#### **Memory Allocation (fin)** Computer Operating Systems, Spring 2025

Instructors: Joel Ramirez Travis McGaha

**Head TAs:** Ash Fujiyama Emily Shen Maya Huizar

#### TAs:

Ahmed Abdellah	Bo Sun	Joy Liu	Susan Zhang	Zihao Zhou
Akash Kaukuntla	Connor Cummings	Khush Gupta	Vedansh Goenka	
Alexander Cho	Eric Zou	Kyrie Dowling	Vivi Li	
Alicia Sun	Haoyun Qin	Rafael Sakamoto	Yousef AlRabiah	
August Fu	Jonathan Hong	Sarah Zhang	Yu Cao	



pollev.com/tqm

How is PennOS Going? Any Questions about Memory Allocation?

- Apply to be a CIS 4480/5480 TA!!!!! (and/or CIS 3990 TA)
  - Do you have an interest in teaching?
  - Did you enjoy the material in this course?
  - Were you super grateful to the TAs who helped/are helping you through this course?
  - Are you interested in supporting your fellow students who will take a course?
  - Looking to either brush up on the material or share your passion for it?
  - Ready to become that "goated" TA?
- If you apply, you should hear from us by Monday of next week...
  - It will be taught by Kyrie and Joel!
- LINK HERE also: <u>https://www.cis.upenn.edu/ta-information/</u>

- PennOS
  - PennOS Final Demo posted sometime tomorrow today?
    - You will have to demo FAT here even if you pass the autograder.
  - Integration can be a big pain, make sure you allocate enough time to it!
- Some notes:
  - DO NOT mmap the entire File System. Only mmap the Allocation Table, the rest of the file system needs to be handled with lseek/write.
    - Do not keep the contents of the file in memory, it should be stored in the file
    - If your PennFat is killed with kill -9, your file contents should still be saved in disk
  - Advice for using gdb to debug
    - handle SIGUSR1 noprint nostop

Makes it so that gdb doesn't report every time SIGUSR1 goes and interrupts you

(more on next slide)

- Some notes:
  - Reminder, you instead of just doing:

you may need to do:

lseek(FAT\_FD, offset, SEEK\_SET);
write(FAT\_FD, contents, size);

lseek(FAT\_FD, offset, SEEK\_SET);
write(FAT\_FD, contents, size);
lseek(FAT\_FD, offset, SEEK\_SET);
write(FAT\_FD, contents, size);

- With the description of setitimer(), it just says that sigalarm is delivered to the process, not necessarily the calling thread. To make sure siglaram goes to the scheduler, you may want to make it so that all threads (spthread or otherwise) that aren't the scheduler call something like: pthread\_sigmask(SIG\_BLOCK, SIGALARM)
  - Which will block SIGALARM in that thread.

- If you are having issues with the scheduler not running you can try running
  - strace -e 'trace=!all' ./bin/pennos
  - You may have to install strace: sudo apt install strace
  - This will print out every time a signal is sent to your pennos
  - (Usual fix is the pthread\_sigmask thing on the previous slide)

# **Lecture Outline**

- Garbage collection
- Arena Allocation
- Buddy Algorithm
- Slab/Slub Allocator

# **Memory Leaks**

- The most common Memory Pitfall
- What happens if we allocate something, but don't delete it?
  - That block of memory cannot be reallocated, even if we don't use it anymore, until it is delete-d
- Garbage Collection
  - Automatically "frees" anything once the program has lost all references to it
  - Affects performance, but avoid memory leaks
  - Java and other "high level" languages
- RAII (Resource Acquisition Is Initialization)
  - C++ and Rust have this, it is VERY GOOD

# Garbage Collection

- When memory is automatically deallocated for us, so we do not need to explicitly free memory
- Very common in higher level languages:
  - Java, C#, Python, Javascript, Ruby, Lisp, Erlang, Racket, Haskell, Scala, Dart, etc.
- Big difference between these languages and languages like C / C++ / Rust:
  - Many of these languages are not run directly on your hardware.
  - Java (for example) runs on the JVM (Java Virtual Machine) which then runs on your computer
  - Garbage collection requires some help from the "runtime" environment" and/or the compiler to keep track of pointers, memory allocations etc and decide when to free them

# **Garbage Collection**

- With the aid of the runtime and compiler we can keep track of all memory allocation and represent it as a directed graph
  - Each allocation is a node in the graph
  - Each pointer is an edge in the graph
  - If an object contains a pointer to another object we draw an edge from that node to the other.



