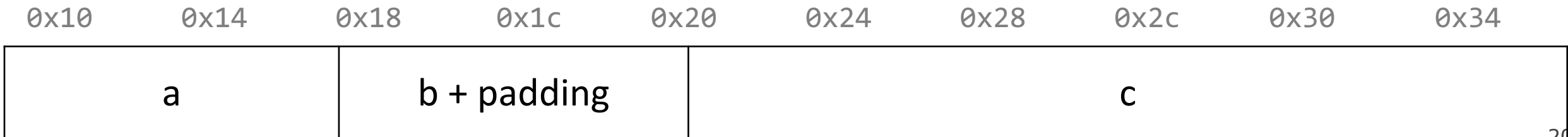
Variable	Value
a	0x10
b	0x18
c	0x20



Bump Allocator

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

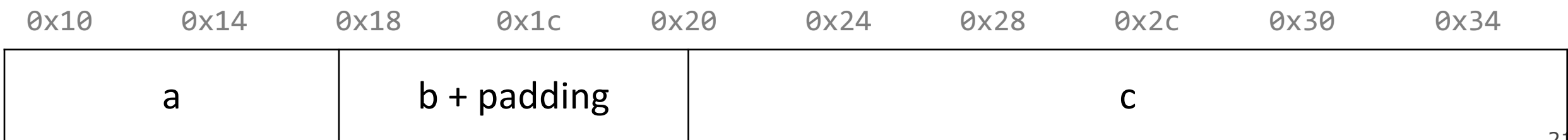
Variable	Value
a	0x10
b	0x18
c	0x20



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

- **Recap:** heap allocators so far
- 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 at the start of each block for a **header** storing its payload size and whether it is allocated or free. (Payload=client data space)
- 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). Overhead!
- 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);
```

“Header”



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);
```

“Payload”

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);
```

Block = header + payload

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);
```

