

TCSS 422: OPERATING SYSTEMS

Free Space Management

Wes J. Lloyd
Institute of Technology
University of Washington - Tacoma



OBJECTIVES

- Chapter 17
 - Fragmentation
 - Free List
 - Memory header
 - Free list operations

February 22, 2017

TCSS422: Operating Systems [Winter 2017]
Institute of Technology, University of Washington - Tacoma

L14.2

FREE SPACE MANAGEMENT

- Management of memory using
 - Only fixed-sized units
 - Easy: keep a list
 - Memory request → return first free entry
 - Simple search
 - With variable sized units
 - More challenging
 - Results from variable sized malloc requests
 - Leads to fragmentation

February 22, 2017

TCSS422: Operating Systems [Winter 2017]
Institute of Technology, University of Washington - Tacoma

L14.3

FRAGMENTATION

- Consider a 30-byte heap

30-byte heap:

free	used	free	
0	10	20	30
- Request for 15-bytes

free list: head →

addr:0
len:10

addr:20
len:10

→ NULL
- Free space: 20 bytes
- No available contiguous chunk → return NULL

February 22, 2017

TCSS422: Operating Systems [Winter 2017]
Institute of Technology, University of Washington - Tacoma

L14.4

FRAGMENTATION - 2

- External:** we can compact
 - Example: Client asks for 100 bytes: malloc(100)
 - OS: No 100byte contiguous chunk(s) available: returns NULL
 - Memory is externally fragmented - - Compaction can fix!
- Internal:** lost space – can't compact
 - OS returns memory units that are too large
 - Example: Client asks for 100 bytes: malloc(100)
 - OS: Returns 125 byte chunk
 - Fragmentation is *in* the allocated chunk
 - Memory is lost, and unaccounted for – can't compact

February 22, 2017

TCSS422: Operating Systems [Winter 2017]
Institute of Technology, University of Washington - Tacoma

L14.5

ALLOCATION STRATEGY: SPLITTING

- Request for 1 byte of memory: malloc(1)

30-byte heap:

free	used	free	
0	10	20	30

free list: head →

addr:0
len:10

addr:20
len:10

→ NULL
- OS locates a free chunk to satisfy request
- Splits chunk into two, returns first chunk

30-byte heap:

free	used	free	
0	10	21	30

free list: head →

addr:0
len:10

addr:21
len:9

→ NULL

February 22, 2017

TCSS422: Operating Systems [Winter 2017]
Institute of Technology, University of Washington - Tacoma

L14.6

ALLOCATION STRATEGY: COALESCING

- Consider 30-byte heap
- Free() frees all 10 bytes segments (list of 3-free 10-byte chunks)



- Request arrives: malloc(30)
- No contiguous 30-byte chunk exists
- Coalescing regroups chunks into contiguous chunk



- Allocation can now proceed

February 22, 2017

TCS5422: Operating Systems [Winter 2017]
Institute of Technology, University of Washington - Tacoma

L14.7

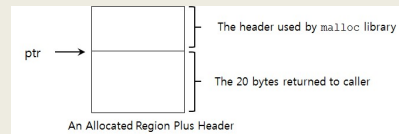
MEMORY HEADERS

- free(void *ptr): Does not require a size parameter

- How does the OS know how much memory to free?

- Header block

- Small descriptive block of memory at start of chunk

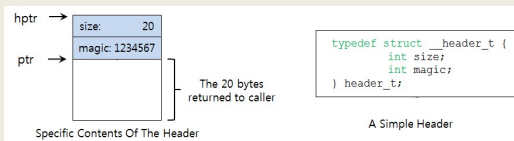


February 22, 2017

TCS5422: Operating Systems [Winter 2017]
Institute of Technology, University of Washington - Tacoma

L14.8

MEMORY HEADERS - 2



- Contains size
- Pointers: for faster memory access
- Magic number: integrity checking

February 22, 2017

TCS5422: Operating Systems [Winter 2017]
Institute of Technology, University of Washington - Tacoma

L14.9

MEMORY HEADERS - 3

- Size of memory chunk is:
- Header size + user malloc size
- N bytes + sizeof(header)

