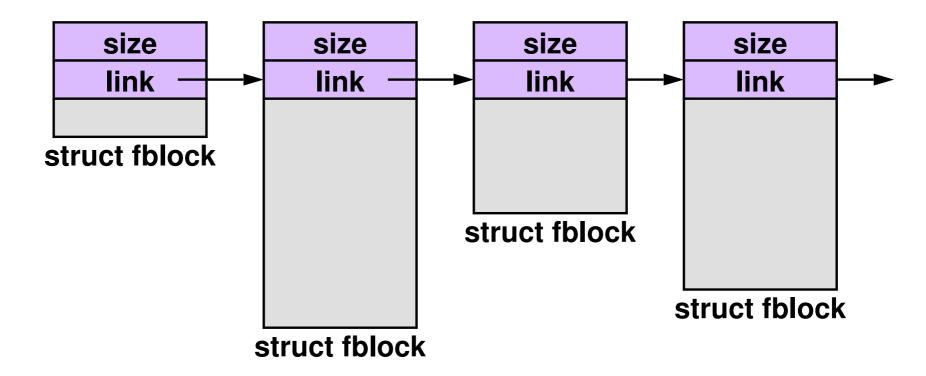
Implementing First Fit: Data Structures



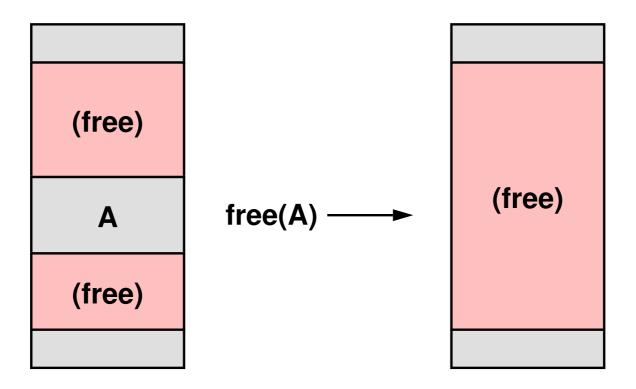


Free list: a linked list of free blocks

- sorted according to block addresses
 - no 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



Liberation of Storage

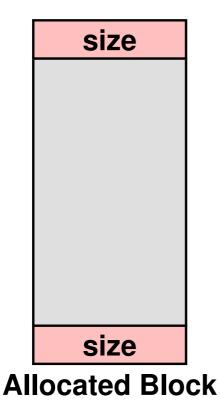


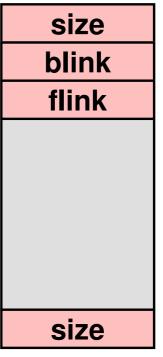


This is known as *coalescing*

- 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





Free Block



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

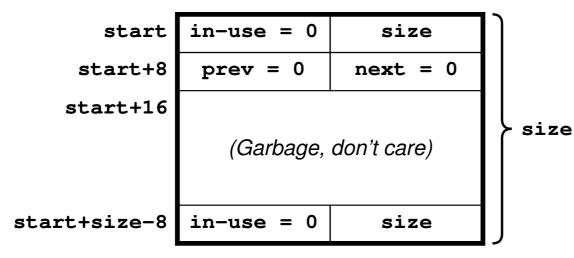
Detailed Examples

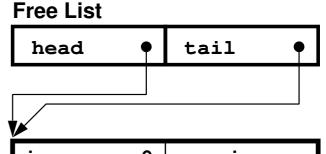


Free block



Free list

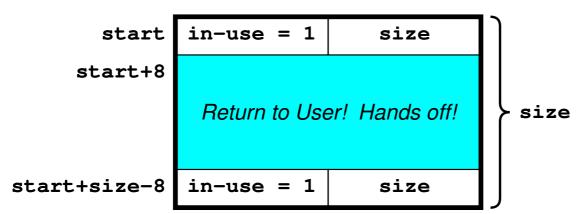


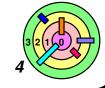


in-use = 0 size prev = 0 next = 0

(Garbage, don't care)