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

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

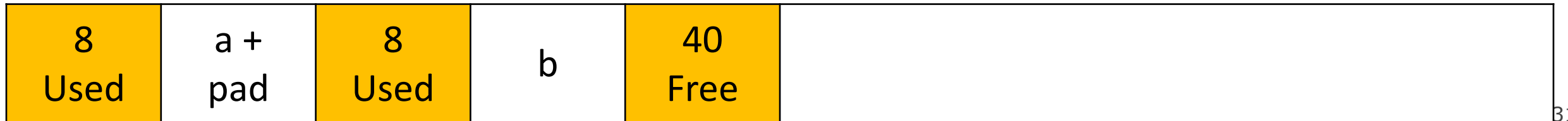


Implicit Free List Allocator

```
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

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



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

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



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

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



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



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



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



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



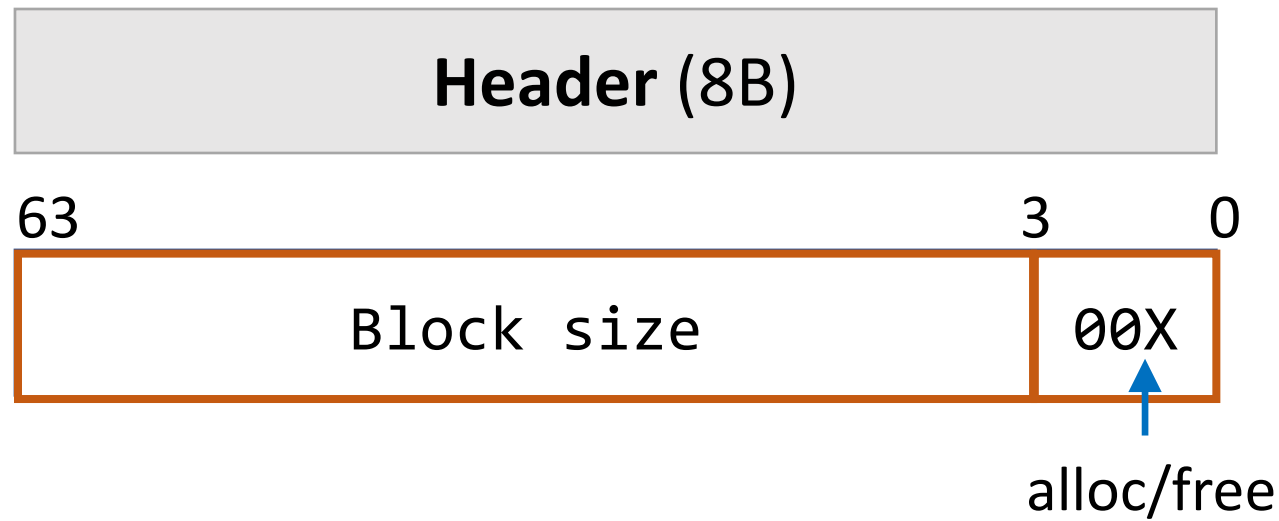
Representing Headers

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



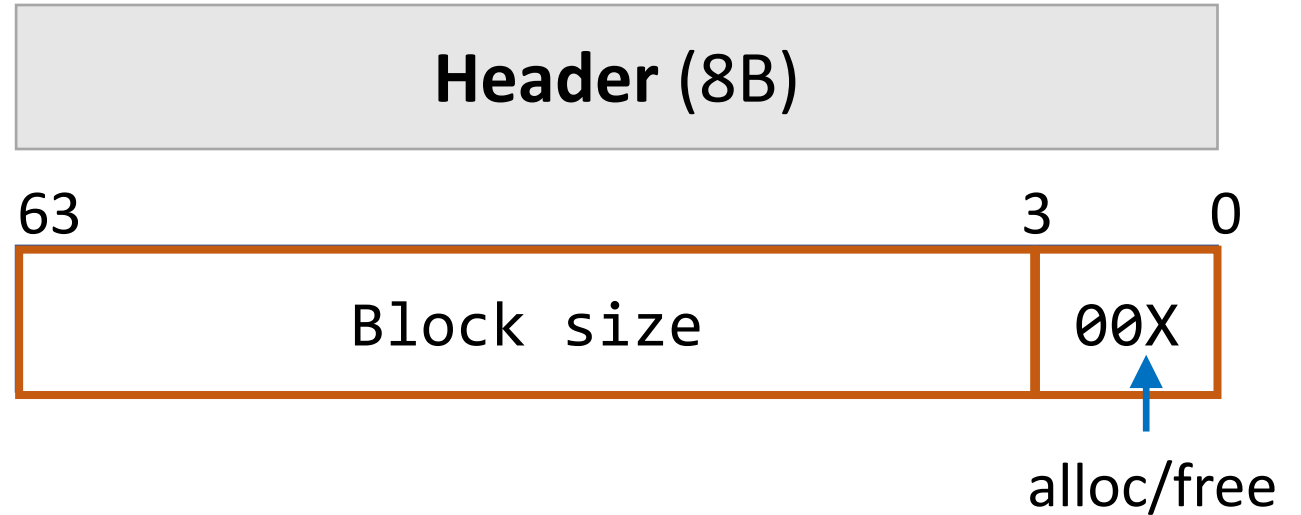
Implicit Free List Allocator

- 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?

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 😊

Up to you!

Implicit free list header design

Should we store the **block size** as

- (A) payload size, or
- (B) header + payload size?

Up to you!

Your decision affects how you traverse the list (be careful of off-by-one)

Up to you!

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?**

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



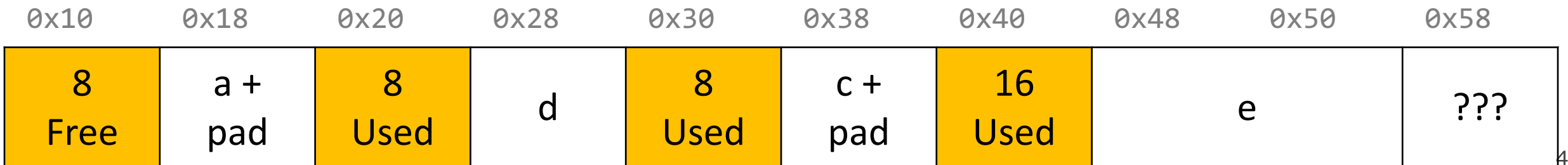
Up to you!

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?**



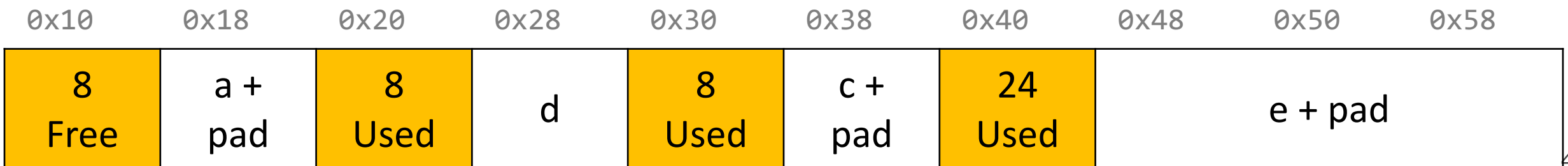
Splitting Policy

...

