Memory Allocation I

CSE 351 Summer 2024

Instructor: Ellis Haker

Teaching Assistants:

Naama Amiel Micah Chang Shananda Dokka Nikolas McNamee Jiawei Huang My C program exit without freeing allocated memory OS:



Administrivia

- Today
 - HW17 due (11:59pm)
- Friday, 8/2
 - RD20 due (1pm)
 - \circ No HW :)
 - Lab 5 released (due 8/14)
- Monday, 8/5
 - RD21 due (1pm)
 - HW19 due (11:59pm)

Topic Group 3: Scale & Coherence

- How do we make memory accesses faster?
- How do programs manage large amounts of memory?
 - How can we allocate memory dynamically (i.e. at runtime)
- How does your computer run multiple programs at once?



Dynamic Memory Allocation

- Overview
 - Introduction and goals
 - Allocation and deallocation (free)
- Malloc and free in C
 - Common memory-related bugs in C
- Fragmentation
- Explicit allocation implementation
 - Implicit free lists
 - Explicit free lists (Lab 5)
 - Splitting and coalescing
- Implicit deallocation: garbage collection

Multiple Ways to Store Program Data



Dynamic Memory Allocation

- Programmers use **dynamic memory allocators** to acquire memory at runtime
 - For data structures whose size or lifetime is known only at runtime
 - Stored in heap segment of memory
- Types of allocators
 - Implicit: language handles garbage collection, programmer only needs to allocate the data
 - <u>Ex</u>: new in Java
 - Explicit: programmer needs to allocate *and* deallocate memory
 - <u>Ex</u>: malloc and free in C

Dynamic Memory Allocation (pt 2)

• Allocator organizes data into variable-sized **blocks**, which can be allocated or free

- What happens when the heap runs out of space?
 - Allocator asks the OS for more!
 - sbrk in Unix



Allocating Memory in C

- Need to #include <stdlib.h> \bullet
- void* malloc(size_t size),
- colles returned from multiple colles to malloc are 1 not grasanteed to be contiguous with each other Allocates a **contiguous** block of size bytes of **uninitialized** memory Ο
 - In C, NULL always size_t?! Simple typedef for an unsigned 8-byte integer Ο
 - Returns a pointer to the beginning of the allocated block Ο
 - Returns NULL if allocation failed, or size=0
 - Pointer to an invalid address (represented as address 0) in x 84, wh true in all
 - Blocks typically aligned to 8 or 16 bytes Ο
- Other versions:
 - calloc: initializes memory to 0 Ο
 - realloc: moves existing block to a larger one Ο

Trouslates to 0

ISAS

Freeing Memory in C

- void free(void* p)
 - Releases block pointed to by p back to the pool of available memory
 - Pointer p must be the address originally returned by malloc/calloc/realloc (i.e., beginning of the block)
 - Don't call free on a block that has already been released!
 - Undefined behavior can even introduce buffer overflow vulnerabilities!
 - No action occurs if you call free(NULL)

Best Practices for malloc and free

- malloc
 - Using sizeof makes code more portable (ints aren't 4B on all machines...)
 - void* is implicitly cast, but explicitly casting will help you catch errors
 - o Ex: int* ptr = (int*) malloc(n * sizeof(int));

• free

- After calling free, set the pointer to NULL
 - Avoids double-free errors

```
free(ptr);
ptr = NULL;
```

malloc and free Example

```
int i, *p;

week p = (int*) malloc(n*sizeof(int)); // allocate space for array of n ints

for if (p == NULL) {

    perror("malloc");    print error

    exit(0);
    wessege

         for (i=0; i<n; i++)
                                                                                 // initialize int array
                p[i] = i;
         free(p);
                                                                                             // free p
                                                                                              // good practice to
         p = NULL;
   set to NULL after free
```

malloc and free Interface

- Applications
 - Must never access memory not currently allocated
 - Must never free memory not currently allocated
 - Must only use free with previously malloc'ed blocks
- Allocators
 - Must respond immediately to malloc
 - Must allocate blocks from <u>free</u> memory
 - Must align blocks so they satisfy all alignment requirements
 - Can't move the allocated blocks

Heap Management in C

- Likely very different from what you're used to
- Programmer has to remember to free data when they're done with it
 - Otherwise, causes a memory leak
- Requires keeping track of where your data is stored
- Bugs are common!



