

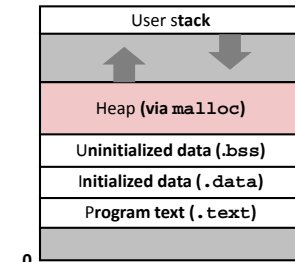
Dynamic Memory Allocation: Basic Concepts

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 start our discussion on explicit memory allocation

3

The malloc Package

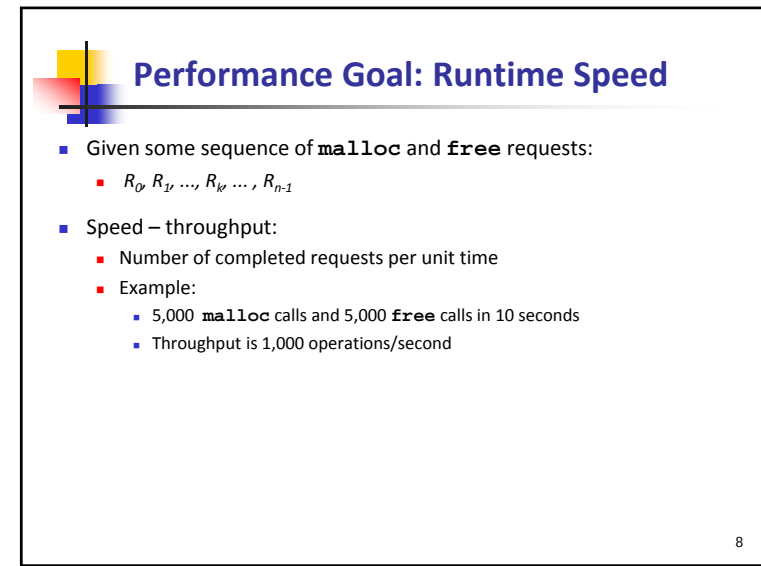
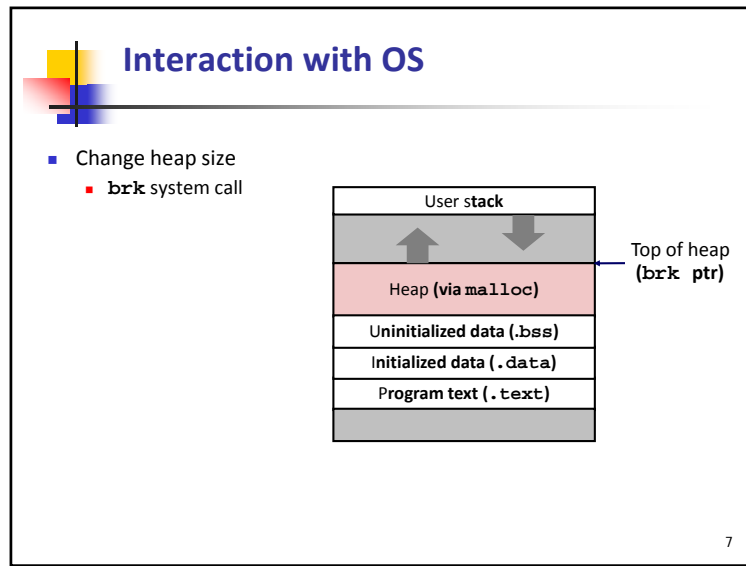
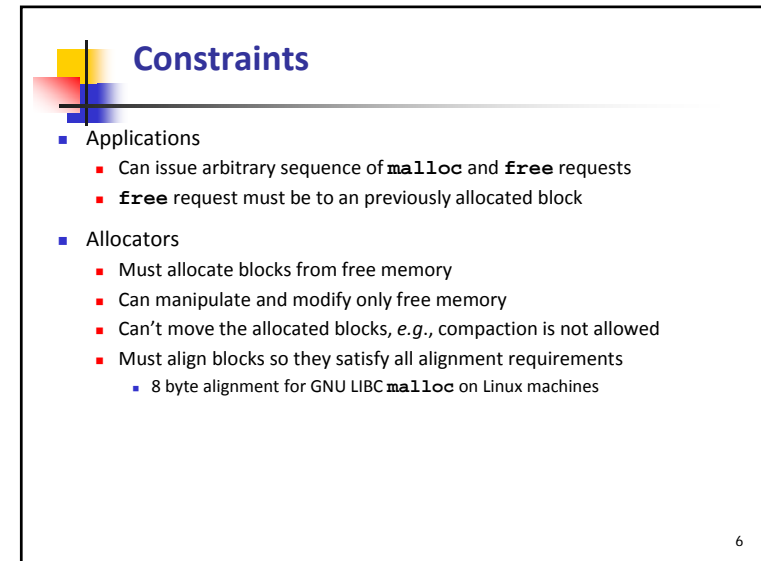
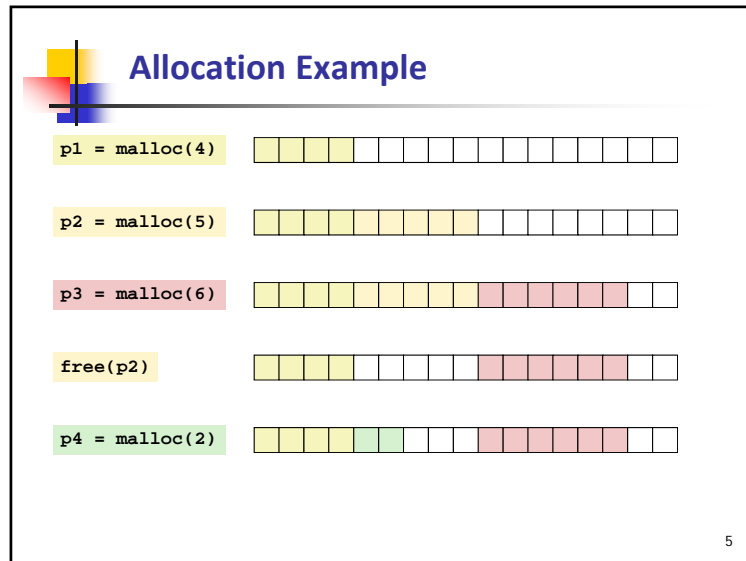
```
#include <stdlib.h>

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 memory allocation call like malloc

    Part of the C library, occasionally invokes a system call to change the
    heap size
```