```
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

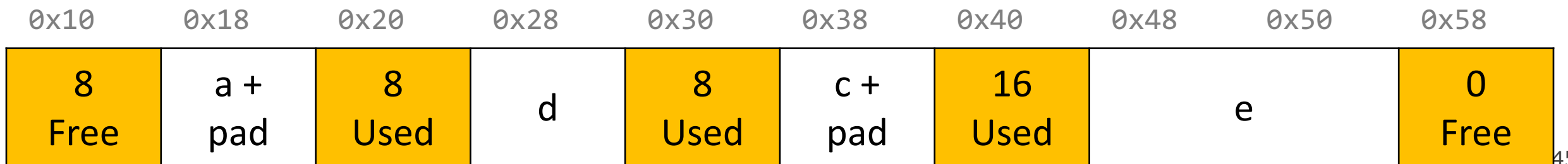
...

```
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!**

Practice 1: Implicit (first-fit)

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?

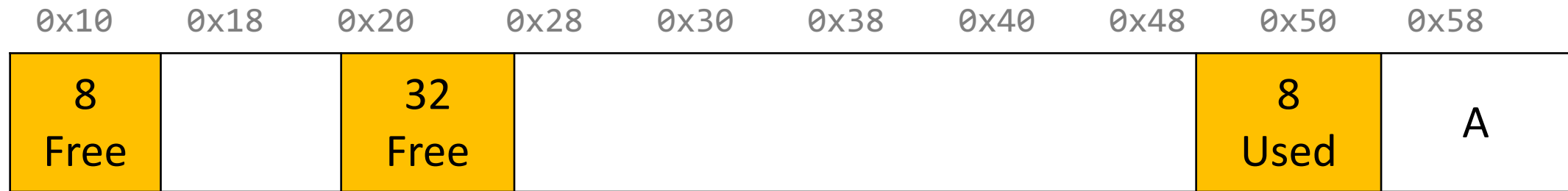


```
void *b = malloc(8);
```

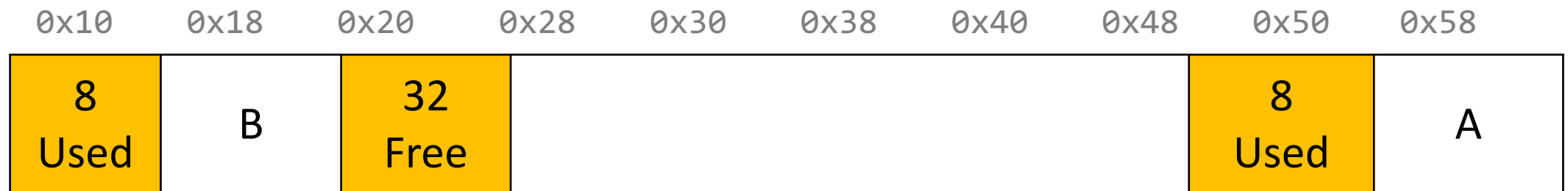


Practice 1: Implicit (first-fit)

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?

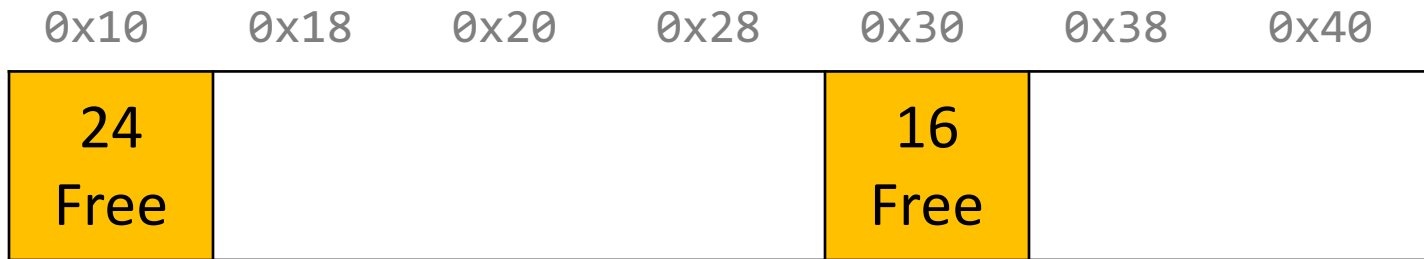