n-use block

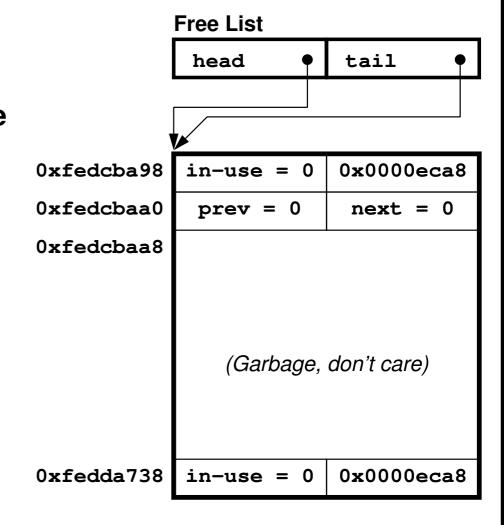


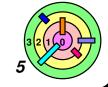


malloc() Example

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:

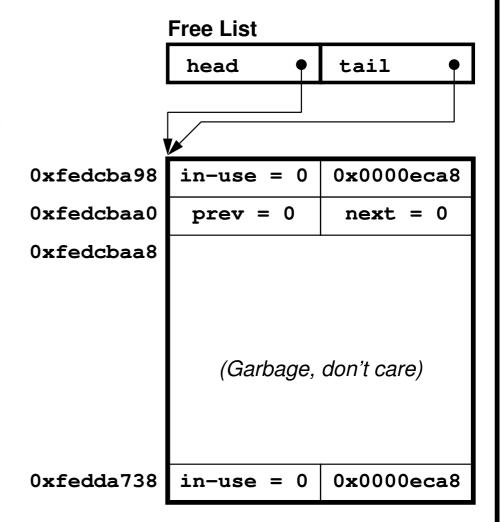




malloc() Example

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:





Ex: Request block size is 100

- split the block into two
- busy block size is 116
- remaining free block size is 60584-116 =60468=0xec34



malloc() Example

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: Free List

head • tail •

0xfedcba98 in

0xfedcbaa0

...

0xfedcbb04 in

0xfedcbb0c in

0xfedcbb14 p

0xfedda738

in-use = 1 0x00000074

Return to user! Hands off!

in-use = 1 0x00000074

in-use = 0 0x0000ec34

prev = 0 next = 0

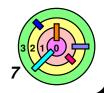
(Garbage, don't care)

in-use = 0 0x0000ec34



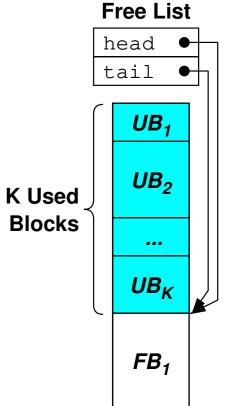
Ex: Request block size is 100

- split the block into two
- busy block size is 116
- remaining free block size is 60584-116 =60468=0xec34



After *K* blocks of memory have been allocated (and assume that none of them have been deallocated)

in the memory layout, the first K blocks are used block, followed by one free block







After *K* blocks of memory have been allocated (and assume that none of them have been deallocated)

in the memory layout, the first K blocks are used block, followed by one free block



Memory blocks can be freed in any order

when a memory block is freed, we need to check if the blocks before it and after it are also free

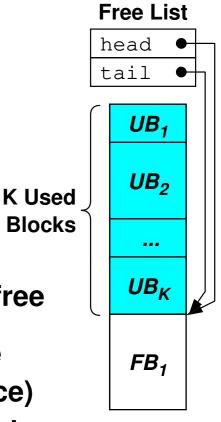


If neither of them are free, we just need to insert the newly freed block into the Free List (at the right place)

- need to search the Free List to find insertion point
- searching through a linear list is "slow", O(n)