- Easy to determine address of header

```
void free(void *ptr) {
    header_t *hptr = (void *)ptr - sizeof(header_t);
}
```

February 22, 2017

TCS5422: Operating Systems [Winter 2017]
Institute of Technology, University of Washington - Tacoma

L14.10

THE FREE LIST

- Simple free list struct

```
typedef struct __node_t {
    int size;
    struct __node_t *next;
} node_t;
```

- Use mmap to create free list
- 4kb heap, 4 byte header, one contiguous free chunk

```
// mmap() returns a pointer to a chunk of free space
node_t *head = mmap(NULL, 4096, PROT_READ|PROT_WRITE,
    MAP_ANON|MAP_PRIVATE, -1, 0);
head->size = 4096 - sizeof(node_t);
head->next = NULL;
```

February 22, 2017

TCS5422: Operating Systems [Winter 2017]
Institute of Technology, University of Washington - Tacoma

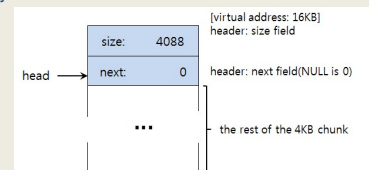
L14.11

FREE LIST - 2

- Create and initialize free-list "heap"

```
// mmap() returns a pointer to a chunk of free space
node_t *head = mmap(NULL, 4096, PROT_READ|PROT_WRITE,
    MAP_ANON|MAP_PRIVATE, -1, 0);
head->size = 4096 - sizeof(node_t);
head->next = NULL;
```

- Heap layout:



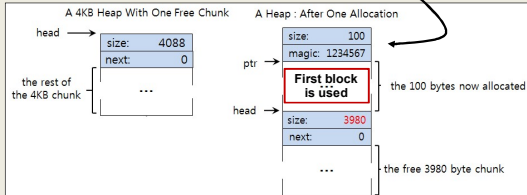
February 22, 2017

TCS5422: Operating Systems [Winter 2017]
Institute of Technology, University of Washington - Tacoma

L14.12

FREE LIST: MALLOC() CALL

- Consider a request for a 100 bytes: `malloc(100)`
- Header block requires 8 bytes
 - 4 bytes for size, 4 bytes for magic number
- Split the heap - header goes with each block



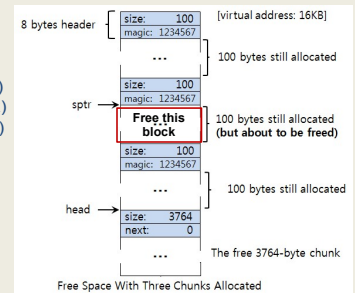
February 22, 2017

TCCS422: Operating Systems [Winter 2017]
Institute of Technology, University of Washington - Tacoma

L14.13

FREE LIST: FREE() CALL

- Addresses of chunks
- Start=16384
+ 108 (end of 1st chunk)
+ 108 (end of 2nd chunk)
+ 108 (end of 3rd chunk)
= 16708



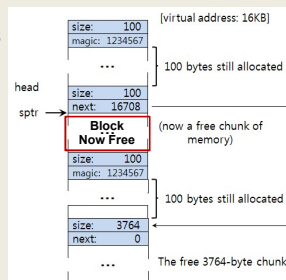
February 22, 2017

TCCS422: Operating Systems [Winter 2017]
Institute of Technology, University of Washington - Tacoma

L14.14

FREE LIST: FREE() CHUNK #2

- `Free(sptr)`
- Our 3 chunks start at 16 KB (@ 16,384 bytes)
- Free chunk #2 - `sptr`
- `Sptr = 16500`
 - `addr - sizeof(node_t)`
- Actual start of chunk #2
 - 16492



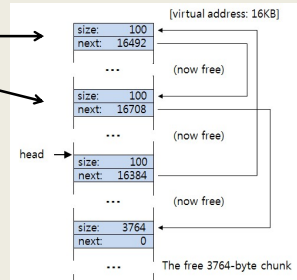
February 22, 2017

TCCS422: Operating Systems [Winter 2017]
Institute of Technology, University of Washington - Tacoma