```
void *b = malloc(8);
```



Practice 2: Implicit (first-fit)

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?



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

Respond on PollEv: pollev.com/cs107
or text CS107 to 22333 once to join.

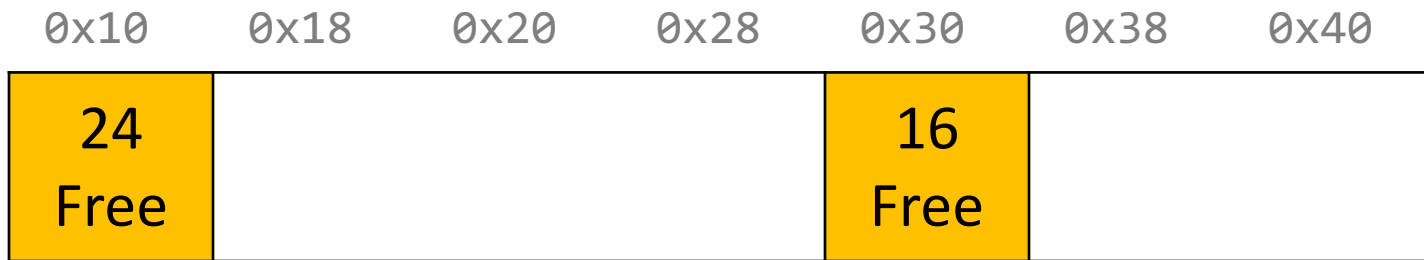


Which image represents what the heap would look like following the request?

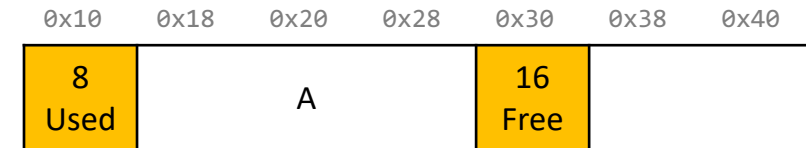
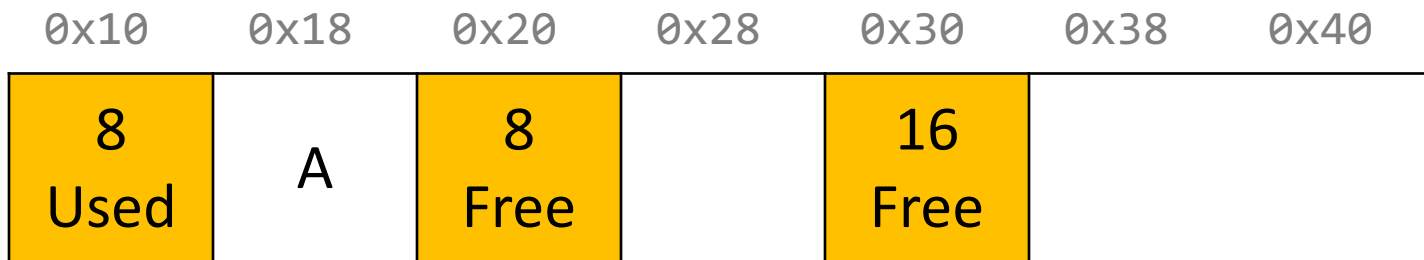


Practice 2: Implicit (first-fit)

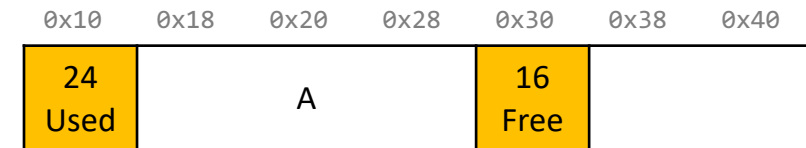
For the following heap layout, what would the heap look like after the following request is made, assuming we are using an **implicit** free list allocator with a **first-fit** approach?



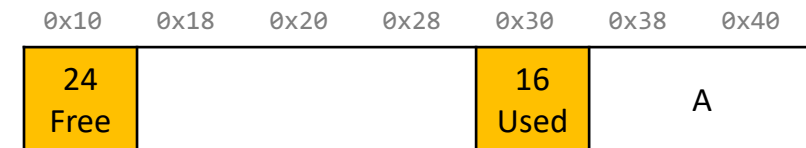
```
void *a = malloc(8);
```



✗ *Space not tracked correctly*



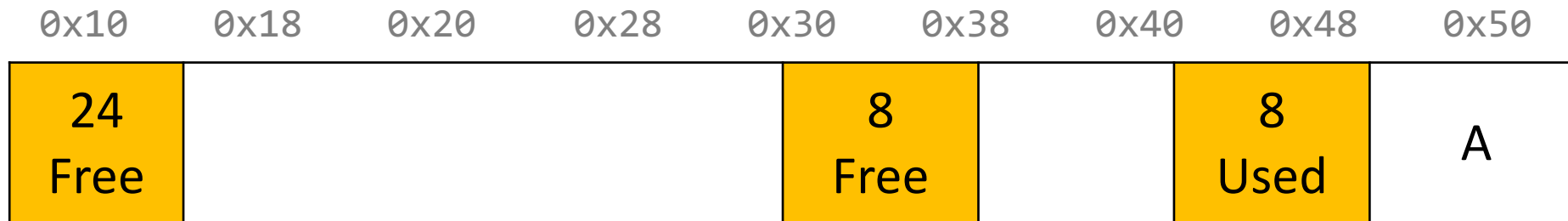
✗ *We can save extra for later*



✗ *First fit chooses first available*