Otherwise, we can *merge/coalesce* the block in question with neighboring free block(s)







Ex: free(Y)

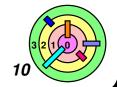
- Y-8-(*(Y-12))
- Y-16 tells you if the *previous* block is free or not
- Y-8+z tells you if the next block is free or not
 - where z is what's in Y-4



Coalescing:

need to make sure that everything is consistent

-		
-12))	in-use=?	size
		?
Y-16	in-use=?	size
Y-8	in-use=1	size=Z
Y	Return to use	er! Hands off!
	in-use=1	size=Z
Y-8+Z	in-use=?	size
		?
	in-use=?	size

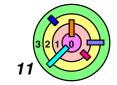




Ex: free (Y) and previous block is free and next block is busy

- i.e., Y-16 is 0 and Y-8+z is 1
 - where z is what's in Y-4 and w is what's in Y-12
- **-** furthermore, Y-8-₩ is on the Free List
- coalesce this block and the previous block

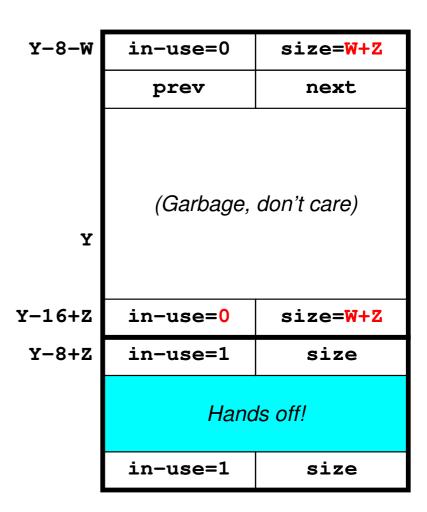
_		
W-8-Y	in-use=0	size=W
	prev	next
	(Garbage,	don't care)
Y-16	in-use=0	size=W
Y-8	in-use=1	size=Z
Y	Return to use	er! Hands off!
Y-16+Z	in-use=1	size=Z
Y-8+Z	in-use=1	size
	Hands off!	
	in-use=1	size

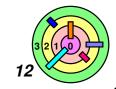




Ex: free (Y) and previous block is free and next block is busy

- i.e., Y-16 is 0 and Y-8+z is 1
 - where z is what's in Y-4 and w is what's in Y-12
- furthermore, Y-8-₩ is on the
 Free List
- coalesce this block and the previous block
 - easy!
 - just change Y-12+z and Y-4-W to W+z and Y-16+z to 0
 - o don't even need to change prev and next!



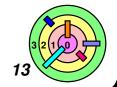




Ex: free (Y) and previous block is busy and next block is free

- i.e., Y-16 is 1 and Y-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

W-8-Y	in-use=1	size=W
	Hand	ls off!
Y-16	in-use=1	size=W
Y-8	in-use=1	size=Z
Y	Return to use	er! Hands off!
	in-use=1	size=Z
Y-8+Z	in-use=0	size=X
	prev	next
	(Garbage,	don't care)
Y-16+Z+X	in-use=0	size=X

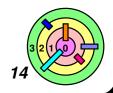




Ex: free (Y) and previous block is busy and next block is free

- i.e., Y-16 is 1 and Y-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

- W-8-Y in-use=1 size=W Hands off! Y-16 in-use=1size=W Y-8 in-use=0size=Z+Xprev Y next (Garbage, don't care) Y-16+Z+X in-use=0 size=Z+X
- adjust next field in previous block in Free List
- adjust prev field in next block in Free List
- may need to update where Free List points

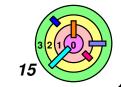




Ex: free (Y) and previous block is free and next block is also free

- i.e., Y-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 Y-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

