

Copyright © William C. Cheng

11

➤ **First-fit vs. best-fit** allocation

- studies have shown that first-fit works better
- best-fit tends to leave behind a large number of regions of memory that are too small to be useful
- best-fit tends to create smallest left-over blocks!
- this is the general problem of **fragmentation**
- **internal fragmentation**: unusable memory is contained within an allocated region (e.g., buddy system)
- **external fragmentation**: unusable memory is separated into small blocks and is interspersed by allocated memory (e.g., best-fit)

Fragmentation

Copyright © William C. Cheng

12

- a linked list of **free blocks**
- don't need to manage allocated blocks
- use a doubly-linked list
- insertion and deletion are fast, i.e., **$O(1)$** , once you know where to insert or delete

struct tblock

size
link

struct tblock

size
link

struct tblock

size
link

Implementing First Fit: Data Structures

Copyright © William C. Cheng

9

Allocate 1000 bytes:

First Fit

Best Fit

Allocate 1100 bytes:

First Fit

Best Fit

Allocate 250 bytes:

First Fit

Best Fit

Allocation Example

Copyright © William C. Cheng

10

Allocate 1000 bytes:

First Fit

Best Fit

Allocate 1100 bytes:

First Fit

Best Fit

Allocate 250 bytes:

First Fit

Best Fit

Allocation Example

Copyright © William C. Cheng

7

Allocate 1000 bytes:

First Fit

Best Fit

Allocate 1100 bytes:

First Fit

Best Fit

Allocate 250 bytes:

First Fit

Best Fit

Allocation Example

Copyright © William C. Cheng

8

Allocate 1000 bytes:

First Fit

Best Fit

Allocate 1100 bytes:

First Fit

Best Fit

Allocate 250 bytes:

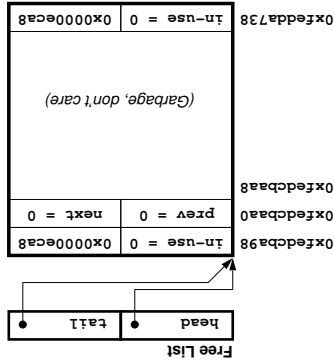
First Fit

Best Fit

Allocation Example



- Ex: Heap starts at 0xfedcba98 and size of the heap is 0x0000e0ca8 (60,584) bytes
- the Free List contains one free block and it looks like this:



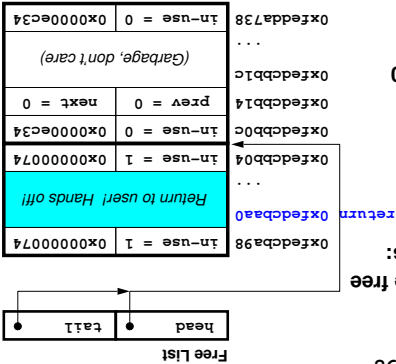
malloc()



- Ex: Heap starts at 0xfedcba98 and size of the heap is 0x0000eca8 (60,584) bytes
- the Free List contains one free block and it looks like this:

```
return
```

- Ex: Request block size is 100
- split the block into two
- busy block size is 116
- remaining free block size is 605



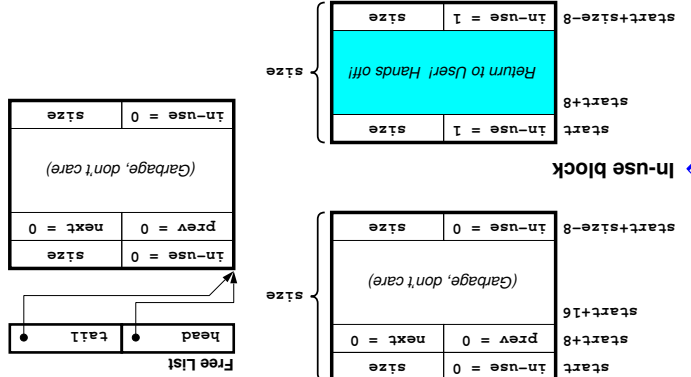
malloc()



Free block

start
start+8
start+16
start+size-8

In-use block

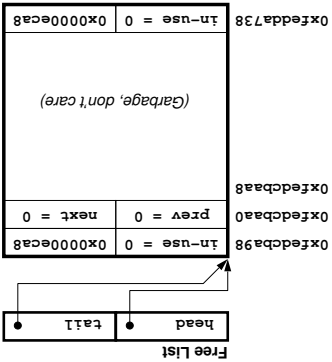


Detailed Examples



Ex: Heap starts at 0xfedcba98 and size of the heap is 0x0000eca8 (60,584) bytes

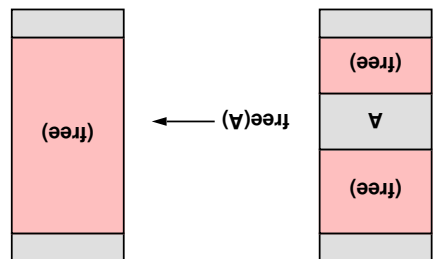
- the Free List contains one free block and it looks like this:



malloc()