Practice 3: Implicit (best-fit)

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?



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

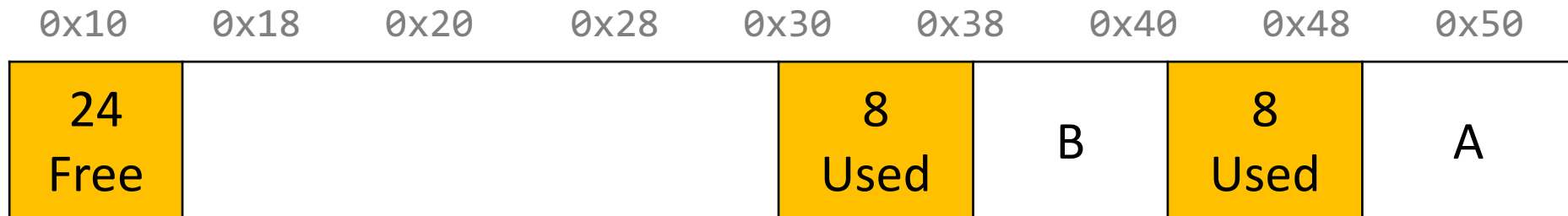


Practice 3: Implicit (best-fit)

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?



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



Final Assignment: 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).

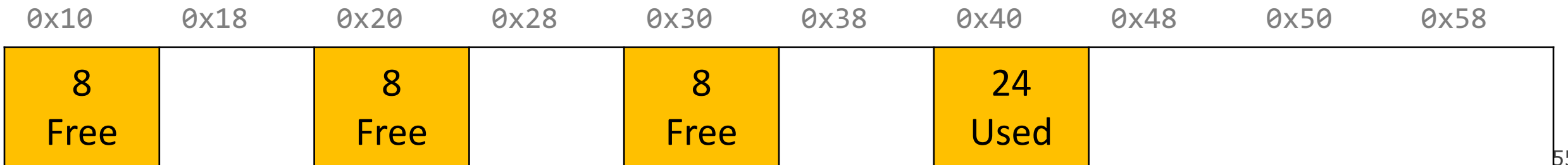
- **Does not need to** have coalescing of free blocks
- **Does not need to** support in-place realloc

(Note: these could be part of an implicit allocator, it's just not a requirement for this assignment)

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



In-Place Realloc

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

Variable	Value
a	0x18

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

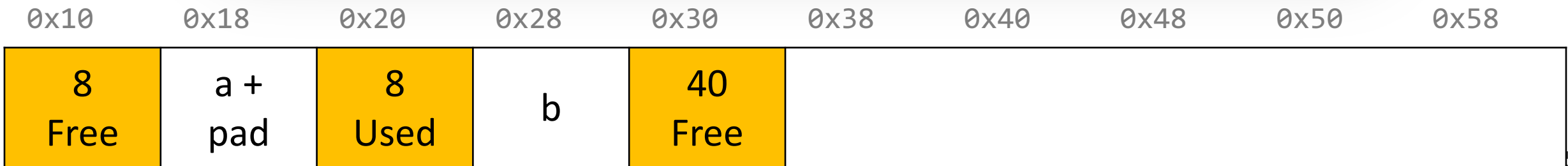


In-Place Realloc