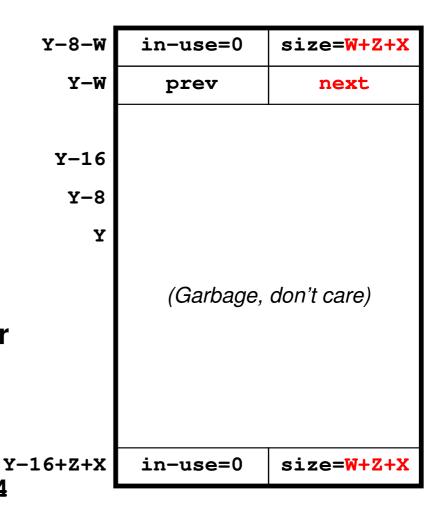
W-8-Y	in-use=0	size=W
Y-W	prev	Y-8+Z
	(Garbage,	don't care)
Y-16	in-use=0	size=W
Y-8	in-use=1	size=Z
Y	Return to use	er! Hands off!
	in-use=1	size=Z
Y-8+Z	in-use=0	size=X
Y+Z	W-8-Y	next
	(Garbage,	don't care)
	in-use=0	size=X

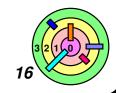




Ex: free (Y) and previous block is free and next block is also free

- i.e., Y-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 Y-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 Y-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 Y-8-W
 - may need to update where Free List points





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 any memory allocator data structure can break the memory allocator code

- if you write into a boundary tag, your program may die later 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 debug



First-fit Algorithm



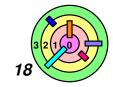
Let *n* be the number of free blocks on the free list

- \rightarrow in the worst case, malloc() is O(n)
- in the worst case, free (ptr) is O(n)
 - occurs when the blocks around the block containing ptr are both in-use



Such performance in unacceptable in the kernel

it is desirable that the kernel's worst-case performance has a bound



3.3 Dynamic Storage Allocation



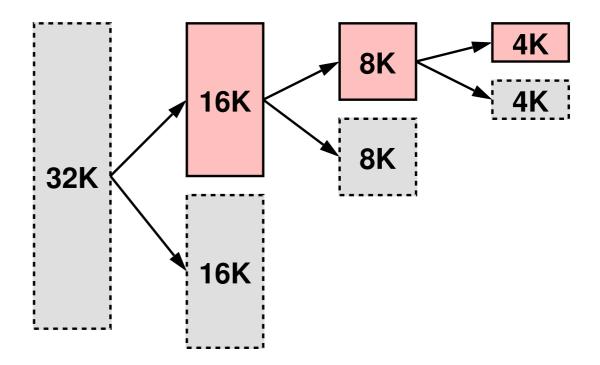
Buddy System

Slab Allocation

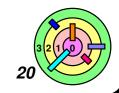


Buddy Lists

Ex: malloc(4000)



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

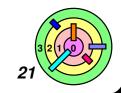


Buddy Systems



Faster memory allocation system (at the cost of more fragmentation, internal fragmentation)

- restrict block size to be a power of 2
 - 1) all blocks of size 2^k start at location x where $x \mod 2^k = 0$
 - 2) given a block starting at location x such that $x \mod 2^k = 0$
 - \Rightarrow BUDDY_k(x) = x+2^k if x mod 2^{k+1}=0
 - \Rightarrow BUDDY_k(x) =x-2^k if x mod 2^{k+1}=2^k
 - \Rightarrow Ex: BUDDY₂(1010100) = 1010000
 - 3) only buddies can be merged
 - 4) try to coalesce buddies when storage is deallocated
 - k different available block lists, one for each block size
 - When request a block of size 2^k and none is available:
 - 1) split smallest block $2^{j} > 2^{k}$ into a pair of blocks of size 2^{j-1}
 - 2) place block on appropriate free list and try again



Buddy Systems

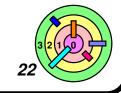


Data Structure

- 1) 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
- 2) each block contains
 - in−use bit
 - size
 - **⋄** NEXT and PREV links for FREE list
- lots of details
 - read weenix source code for its "page allocator"