4



Performance Goal: Memory Utilization

- Given some sequence of **malloc** and **free** requests:
 - $R_0, R_1, \dots, R_k, \dots, R_{n-1}$
- Def:** Aggregate payload P_k
 - malloc**(p) results in a block with a **payload** of p bytes
 - After request R_k has completed, the **aggregate payload** P_k is the sum of currently allocated payloads
- Def:** Heap size H_k
 - Maximum heap size so far
- Def:** Peak memory utilization after k requests
 - $U_k = (\max_{i \leq k} P_i) / H_k$
- Goals: achieve high runtime speed and memory utilization
 - Somewhat conflicting

9

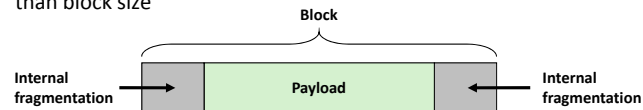
Fragmentation

- Poor memory utilization caused by **fragmentation**
 - internal** fragmentation
 - external** fragmentation

10

Internal Fragmentation

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

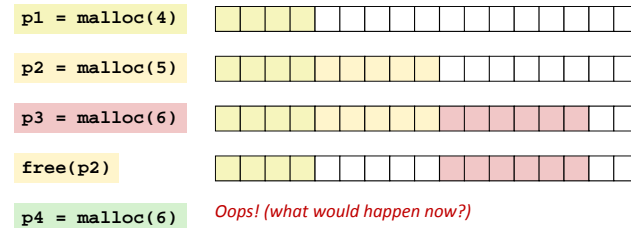


- Caused by
 - Overhead of maintaining heap data structures
 - Padding for alignment purposes
 - ...

11

External Fragmentation

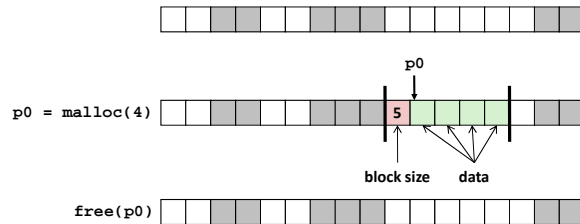
- Occurs when there is enough aggregate heap memory, but no single free block is large enough



12

Structure of Allocated Block

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



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?
- How do we reinsert freed block?