```
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?

Checkpoint Review

Heap allocator terminology: What do the below terms mean/imply?

- Payload, Header, Free/Used(Allocated) status
- Splitting policy
- Memory utilization vs Throughput
- Bump allocator, Implicit free list Allocator
- First-fit approach, Best-fit approach
- Coalescing
- Realloc in place
- Fragmentation

Lecture Plan

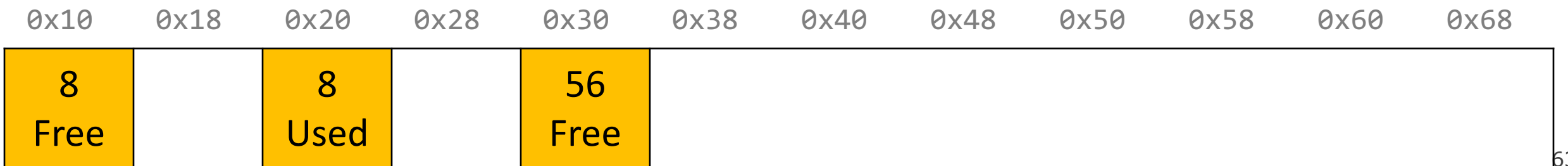
- **Recap:** heap allocators so far
- Method 0: Bump Allocator
- Method 1: Implicit Free List Allocator
- **Method 2: Explicit Free List Allocator**

Lecture Plan

- **Recap:** heap allocators so far
- Method 0: Bump Allocator
- Method 1: Implicit Free List Allocator
- **Method 2: Explicit Free List Allocator**
 - **Explicit Allocator**
 - Coalescing
 - In-place realloc

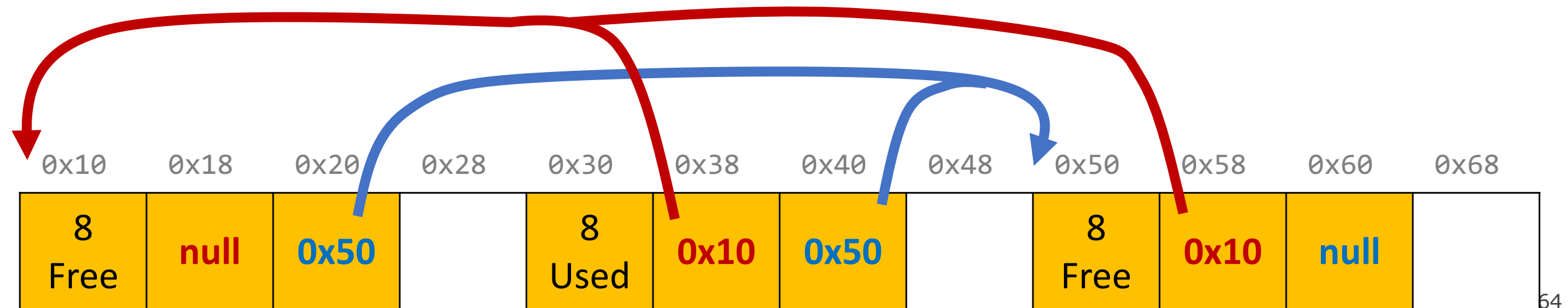
Can We Do Better?