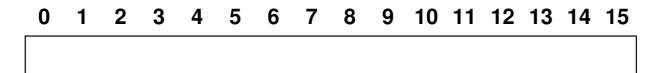


ratio of successive block sizes is 2/3 instead of 1/2

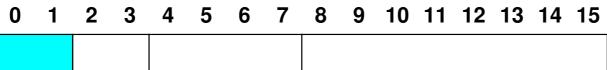


High-level Example of Buddy Algorithm

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

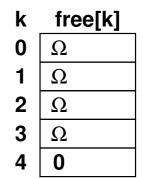


1) allocate a block of size 2



2) allocate a block of size 4





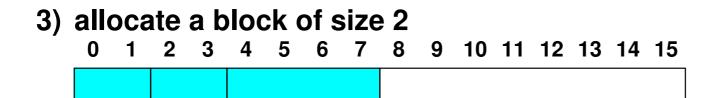
k	free[k]	
0	Ω	
1	X 2	
2	X 4	
3	№ 8	
4	X O	

k	free[k]
0	Ω
1	X 2
2	Ω X X
3	8
4	χ Ω



High-level Example of Buddy Algorithm







k	free[k]
0	Ω
1	Ω Χ Ω
2	Ω X Ω
3	X 8
4	Θ Ω

K	_free[k]
0	Ω
1	X X 10
2	X X 12
3	ΩΧΧΩ
4	Θ (Ω

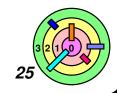


3.3 Dynamic Storage Allocation



Buddy System

Slab Allocation



Slab Allocation



Objects are allocated and freed frequently

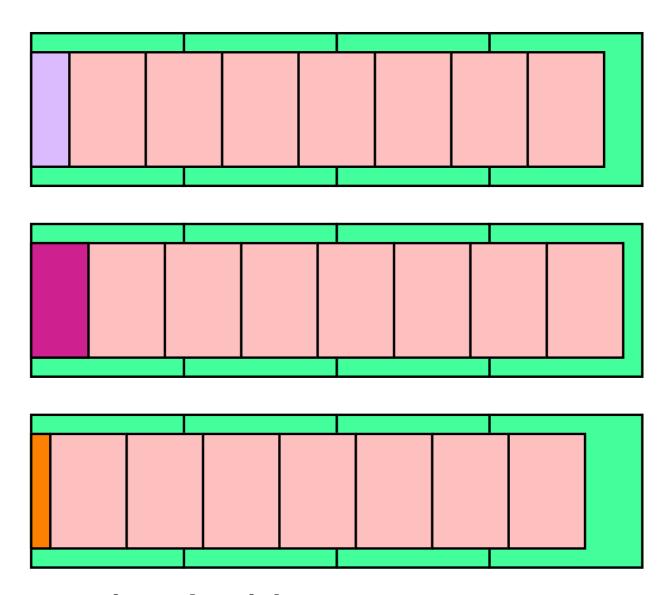
- allocation involves
 - finding an appropriate-sized storage
 - initialize it
 - pointers need to point at the right places
 - may even need to initialize synchronization data structures
- deallocation involves
 - tearing down the data structures
 - freeing the storage
- lots of "overhead"



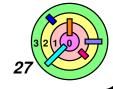
Difficulties with dynamic storage allocation

- 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 at a time
 - return one of them from a slab

Slab Allocation



see weenix kernel code!



Slab Allocation



Slab Allocation

- sets up a separate cache for each type of object to be managed
- contiguous sets of pages called slabs, allocated to hold objects
 - we will cover "pages" later, won't get into too much detail now



- this is where you pay for initialization, but it's done in a batch
- As *objects* are being allocated, they are taken from the set of existing slabs in the cache
 - objects are considered "preallocated" since they have all been initialized already
- As *objects* are being freed, they are simply marked as free
 - don't have to free up storage
 - when appropriate can free up an entire slab

