

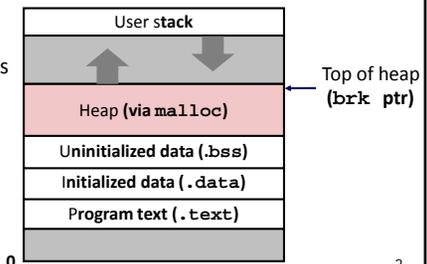
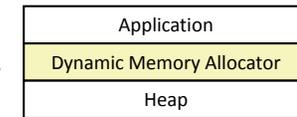
Dynamic Memory Allocation: Basic Concepts

Kai Shen

1

Dynamic Memory Allocation

- Programmers use *dynamic memory allocators* (such as `malloc`) to acquire memory at run time.
 - For data structures whose size is only known at runtime.
- Dynamic memory allocators manage an area of process virtual memory known as the *heap*.



2

Dynamic Memory Allocation

- Allocator maintains heap as collection of variable sized *blocks*, which are either *allocated* or *free*
- Types of allocators
 - Explicit allocator:** application allocates and frees space
 - E.g., `malloc` and `free` in C
 - Implicit allocator:** application allocates, but does not free space
 - E.g. garbage collection in Java, ML, and Lisp
- We discuss explicit memory allocation today

3

The malloc Package

```
#include <stdlib.h>

void *malloc(size_t size)
    Successful:
        Returns a pointer to a memory block of at least size bytes
        (typically) aligned to 8-byte boundary
        If size=0, returns NULL
    Unsuccessful: returns NULL (0) and sets errno

void free(void *p)
    Returns the block pointed at by p to pool of available memory
    p must come from a previous call to malloc or realloc
```

4

Allocation Example

The diagram illustrates a sequence of memory allocation and deallocation operations on a 20-byte heap:

- `p1 = malloc(4)`: Allocates 4 bytes (yellow).
- `p2 = malloc(5)`: Allocates 5 bytes (yellow).
- `p3 = malloc(6)`: Allocates 6 bytes (yellow).
- `free(p2)`: Frees the 5 bytes allocated to p2 (yellow).
- `p4 = malloc(2)`: Allocates 2 bytes (green).

Final state: p1 (4 bytes), p4 (2 bytes), and p3 (6 bytes) are allocated. The freed memory (5 bytes) is now free.

5

Constraints

- Applications
 - Can issue arbitrary sequence of `malloc` and `free` requests
 - `free` request must be to a `malloc`'d block
- Allocators
 - Must respond immediately to `malloc` requests
 - Can manipulate and modify only free memory
 - Must allocate blocks from free memory
 - Must align blocks so they satisfy all alignment requirements
 - 8 byte alignment for GNU `malloc` (`libc malloc`) on Linux boxes
 - Can't move the allocated blocks once they are `malloc`'d
 - i.e.*, compaction is not allowed

6

Interaction with OS

- Change heap size
 - `brk` system call

The diagram shows the memory layout from top to bottom:

- User stack (grows down)
- Heap (via `malloc`) (grows up, bounded by `brk ptr`)
- Uninitialized data (`.bss`)
- Initialized data (`.data`)
- Program text (`.text`)

7

Performance Goal: Runtime Speed

- Given some sequence of `malloc` and `free` requests:
 - $R_0, R_1, \dots, R_k, \dots, R_{n-1}$
- Speed – throughput:
 - Number of completed requests per unit time
 - Example:
 - 5,000 `malloc` calls and 5,000 `free` calls in 10 seconds
 - Throughput is 1,000 operations/second

8

Structure of Allocated Block

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

`p0 = malloc(4)`

`free(p0)`

13

Free Block Management

- How to keep track of free blocks?
- Allocate from free blocks:
 - How do we pick a block to use for allocation – 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 freed block?