- It would be nice if we could jump *just between free blocks*, rather than all blocks, to find a block to reuse.
- **Idea:** let's modify each header to add a pointer to the previous free block and a pointer to the next free block.



Can We Do Better?

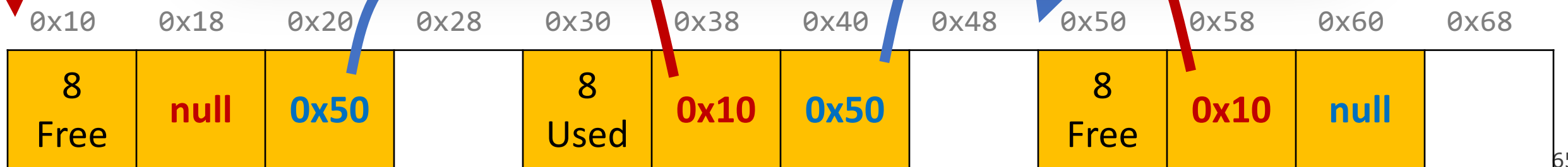
- It would be nice if we could jump *just between free blocks*, rather than all blocks, to find a block to reuse.
- **Idea:** let's modify each header to add a pointer to the **previous** free block and a pointer to the **next** free block.



Can We Do Better?

- It would be nice if we could jump *just between free blocks*, rather than all blocks, to find a block to reuse.
- **Idea:** let's modify each header to add a pointer to the **previous** free block and a pointer to the **next** free block.

This is inefficient – it triples the size of *every* header, when we just need to jump from one free block to another. And even if we just made free headers bigger, it's complicated to have *two* different header sizes.



Can We Do Better?

- It would be nice if we could jump *just between free blocks*, rather than all blocks, to find a block to reuse.
- **Idea:** let's modify each header to add a pointer to the previous free block and a pointer to the next free block. *This is inefficient / complicated.*
- **Where can we put these pointers to the next/previous free block?**
- **Idea:** In a separate data structure?

Can We Do Better?

- It would be nice if we could jump *just between free blocks*, rather than all blocks, to find a block to reuse.
- **Idea:** let's modify each header to add a pointer to the previous free block and a pointer to the next free block. *This is inefficient / complicated.*
- **Where can we put these pointers to the next/previous free block?**
- **Idea:** In a separate data structure? *More difficult to access in a separate place – prefer storing near blocks on the heap itself.*

Can We Do Better?

- **Key Insight:** the payloads of the free blocks aren't being used, because they're free.
- **Idea:** since we only need to store these pointers for free blocks, let's store them in the first 16 bytes of each free block's payload!