#### **Garbage Collection**

- Each allocation is a node in the graph
- Each pointer is an edge in the graph
- If an object contains a pointer to another object we draw an edge from that node to the other.
- We also keep track of which pointers are held by local variables (pointers that are not on the heap, but point to the heap). These are "roots"

Nodes that are "reachable" from a root are safe if it can't be reached from a root, then it is garbage

# **Lecture Outline**

- Garbage collection
- Arena Allocation
- Buddy Algorithm
- Slab/Slub Allocator

#### **Arena Allocator**

- In some instances, we want to allocate a lot of items and limit those allocations to one scope. We call our allocator a "arena allocator". It allocates things within the same "arena"/region/pool of the same scope
- ✤ For example, Consider we start with:



- Note that there is a little bit more metadata than just these two pointers
- Then we allocate 4 bytes



#### Arena Allocator: Alloc

For example, Consider we start with:



- Note that there is no metadata, just these two pointers
- Then we allocate 4 bytes

Then we allocate 16 bytes



\*Alignment could be a thing that affects how we allocate things, but we are leaving that out

#### **Arena Allocator: Free**

 Once we are done with our temporaries, we free the all allocations, and we can then use it again as if "fresh"



- Looks the same as when we started!
- \* That is the API
- Example usage:



 Instead of being scoped to a function, an arena allocator can also be scoped to an "object" or the lifetime of some "task"

#### **Arena Allocator: Growing**

- This simple arena allocator we are showing can also be called a "bump allocator" since to allocate we just "bump" the pointer
- All the memory for an arena allocator is allocated before hand, typically there is a good guess for the memory that a given scope will need, so we can just allocate that many pages or bytes
- If we want to handle growing an arena allocator, it may handle multiple "arenas" and simply allocate a new arena whenever one is requested.
  - Can allocate new pages by using mmap() to create "anonymous" mappings (anonymous = pages aren't mapped to a file)



pollev.com/tqm

- How fast is our arena allocator at allocating things on average? At freeing things?
- What does the internal and external fragmentation look like with our arena allocator?
- Why can't we use this as a replacement for malloc maintaining lists of allocated & freed memory?

# **Lecture Outline**

- Garbage collection
- Arena Allocation
- Buddy Algorithm
- Slab/Slub Allocator

# **Buddy Algorithm**

- Keeps in mind that there is some "maximum" amount of memory and divides memory into partitions that are powers of 2.
  - Power of 2 allows for compact allocation tracking and makes coalescing memory quick.
  - Usually with the smallest unit being 1 page, 4096 bytes.

- Modified implementation of the buddy system is one of many things used by the Linux kernel and the others (like a version of malloc called jemalloc)
  - Linux Kernel uses the buddy algorithm to allocate physical pages to the kernel

0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
							24 pa	ages							

- We start with the full pool of memory, in this example, 2<sup>4</sup>
   pages (usually a higher cap than this, this is for example)
- What happens if someone asks to allocate 1 page?