L14.15

FREE LIST- FREE ALL CHUNKS

- Now free remaining chunks:
 - `Free(16392)`
 - `Free(16608)`
- Walk back 8 bytes for actual start of chunk
- External fragmentation
- Free chunk pointers out of order
- Coalescing of next pointers is needed



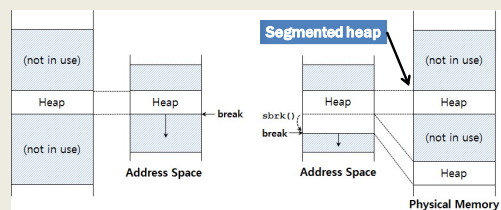
February 22, 2017

TCCS422: Operating Systems [Winter 2017]
Institute of Technology, University of Washington - Tacoma

L14.16

GROWING THE HEAP

- Start with small sized heap
- Request more memory when full
- `sbrk()`, `brk()`



February 22, 2017

TCCS422: Operating Systems [Winter 2017]
Institute of Technology, University of Washington - Tacoma

L14.17

MEMORY ALLOCATION STRATEGIES

- Best fit**
 - Traverse free list
 - Identify all candidate free chunks
 - Note which is smallest (has best fit)
 - When splitting, "leftover" pieces are small (and potentially less useful -- fragmented)
- Worst fit**
 - Traverse free list
 - Identify largest free chunk
 - Split largest free chunk, leaving a still large free chunk

February 22, 2017

TCCS422: Operating Systems [Winter 2017]
Institute of Technology, University of Washington - Tacoma

L14.18

EXAMPLES

- Allocation request for 15 bytes

head → 10 → 30 → 20 → NULL

- Result of Best Fit

head → 10 → 30 → 5 → NULL

- Result of Worst Fit

head → 10 → 15 → 20 → NULL

February 22, 2017

TCCS422: Operating Systems [Winter 2017]
Institute of Technology, University of Washington - Tacoma

L14.19

MEMORY ALLOCATION STRATEGIES - 2

- First fit**
 - Start search at beginning of free list
 - Find first chunk large enough for request
 - Split chunk, returning a "fit" chunk, saving the remainder
 - Avoids full free list traversal of best and worst fit
- Next fit**
 - Similar to first fit, but start search at last search location
 - Maintain a pointer that "cycles" through the list
 - Helps balance chunk distribution vs. first fit
 - Find first chunk, that is large enough for the request, and split
 - Avoids full free list traversal

February 22, 2017

TCCS422: Operating Systems [Winter 2017]
Institute of Technology, University of Washington - Tacoma

L14.20

SEGREGATED LISTS

- For popular sized requests
e.g. for kernel objects such as locks, inodes, etc.
- Manage as segregated free lists
- Provide object caches: stores pre-initialized objects

- How much memory should be dedicated for specialized requests (object caches)?
- If a given cache is low in memory, can request "slabs" of memory from the general allocator for caches.
- General allocator will reclaim slabs when not used

February 22, 2017

TCCS422: Operating Systems [Winter 2017]
Institute of Technology, University of Washington - Tacoma

L14.21

BUDDY ALLOCATION

- Binary buddy allocation**
 - Divides free space by two to find a block that is big enough to accommodate the request; the next split is too small...
- Consider a 7KB request

64 KB

32 KB32 KB

16 KB16 KB

8 KB8 KB

64KB free space for 7KB request

February 22, 2017

TCCS422: Operating Systems [Winter 2017]
Institute of Technology, University of Washington - Tacoma

L14.22

BUDDY ALLOCATION - 2

- Buddy allocation: suffers from internal fragmentation
- Allocated fragments, typically too large
- Coalescing is simple
 - Two adjacent blocks are promoted up

February 22, 2017

TCCS422: Operating Systems [Winter 2017]
Institute of Technology, University of Washington - Tacoma

L14.23

INTRODUCTION TO PAGING

February 22, 2017

TCCS422: Operating Systems [Winter 2017]
Institute of Technology, University of Washington - Tacoma

L15.24

OBJECTIVES

- Chapter 18
 - Paging
 - Address translation
 - Paging questions

February 22, 2017

TCCS422: Operating Systems [Winter 2017]
Institute of Technology, University of Washington - Tacoma