14

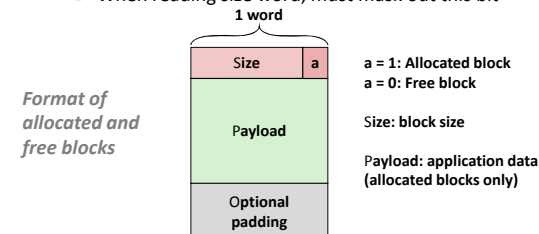
Keeping Track of Free Blocks

- 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

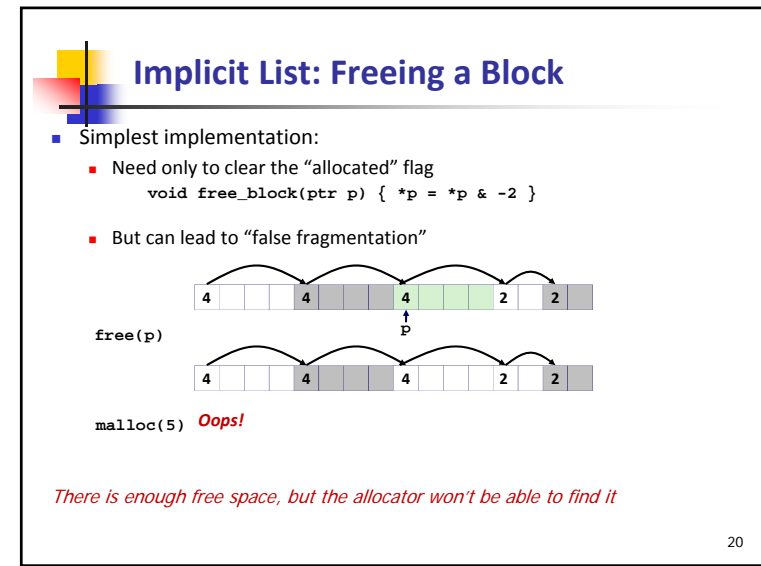
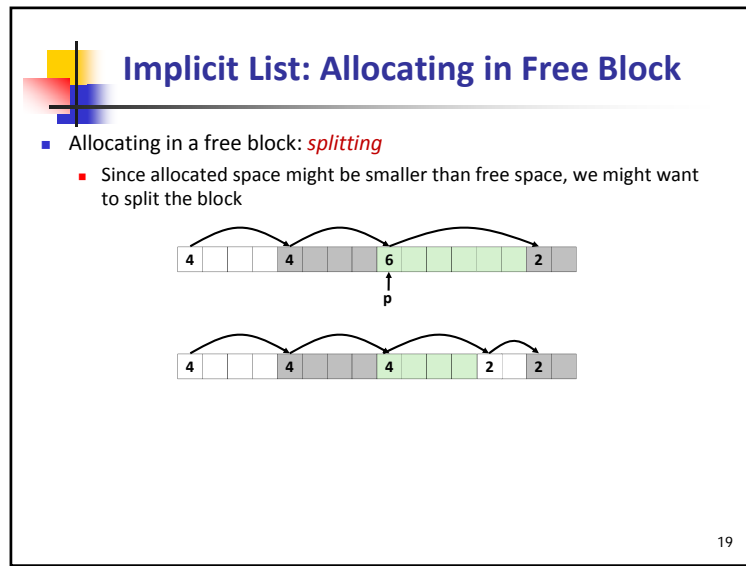
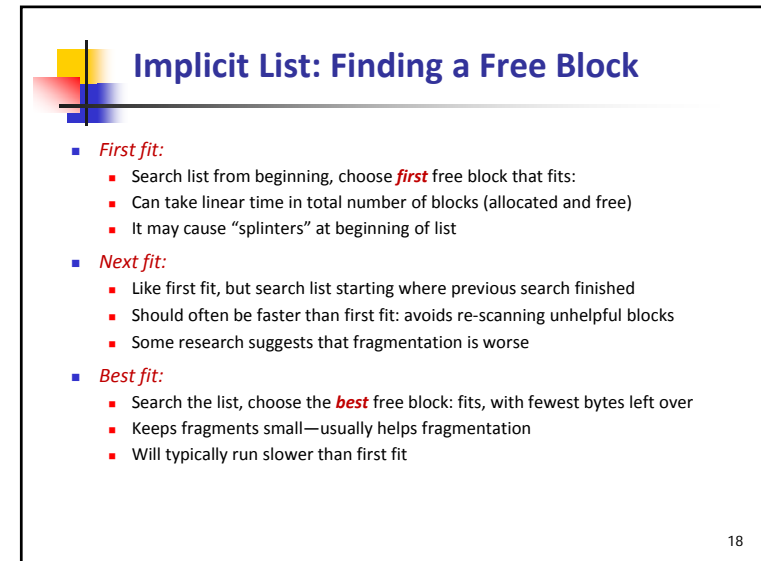
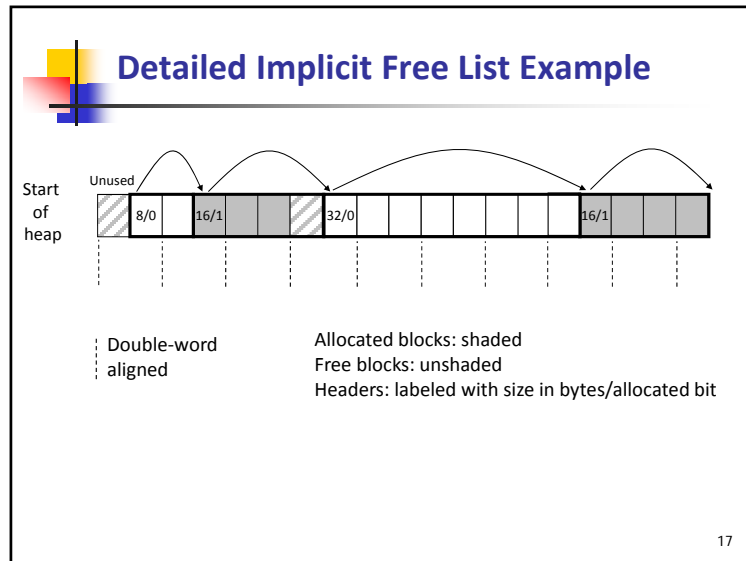
15