0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
			2 <sup>3</sup> pa	ages							2 <sup>3</sup> p	ages			

- What happens if someone asks to allocate 1 page?
  - Split page chunks into half until we have enough

0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
	2 <sup>2</sup> pa	ages			2 <sup>2</sup> pa	ages					2 <sup>3</sup> p	ages			

- What happens if someone asks to allocate 1 page?
  - Split page chunks into half until we have enough

0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
<b>2</b> <sup>1</sup> <b>p</b>	ages	21 pa	ages		2² pa	ages					2 <sup>3</sup> p	ages			

- What happens if someone asks to allocate 1 page?
  - Split page chunks into half until we have enough

0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
1	1	21 p	ages		2² p	ages					2 <sup>3</sup> p	ages			

L23: Memory Allocation

- What happens if someone asks to allocate 1 page?
  - Split page chunks into half until we have enough

0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
1	1	21 p	ages		2² pa	ages				-	2 <sup>3</sup> p	ages			
Α		•		-											

- What happens if someone asks to allocate 1 page?
  - Split page chunks into half until we have enough
- Can mark the one page as being used by allocation A

0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
1	1	21 pa	ages		2² pa	ages					2 <sup>3</sup> p	ages			
A				-											

Now someone requests 2 pages, what happens?

0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
1	1	21 pa	ages		2² pa	ages					2 <sup>3</sup> p	ages			
A		В													

- Now someone requests 2 pages, what happens?
- We can claim the 2<sup>1</sup>-page chunk and mark it as being used by allocation B

0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
1	1	21 pa	ages	2 <sup>2</sup> pages							2 <sup>3</sup> p	ages			
A		В													

Now someone requests 3 pages, what happens?



- Now someone requests 3 pages, what happens?
- Buddy ONLY deals with powers of 2, this gets rounded up to 2<sup>2</sup> pages (4 pages)
- We can claim the 2<sup>2</sup>-page chunk and mark it as being used by allocation C



- Last allocation: someone allocates 1 page, what happens?
- We can claim the 1-page chunk and mark it as being used by allocation D



- Let's walk through the freeing process
- First, allocation D is done and frees its page



- Let's walk through the freeing process
- First, allocation D is done and frees its page
- To free the page, we just mark it as no longer being allocated. Nothing we can coalesce (yet)



- Let's walk through the freeing process
- Second, allocation A is done and frees its page
- To start, we just mark it as no longer being allocated.



- Let's walk through the freeing process
- Second, allocation A is done and frees its page
- To start, we just mark it as no longer being allocated.
- Then we can coalesce!
- Each "chunk" has a "buddy", the buddy being the its "twin" created while spitting chunks in half.
- ✤ If both buddies are free, they can be combined ☺



- Let's walk through the freeing process
- Second, allocation A is done and frees its page
- To start, we just mark it as no longer being allocated.
- Then we can coalesce!
- Each "chunk" has a "buddy", the buddy being the its "twin" created while spitting chunks in half.
- ✤ If both buddies are free, they can be combined ☺



- Let's walk through the freeing process
- Third, allocation C is done and frees its pages

0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
21 pa	ages	21 pa	ages		2² pa	ages					2 <sup>3</sup> p	ages			
		В													

- Let's walk through the freeing process
- Third, allocation C is done and frees its pages
- Can't coalesce since its buddy is not completely free

0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
21 p	ages	21 p	ages		2² pa	ages					2 <sup>3</sup> p	ages			

- Let's walk through the freeing process
- lastly, allocation B is done and frees its pages

0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	
	2 <sup>2</sup> pa	ages		2 <sup>2</sup> pages				2 <sup>3</sup> pages								

- Let's walk through the freeing process
- lastly, allocation B is done and frees its pages
- Its buddy is free so we can coalesce!

0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
2 <sup>3</sup> pages											2 <sup>3</sup> p	ages			

- Let's walk through the freeing process
- lastly, allocation B is done and frees its pages
- Its buddy is free so we can coalesce!
- The newly coalesced chunk can be further coalesced!

0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
	2 <sup>4</sup> pages														

- Let's walk through the freeing process
- lastly, allocation B is done and frees its pages
- Its buddy is free so we can coalesce!
- The newly coalesced chunk can be further coalesced!
- The newly coalesced chunk can be further coalesced!