L15.25

PAGING

- Split up address space of process into *fixed sized pieces* called **pages**
- Alternative to *variable sized pieces* (Segmentation) which suffers from significant fragmentation
- Physical memory is split up into an array of fixed-size slots called **page frames**.
- Each process has a **page table** which translates virtual addresses to physical addresses

February 22, 2017

TCCS422: Operating Systems [Winter 2017]
Institute of Technology, University of Washington - Tacoma

L15.26

ADVANTAGES OF PAGING

- Flexibility
 - Abstracts the process address space into pages
 - No need to track direction of HEAP / STACK growth
 - No need to store unused space
- Simplicity
 - Pages and page frames are the same size
 - Easy to allocate and keep a free list of pages

February 22, 2017

TCCS422: Operating Systems [Winter 2017]
Institute of Technology, University of Washington - Tacoma

L15.27

PAGING: EXAMPLE

Page Table:
VP0 → PF3
VP1 → PF7
VP2 → PF5
VP3 → PF2

- Consider a 128 byte address space with 16-byte pages
- Consider a 64-byte program address space

A Simple 64-byte Address Space

64-Byte Address Space Placed In Physical Memory

February 22, 2017

TCCS422: Operating Systems [Winter 2017]
Institute of Technology, University of Washington - Tacoma

L15.28

PAGING: ADDRESS TRANSLATION

- PAGE: Has two address components
 - VPN: Virtual Page Number
 - Offset: Offset within a Page

- Example:
Page Size: 16-bytes, Address Space: 64-bytes

Here there are just four pages...

February 22, 2017

TCCS422: Operating Systems [Winter 2017]
Institute of Technology, University of Washington - Tacoma

L15.29

EXAMPLE:
PAGING ADDRESS TRANSLATION

- Consider a 64-byte program address space (4 pages)
- Stored in 128-byte physical memory (8 frames)
- Offset is preserved
- VPN is looked up

Page Table:
VP0 → PF3
VP1 → PF7
VP2 → PF5
VP3 → PF2

February 22, 2017

TCCS422: Operating Systems [Winter 2017]
Institute of Technology, University of Washington - Tacoma

L15.30

PAGING DESIGN QUESTIONS

- Where are page tables stored?
- What are the typical contents of the page table?
- How big are page tables?
- Does paging make the system too slow?

February 22, 2017

TCCS422: Operating Systems [Winter 2017]
Institute of Technology, University of Washington - Tacoma

L15.31

WHERE ARE PAGE TABLES STORED?

- Example:
 - Consider a 32-bit process address space (up to 4GB)
 - With 4 KB pages
 - 20 bits for VPN (2^{20} pages)
 - 12 bits for the page offset (2^{12} unique bytes in a page)
- Page tables for each process are stored in RAM
 - Support potential storage of 2^{20} translations = 1,048,576 pages per process
 - Each page has a page table entry size of 4 bytes

February 22, 2017

TCCS422: Operating Systems [Winter 2017]
Institute of Technology, University of Washington - Tacoma

L15.32

PAGE TABLE EXAMPLE

- With 2^{20} slots in our page table for a single process
- Each slot dereferences a VPN
- Provides physical frame number
- Each slot requires 4 bytes (32 bits)
 - 20 for the PFN on a 4GB system with 4KB pages
 - 12 for the offset which is preserved
 - (note we have no status bits, so this is unrealistically small)
- How much memory to store page table for 1 process?
 - 4,194,304 bytes (or 4MB) to index one process

VPN ₀
VPN ₁
VPN ₂
...
...
VPN ₁₀₄₈₅₇₆

February 22, 2017

TCCS422: Operating Systems [Winter 2017]
Institute of Technology, University of Washington - Tacoma

L15.33

NOW FOR AN ENTIRE OS

- If 4 MB is required to store one process
- Consider how much memory is required for an entire OS?
 - With for example 100 processes...
- Page table memory requirement is now 4MB x 100 = 400MB
- If computer has 4GB memory (maximum for 32-bits), the page table consumes 10% of memory

400 MB / 4000 GB

- Is this efficient?