14

Keeping Track of Free Blocks

- Method 1: *Implicit list* using length—links all blocks
- Method 2: *Explicit list* among the free blocks using pointers
- Method 3: *Segregated free list*
 - Different free lists for different size classes
- Method 4: *Blocks sorted by size*
 - Can use a balanced tree (e.g. Red-Black tree) with pointers within each free block, and the length used as a key

15

Method 1: Implicit Free Block List

- For each block we need both size and allocation status
 - Could store this information in two words: wasteful!
- Standard trick
 - If blocks are aligned, some low-order address bits are always 0
 - Instead of storing an always-0 bit, use it as a allocated/free flag
 - When reading size word, must mask out this bit

Format of allocated and free blocks

1 word

Size a

a = 1: Allocated block
a = 0: Free block

Size: block size

Payload: application data (allocated blocks only)

Optional padding

16

Detailed Implicit Free List Example

Start of heap

Unused

8/0 16/1 32/0 16/1

Double-word aligned

Allocated blocks: shaded
Free blocks: unshaded
Headers: labeled with size in bytes/allocated bit

17

Implicit List: Finding a Free Block

- **First fit:**
 - Search list from beginning, choose **first** free block that fits:
 - Can take linear time in total number of blocks (allocated and free)
 - It may 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: fits, with fewest bytes left over
 - Keeps fragments small—usually helps fragmentation
 - Will typically run slower than first fit

18

Implicit List: Allocating in Free Block

- Allocating in a free block: **splitting**
 - Since allocated space might be smaller than free space, we might want to split the block

19

Implicit List: Freeing a Block

- Simplest implementation:
 - Need only clear the "allocated" flag

```
void free_block(ptr p) { *p = *p & ~2 }
```
- But can lead to "false fragmentation"


```
free(p)
```

```
malloc(5) Oops!
```

There is enough free space, but the allocator won't be able to find it

20

Implicit List: Coalescing

- Join (*coalesce*) with next/previous blocks, if they are free
 - Coalescing with next block

free(p)

- But how do we coalesce with *previous* block?

21

Implicit List: Bidirectional Coalescing

- Boundary tags** [Knuth73]
 - Replicate size/allocated word at "bottom" (end) of free blocks
 - Allows us to traverse the "list" backwards, but requires extra space
 - Important and general technique!

Header →

Size	a
Payload and padding	
Boundary tag (footer)	Size a

a = 1: Allocated block
a = 0: Free block

Size: Total block size
Payload: Application data (allocated blocks only)

Format of allocated and free blocks

22

Constant Time Coalescing

Block being freed →

Case 1	Case 2	Case 3	Case 4
Allocated	Allocated	Free	Free
Free	Free	Allocated	Free
Allocated	Free	Free	Free

23

Constant Time Coalescing (Case 1)

m1	1
m1	1
n	1
n	1
m2	1
m2	1

→

m1	1
m1	1
n	0
n	0
m2	1
m2	1

24

Constant Time Coalescing (Case 2)

m1	1
m1	1
n	1
n	1
m2	0
m2	0

m1	1
m1	1
n+m2	0
n+m2	0

25

Constant Time Coalescing (Case 3)

m1	0
m1	0
n	1
n	1
m2	1
m2	1

n+m1	0
n+m1	0
m2	1
m2	1

26

Constant Time Coalescing (Case 4)

m1	0
m1	0
n	1
n	1
m2	0
m2	0

n+m1+m2	0
n+m1+m2	0

27

Space use of footers

- Additional space consumption
- Can it be optimized?
 - Footer tag is free in free blocks
 - No need to coalesce in allocated blocks

28



Implicit Free Block Lists: Summary

- Implementation: very simple
- Allocate cost:
 - linear time worst case
- Free cost:
 - constant time worst case
 - even with coalescing
- Memory usage:
 - Depend on placement policy: first-fit, next-fit or best-fit

- Not used in general allocator because of linear-time allocation
- Segregated free lists (next lecture) approximate a best fit placement policy without having to search entire free list

29



Disclaimer

These slides were adapted from the CMU course slides provided along with the textbook of "Computer Systems: A programmer's Perspective" by Bryant and O'Hallaron. The slides are intended for the sole purpose of teaching the computer organization course at the University of Rochester.

30