- This is known as *coalescing*
- in order to make coalescing possible, you need to know that *size* of the blocks above and below the block being freed
- you also need to know if they are *allocated* or *free*

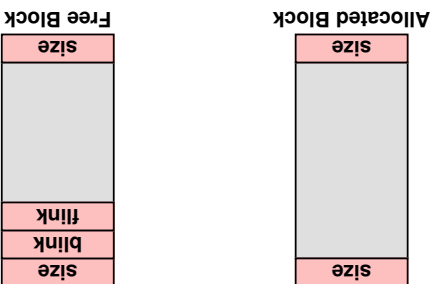


Liberation of Storage



← This is known as **coalescing**

- = in order to make coalescing possible, you need to know that **size** of the blocks above and below the block being freed
- you also need to know if they are **allocated** or **free**



Boundary Tags



First-fit Algorithm

- Let n be the number of tree blocks on the free list
 - $\text{malloc}()$ is $O(n)$
 - $\text{free}(\text{ptr})$ is $O(n)$
 - occurs when the blocks around the block containing `ptr` are both in-use
- Such performance is unacceptable in the kernel



free() Example

- Ex: free (x) and previous block is also free**

 - i.e., x-16 is 0 and y-8+z is 0
 - where z is what's in y-4,
 - x is what's in y-4+z, and w is what's in x-12
 - blocks starting at y-8-w and y-8+z are both on the Free List and next to and point at each other
 - coalesce all 3 blocks
 - just change y-4-w and x-12+z+x to w+z+X
 - copy next from y+z+4 to y-w+4
 - adjust prev field in the new next block in Free List to point to x-8-w
 - may need to update where Free List points

prev	size

$X+Z+9$	$0 \leftarrow \text{use} = 0$	$\text{size} = 0$
X	<p>(Garbage, don't care)</p>	
$8-X$		
$16-X$		
$M-X$	prev	next
$M-8-X$	$0 \leftarrow \text{use} = 0$	$\text{size} = 8+Z+X$



Example () free

- Ex: free (x) and *previous block* is busy and *next block* is free

 - i.e., $x-16$ is 1 and $x-8+z$ is 0
 - where z is what's in $y-4$
 - and x is what's in $y-4+z$

furthermore, $y-8+z$ is on the Free List

⇒ coalesce this block and the *next* block

 - just change $y-4$ and $y-12+z+x$ to $z+x$ and $y-8$ to 0
 - move prev and next pointers

0	8	16	24
1	2	3	4
→ 8	→ 16	→ 24	→ null

[illegible]

3.3 Dynamic Storage Allocation

- ➡ Best-fit & First-fit Algorithms
- ➡ *Buddy System*
- ➡ Slab Allocation



First-fit & Best-fit Algorithms

- ▶ Memory allocator must run fast
 - ▶ it does not check if the free list is in a consistent state
 - ▶ just like our warmup 1 assignment
- ▶ One bad bit in the memory allocator data structure and it can break the memory allocator code
- ▶ if you write into a *boundary tag*, your program may die in malloc() or free()
 - ▶ what would happen if you call free() twice on the same address?
 - ▶ user/application code can *corrupt the memory allocation chain* easily
 - ▶ the result can lead to *segmentation faults*
 - ▶ unfortunately, the corruption can *stay hidden* for a long time and *eventually* lead to a segmentation fault
 - ▶ memory corruption bugs are very difficult to squash



Example free()

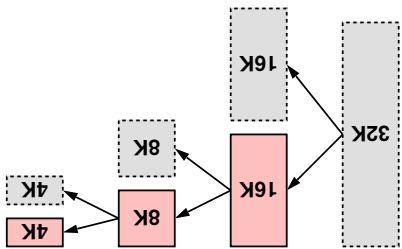
- Ex: free (x) and *previous block*
- ie., $y-16$ is 0 and $x-8$ is 0
- where z is what's in $y-4$,
- x is what's in $y-4+z$, and
- w is what's in $y-12$
- blocks starting at $y-8-w$ and $x-8+z$ are both on the Free List
- and next to and point at each other
- coalesce all 3 blocks

M-8-X	in-use=0	size=M	prev	X-8-Z	size=M
M-X	prev	X-8-Z	size=M	(Garbage, don't care)	in-use=0
91-X	in-use=0	size=M	in-use=1	size=Z	
8-X	in-use=1	size=Z	Return to user! Hands off!		
X					
Z+8-X	in-use=0	size=X			
Z+X	X-8-M	next	(Garbage, don't care)		
	in-use=0	size=X			



Copyright © William C. Cheng

- blocks get evenly divided into two blocks that are buddies with each other
- can only merge with your buddy if your buddy is also free
- internal fragmentation**
- Ex: malloc(4000)
- return a 4K block