February 22, 2017

TCCS422: Operating Systems [Winter 2017]
Institute of Technology, University of Washington - Tacoma

L15.34

WHAT'S ACTUALLY IN THE PAGE TABLE

- Page table is data structure used to map virtual page numbers (VPN) to the physical address (Physical Frame Number PFN)
 - Linear page table → simple array
- Page-table entry
 - 32 bits for capturing state

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0
PFN																				U	S	D	A	P	R	W					

An x86 Page Table Entry(PTE)

February 22, 2017

TCCS422: Operating Systems [Winter 2017]
Institute of Technology, University of Washington - Tacoma

L15.35

PAGE TABLE ENTRY

- P: present
- R/W: read/write bit
- U/S: supervisor
- A: accessed bit
- D: dirty bit
- PFN: the page frame number

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0
PFN																				U	S	D	A	P	R	W					

An x86 Page Table Entry(PTE)

February 22, 2017

TCCS422: Operating Systems [Winter 2017]
Institute of Technology, University of Washington - Tacoma

L15.36

PAGE TABLE ENTRY - 2

- Common flags:
- Valid Bit:** Indicating whether the particular translation is valid.
- Protection Bit:** Indicating whether the page could be read from, written to, or executed from
- Present Bit:** Indicating whether this page is in physical memory or on disk(swapped out)
- Dirty Bit:** Indicating whether the page has been modified since it was brought into memory
- Reference Bit(Accessed Bit):** Indicating that a page has been accessed

February 22, 2017

TCCS422: Operating Systems [Winter 2017]
Institute of Technology, University of Washington - Tacoma

L15.37

HOW BIG ARE PAGE TABLES?

- Page tables are too big to store on the CPU
- Page tables are stored using physical memory
- Paging supports efficiently storing a sparsely populated address space
 - Reduced memory requirement
Compared to base and bounds, and segments

February 22, 2017

TCCS422: Operating Systems [Winter 2017]
Institute of Technology, University of Washington - Tacoma

L15.38

DOES PAGING MAKE THE SYSTEM TOO SLOW?

- Translation
- Issue #1:** Starting location of the page table is needed
 - HW Support: Page-table base register
 - stores active process
 - Facilitates translation
- Issue #2:** Each memory address translation for paging requires an extra memory reference
 - HW Support: TLBs (Chapter 19)

February 22, 2017

TCCS422: Operating Systems [Winter 2017]
Institute of Technology, University of Washington - Tacoma

L15.39

PAGING MEMORY ACCESS

```
1. // Extract the VPN from the virtual address
2. VPN = (VirtualAddress & VPN_MASK) >> SHIFT
3.
4. // Form the address of the page-table entry (PTE)
5. PTEAddr = PTBR + (VPN * sizeof(PTE))
6.
7. // Fetch the PTE
8. PTE = AccessMemory(PTEAddr)
9.
10. // Check if process can access the page
11. if (PTE.Valid == False)
12.     RaiseException(SEGMENTATION_FAULT)
13. else if (CanAccess(PTE.ProtectBits) == False)
14.     RaiseException(PROTECTION_FAULT)
15. else
16.     // Access is OK: form physical address and fetch it
17.     offset = VirtualAddress & OFFSET_MASK
18.     PhysAddr = (PTE.PFN << PFN_SHIFT) | offset
19.     Register = AccessMemory(PhysAddr)
```

February 22, 2017

TCCS422: Operating Systems [Winter 2017]
Institute of Technology, University of Washington - Tacoma

L15.40

COUNTING MEMORY ACCESSES

- Example: Use this Array initialization Code

```
int array[1000];
...
for (i = 0; i < 1000; i++)
    array[i] = 0;
```

- Assembly equivalent:

```
0x1024 movl $0x0, (%edi,%eax,4)
0x1028 incl %eax
0x102c cmpl $0x03e8,%eax
0x1030 jne 0x1024
```

February 22, 2017

TCCS422: Operating Systems [Winter 2017]
Institute of Technology, University of Washington - Tacoma

L15.41

WHERE THE MEMORY ACCESSES ARE:
FOR THE FIRST 5 LOOP ITERATIONS

- Locations:
 - Page table
 - Array
 - Code
- 50 accesses for 5 loop iterations

February 22, 2017

TCCS422: Operating Systems [Winter 2017]
Institute of Technology, University of Washington - Tacoma

L15.42