0x10 0x18 0x20 0x28 0x30 0x38 0x40 0x48 0x50 0x58 0x60 0x68

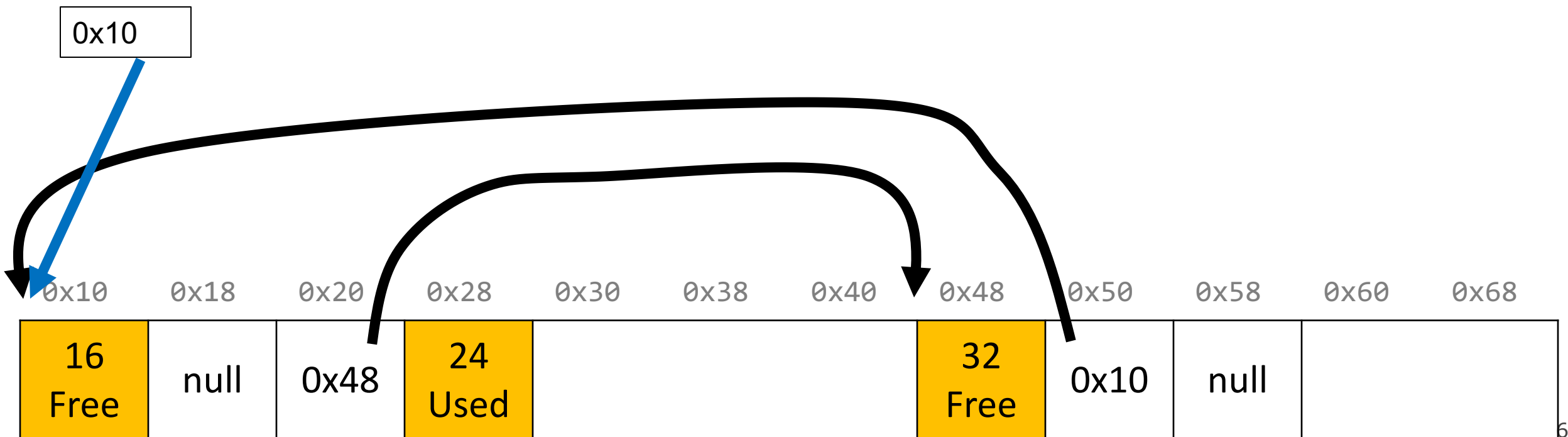


Can We Do Better?

- **Key Insight:** the payloads of the free blocks aren't being used, because they're free.
- **Idea:** since we only need to store these pointers for free blocks, let's store them in the first 16 bytes of each free block's payload!

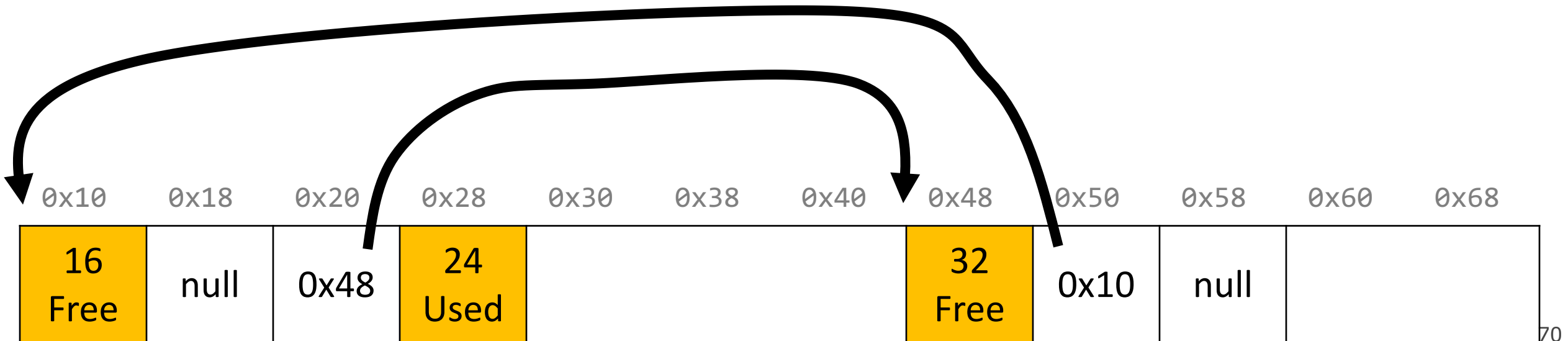
First free block

0x10



Can We Do Better?

- **Key Insight:** the payloads of the free blocks aren't being used, because they're free.
- **Idea:** since we only need to store these pointers for free blocks, let's store them in the first 16 bytes of each free block's payload!
- This means each payload must be big enough to store 2 pointers (16 bytes). So we must require that for every block, free and allocated. (why?)



Explicit Free List Allocator

- This design builds on the implicit allocator, but also stores pointers to the next and previous free block inside each free block's payload.
- When we allocate a block, we look through just the free blocks using our linked list to find a free one, and we update its header and the linked list 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 and update the linked list.

This **explicit** list of free blocks increases request throughput, with some costs (design and internal fragmentation)

Recap

- **Recap:** heap allocators so far
- Method 0: Bump Allocator
- Method 1: Implicit Free List Allocator
- Method 2: Explicit Free List Allocator

Lecture 23 takeaway: Bump, implicit free list and explicit free list are 3 heap allocator designs, each with their own tradeoffs. The implicit free list and explicit free list designs use headers to keep track of blocks. The explicit free list allocator adds an embedded doubly-linked list between free blocks, stored in the free block payloads.

Next time: more about explicit allocators, and optimization