Buddy Lists

Operating Systems - CSCI 402



Copyright © William C. Cheng

- Data Structure**
 - doubly-linked list (not circular) FREE list indexed by k
 - links stored in actual blocks
 - FREE[k] points to first available block of size 2^k
- each block contains
 - in-use bit
 - size
 - NEXT and PREV links for FREE list
- lots of details
 - read `wenix` source code for its "page allocator"
- Can get greater variety in block sizes using Fibonacci sequence of block sizes so $b_i = b_{i-1} + b_{i-2}$
- ratio of successive block sizes is $2/3$ instead of $1/2$

Buddy Systems

Operating Systems - CSCI 402



Copyright © William C. Cheng

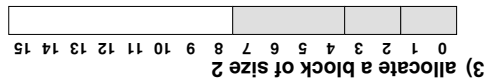
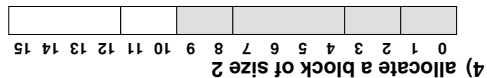
- Faster memory allocation system (at the cost of more fragmentation, internal fragmentation)**
 - restrict block size to be a power of 2
 - all blocks of size 2^k start at location x where $x \bmod 2^k = 0$
 - given a block starting at location x such that $x \bmod 2^k = 0$
 - only buddies can be merged
 - $BUDDX_k(x) = x + 2^k$ if $x \bmod 2^{k+1} = 0$
 - $BUDDX_k(x) = x - 2^k$ if $x \bmod 2^{k+1} = 2^k$
 - Ex: $BUDDX_2(1010100) = 1010000$
 - try to coalesce buddies when storage is deallocated
 - different available block lists, one for each block size
 - When request a block of size 2^k and none is available:
 - split smallest block $2^l > 2^k$ into a pair of blocks of size 2^{l-1}
 - place block on appropriate free list and try again

Buddy Systems

Operating Systems - CSCI 402



Copyright © William C. Cheng



Ex: 16 "pages" (minimum allocation is 1 page)

High-level Example of Buddy Algorithm

Operating Systems - CSCI 402

free[k]	0
1	0
2	0
3	0
4	0

free[k]	0
1	0
2	0
3	0
4	0

3.3 Dynamic Storage Allocation

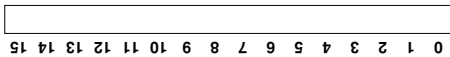
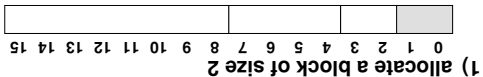
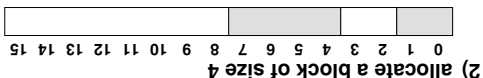
- Best-fit & First-fit Algorithms
- Buddy System
- Slab Allocation



Copyright © William C. Cheng



Copyright © William C. Cheng



Ex: 16 "pages" (minimum allocation is 1 page)

High-level Example of Buddy Algorithm

Operating Systems - CSCI 402

free[k]	0
1	0
2	0
3	0
4	0

free[k]	0
1	0
2	0
3	0
4	0

free[k]	0
1	0
2	0
3	0
4	0



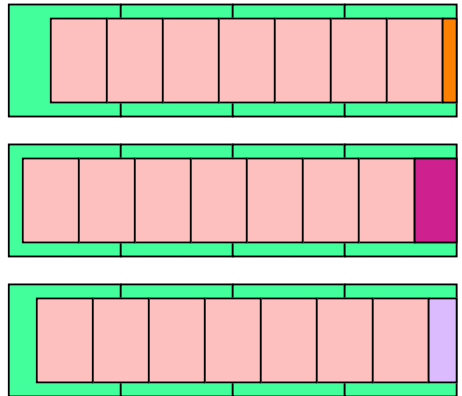
Slab Allocation

- we will cover "pages" later, won't get into too much detail now
 - contiguous sets of pages called **slabs**, allocated to hold objects
 - sets up a separate cache for each type of object to be managed
- Slab Allocation**
- As **objects** are being freed, they are simply marked as free
 - don't have to free up storage
 - As **objects** are being freed, they are simply marked as free
 - objects are considered "preallocated" since they have all been initialized already
 - existing slabs in the cache
 - As **objects** are being allocated, they are taken from the set of existing slabs in the cache
 - this is where you pay for initialization, but it's done in a **batch**
 - the objects it hold
 - Whenever a **slab** is allocated, a constructor is called to initialize all the objects it hold



Slab Allocation

- see weenix kernel code!



Slab Allocation

- ▶ Objects are allocated and freed frequently
 - allocation involves
 - finding an appropriate-sized storage
 - initialize it
 - ◆ pointers need to point at the right places
 - ◆ may even need to initialize synchronization data
 - deallocation involves
 - tearing down the data structures
 - freeing the storage
 - lots of "overhead"
- ▶ Difficulties with dynamic storage allocation
 - you cannot predict what an application will ask for
 - but it's not true for the kernel
 - e.g., can allocate a slab of process control blocks and return one of them from a slab