# **Buddy Algorithm Implementation**

Buddy Algorithm can be maintained with a binary search tree



#### **Buddy Algorithm Implementation**

 Since Buddy has a known max size, we can represent the tree in an array or bitmap. (example shows up to 2<sup>2</sup> for space on the slide)

$$2^2$$
  $2^1$   $2^1$   $2^0$   $2^0$   $2^0$   $2^0$ 



(alternate way to show the array, may make the connection between array and tree easier to see). Indexes go Left -> Right, top to bottom



#### **Buddy Algorithm Implementation**

- The tree (array representation) is useful for coalescing, but we can make algorithm faster by keeping track of several free lists, roughly one list per size
  - Quicker lookup for memory allocation
  - Coalescing is still fast since we can maintain a bitmap and easily find the location of a "buddy". If an allocation's "Buddy" is free it should be 2<sup>k</sup> bytes before/after it.





pollev.com/tqm

How does the fragmentation for the buddy algorithm look?

# **Buddy Algorithm**

- A bit restrictive in the interface, must be a power of 2
  - Internal fragmentation can be a lot  $\ensuremath{\mathfrak{S}}$
  - If someone needs 2<sup>4</sup> +1 pages, buddy algorithm will allocate 2<sup>5</sup> pages, 2<sup>4</sup> 1 pages of fragmentation
- External fragmentation is generally kept pretty small
- Small allocations don't really work for this

# **Lecture Outline**

- Garbage collection
- Arena Allocation
- Buddy Algorithm
- Slab/Slub Allocator

# **Slab Allocator\***

- What if we restrict the API to a *single* size that can be allocated or freed?
- First, you need to allocate the thing you will allocate from
  - When you create it, you specify a name and some other information
  - The thing we care about is that you specify the size of the objects that the slab allocator will allocate from

#### We are simplifying this allocator a good bit

#### **Slab Allocator High Level**

- In the context of a slab allocator
  - Object: the thing we want to allocate, some fixed size memory that we want to allocate NOT the same as a java object
  - slab: a chunk of memory containing the "objects"
  - A cache maintains lists of slabs noting which slabs are full/empty/partially in use



## **Slab Allocator High Level**

There can be multiple slabs that are partial/empty free



# **Slab Allocator High Level: Alloc**

- Each slab maintains a pointer to an element that is free in the slab.
   (This pointer is stored in some metadata somewhere.)
- Each free object contains a pointer to the next free object in the slab
- When we allocate from the cache, we get a pointer to the first element that is free



# **Slab Allocator High Level: Alloc**

- Each slab maintains a pointer to an element that is free in the slab.
   (This pointer is stored in some metadata somewhere.)
- Each free object contains a pointer to the next free object in the slab
- When we allocate from the cache, we get a pointer to the first element that is free



#### **Slab Allocator High Level: Free**

- When we want to free something, we are given the pointer to that object
   So we can do math on the address to calculate the page (and thus which slab it goes to)
- From there we can just "push it to the front of the free list"



# **Slab Allocator High Level: Free**

- When we want to free something, we are given the pointer to that object
   So we can do math on the address to calculate the page (and thus which slab it goes to)
- From there we can just "push it to the front of the free list"





pollev.com/tqm

What is the runtime for slab?

How does the fragmentation look?

# **Slab Allocator Analysis**

- ✤ Slab allocator is very useful for minimizing overhead for allocating and freeing.
- Can be minimal internal and external fragmentation (gets more complicated when you account for alignment and buddy algo requirements)

#### **Slab Allocator Usage**

- Used on top of the buddy algorithm in the kernel.
  - This allows us to use the buddy algorithm still, but can quickly allocate smaller sized "objects" within the *slabs* of memory returned by the buddy algorithm
- General Memory allocators may use something like this, allocate many slabs of various sizes and try to mostly use those for allocation
  - The generic "kmalloc" (kernel malloc) is backed by the slab allocator. When it asks for N bytes it allocates from a slab that will best fit that allocation size.