Find that Bug!

heef

Find that Bug! (pt 2)

```
x = (int*)malloc(N * sizeof(int));
// manipulate x
free(x);
...
y = (int*)malloc(N * sizeof(int));
// manipulate y
free(x);
```

Find that Bug! (pt 3)

Une after free

Find that Bug! (pt 4)

```
allocate
typedef struct L {
    int val;
    struct L* next;
} list; // linked list node
                                                        torgot
void main() {
    list* head = (list*) malloc(sizeof(list));
                                                       free
    head->val = 0;
    head->next = NULL;
    // create and manipulate the rest of the list
    . . .
                           call vigloc more
   _free(head);
                        times to make more
vokes (not shown)
    return:
```

Dynamic Memory Allocation

- Overview
 - Introduction and goals
 - Allocation and deallocation (free)
- Malloc and free in C
 - Common memory-related bugs in C
- Fragmentation
- Explicit allocation implementation
 - Implicit free lists
 - Explicit free lists (Lab 5)
 - Splitting and coalescing
- Implicit deallocation: garbage collection

Notation

- We will draw memory divided into words
 - \circ Each word is 64 bit = 8 bytes
 - Note: textbook and old videos still use 32-bit words
- Allocations will be aligned and in sizes that are a multiple of words



Memory Allocation Example



Performance Goals

- Given some sequence of malloc and free requests $R_0, R_1, ..., R_k, ..., R_{n-1}$, maximize throughput and peak memory utilization
 - Often conflicting goals!

1) Throughput:

- Number of completed requests per some unit of time
- Example:
 - If 5,000 malloc calls and 5,000 free calls completed in 10 seconds, then throughput is 1,000 operations/second

Performance Goals (pt 2)

- Aggregate Payload (P_k)
 - malloc(p) returns a payload of p bytes
 - After request R_k has completed, P_k = the sum of currently allocated payloads
- Current Heap Size (*H_k*)
 - Allocator can increase heap size using sbrk
 - Assume allocator cannot decrease heap size
- 2) **Peak Memory Utilization** (U_k) :
 - Memory utilization = aggregate payload ÷ heap size
 - i.e., how much of our heap contains payload data
 - $U_k = (\max_{i \le k} P_k) \div H_k$ after *k*+1 requests
 - i.e., the maximum utilization at any point up through when R_k has completed



Fragmentation

- Recall fragmentation in structs
 - Used to preserve alignment
 - Internal: unused space inside of a struct (between fields)
 - **External**: unused space *between* struct instances (after fields)

- The heap is similar
 - Internal fragmentation: extra space inside an allocated block
 - External fragmentation: extra space between allocated blocks
 Con he for alignment, or other reasons

Internal Fragmentation

• Additional space inside of an allocated block not used to store a payload



• Causes

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

External Fragmentation

- Occurs when allocation/free pattern leaves "holes" between blocks
 - Can cause situations where there is enough aggregate heap space to satisfy a request, but no single free block is large enough



Polling Question



- A) Temporary arrays should not be allocated on the Heap
- B) malloc returns an address of a block that is filled with mystery data
 - Peak memory utilization is a measure of both internal and external fragmentation
 -) An allocation failure will cause your program to stop

```
False,
returns NULL
on failure
```

Utilization = poul

uselly Letter te stock I put ou the stock

Dynamic Memory Allocation

- Overview
 - Introduction and goals
 - Allocation and deallocation (free)
- Malloc and free in C
 - Common memory-related bugs in C
- Fragmentation
- Explicit allocation implementation
 - Implicit free lists
 - Explicit free lists (Lab 5)
 - Splitting and coalescing
- Implicit deallocation: garbage collection

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 reclaim free blocks?

Knowing how much to free

- Store the length of the block in a header
 - The word preceding the data
 - Also may contain other info, like whether the block is allocated



Keeping track of free blocks

Implicit Free List

Use pointer arithmetic to traverse through entire heap until we find a free bock Ο pros' sover memory



Explicit Free List

Free block stores pointer to the next free block, forming a linked list Ο



- Others not covered in this class
 - Segregated free lists (different lists for each object type), sorting blocks by size Ο

coust finding free block

Implicit Free Lists (i.e. no pointers)

- plicit Free Lists (i.e. no pointers)
 For each block, we need to store: size, is_allocated
 O Could use two works, but kinda wasteful... •
- Recap: if size is a multiple of 2^n , then lowest *n* bits of the size are always 0
 - Use the lowest bit of header word to store **is_allocated** flag
 - When reading **size**, mask this bit out

Block format:	size a	a = 1 if allocated, 0 if free	If the header value is h:
	payload	size = total block size in bytes	$h = size a$ $a = h \& 1$ $size = b \& \sim 1$
	optional padding	allocated blocks only	=061116, US out Lovest 614

Header Questions

 How many "flags" (boolean values) could we fit in our header if our allocator uses 16-byte alignment?

16 = 2", lowest 4 hits will be 0

If we placed a new flag in the second least significant bit, write out a C expression that will extract this new flag from the header!
 2=0,00 ... 010

Summary

- The heap is a segment of memory used to dynamically allocate data
 - Useful when we don't know the size or when it can be freed until runtime
- **Fragmentation** is space in the heap that is not used to store payloads
 - Internal: inside of a block
 - External: between blocks
- C uses an **explicit allocator**, meaning the programmer decides when heap data is freed
- Different heap implementations
 - Implicit free list: only store header and payload
 - Explicit free list: free blocks store pointers to the next free block (next lecture)