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 bits in block size (e.g., 3 bits for 8-byte alignment) are always 0
 - Instead of storing an always-0 bit, use it as an allocated/free flag
 - When reading size word, must mask out this bit



16



Implicit List: Coalescing

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

free(p)

logically gone

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

21

Implicit List: Bidirectional Coalescing

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

Header

Format of allocated and free blocks

Boundary tag (footer)

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

Size: Total block size

Payload: Application data (allocated blocks only)

22

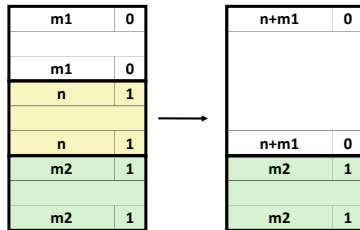
Constant Time Coalescing (Case 1)

23

Constant Time Coalescing (Case 2)

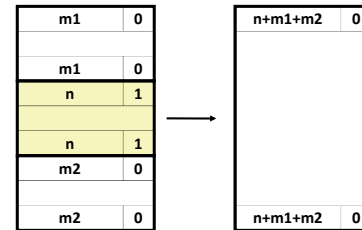
24

Constant Time Coalescing (Case 3)



25

Constant Time Coalescing (Case 4)



26

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
- Solution to eliminate the footer's space cost:
 - Footer only in free blocks
 - Use another bit in each allocated block's header (remember that we have 3 free bits) to indicate whether the preceding block is allocated or free
 - But require more maintenance
 - Every free must update next allocated block's preceding-block-allocation bit
 - Every allocation must know whether the preceding block is free or allocated

27

Implicit Free Block Lists: Summary

- Implementation: very simple
- Free cost:
 - constant time worst case
 - even with coalescing
- Allocation cost:
 - linear time worst case
- Memory usage:
 - Depend on placement policy: first-fit, next-fit or best-fit
- Not used in general allocator because of slow (linear-time) allocation

28



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.

29